

Low-Rank Adaptation Approach for Vietnamese-Bahnaric Lexical Mapping from Non-Parallel Corpora

La Cam Huy^{1,2}, Le Quang Minh^{1,2}, Tran Ngoc Oanh^{1,2}, Le Due Dong^{1,2}, Duc Q. Nguyen^{1,2}, Nguyen Tan Sang^{1,2}, Tran Quan^{1,2}, Tho Quan^{1,2,*}



Use your smartphone to scan this QR code and download this article

¹Ho Chi Minh City University of Technology (HCMUT), 268 Ly Thuong Kiet Street, District 10, Ho Chi Minh City, Vietnam

²Vietnam National University Ho Chi Minh City, Linh Trung Ward, Thu Due City, Ho Chi Minh City, Vietnam

Correspondence

Tho Quan, Ho Chi Minh City University of Technology (HCMUT), 268 Ly Thuong Kiet Street, District 10, Ho Chi Minh City, Vietnam

Vietnam National University Ho Chi Minh City, Linh Trung Ward, Thu Due City, Ho Chi Minh City, Vietnam

Email: qttho@hcmut.edu.vn

History

- Received: 7-9-2023
- Accepted: 26-4-2024
- Published Online:

DOI :



Copyright

© VNUHCM Press. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license.



ABSTRACT

Bilingual dictionaries are vital tools for automated machine translation. Leveraging advanced machine learning techniques, it is possible to construct bilingual dictionaries by automatically learning lexical mappings from bilingual corpora. However, procuring extensive bilingual corpora for low-resource languages, such as Bahnaric, poses a significant challenge. Recent studies suggest that non-parallel corpora, supplemented with a handful of anchor words, can aid in the learning of these mappings, which contain parameters for automated translation between source and target languages. The prevailing methodology involves using Generative Adversarial Networks (GANs) and solving the Procrustes orthogonal problem to generate this mapping. This approach, while innovative, exhibits instability and demands substantial computational resources, posing potential issues in rural regions where Bahnaric is spoken natively. To mitigate this, we propose a low-rank adaptation strategy, where the limitations of GANs can be circumvented by directly calculating the rigid transformation between the source and target languages. We evaluated our approach using the French-English dataset, and a low-resource dataset, Vietnamese-Bahnaric. Notably, the Vietnamese-Bahnaric lexical mapping produced by our method is valuable not only to the field of computer science, but also contributes significantly to the preservation of Bahnaric cultural heritage within Vietnam's ethnic minority communities.

Key words: Low-rank adaptation, lexical mapping, low-resource language, Kabsch algorithm

1 INTRODUCTION

The construction of bilingual dictionaries represents a valuable endeavor for both the computational linguistics and computer science communities. This process necessitates the accumulation, classification, and presentation of word pairs and their corresponding translations in two languages¹. Historically, this task has entailed the use of reliable linguistic resources, bilingual documents, and consultations with native speakers to ensure precision. However, with recent developments in Artificial Intelligence (AI), it is now feasible to apply machine learning algorithms to train language models capable of comprehending and generating translations between two languages². Such advancements demonstrate the intersection of AI and linguistics, revolutionizing the way we approach bilingual dictionary construction. However, machine translation methods utilizing machine learning techniques typically rely heavily on a significant volume of parallel bilingual corpora for training, especially in the context of deep learning models³. This poses a substantial challenge, particularly for low-resource languages such as Bahnaric, where

obtaining such parallel language data is notably difficult. Recent research proposes the construction of a lexical mapping between the source and target languages without the necessity for extensive parallel corpora. This is achieved by learning the mapping between language embedding spaces with the aid of selected anchor words. These anchor words can be automatically extracted or manually designated by linguistic specialists. Figure 1 illustrates the approach at a theoretical level. It begins with two language embedding spaces, one for English and the other for French, each with arbitrary shapes. The mapping process endeavors to convert the embedding space of the source language into that of the target language. Subsequently, adjustments are made to minimize the disparity between the shapes of these two spaces. To isolate the problem of finding the mapping, current state-of-the-art (SOTA) approach⁴ presupposes that the two languages under consideration possess analogous structures. Consequently, after training two distinct embedding models, their embedding point cloud shapes are similar⁵. With this assumption, Generative Adversarial Networks (GANs) are then employed to compute the linear mapping matrix R

Cite this article : Huy L C, Minh L Q, Oanh T N, Dong L D, Nguyen D Q, Sang N T, Quan T, Quan T. **Low-Rank Adaptation Approach for Vietnamese-Bahnaric Lexical Mapping from Non-Parallel Corpora.** *Sci. Tech. Dev. J. – Engineering and Technology* 2024; (1):1-13.

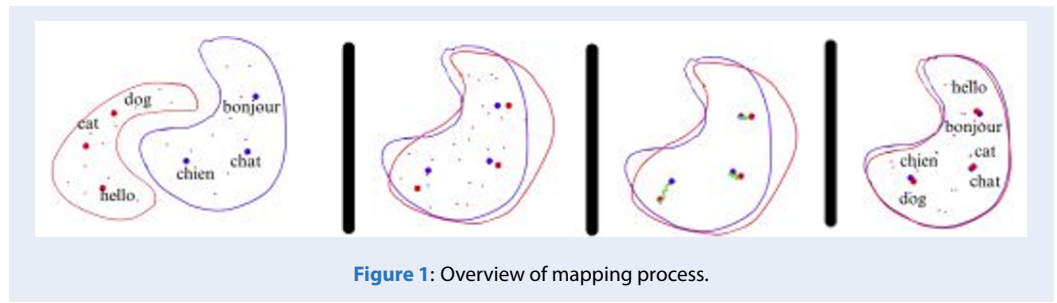


Figure 1: Overview of mapping process.

48 $\in \mathbb{R}^{n \times n}$. During the refinement phase, this method
 49 constructs a synthetic bilingual dictionary contain-
 50 ing only high-frequency words, serving as anchors
 51 to compute the refined mapping matrix $\mathbf{R} \in \mathbb{R}^{n \times n}$.
 52 However, this method exhibits three primary disad-
 53 vantages, both theoretically and practically. From
 54 a theoretical standpoint, assuming similar embed-
 55 ding point cloud shapes and according to the geo-
 56 metric transformation theories^{6,7}, the transformation
 57 a between point clouds must operate within the n -
 58 dimensional special Euclidean group ($SE(n)$ group)⁸,
 59 $a \in SE(n)$. Additionally, based on the theory of special
 60 Euclidean group,

$$SE(n) = T(n) \times SO(n) \quad (1)$$

61 Without any enforcement, $\mathbf{R}, \mathbf{R}' \in O(n)$, leading to
 62 embedding points of corresponding words in two lan-
 63 guages failing to align after transformation (Figure 2).
 64 This stems from the group $O(n)$ containing reflection
 65 and omitted translation actions within the group $T(n)$.
 66 From a practical perspective, constructing bilingual
 67 dictionaries with less than 100 words in low-resource
 68 languages is conceivable⁹, rendering automatic iden-
 69 tification of anchor words unnecessary in general use-
 70 cases. In certain instances, should the automatically
 71 detected anchors deviate from the correct mapping,
 72 the resultant computation of the transformation may
 73 yield incorrect or erroneous results, as illustrated in
 74 Figure 3. Additionally, the adversarial training pro-
 75 cess in GANs may be unstable¹⁰, resulting in poten-
 76 tial model collapse.

77 Another challenge associated with low-resource lan-
 78 guages is the scarcity of available documents. Without
 79 sufficient data, deep learning-based embedding mod-
 80 els are not well learned, which may contradict our as-
 81 sumption. To mitigate this, without the need of par-
 82 allel corpora, data augmentation, via modern tech-
 83 niques, can foster robust embedding models without
 84 any further data collection costs¹¹.

85 In this study, we propose an effective method known
 86 as Augmenting and Sampling with Kabsch (ASK) to

87 address the data scarcity in low-resource languages
 88 and the aforementioned issues of the SOTA approach.
 89 By augmenting the available low-resource language
 90 data and utilizing the Kabsch algorithm¹² to fine-tune
 91 embedding models with randomly sampled anchor
 92 words, we create the transformation $\alpha \in SE(n)$ to map
 93 the source embedding space to the target one. Our
 94 contributions are outlined as follows.

- Implementation of contemporary data augmen-
 95 tation techniques, including sentence boundary
 96 augmentation and multitask learning data aug-
 97 mentation, to enhance low- resource language
 98 data, thus improving the performance of em-
 99 bedding model. 100
- Adaptation of the Kabsch algorithm with ran-
 101 domly sampled anchors to fine-tune and com-
 102 pute the mapping of two language embedding
 103 spaces. 104
- Execution of experiments to assess the efficacy
 105 of our proposed method across various set-
 106 tings, including the well-known French-English
 107 dictionary and the low- resource Vietnamese-
 108 Bahnaric dictionary, underlines the importance
 109 of data augmentation and demonstrates the cor-
 110 rectness and efficiency of our approach. 111

112 RELATED WORKS

113 A. Similarity between embedding spaces 114 across languages

115 Recent advancements in the field of language repre-
 116 sentation have unveiled compelling insights into the
 117 structural similarities that exist across various lan-
 118 guages. A study by¹³⁻¹⁵ reveals that languages sharing
 119 a similar grammatical structure tend to exhibit cor-
 120 responding shapes within their embedding point clouds
 121 when analyzed using identical embedding models.
 122 This congruence between different language spaces is
 123 not merely coincidental but is likely indicative of un-
 124 derlying linguistic parallels that manifest in the syn-
 125 tactic and semantic dimensions of the languages. The

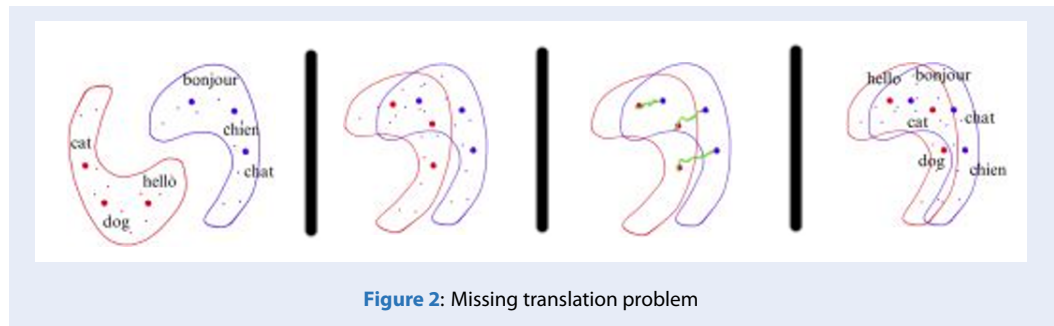


Figure 2: Missing translation problem

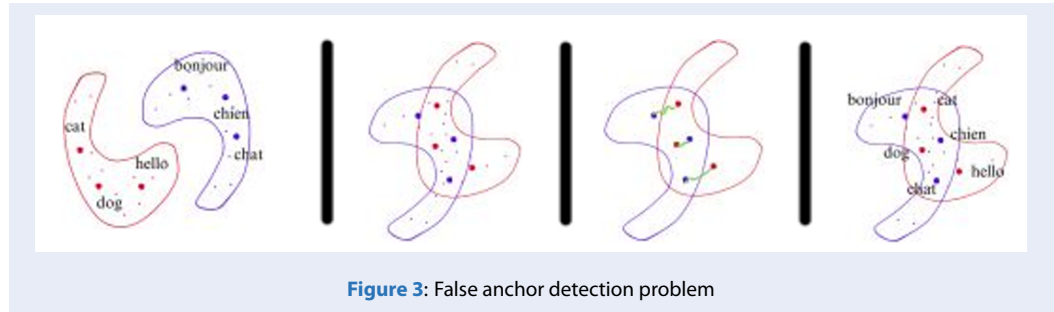


Figure 3: False anchor detection problem

126 discovery has profound implications for cross-lingual
 127 modeling and machine translation, as it could lead to
 128 more efficient algorithms for mapping between differ-
 129 ent language spaces¹⁵. However, the correctness of
 130 an embedding model strongly depends on the train-
 131 ing dataset. In case the two languages have analogous
 132 structures, if one of them does not have rich and
 133 diverse data, their embedding point clouds could be signif-
 134 icantly different.

135 B. Lexical mapping for low-resource lan- 136 guages

137 Lexical mapping, the computational process of align-
 138 ing words or phrases across different languages, rep-
 139 resents an active area of research with critical impli-
 140 cations for the creation of bilingual dictionaries, espe-
 141 cially for low-resource languages such as those spoken
 142 by ethnic minority groups. This research is essential
 143 for the enhancement of machine translation systems
 144 that rely on these dictionaries. Lexical mapping solu-
 145 tions can be broadly divided into three categories: (i)
 146 methods requiring parallel data; (ii) methods neces-
 147 sitating only a few parallel anchors; and (iii) methods
 148 operating with non-parallel data.

149 Approaches utilizing parallel data typically exhibit su-
 150 perior performance, with techniques ranging from
 151 the normalization and application of orthogonal
 152 mapping for translation¹⁶ to the development of ex-
 153 tensive multilingual word embeddings¹⁷. However,

obtaining sufficient parallel data for low- resource lan- 154
 guages remains a significant challenge, limiting the ef- 155
 fective deployment of deep learning-based methods 156
 in practical applications. 157

In response to this limitation, research has explored 158
 solutions that do not require parallel data. A recent 159
 example involves the utilization of adversarial training 160
 to automatically identify anchor words, which 161
 are then used to compute transformations between 162
 embedding spaces⁴. Though this approach circum- 163
 vents the need for parallel corpora and achieves SOTA 164
 performance among non-parallel data approaches, its 165
 performance remains markedly below that of meth- 166
 ods relying on parallel corpora. 167

It is worth noting that the construction of a small 168
 bilingual dictionary is often feasible, making meth- 169
 ods that use such dictionaries as anchors particularly 170
 promising. These approaches are designed to strike a 171
 balance between data requirements and methodolog- 172
 ical performance, addressing a critical tradeoff in the 173
 quest to automate the process of bilingual dictionary 174
 creation and enhance machine translation capabili- 175
 ties. 176

177 C. Rigid transformation and Special Eu- 178 clidean Group

A rigid transformation, also known as a Euclidean 179
 transformation or isometry¹⁸, is a geometric trans- 180
 formation that preserves distance between every pair 181

of points. In more formal terms, a transformation α is considered rigid if for any two points A and B , the distance between A and B is the same as the distance between a $\alpha(A)$ and $\alpha(B)$. The Euclidean group¹⁹, denoted as $E(n)$, is the group of all Euclidean transformations in n -dimensional Euclidean space. It is a mathematical structure that encodes the geometry of Euclidean space and captures the ways objects can be moved around without changing their shape or size. Transformations in $E(n)$ group can be decomposed into components in two subgroups which are rotation ($O(n)$) and translation ($T(n)$) groups (Equation 2).

$$E(n) = T(n) \times O(n) \tag{2}$$

In linear algebra, transformation in $E(n)$ can be also defined as Equation 3.

$$E(n) = \{A|A = \begin{bmatrix} R & t \\ O_{1 \times n} & 1 \end{bmatrix}, \\ R \in R^{n \times n}, \\ t \in R^n, R^T R = RR^T = I\}$$

Assuming that X is a point in a n -dimensional Euclidean space, the transformation a can be expressed as

$$\alpha(x) = R_x + t \tag{4}$$

However, in $(n > 2)$ -dimensional spaces, the transformation can include reflections, which is unnecessary in some usecases such as moving aerospace rocket in spaces. Therefore, theoretically, we do have a subgroup known as special Euclidean group ($SE(n)$) which includes only the isometries that preserve orientation. This means it consists of translations and rotations, but excludes reflections. The term “special” in the name refers to the preservation of orientation. Formal definition of $SE(n)$ in linear algebra is illustrated in (5).

$$E(n) = \{A|A = \begin{bmatrix} R & t \\ O_{1 \times n} & 1 \end{bmatrix}, \\ R \in R^{n \times n}, \\ t \in R^n, R^T R = RR^T = I, |R| = 1\} \tag{5}$$

In $SE(n)$ group, the movement of a rigid body B in Figure 4 can be explained by reference frame $\{A\}$ by creating another reference frame $\{B\}$ on B and describing the position and direction of B in relation to A using a homogeneous transformation matrix¹⁹.

$$A_{A_B} \begin{bmatrix} A_{R_B} & A_{t_{O'}} \\ O_{1 \times n} & 1 \end{bmatrix} \tag{6}$$

where $A_{t_{O'}}$ is the translation vector of the origin O' of $\{B\}$ in the reference frame $\{A\}$, and A_{R_B} is a rotation

matrix that transforms the components of vectors in $\{B\}$ into components in $\{A\}$. Figure 4 presents an example of transformation from

B to A which can be written as $A_{t_{O'}} = A_{R_B} t_{O'} + A_{t_{O'}}$ in 3-dimensional Euclidean space. Moreover, the composition of two displacements, from $\{A\}$ to $\{B\}$, and from $\{B\}$ to $\{C\}$, is equal to the matrix multiplication of A_{A_B} and B_{A_C} . Equation 7 illustrates the decomposition of the transformation $\{C\}$ to $\{A\}$ into two sub-transformations $\{C\}$ to $\{B\}$ and $\{B\}$ to $\{A\}$.

$$A_{A_C} = \begin{bmatrix} A_{R_C} & A_{t_{O'}} \\ O_{1 \times n} & 1 \end{bmatrix} \\ = \begin{bmatrix} A_{R_B} & A_{t_{O'}} \\ O_{1 \times n} & 1 \end{bmatrix} \times \begin{bmatrix} B_{R_C} & B_{t_{O'}} \\ O_{1 \times n} & 1 \end{bmatrix} \\ = \begin{bmatrix} A_{R_B} \times B_{R_C} & A_{R_B} \times B_{t_{O'}} + A_{t_{O'}} \\ O_{1 \times n} & 1 \end{bmatrix} \tag{7}$$

It is evident from (7) that the transformation is reversible, meaning we can aggregate multiple transformations into one. Due to this property, assuming that the transformation AAB consists of a single rotation followed by a single translation, then $\exists^A A'B \in SE(n) \Rightarrow^A A'B = A_{A_B}$.

METHODOLOGY

A. Overview of pipeline

Assume the task at hand is to identify the lexical mapping between two languages: a low-resource language and another language with a grammatical structure that exhibits similarity. In this context, the proposed method, referred to as ASK, functions as a comprehensive, end-to-end pipeline designed specifically to discover the mapping between the embedding spaces of the two languages. The ASK method is articulated into two primary phases, detailed as follows.

- 1. Embedding Model Construction:** The initial phase involves constructing a unique embedding model for each language. For the low-resource language, two specific data augmentation techniques are employed to enhance the modeling process: Sentence Boundary Augmentation (SB)²⁰ and Multitask Learning Data Augmentation (MD)²¹. These techniques aim to improve the representational capacity of the embeddings, especially when dealing with limited data availability.
- 2. Fine-tuning and Mapping Computation:** In the subsequent phase, the focus shifts to fine-tuning embedding models and computing the mapping between the embedding spaces of the

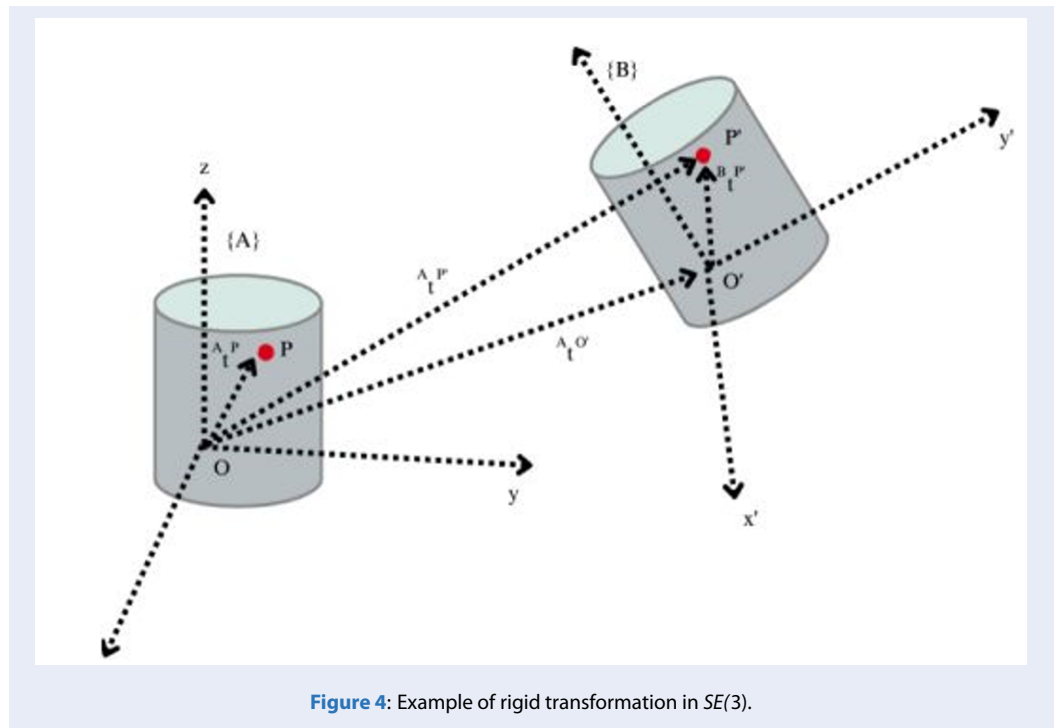


Figure 4: Example of rigid transformation in SE(3).

259 two languages. A set of parallel words is randomly sampled from the collected bilingual dictionary and designated as anchor points. Utilizing the Kabsch algorithm, we fine-tune two embedding models for anchors to be aligned. Then, these anchors are employed to calculate the n-dimensional rigid transformation between the embedding spaces. This rigorous approach leverages the intrinsic geometric properties of the data, ensuring an accurate alignment of the linguistic structures.

270 **B. Embedding model construction**

271 In this study, we applied two below techniques to deal with data shortage of low-resource languages.

- 273 1. **Sentence Boundary Augmentation** is a noise-based approach at the sentence level. By truncating parts of sentences and then combining them, it can remove context from the first sentence, add context from the second sentence, and combine them into a single training example. The proportion of the sentences is governed by a hyperparameter.²⁰
- 274
- 275
- 276
- 277
- 278
- 279
- 280
- 281 2. **Multitask Learning Data Augmentation** combines a set of simple data augmentation methods including Word Swap, Reverse, Semantic Embedding²², Exploratory Data Analysis (EDA)²³ to produce synthetic sentences.
- 282
- 283
- 284
- 285

By adding noise to the text in this way, the embedding model can learn different embeddings for words based on the combination of sentences. These generated sentences along with the original ones are then used as the training data for learning monolingual embedding model^{24,25}.

292 **C. Fine-tuning and mapping computation with Kabsch algorithm**

293 Firstly, we denote the real mapping between two languages as $f^*(\cdot)$ and the set of anchor words of these languages as $W_A = \{w_i^A\}_{i=1}^N$ and $W_B = \{w_i^B\}_{i=1}^N$ where $w_i^A = f^*(w_i^B)$. Considering the original embedding models for two languages are M_A and M_B . We add linear transformations to the end of each model, thus, the embedding model should become M_A^θ, M_B^γ where θ and γ are learnable parameters. Then the vector sets of anchor words can be expressed as (8).

$$\begin{aligned}
 X^\gamma &= \{x_i = M_B^\gamma(w_i^B) \in \mathbb{R}^n\}_{i=1}^N \\
 Y^\theta &= \{y_i = M_A^\theta(w_i^A) \in \mathbb{R}^n\}_{i=1}^N
 \end{aligned}
 \tag{8}$$

In this study, we treat the problem of finding mapping between two embedding spaces as Procrustes superimposition problem²⁶. Therefore, we utilize the Kabsch algorithm to find the mapping or the transformation between two embedding point cloud, mathematically speaking. The objective of Kabsch algorithm is

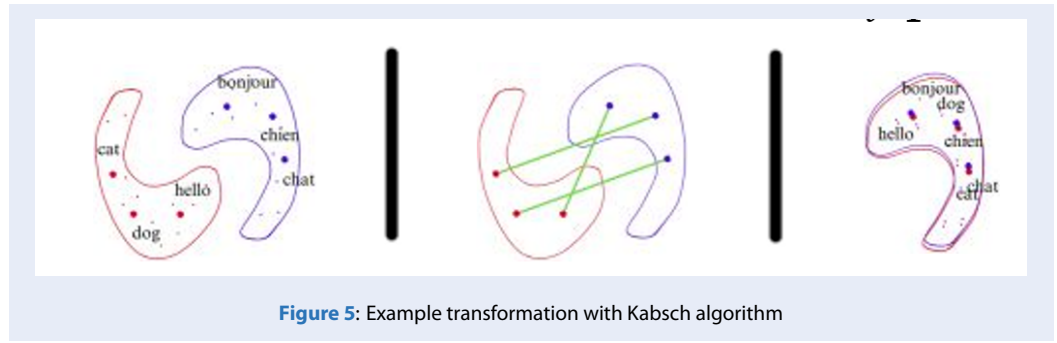


Figure 5: Example transformation with Kabsch algorithm

310 computing an approximation $f(\cdot)$ of the mapping $f^*(\cdot)$
 311 to optimize the objective function in (9).

$$f = \underset{f}{\operatorname{argmin}} E_{\substack{X \sim B \\ Y \sim A}} [\|f(X) - Y\|^2] \quad (9)$$

312 However, we can not directly optimize (9), so that we
 313 reparameterize it with θ and γ . The new objective
 314 function is then become (10). This objective func-
 315 tion is also the loss function for fine-tuning embed-
 316 ding models.

$$L = \underset{\theta, \gamma}{\operatorname{argmin}} E_{\substack{X^\gamma \sim B \\ Y^\theta \sim A}} [\|f(X^\gamma) - Y^\theta\|^2] \quad (10)$$

317 Base on the theory of $SE(n)$ group, the $f(\cdot)$ repre-
 318 sents an affine linear function: $R^n \rightarrow R^n$, which corre-
 319 sponds to a rigid motion in R^n . Under the perspective
 320 of linear algebra, $f(x)=Rx+t$ with $X \in R^n$, where $R \in$
 321 $R^{n \times n}$, $|R|=1$, and $t \in R^n$. Nextly, we denote the cen-
 322 troid if point cloud X and Y in Equation 11.

$$\begin{aligned} \mu_X &= \frac{1}{N} \sum_{x_i \in X} x_i \\ \mu_Y &= \frac{1}{N} \sum_{y_i \in Y} y_i \end{aligned} \quad (11)$$

323 The Kabsch algorithm is summarized in Table 1. Fig-
 324 ure 5 illustrates the transformation with Kabsch algo-
 325 rithm.

326 After the embedding models are fine-tuning, we cal-
 327 culate the approximate mapping function using the
 328 same procedure. Consequently, the process of iden-
 329 tifying the mapping of a source language word in the
 330 target language involves ranking the neighboring em-
 331 bedding points based on cosine similarity. Cosine
 332 similarity is a widely used metric in natural language
 333 processing that measures the similarity between two
 334 vectors in a high-dimensional space. By employing
 335 this approach, we can effectively determine the closest
 336 matching target language word or its nearest neigh-
 337 bors in the embedding space.

Next, we present the proof of better performance of
 the Kabsch algorithm in n -dimensional space in com-
 parison to the original Procrustes problem and the
 SOTA approach.

a) *Ensuring rigid transformation*: Assuming that the
 objective of Procrustes problem is hold, denoted as
 (12).

$$\begin{aligned} g &= \underset{g}{\operatorname{argmin}} E_{\substack{X \sim A \\ Y \sim B}} [\|g(X) - Y\|^2], \quad g \in O(n) \\ &= \underset{g}{\operatorname{argmin}} E_{\substack{X \sim A \\ Y \sim B}} \left[\operatorname{tr} \left((RX - Y)^T (RX - Y) \right) \right], \quad g \in O(n) \\ &= \underset{g}{\operatorname{argmin}} E_{\substack{X \sim A \\ Y \sim B}} \left[\operatorname{tr} (X^T X) + \operatorname{tr} (Y^T Y) - 2 \operatorname{tr} (Y^T R X) \right] \\ &= \underset{g}{\operatorname{argmin}} E_{\substack{X \sim A \\ Y \sim B}} \left[\operatorname{tr} (Y^T R X) \right] \quad (12) \end{aligned}$$

Let $C = XY^T = U\Sigma V^T$, since $V^T R U$ is orthogonal,
 then

$$\begin{aligned} \operatorname{tr}(RC) &= \operatorname{tr}(R U \Sigma V^T) \\ &= \operatorname{tr}(V^T R U \Sigma) \leq \operatorname{tr}(\Sigma) = \sum_{j=1}^n \sigma_j \end{aligned} \quad (13)$$

The euqation holds if $R = V U^T$ and $|V U^T| > 0$. How-
 ever, in case $|V U^T| < 0$, the (13) becomes (14).

$$\begin{aligned} \operatorname{tr}(RC) &= \operatorname{tr}(R U \Sigma V^T) \\ &= \operatorname{tr}(V^T R U \Sigma) \leq \sum_{j=1}^n (\sigma_j - \sigma_n) \end{aligned} \quad (14)$$

If we keep $|R| = |V U^T|$, we still achieve the equality
 but $|R| = -1$ which causes the reflections in the origi-
 nal point cloud, which is not what we expect since
 we assume that the two sets of point cloud have the
 same shape. The Kabsch algorithm resolves this issue
 and get $g \in SO(n)$ by choosing $R = V \Sigma' U^T$ where
 $\Sigma' = \{\sigma_{i < n} = 1, \sigma_n = -1\}$.

b) *Tackling translation in high-dimensional space*: As-
 suming that we already solve the original Procrustes
 problem and get the mapping function $g(\cdot)$, we define
 our mapping function $f(\cdot)$ as (15).

$$f(X) = g(X) - g(\mu_X) + \mu_Y \quad (15)$$

Table 1: Kabsch algorithm

Algorithm 1 Kabsch Algorithm
Input: Point cloud set $X, Y \in \mathbb{R}^n$
Output: $R \in \mathbb{R}^{n \times n}, t \in \mathbb{R}^n$
$C = XY^T$
Perform SVD: $C = U\Sigma V^T$
$\Sigma' = \{\sigma_i\}_{i=1}^n$, where $\sigma_{i < n}$ and $\sigma_n = \text{sign}(VU^T)$
$R = V\Sigma'U^T$
$t = \mu_Y - R\mu_X$
return R, t

Table 2: Number of sentences in Vietnamese and Bahnaric corpora

Dataset	Original	Augmented
Vietnamese	16105	78307
Bahnaric	16105	78307

Table 3: Examples of French-to-English on 10000 anchors

Source Word	Top1	Top2	Top3	Top4
soins	care	deal	fear	attention
fin	end	close	goal	stop
chaque	each	apiece	vice	canso
position	position	place	emplacement	location
accès	access	accession	approach	admission
ouest	west	westward	eastern	easterly
période	period	stop	point	flow
emplois	jobs	job	subcontract	line
impôt	tax	taxes	taxation	assess
rôle	role	persona	character	function

Table 4: Examples of Bahnaric-to-Vietnamese on 500 anchors

Source	Top1	Top2	Top3	Top4
máu	pham	thăm	chăn	tâng_kojung
sũa	Đak_toh	dư_dư	bek_bỏ	hla_piết_yêr
yên	an	krũ	kopung	areh
trôi	đơng	pơdrăn	prah	kônăr
gì	kio	kỏjong	totuanh	bok_y
VỔ	pơchah	brôm	apăl_asol	kokóch
bay	apăl	srang	búp_búp	long_wắk
tỏa	toprah	chă_hming	hla_piết_yêr	bluh_lêch
thiếu	bĩ_mah	mớng_kotang	ping_ngil	hmingji
công	kowơng	dử_dử	yêr_tomông	nguk_ich

360 Considering the difference between original solution
 361 and Kab- sch algorithm as in (16), we observe that
 362 when $g(\mu_X) \neq \mu_Y$, the Kabsch algorithm, that takes
 363 translation into account will be convergence to the
 364 maxima while the original one can not.

$$\begin{aligned} \Delta &= \|g(X) - Y\|^2 - \|f(X) - Y\|^2 \\ &= \sum_{i=1}^n (Rx_i - y_i) - \sum_{i=1}^n (Rx_i - R\mu_X + \mu_Y - y_i) \\ &= \sum_{i=1}^n (R\mu_X - \mu_Y) \\ &= n \|g(\mu_X) - \mu_Y\|^2 \geq 0 \quad (16) \end{aligned}$$

365 EXPERIMENTS

366 In this section, we conduct a comprehensive com-
 367 parison of our proposed approach with other base-
 368 line methods across various benchmarks. Our ex-
 369 perimental analysis consists of two distinct phases.
 370 Firstly, we concentrate on well-resourced language
 371 pairs, particularly French-English, to showcase the ef-
 372 fectiveness and efficiency of our method. Secondly,
 373 we extend our evaluation to the Vietnamese-Bahnaric
 374 language pair, strategically chosen to assess and verify
 375 our method performance in a setting with limited lin-
 376 guistic resources. This two-phase evaluation enables a
 377 robust examination of the generalizability and adapt-
 378 ability of our approach across different language sce-
 379 narios, contributing to a deeper understanding of its
 380 capabilities and limitations.

381 A. Experimental setups

382 Toward experiments on rich-resource datasets,
 383 French- English, we uses a French-English corpus
 384 containing 53,241 words. We will trã embeddings
 385 with three options:

- 386 1. 1,000 anchor words along with 52,241 test
 387 words.
- 388 2. 10,000 anchor words along with 43,241 test
 389 words.
- 390 3. 50,000 anchor words along with 3,241 test
 391 words.

392 For a fair model comparison, we use the rich-resource
 393 dataset without augmentation. Synonyms of English
 394 words are found using WordNet from Princeton Uni-
 395 versity²⁷ and implemented by NLTK^a for evaluation.
 396 Furthermore, we will assess the impact of data aug-
 397 mentation on our low-resource datasets through two
 398 different tests:

- 399 1. Evaluation using the original datasets.

^a<https://www.nltk.org/>

2. Evaluation using augmented data from the orig- 400
 401 inal dataset, which includes sentences with sen-
 402 tence boundaries, EDA, and semantic embed-
 403 ding augmentation combined with the original
 404 datasets.

The dataset information, comprising both the original 405
 data and its augmented counterpart, is provided in Ta- 406
 ble 2. The original dataset is represented in the ‘Orig- 407
 inal’ column, while the augmented dataset is found in 408
 the ‘DA’ column. 409

The embeddings will be trained with three options: 410

1. 100 anchor words along with the rest being test 411
 412 words.
2. 500 anchor words along with the rest being test 413
 414 words.
3. 1000 anchor words along with the rest being 415
 416 test words. During training, ASK utilizes Singu-
 417 lar Value Decomposition (SVD) for learning the
 418 mapping, and no hyperparameters are required. 419
 420 However, the word embeddings also play a criti-
 421 cal role. After conducting multiple experiments,
 422 we selected the Skip-gram model to learn the
 423 word embeddings with the following settings: 424
 425 the hidden dimension is 100, the window size
 is 5, and words whose frequency less than 2 are
 ignored.

We have employed two commonly used metrics 426
 which are listed in the followings to evaluate the rank- 427
 ing performance of our model. 428

1. Mean Reciprocal Rank (MRR): This metric in- 429
 430 corporates synonyms in addition to exact word
 431 matching. By considering synonyms, we obtain
 432 a more comprehensive evaluation of the map-
 433 ping quality. To evaluate the model, we compute
 434 the mean MRR across all testing words.
2. Top-K accuracy (Top-KAcc): This metric eval- 435
 436 uates the model performance by examining the
 437 Top-A ranked results and assessing the position
 438 of the correct word.
3. Runtime: This metric quantifies the elapsed 439
 440 time taken by the model to identify the mapping
 441 function responsible for translating source lan-
 442 guage words to their corresponding target lan-
 443 guage words.

To improve performance on low-resource datasets, 444
 we employ a fine-tuning strategy. Our model con- 445
 sists of three linear layers that project the original em- 446
 beddings into a shared space, ensuring that both the 447
 source and target mapped embeddings have the same 448

Table 5: The comparison between Kabsch and the other supervised models on French-English

Method	Top-1Acc(%)↑	Top-5Acc(%)↑	Top-10Acc(%)↑	MRR↑	Runtime(ms) ↓
1000 anchor words					
Arttxem	3.678 ± 0.289	7.094 ± 0.419	8.842 ± 0.461	0.05478 ± 0.0034	3819.0170 ± 99.1973
Dino	1.386 ± 0.198	3.353 ± 0.387	4.595 ± 0.474	0.02542 ± 0.00285	7.7471 ± 4.1768
Mikolov	1.388 ± 0.196	3.342 ± 0.385	4.584 ± 0.462	0.02537 ± 0.00283	2535.6096 ± 57.1759
Kabsch	3.984 ± 0.251	7.41 ± 0.375	9.066 ± 0.406	0.05779 ± 0.00301	1.3448 ± 0.2295
ASK	19.14 ± 0.217	25.32 ± 0.323	27.13 ± 0.354	0.2056 ± 0.001	1.4288 ± 0.2135
10000 anchor words					
Arttxem	7.812 ± 0.194	11.926 ± 0.295	13.679 ± 0.329	0.09909 ± 0.00238	3972.3056 ± 199.4734
Dino	1.886 ± 0.122	4.25 ± 0.195	5.674 ± 0.227	0.03233 ± 0.00157	25.0557 ± 0.5525
Mikolov	1.887 ± 0.109	4.256 ± 0.175	5.686 ± 0.198	0.03239 ± 0.00136	2524.3594 ± 23.2422
Kabsch	9.088 ± 0.054	13.547 ± 0.069	15.438 ± 0.082	0.1135 ± 0.0006	2.7956 ± 0.2618
ASK	46.25 ± 0.032	53.19 ± 0.035	55.5 ± 0.042	0.4787 ± 0.0014	3.2143 ± 0.0.1538
50000 anchor words					
Arttxem	10.361 ± 0.508	15.369 ± 0.592	17.603 ± 0.431	0.1294 ± 0.00481	4538.6975 ± 53.7677
Dino	1.867 ± 0.195	4.141 ± 0.307	5.564 ± 0.302	0.03185 ± 0.0016	259.4690 ± 2.9341
Mikolov	1.922 ± 0.179	4.172 ± 0.306	5.5659 ± 0.288	0.03209 ± 0.00156	13438.7266 ± 93.4700
Kabsch	9.719 ± 0.414	14.07 ± 0.476	15.89 ± 0.51	0.11926 ± 0.0043	9.6051 ± 0.9784
ASK	61.71 ± 0.396	66.34 ± 0.413	69.423 ± 0.442	0.6312 ± 0.0036	10.1524 ± 0.1.1226

shape. We use hidden state dimensions are set to 1024 and 2048 and activate these layers using Relu and Tanh functions, as they yielded the best results during experimentation. The training process maintains a constant learning rate of 10^{-3} across dataset sizes (100, 500, 1000) but extends the number of epochs (20000, 40000, 80000) for enhanced optimization. Our chosen optimization method is Stochastic Gradient Descent (SGD).

B. Baselines

The study of Mikolov¹³ utilizes skip-gram word embedding to learn high-quality word embeddings, optimizing for a rotation matrix that minimizes the loss function $\sum_{i=1}^n \|Wx_i - z_i\|^2$. By employing gradient descent, they find optimal values for the matrix w , enabling seamless mapping between the word spaces of source and target languages without constraints. The authors then identify the target language word with the highest cosine similarity to z , establishing meaningful associations between words in different languages for crosslingual tasks like translation and word alignment.

The Mikolov model¹³ lacks constraints, which may lead to overfitting and underutilization of word embedding features. To address this, the Dinu model²⁸ introduces regularization to prevent specific words from being consistently mapped to particular targets. Additionally, they modify the method for selecting the correct word after mapping the source language word using the matrix w . This change is necessary because cosine similarity, commonly used for this task, encounters the Hubness problem—an inherent challenge in high-dimensional spaces²⁹ and a recognized issue for word-based vectors²⁹. As a result, the focus lies on proposing a straightforward and efficient solution to handle this problem by adjusting the similarity matrix post-mapping process.

And the last model which we use for comparing our result is Artetxe model³⁰. The Artetxe method is remarkable for its effectiveness even with just 25 word pairs, a departure from previous methods that often require thousands of words for satisfactory performance. They emphasize the adaptability of the approach with low-dimensional pre-trained word embeddings. For inducing bilingual lexicons, a common evaluation task, they use a small train set (seed dictionary) to learn an initial mapping, leading to a larger and potentially enhanced dictionary. In the second step, they train the model to refine the source-to-target language mapping, aiming for improvements over the input dictionary. This iterative process allows for continuous refinement until a convergence criterion is met.

RESULT

A. Evaluations using rich-resource datasets

This experiment assesses the effectiveness of Kabsch algorithm, in finding language mappings between French and English datasets (rich-resource datasets) with similar point cloud shapes. The analysis (Table 5) demonstrates that Kabsch outperforms most other methods when utilizing 1000 and 10,000 anchor points. However, Our ASK model outperforms other methods due to its fine-tuned embedding, which aligns the shapes of the source and target language embeddings.

B. Evaluations using low-resource datasets

In this scenario, we executed full pipeline of ASK including data augmentation, fine-tuning embedding models and computing mapping with Kabsch. We compared our method with its ablated versions and other supervised learning models in terms of Top-K Accuracy and mapping computation runtime for Vietnamese-Bahnaric in the Table 6.

DISCUSSION

In rich-resource dataset, Kabsch consistently achieves favorable results across all cases, maintaining a relatively lower runtime compared to other methods. Kabsch exhibits the lowest runtime among the tested models, making it a promising approach for efficient and accurate language mapping tasks. To showcase the mapping process, we have randomly chosen 10 words, which are presented in Table 3. Each “Top i ” column representing the i th target word with the highest similarity score.

In low-resource dataset, Kabsch algorithm’s result tends to be slightly less impressive compared to alternative models. This can be traced back to the data’s limited scale. Since the dataset is small, it might fail to meet the criteria for the embedding shapes to match exactly, resulting in a decline in accuracy. However, by implementing Finetuning on Kabsch. It’s important to highlight that Kabsch’s runtime has been notably performer in terms of execution speed.

Following the application of various augmentation techniques, such as sentence boundary augmentation, EDA, and word2vec, on the initial dataset, we significantly expanded its size. Consequently, we observed a substantial improvement in performance compared to evaluating the model on the original low-resource dataset. This enhancement stems from the model’s enhanced capability to learn the underlying distribution of the data. Notably, our proposed method achieves higher Top-1 Accuracy and MRR

Table 6: The comparison between our method, its ablated versions (with fine-tuning (FT) and data augmentation (DA)) and the other supervised models on Vietnamese-Bahnaric

Method	Top-1Acc(%)↑	Top-5 Acc(%)↑	Top-10Acc(%)↑	MRR↑	Time(ms) ↓
100 anchor words					
Artetxem	0.8 ±1.5	1.2 ±1.7	1.5 ±1.7	0.012 ±0.016	0.641 ±0.122
Dino	0.1 ±0.1	0.6 ±0.2	1.3 ±0.3	0.008 ± 0.001	0.015 ± 0.003
Mikolov	4.9 ±0.1	5.3 ±0.1	5.5 ±0.2	0.054 ± 0.001	2.450 ±0.122
Kabsch	1.3 ±0.1	2.2 ±0.2	3.1 ±0.3	0.022 ± 0.001	0.001 ± 0.0003
Kabsch + FT	5.0 ±0.1	5.3 ±0.1	5.6 ±0.3	0.054 ± 0.001	0.0033 ± 0.0001
Kabsch + DA	2.9 ±0.2	3.8 ±0.3	4.4 ±0.5	0.037 ± 0.003	0.001 ± 0.0003
ASK	6.0 ±0.04	6.1 ±0.05	6.4 ±0.1	0.063 ± 0.0004	0.0035 ± 0.016
500 anchor words					
Artetxem	9.9 ±0.5	13.3 ±0.5	15.2 ±0.5	0.119 ±0.004	0.679 ±0.127
Dino	0.2 ±0.1	1.0 ±0.3	2.0 ±0.4	0.012 ±0.001	0.015 ± 0.004
Mikolov	6.0 ±0.3	7.5 ±0.3	8.4 ±0.3	0.071 ± 0.002	2.567 ±0.054
Kabsch	3.2 ±0.5	4.8 ±0.6	6.0 ±0.9	0.044 ± 0.005	0.001 ± 0.0001
Kabsch + FT	30.4 ±0.1	30.6 ±0.2	30.8 ±0.2	0.306 ± 0.001	0.0037 ± 0.0001
Kabsch + DA	8.3 ±0.4	10.4 ±0.6	11.8 ±0.8	0.097 ± 0.004	0.001 ± 0.0001
ASK	33.3 ±0.1	33.5 ±0.1	33.7 ±0.1	0.336 ±0.001	0.0038 ± 0.0001
1000 anchor words					
Artetxem	11.9 ±0.8	16.7 ±0.7	19.0 ±0.6	0.145 ±0.007	0.685 ±0.164
Dino	0.2 ±0.1	0.99 ±0.3	1.9 ±0.4	0.012 ± 0.002	0.017 ±0.003
Mikolov	8.2 ±0.8	9.5 ±0.6	10.5 ±0.7	0.093 ± 0.007	2.540 ± 0.028
Kabsch	4.6 ±0.5	6.7 ±0.6	8.1 ±0.5	0.061 ± 0.005	0.001 ± 0.0001
Kabsch + FT	59.1 ±0.8	60.0 ±0.3	60.4 ±0.3	0.597 ± 0.005	0.004 ± 0.0001
Kabsch + DA	11.1 ±0.6	14.1 ±0.8	16.2 ±1.0	0.131 ±0.007	0.001 ± 0.0002
ASK	64.9 ±0.2	65.0 ±0.1	65.0 ±0.1	0.650 ±0.001	0.004 ± 0.00035

552 scores compared to alternative approaches. This ob-
 553 servation underscores the advantage of employing a
 554 larger dataset and highlights the fulfillment of the
 555 underlying assumption, contributing to the superior
 556 performance of our approach over other methods
 557 while maintaining a lower runtime. Additionally, we
 558 randomly selected 10 Vietnamese words to illustrate
 559 the mapping process. These words are presented in
 560 Table 4, with the same column meanings as in Table 3.
 561

562 CONCLUSION

563 This paper introduces a novel approach for word
 564 alignment based on distribution representations.

Leveraging two monolingual language corpora and
 an initial dictionary, our method effectively learns
 a meaningful transformation for individual words.
 The experimental results reveal the efficacy of our
 approach on rich-resource datasets, exhibiting supe-
 rior training time compared to alternative methods.
 Additionally, promising performance is observed on
 low-resource datasets, highlighting the potential for
 broader applicability.

In the future, we intend to conduct further investi-
 gations in this direction, aiming to refine and opti-
 mize our method to ensure a more coherent shape for
 word embeddings from two monolingual language
 corpora. This enhancement will facilitate more effi-

579 cient alignment between the corpora, ultimately lead-
 580 ing to improved alignment accuracy and precision.
 581 Our ongoing research aims to enhance the practical-
 582 ity and versatility of our approach, enabling cross-
 583 lingual language processing and effective multilingual
 584 resource alignment.

585 ACKNOWLEDGMENT

586 This research is funded by Ministry of Science and
 587 Technology (MOST) within the framework of the
 588 Program "Supporting research, development and
 589 technology application of Industry 4.0" KC-4.0/19-25
 590 - Project "Development of a Vietnamese- Bahnaric
 591 machine translation and Bahnaric text- to-speech sys-
 592 tem (all dialects)" - KC-4.0-29/19-25.

593 LIST OF ABBREVIATION

- 594 **AI:** Artificial Intelligence
- 595 **GANs:** Generative Adversarial Networks
- 596 **EDA:** Exploratory Data Analysis
- 597 **ASK:** Augmenting and Sampling with Kabsch
- 598 **MRR:** Mean Reciprocal Rank
- 599 **SVD:** Singular Value Decomposition
- 600 **SOTA:** state-of-the-art

601 CONFLICT OF INTEREST

602 The authors hereby declare that there is no conflict of
 603 interest in the publication of this article.

604 AUTHORS' CONTRIBUTION

- 605 • La Cam Huy: Gathering data in English and
 606 French, performing preprocessing on data in
 607 English, French, Vietnamese, and Bahnaric lan-
 608 guages, searching for relevant problem-solving
 609 models, constructing models, comparing re-
 610 sults, and writing research papers.
- 611 • Le Quang Minh: Collecting information in En-
 612 glish and French, organizing information in En-
 613 glish, French, Vietnamese, and Bahnaric lan-
 614 guages, finding problem-solving methods that
 615 are related to the topic, and writing research pa-
 616 pers.
- 617 • Tran Ngoe Oanh: Performing preprocessing on
 618 data in English, French, Vietnamese and Bah-
 619 naric languages. Augmenting the dataset and
 620 writing research papers
- 621 • Le Due Dong: Augmenting the dataset, sup-
 622 porting model construction, writing the re-
 623 search paper
- 624 • Due Q. Nguyen: Come up with ideas for writ-
 625 ing articles, collect data in English, French, Viet-
 626 namese and Bahnaric. Testing models, tutorials
 627 and editing paper.

- Nguyen Tan Sang: Participate in the extending
 data for Vietnamese and Bahnaric. 628
- Tran Quan: Participate in coming up writing
 ideas 630
- Tho Quan: Come up with ideas for writing ar-
 ticles, collecting data in Vietnamese, Bahnaric.
 Providing paper tutorials and editing. 632

635 REFERENCES

1. Zhu W, Zhou Z, Huang S, Lin Z, Zhou X, Tu Y, et al. Improv-
 ing bilingual lexicon induction on distant language pairs. In:
 Huang S, Knight K, editors. Machine Translation. Singapore:
 Springer Singapore; 2019. p. 1-10; Available from: https://doi.org/10.1007/978-981-15-1721-1_1. 636
2. Janiesch C, Zschiech P, Heinrich K. Machine learning and
 deep learning. Electronic Markets. 2021 Apr;31(3):685-
 695; Available from: <https://doi.org/10.1007/s12525-021-00475-2>. 637
3. Mondal SK, Zhang H, Kabir HMD, Ni K, Dai H-N. Machine trans-
 lation and its evaluation: a study. Artificial Intelligence Re-
 view. 2023 Feb;56(9):10137-10226; Available from: <https://doi.org/10.1007/s10462-023-10423-5>. 638
4. Lample G, Conneau A, Ranzato M, Denoyer L, Jegou H. Word
 translation without parallel data. In: International Conference
 on Learning Representations. 2018; 639
5. Tang C, Yang X, Wu B, Han Z, Chang Y. Parts2words: Learn-
 ing joint embedding of point clouds and texts by bidirec-
 tional matching between parts and words. In: 2023 IEEE/CVF
 Conference on Computer Vision and Pattern Recognition
 (CVPR). Los Alamitos, CA, USA: IEEE Computer Society; 2023
 Jun. p. 6884-6893; Available from: <https://doi.org/10.1109/CVPR52729.2023.00665>. 640
6. Yaglom IM. Geometric transformations. Washington: Math-
 ematical Association of America; 1962; Available from: <https://doi.org/10.5948/UPO9780883859254>. 641
7. Guggenheimer HW. Plane geometry and its groups.
 Cambridge University Press; 1968. vol. 11, no. 3,
 p. 508-509; Available from: <https://doi.org/10.1017/S0008439500029660>. 642
8. Satorras VG, Hoogeboom E, Welling M. E(n) equivariant graph
 neural networks. In: Proceedings of the 38th International
 Conference on Machine Learning. Proceedings of Machine
 Learning Research. vol. 139. PMLR; 2021 Jul 18-24. p. 9323-
 9332; 643
9. Rubino R, Marie B, Dabre R, Fujita A, Utiyama M, Sumita
 E. Extremely low-resource neural machine translation for
 Asian languages. Machine Translation. 2020 Dec;34(4):347-
 382; Available from: <https://doi.org/10.1007/s10590-020-09258-6>. 644
10. Li Z, Xia P, Tao R, Niu H, Li B. A new perspective on
 stabilizing GANs training: Direct adversarial training. IEEE
 Transactions on Emerging Topics in Computational Intellig-
 ence. 2023;7(1):178-189; Available from: <https://doi.org/10.1109/TETCI.2022.3193373>. 645
11. Wang S, Yang Y, Wu Z, Qian Y, Yu K. Data augmentation using
 deep generative models for embedding based speaker recog-
 nition. IEEE/ACM Transactions on Audio, Speech, and Lan-
 guage Processing. 2020;28:2598-2609; Available from: <https://doi.org/10.1109/TASLP.2020.3016498>. 646
12. Gupta KK, Sen S, Haque R, Ekbal A, Bhattacharyya P, Way A.
 Augmenting training data with syntactic phrasal-segments
 in low-resource neural machine translation. Machine Transla-
 tion. 2021 Dec;35(4):661-685; Available from: <https://doi.org/10.1007/s10590-021-09290-0>. 647

- 691 13. Mikolov T, Le QV, Sutskever I. Exploiting similarities among
692 languages for machine translation. 2013; 762
- 693 14. Zhou C, Ma X, Wang D, Neubig G. Density matching for
694 bilingual word embedding. In: North American Chapter of the
695 Association for Computational Linguistics. 2019;PMID:
696 31090322. Available from: <https://doi.org/10.18653/v1/N19-1161>. 763
- 697 15. Conneau A, Khandelwal K, Goyal N, Chaudhary V, Wenzek
698 G, Guzman F, et al. Unsupervised cross-lingual representa-
699 tion learning at scale. In: Proceedings of the 58th Annual
700 Meeting of the Association for Computational Linguistics. On-
701 line: Association for Computational Linguistics; 2020 Jul. p.
702 8440-8451; Available from: <https://doi.org/10.18653/v1/2020.acl-main.747>. 764
- 703 16. Xing C, Wang D, Liu C, Lin Y. Normalized word embedding
704 and orthogonal transform for bilingual word translation. In:
705 Proceedings of the 2015 Conference of the North American
706 Chapter of the Association for Computational Linguistics:
707 Human Language Technologies. Denver, Colorado: Associa-
708 tion for Computational Linguistics; 2015 May-Jun. p. 1006-
709 1011; Available from: <https://doi.org/10.3115/v1/N15-1104>. 765
- 710 17. Ammar W, Mulcaire G, Tsvetkov Y, Lample G, Dyer C, Smith NA.
711 Massively multilingual word embeddings. 2016; 766
- 712 18. Artin E. Geometric Algebra. John Wiley & Sons; 2011; 767
- 713 19. Coxeter HSM, Greitzer SL. Geometry Revisited. Mathematical
714 Association of America; 2016;
- 715 20. Li D, I T, Arivazhagan N, Cherry C, Padfield D. Sentence
716 boundary augmentation for neural machine translation ro-
717 bustness. In: ICASSP 2021 - 2021 IEEE International Confer-
718 ence on Acoustics, Speech and Signal Processing (ICASSP).
719 2021. p. 7553-7557; Available from: <https://doi.org/10.1109/ICASSP39728.2021.9413492>.
- 720 21. Meyerson E, Miikkulainen R. Pseudo-task augmentation: From
721 deep multitask learning to intratask sharing-and back. In: Pro-
722 ceedings of the 35th International Conference on Machine
723 Learning. 2018. p. 739-748;
- 724 22. Wang WY, Yang D. That's so annoying! 11: A lexical and frame-
725 semantic embedding based data augmentation approach to
726 automatic categorization of annoying behaviors using #pet-
727 peeve tweets. In: Proceedings of the 2015 Conference on Em-
728 pirical Methods in Natural Language Processing. Lisbon, Por-
729 tugal: Association for Computational Linguistics; 2015 Sep. p.
730 2557-2563; Available from: <https://doi.org/10.18653/v1/D15-1306>. 765
- 731 23. Wei JW, Zou K. EDA: easy data augmentation techniques
732 for boosting performance on text classification tasks. CoRR.
733 2019; Available from: <https://doi.org/10.18653/v1/D19-1670>. 766
- 734 24. Voita E, Sennrich R, Titov I. Analyzing the source and target
735 contributions to predictions in neural machine translation.
736 In: Proceedings of the 59th Annual Meeting of the Associa-
737 tion for Computational Linguistics and the 11th International
738 Joint Conference on Natural Language Processing (Volume 1:
739 Long Papers). Online: Association for Computational Linguis-
740 tics; 2021 Aug. p. 1126-1140; Available from: <https://doi.org/10.18653/v1/2021.acl-long.91>. 767
- 741 25. Dong D, Wu H, He W, Yu D, Wang H. Multi-task learning for
742 multiple language translation. In: Proceedings of the 53rd An-
743 nual Meeting of the Association for Computational Linguistics
744 and the 7th International Joint Conference on Natural Lan-
745 guage Processing (Volume 1: Long Papers). 2015. p. 1723-
746 1732; Available from: <https://doi.org/10.3115/v1/P15-1166>. 765
- 747 26. Kendall DG. A survey of the statistical theory of shape. Statisti-
748 cal Science. 1989;4(2):87-99; Available from: <https://doi.org/10.1214/ss/1177012582>. 766
- 749 27. Princeton University. About WordNet. 2010; 767
- 750 28. Dinu G, Baroni M. Improving zero-shot learning by mitigat-
751 ing the hubness problem. In: 3rd International Conference
752 on Learning Representations, ICLR 2015, San Diego, CA, USA,
753 May 7-9, 2015, Workshop Trade Proceedings. 2015; 765
- 754 29. Radovanovic M, Nanopoulos A, Ivanovic M. Hubs in space:
755 Popular nearest neighbors in high-dimensional data. Journal
756 of Machine Learning Research. 2010;11(sept):2487-2531; 762
- 757 30. Artetxe M, Labaka G, Agirre E. Learning bilingual word em-
758 beddings with (almost) no bilingual data. In: Proceedings
759 of the 55th Annual Meeting of the Association for Computa-
760 tional Linguistics (Volume 1: Long Papers). 2017. p. 451-
761 462; Available from: <https://doi.org/10.18653/v1/P17-1042>. 767

Hướng tiếp cận thu giảm số chiều cho phép ánh xạ từ vựng tiếng Việt sang tiếng Ba Na từ các tập ngữ liệu khổng song song

La Cẩm Huy^{1,2}, Lê Quang Minh^{1,2}, Trần Ngọc Oanh^{1,2}, Lê Đức Đồng^{1,2}, Nguyễn Quang Đức^{1,2}, Nguyễn Tấn Sang^{1,2}, Trần Quân^{1,2}, Quân Thành Thơ^{1,2,*}



Use your smartphone to scan this QR code and download this article

¹Trường Đại học Bách khoa, Đại học Quốc gia Thành phố Hồ Chí Minh, 268 Lý Thường Kiệt, Phường 14, Quận 10, Thành phố Hồ Chí Minh, Việt Nam

²Đại học Quốc gia Thành phố Hồ Chí Minh, Phường Linh Trung, Thành phố Thủ Đức, Thành phố Hồ Chí Minh, Việt Nam

Liên hệ

Quân Thành Thơ, Trường Đại học Bách khoa, Đại học Quốc gia Thành phố Hồ Chí Minh, 268 Lý Thường Kiệt, Phường 14, Quận 10, Thành phố Hồ Chí Minh, Việt Nam

Đại học Quốc gia Thành phố Hồ Chí Minh, Phường Linh Trung, Thành phố Thủ Đức, Thành phố Hồ Chí Minh, Việt Nam

Email: qtttho@hcmut.edu.vn

Lịch sử

- Ngày nhận: 7-9-2023
- Ngày chấp nhận: 26-4-2024
- Ngày đăng:

DOI:



Bản quyền

© ĐHQG Tp.HCM. Đây là bài báo công bố mở được phát hành theo các điều khoản của the Creative Commons Attribution 4.0 International license.



TÓM TẮT

Từ điển song ngữ là công cụ quan trọng cho việc dịch máy tự động. Bằng cách tận dụng các kỹ thuật học máy tiên tiến, chúng ta có thể xây dựng từ điển song ngữ bằng cách tự động học các sự ánh xạ từ vựng từ tập văn bản song ngữ. Tuy nhiên, việc thu thập tập văn bản song ngữ phong phú cho các ngôn ngữ ít tài nguyên, chẳng hạn như ngôn ngữ Ba Na, đặt ra một thách thức đáng kể. Những nghiên cứu gần đây cho thấy rằng các tập văn bản đơn ngữ, kết hợp với từ neo (anchor words), có thể hỗ trợ trong quá trình học các ánh xạ này. Phương pháp thường được áp dụng bao gồm sử dụng Mạng GAN (Generative Adversarial Networks) kết hợp giải quyết vấn đề *trục giao Procrustes* để tạo ra sự ánh xạ này. Phương pháp này thường không ổn định và đòi hỏi tài nguyên tính toán đáng kể, đưa đến những khó khăn tiềm ẩn khi xử lý những ngôn ngữ ít tài nguyên như tiếng Ba Na được thu thập ở vùng sâu vùng xa. Để giảm thiểu điều này, chúng tôi đề xuất một chiến lược điều chỉnh *số chiều thấp* (low-rank), trong đó các hạn chế của GAN có thể được tránh bằng cách tính toán trực tiếp sự biến đổi giữa ngôn ngữ nguồn và ngôn ngữ đích. Chúng tôi đã đánh giá phương pháp của mình bằng cách sử dụng một bộ dữ liệu giàu tài nguyên giữa tiếng Pháp - tiếng Anh và một bộ dữ liệu ít tài nguyên giữa tiếng Việt - tiếng Ba Na. Đáng chú ý, sự ánh xạ từ vựng giữa tiếng Việt - tiếng Ba Na được tạo ra bằng phương pháp của chúng tôi có giá trị không chỉ trong lĩnh vực khoa học máy tính, mà còn đóng góp đáng kể vào việc bảo tồn di sản văn hóa của ngôn ngữ Ba Na trong cộng đồng dân tộc thiểu số của Việt Nam.

Từ khoá: Thu giảm số chiều, ánh xạ từ vựng, ngôn ngữ ít tài nguyên, giải thuật Kabsch

Trích dẫn bài báo này: Huy L C, Minh L Q, Oanh T N, Đồng L D, Đức N Q, Sang N T, Quân T, Thơ Q T. **Hướng tiếp cận thu giảm số chiều cho phép ánh xạ từ vựng tiếng Việt sang tiếng Ba Na từ các tập ngữ liệu khổng song song.** *Sci. Tech. Dev. J. - Eng. Tech.* 2024; ():1-1.