

Brain neuro
Language
Processing

Bilingual Word Embedding with parallel corpus

고려대학교 컴퓨터학과

Brain-neuro Language Processing Lab

이설화

E-mail : whiteldark@korea.ac.kr

Index

2

1

Introduction to
Word representation

2

Word
Embedding

3

Bilingual
Word
Embedding

4

Building a Bilingual
word embedding
space with
parallel corpus

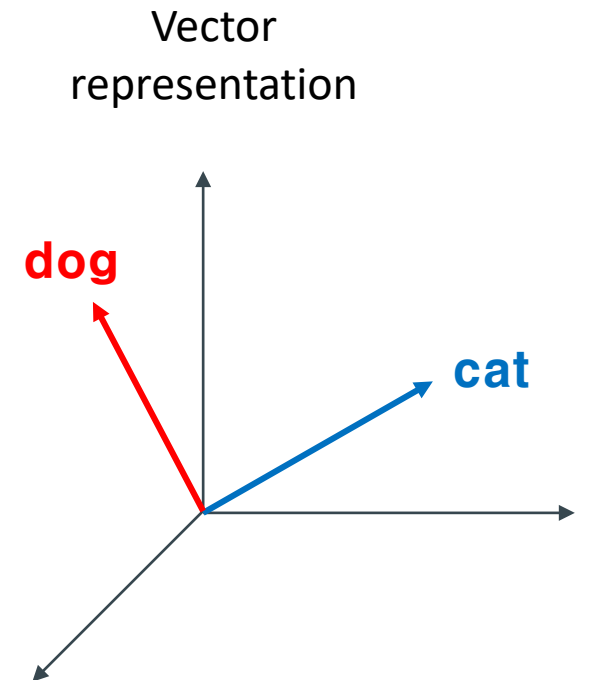


Introduction to Word Representation

Trend of word representation

4

- 단어를 dense 실수 벡터인 vector representation으로 표현하는 방법(word embedding)은 최근 NLP분야에서 인기를 얻고 있음



Traditional word representation

5

- 기존 word representation은 one-hot (or one-of-N) encoding 방식을 사용해왔음

Vocabulary (size=10)

word	id
"the"	1
"dog"	2
"cat"	3
"and"	4
...	...



One-hot vector representation

dog

0 1 0 0 0 0 0 0 0 0

cat

0 0 1 0 0 0 0 0 0 0

Limitation of traditional word representation method (1/2)

6

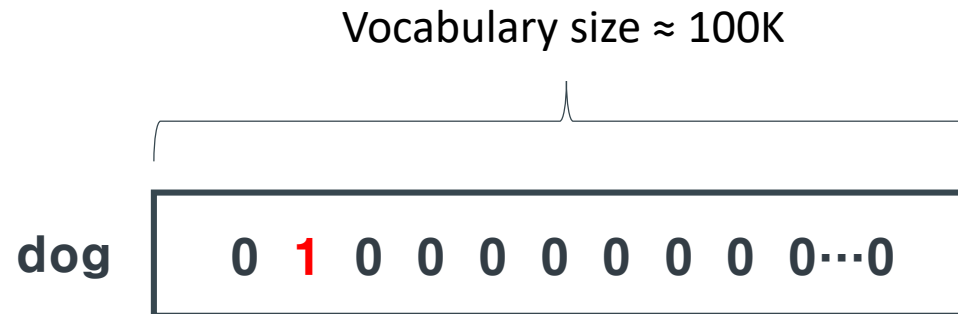
- One-hot encoding은 단어들 간의 관계성을 고려하여 표현하지 않음
 - Ex. Cat : Id3, dog : Id2 => animal



Limitation of traditional word representation method (2/2)

7

- One-hot representation은 매우 높은 dimension을 가지는 문제점이 있음
 - Memory expensive



Word Embedding

The background features a series of overlapping chevron shapes pointing to the right. The colors transition from a dark blue on the left to a lighter blue on the right. A thick, bright cyan line runs diagonally across the image, starting from the top center and extending towards the bottom right, passing over the chevrons.

Word embedding이 란?

9

- 일반적으로 one-hot encoding과 같은 sparse한 방식으로 표현된 단어를 더 낮은 차원의 dense한 실수 벡터 공간에 매핑하는 것
 - 한 단어당 일차원 dimension의 벡터 공간으로 매핑함으로써 lower dimension을 구성함



One-hot vector representation

dog	0	1	0	0	0	0	0	0	0	0
cat	0	0	1	0	0	0	0	0	0	0

sparse



Word embedding

dog	0.3	0.8	0.9	0.1	0.3
cat	0.5	0.1	0.1	0.2	0.6

Dense

Word embedding의 목표

10

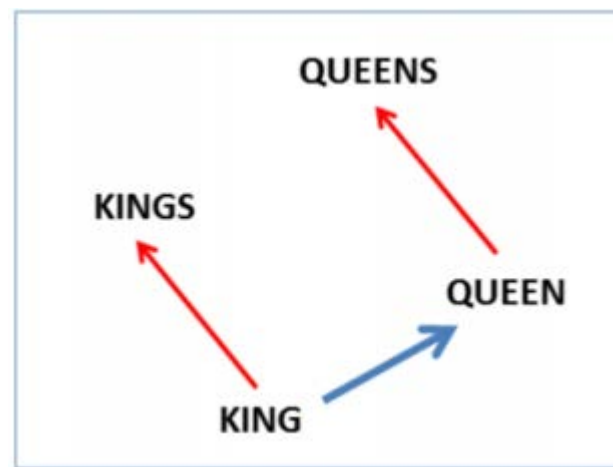
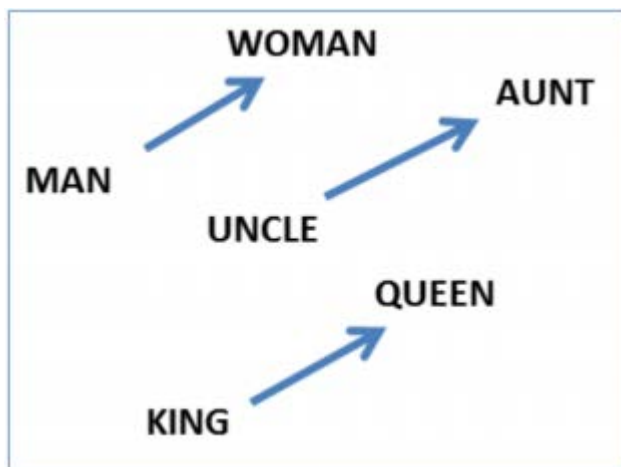
- 단어를 dense한 실수 벡터 공간에 매핑하되,
단어의 의미가 반영되도록 함
 - 유사한 의미의 단어는 벡터 공간 상의 가까운 거리 내에
분포하도록 함




Word embedding의 특성

11

- 단어 의미 간의 상관계도 반영됨
- $king - man + woman \approx ???$
- $king - man + woman \approx queen$
- Example 사이트 : <http://w.elnn.kr/search/>





Bilingual Word Embedding

Bilingual Word Embedding 정의

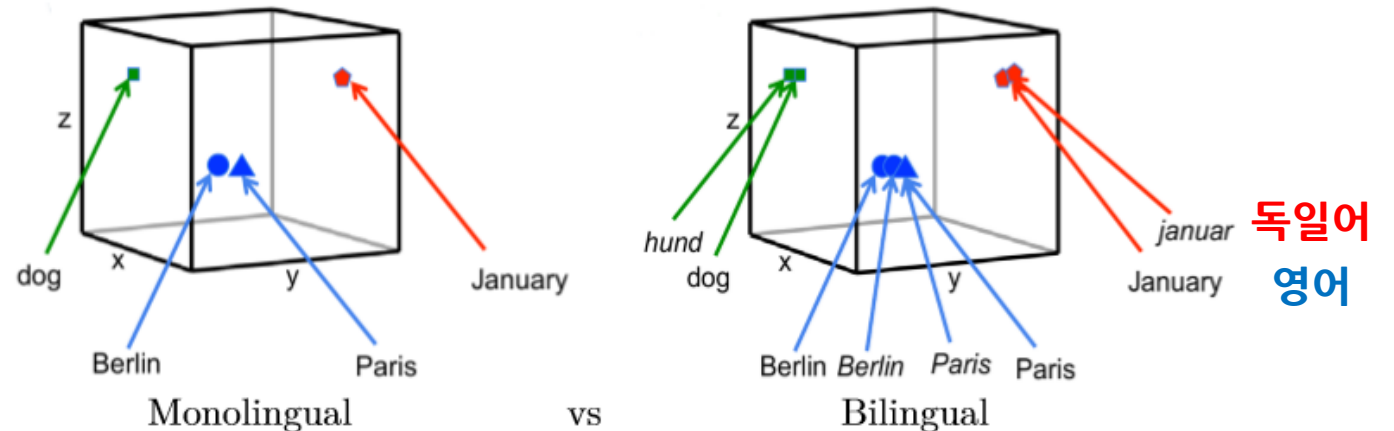
13

Bilingual Word Embedding (BWE) 정의

- 두 개의 다른 언어로부터 하나의 shared space에 단어를 embedding하는 것

BWE goal

- BWE learning model은 서로 다른 두 언어에서 유사한 의미를 가지는 단어가 유사한 공간에 mapping되도록 하는 것

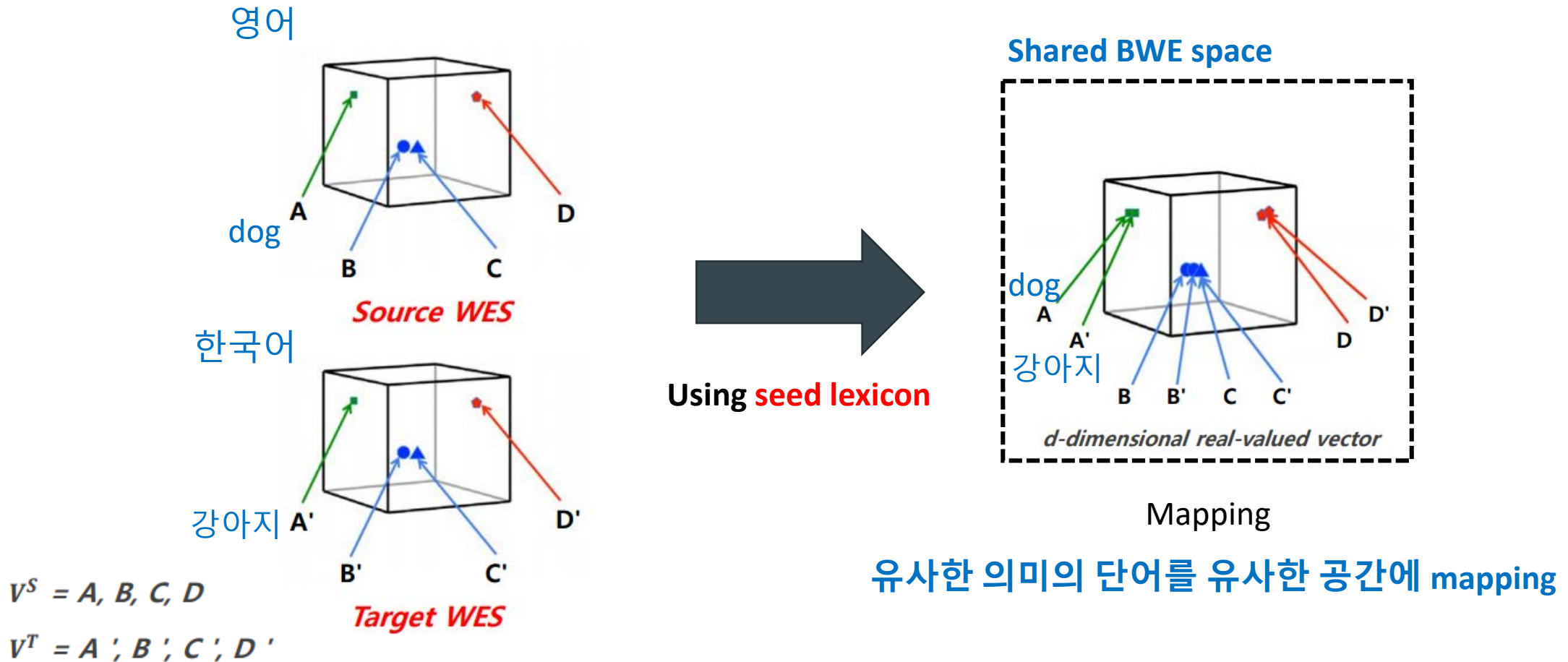


14

- [illegible]

Overview of building a BWE space

15





Building a Bilingual word
embedding space with
parallel corpus

Research Motivation (1/2)

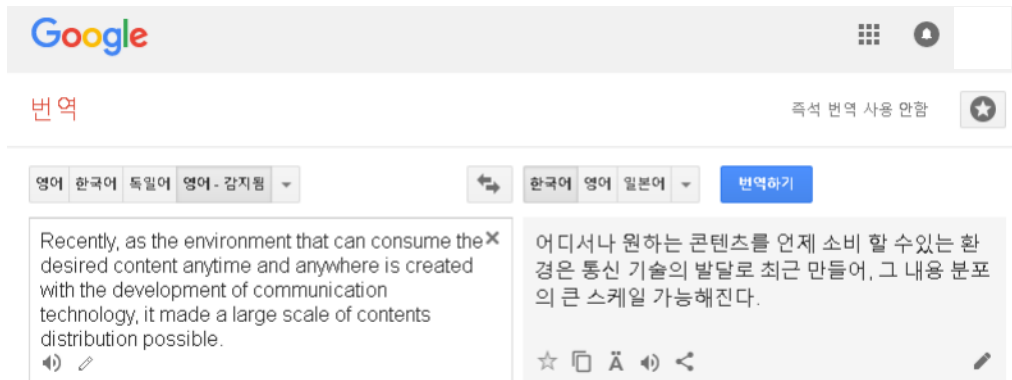
17

- 기존의 BWE 연구에서는 “**High-quality training seed lexicons**”가 존재한다고 가정
 - Vulic, Ivan, and Anna Korhonen. "On the role of seed lexicons in learning bilingual word embeddings." ACL, 2016.
 - Shared Bilingual Word Embedding Space(SBWES)를 유도하는 과정에서 seed lexicons 의 역할과 중요성에 대해 분석하였음

Research Motivation (1/2)

18

High-quality???

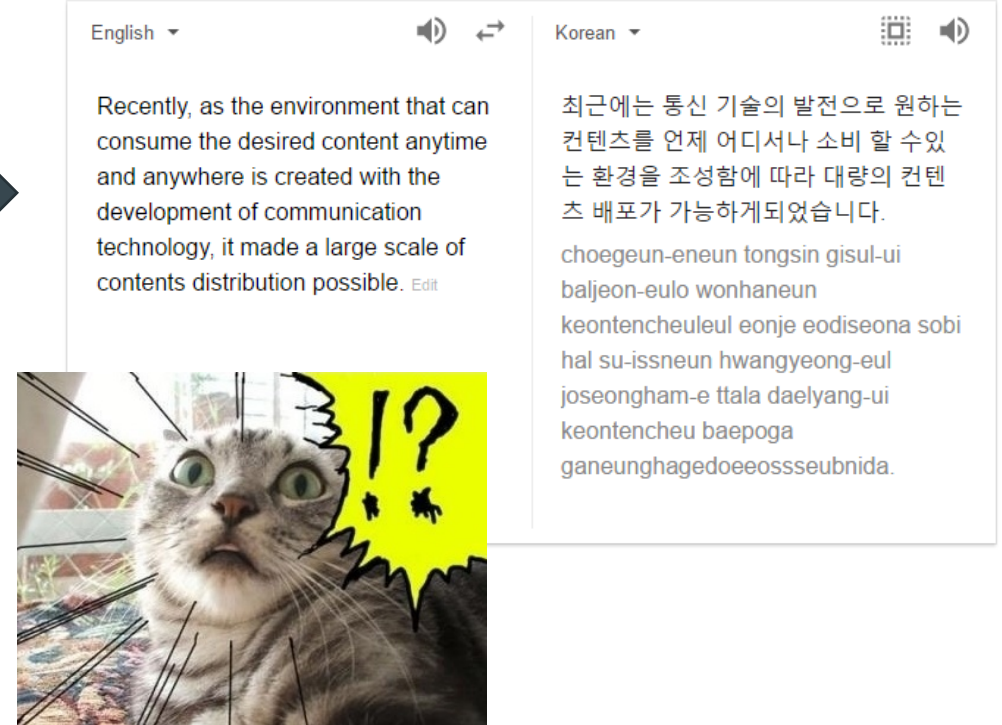


upgrade



Google

High-quality!!!



Research Motivation (2/2)

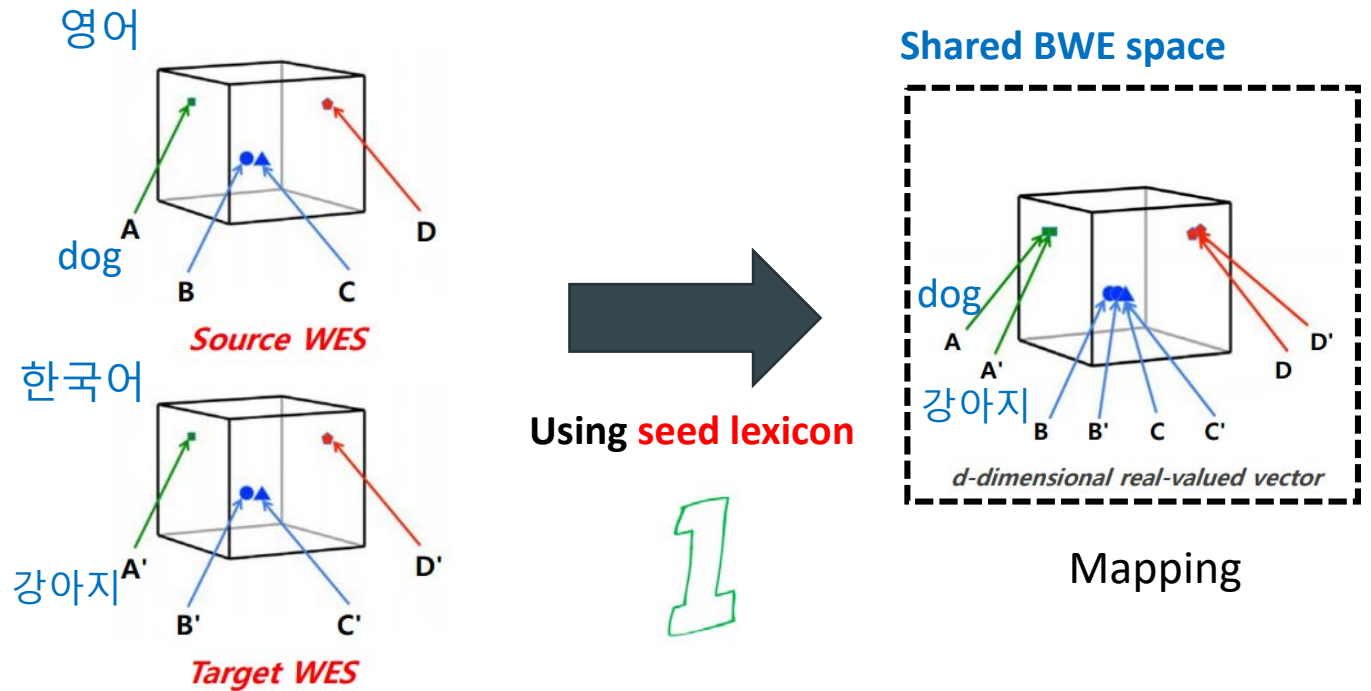
19

SBWES를 생성하는 데는 다음의 두 가지 과정을 거침 (Vulic et al., 2016, ACL)

1

두 개의 별도로 준비된 non-aligned monolingual embedding spaces는 monolingual WE learning model을 사용하여 유도하게 됨

- 비교적 얻기 쉬운 문서 번역쌍 (Wikipedia)을 word translation pairs로 사용하였음
- Our research idea => **Parallel corpus**

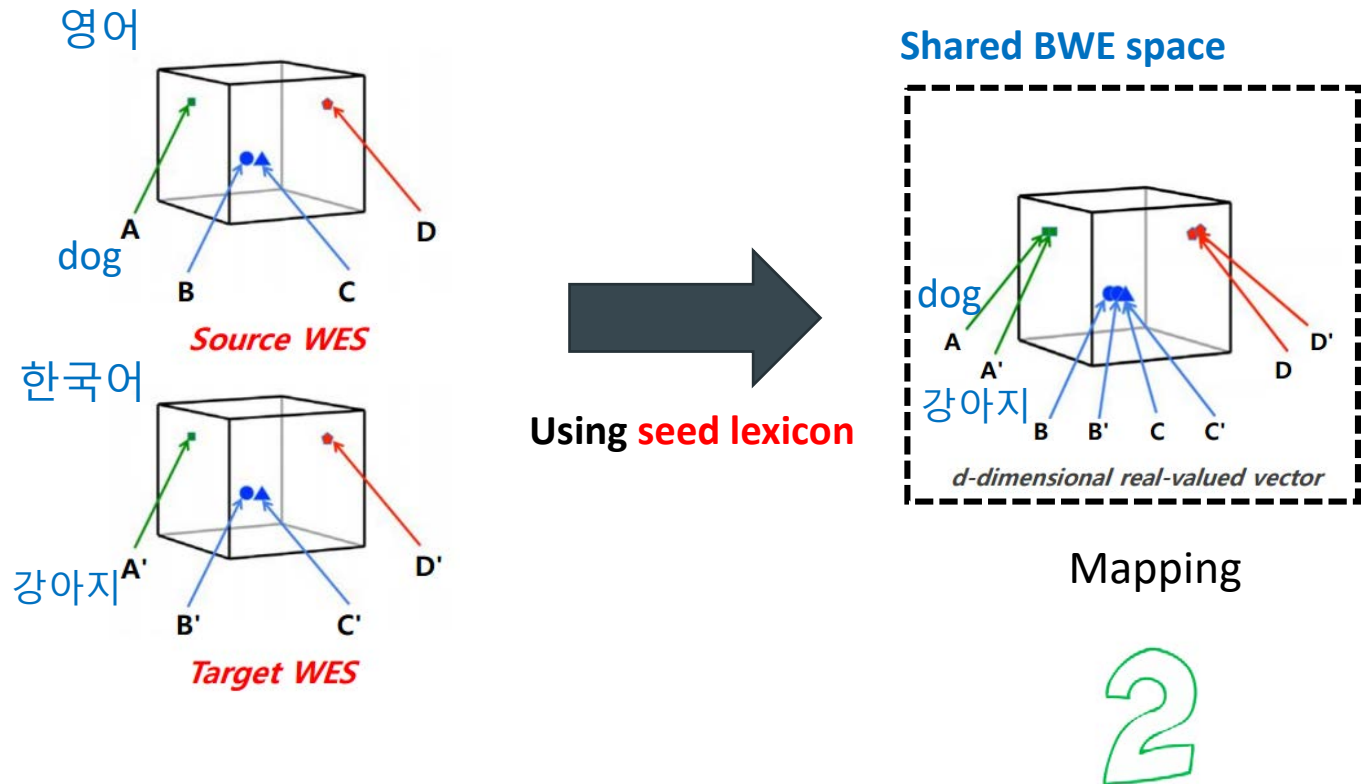


Research Motivation (2/2)

20

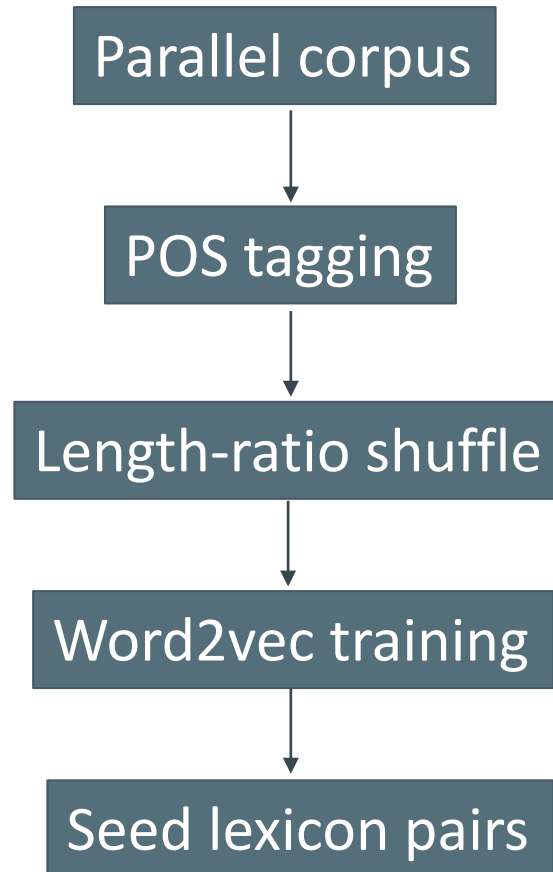
2

단어 번역 쌍을 seed lexicon으로 사용하고, 이를 통해 두 개의 monolingual spaces를 하나로 묶는 mapping function을 학습함으로써 shared BWE space를 구성함



making the seed lexicons

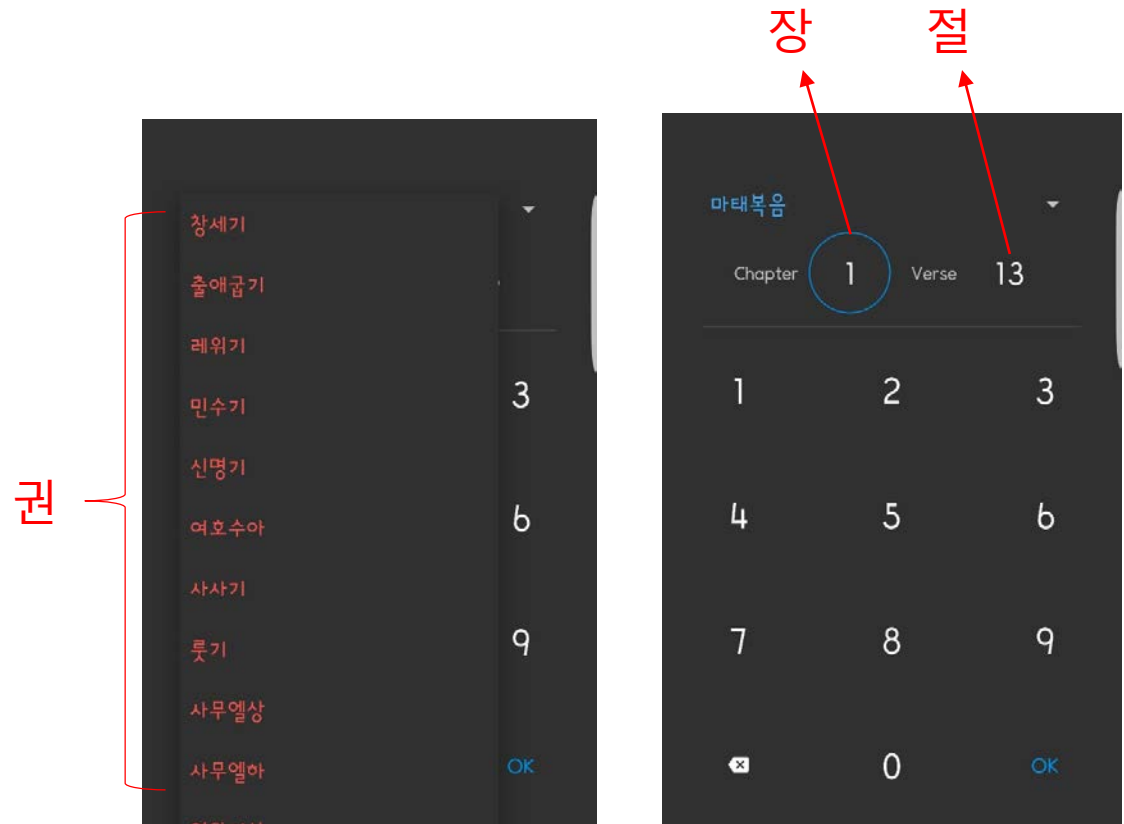
21



Making the seed lexicons

22

- Parallel corpus
 - Sentence-aligned corpus
 - Bible data (총 66권, 1,189장, 31,102 절로 구성)



Making the seed lexicons

23

권, 장, 절

01창 1:1 태초에 하나님이 천지를 창조하시니라
01창 1:2 땅이 혼돈하고 공허하며 흑암이 깊음 위에 있고 하나님의 영은 수면 위에 운행하시니라
01창 1:3 하나님이 이르시되 빛이 있으라 하시니 빛이 있었고
01창 1:4 빛이 하나님이 보시기에 좋았더라 하나님이 빛과 어둠을 나누사
01창 1:5 하나님이 빛을 낮이라 부르시고 어둠을 밤이라 부르시니라 저녁이 되고 아침이 되니 이는 첫째 날이니라
01창 1:6 하나님이 이르시되 물 가운데에 궁창이 있어 물과 물로 나뉘라 하시고
01창 1:7 하나님이 궁창을 만드사 궁창 아래의 물과 궁창 위의 물로 나뉘게 하시니 그대로 되니라
01창 1:8 하나님이 궁창을 하늘이라 부르시니라 저녁이 되고 아침이 되니 이는 둘째 날이니라
01창 1:9 하나님이 이르시되 천하의 물이 한 곳으로 모이고 물이 드러나라 하시니 그대로 되니라

한국어

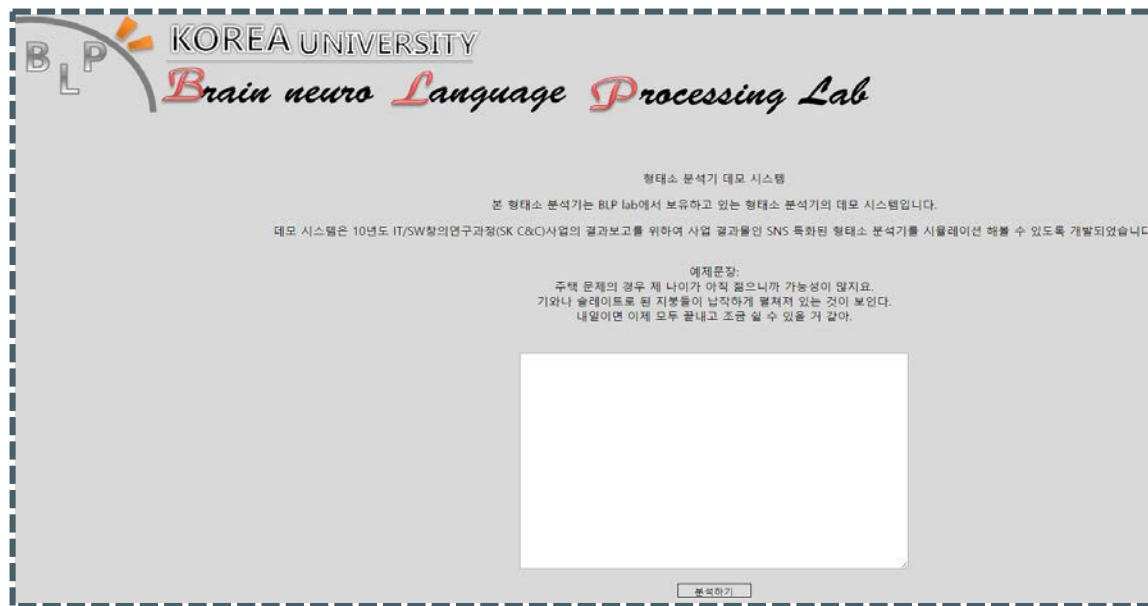
01Gn 1:1 In the beginning God created the heavens and the earth.
01Gn 1:2 Now the earth was formless and empty, darkness was over the surface of the deep, and the Spirit of God was hovering over the waters.
01Gn 1:3 And God said, "Let there be light," and there was light.
01Gn 1:4 God saw that the light was good, and he separated the light from the darkness.
01Gn 1:5 God called the light "day," and the darkness he called "night." And there was evening, and there was morning--the first day.
01Gn 1:6 And God said, "Let there be an expanse between the waters to separate water from water."
01Gn 1:7 So God made the expanse and separated the water under the expanse from the water above it. And it was so.
01Gn 1:8 God called the expanse "sky." And there was evening, and there was morning--the second day.
01Gn 1:9 And God said, "Let the water under the sky be gathered to one place, and let dry ground appear." And it was so.

영어

Making the seed lexicons

Preprocessing

24



■ 한국어 형태소 분석기

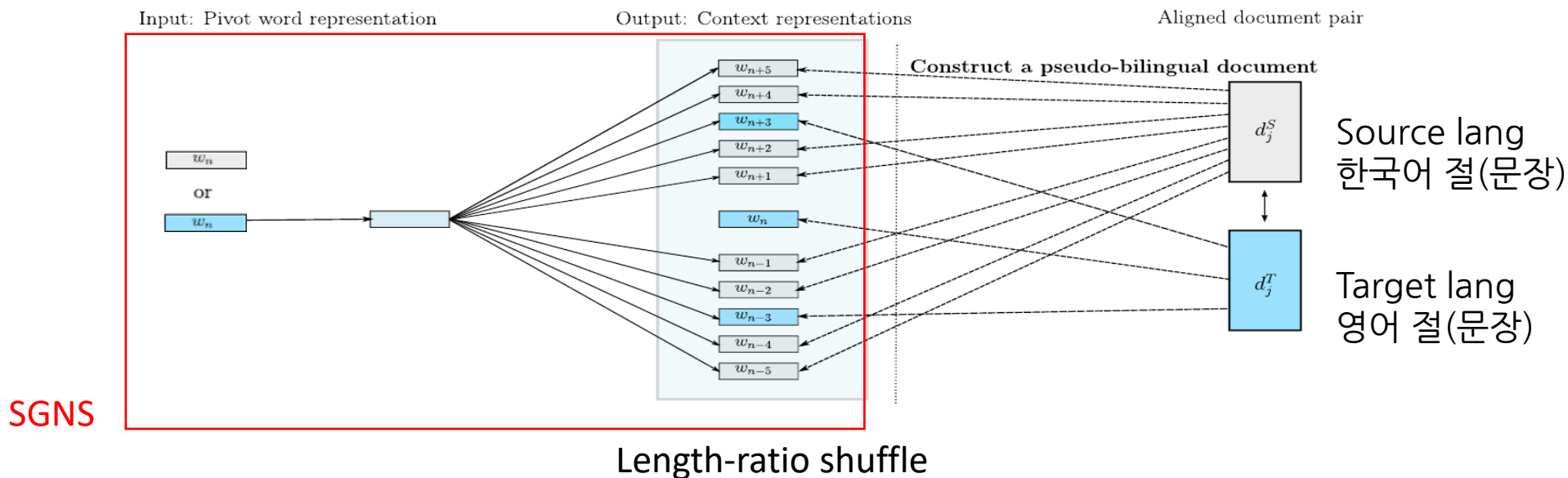
- 기호 제거
- 한국어
 - POS tagger (BLP lab asset)
 - url : <http://blpdemo.korea.ac.kr/MA/>

Making the seed lexicons

Preprocessing

25

- 각 문장별 token수의 비율로 shuffle
 - 예) Ko = {포도, 복숭아}
 - Token number = 2
 - 예) Eng = {carrot, apple, pineapple, egg}
 - Token number = 4
 - Length-ratio = 2 : 1



Making the seed lexicons

Bilingual Word Embedding Skip-Gram (1/2)

26

```
2 01창 1:1 태초에 하나님이 천지를 창조하시니라
3 01창 1:2 땅이 혼돈하고 공허하며 흑암이 깊음 위에 있고 하나님의
4 01창 1:3 하나님이 이르시되 빛이 있으라 하시니 빛이 있었고
5 01창 1:4 빛이 하나님이 보시기에 좋았더라 하나님이 빛과 어둠을
6 01창 1:5 하나님이 빛을 낮이라 부르시고 어둠을 밤이라 부르시니라
7 01창 1:6 하나님이 이르시되 물 가운데에 궁창이 있어 물과 물로 나
8 01창 1:7 하나님이 궁창을 만드사 궁창 아래의 물과 궁창 위의 물로
9 01창 1:8 하나님이 궁창을 하늘이라 부르시니라 저녁이 되고 아침이
```

Korean

```
2 01Gn 1:1 In the beginning God created the heaven
3 01Gn 1:2 Now the earth was formless and empty,
4 01Gn 1:3 And God said, "Let there be light," and
5 01Gn 1:4 God saw that the light was good, and he
6 01Gn 1:5 God called the light "day," and the day
7 01Gn 1:6 And God said, "Let there be an expanse
8 01Gn 1:7 So God made the expanse and separated
9 01Gn 1:8 God called the expanse "sky." And there
10 01Gn 1:9 And God said, "Let the water under the
```

English

하나님
god

- in
- 태초
- the
- 에
- beginning
- god
- 이
- created
- 천지
- the

- in
- 태초
- the
- 에
- beginning
- god
- 이
- created
- 천지
- the

Window size = 5

Length-ratio shuffle

```
1 in 태초 the 에 beginning 하나님 god 이 created 천지 the 를 heavens 창조 and 하 the earth 시
2 땅 now 이 the 혼돈 earth 하 was 고 formless 공허하 and 며 empty 흑암 darkness 이 was 깊 over
3 하나님 and 이 god 이르 시 said 되 let 빛 이 there 있 be 으라 하 light 시 and 니 빛 there 이 was
4 빛 god 이 saw 하나님 that 이 the 보 light 시 기 was 에 good 좋 and 았 he 더라 separated 하나님
5 하나님 god 이 called 빛 을 the 낮 light 이 라 day 부르 and 시 the 고 어둠 darkness 을 he 밤 이
6 하나님 and 이 god 이르 said 시 되 let 물 there 가운데 be 에 an 궁창 이 expanse 있 between 어 th
7 하나님 so 이 god 궁창 made 을 the 만들 expanse 사 and 궁창 separated 아래 the 의 water 물 unde
8 하나님 god 이 called 궁창 을 the 하늘 expanse 이 라 sky 부르 and 시 니라 there 저녁 was 이 even
9 하나님 and 이 god 이르 said 시 let 되 the 천하 water 의 under 물 the 이 sky 한 be 곳 gathered
```

Length-ratio shuffle

Making the seed lexicons

Bilingual Word Embedding Skip-Gram (2/2)

27

- Word2vec package
 - SGNS (Skip-gram Negative Sampling)

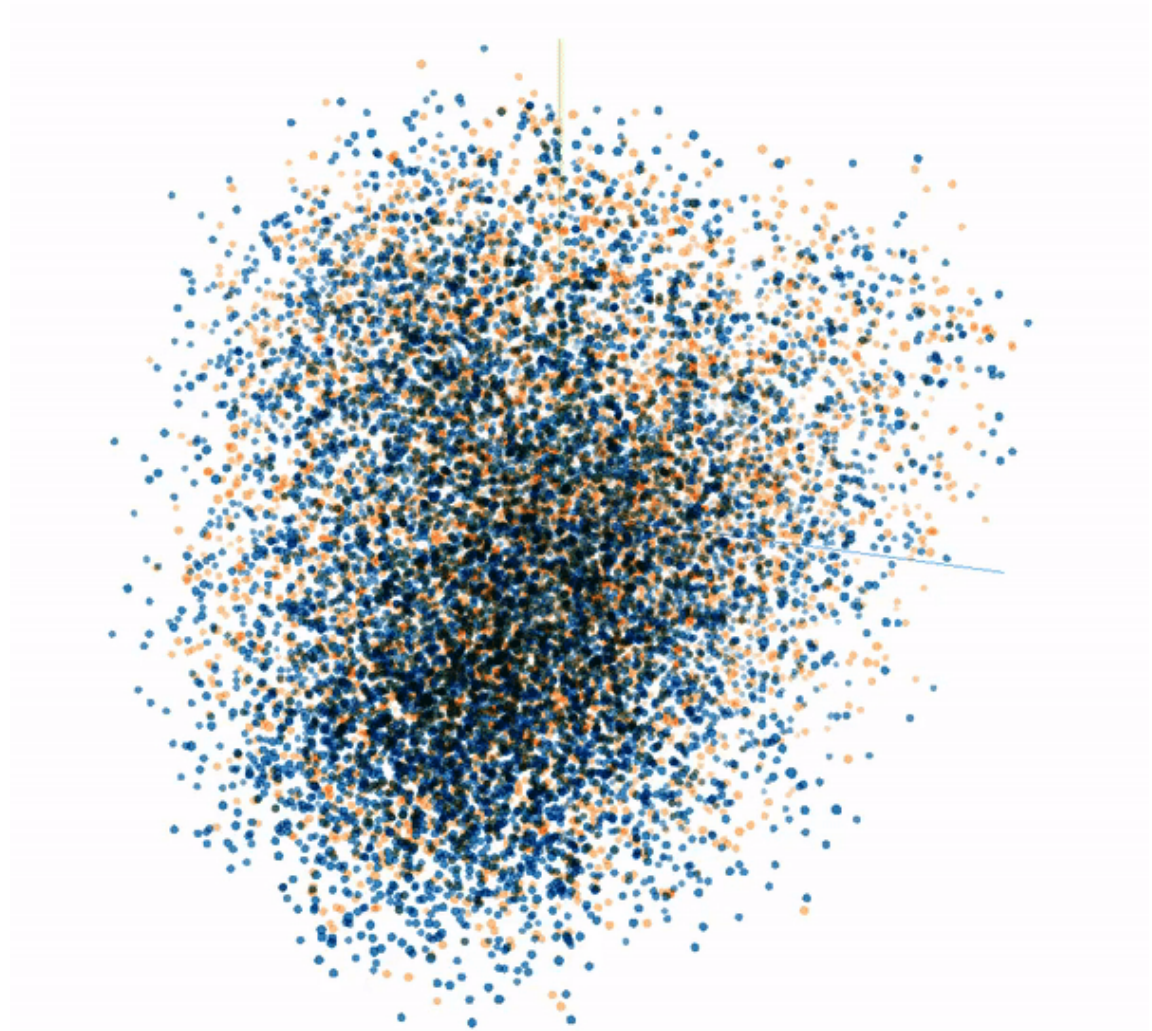


Visualization of seed lexicons

28

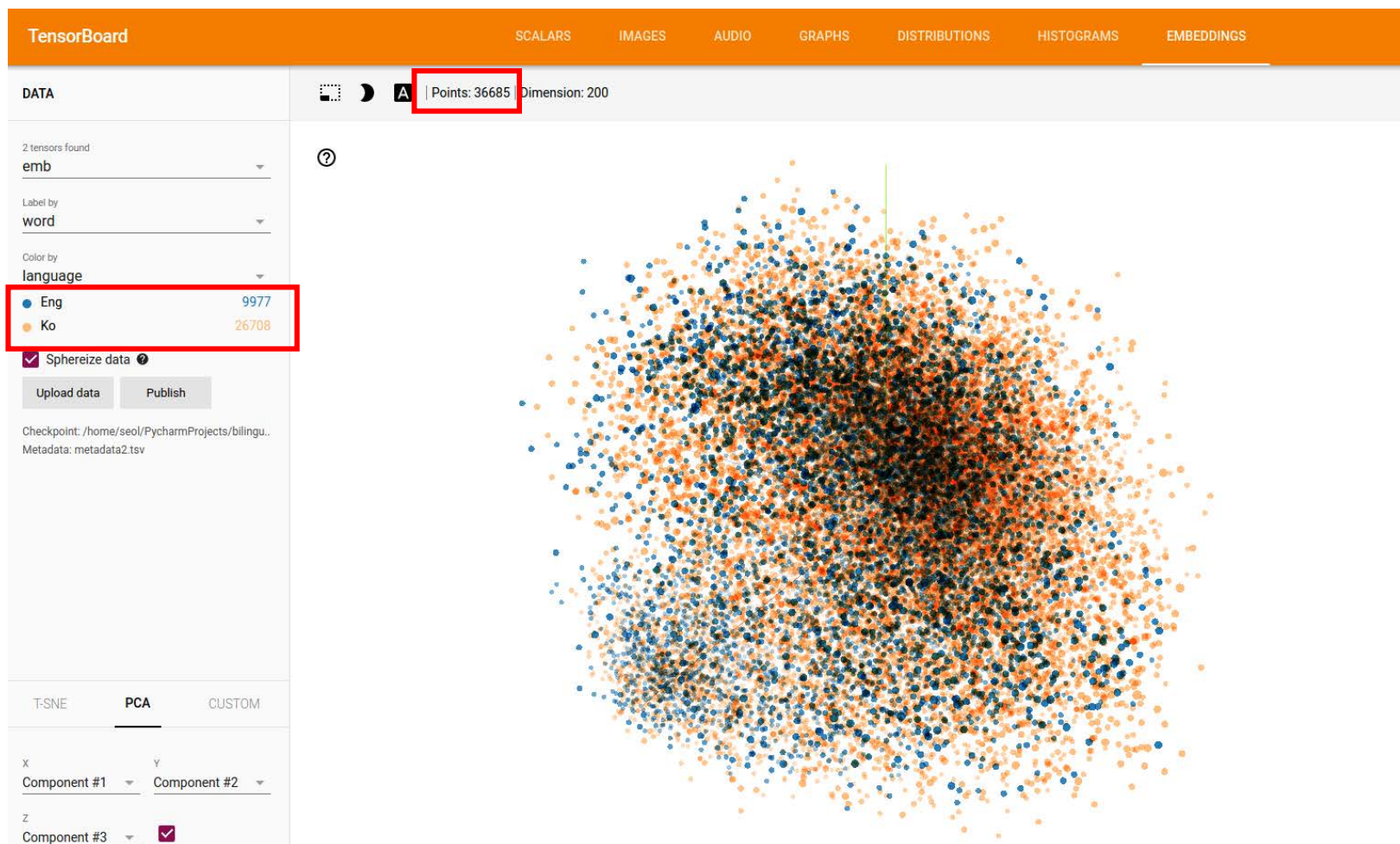
Google

TensorBoard



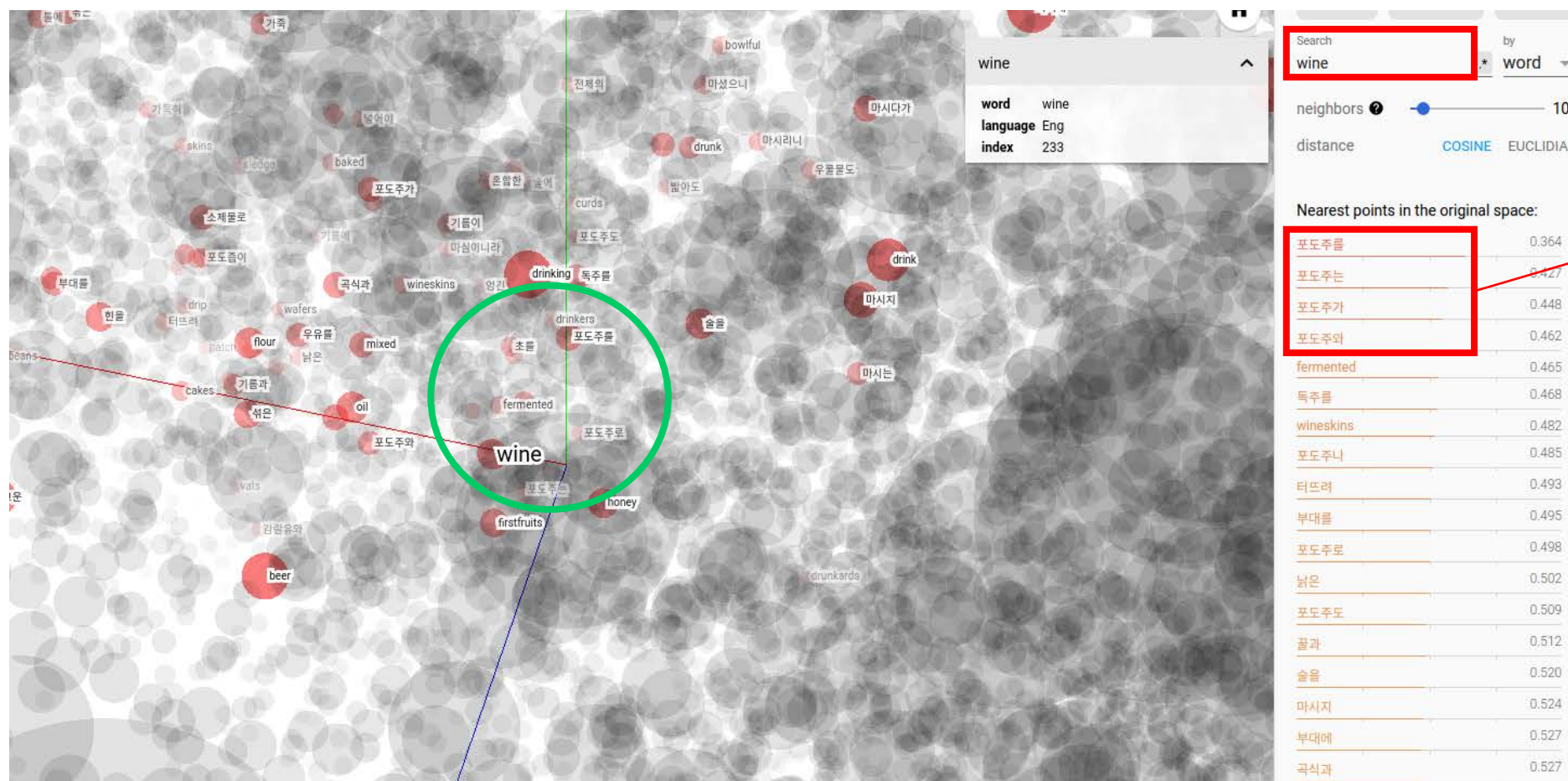
Visualization of seed lexicons_어절단위

29



Visualization of seed lexicons_어절단위

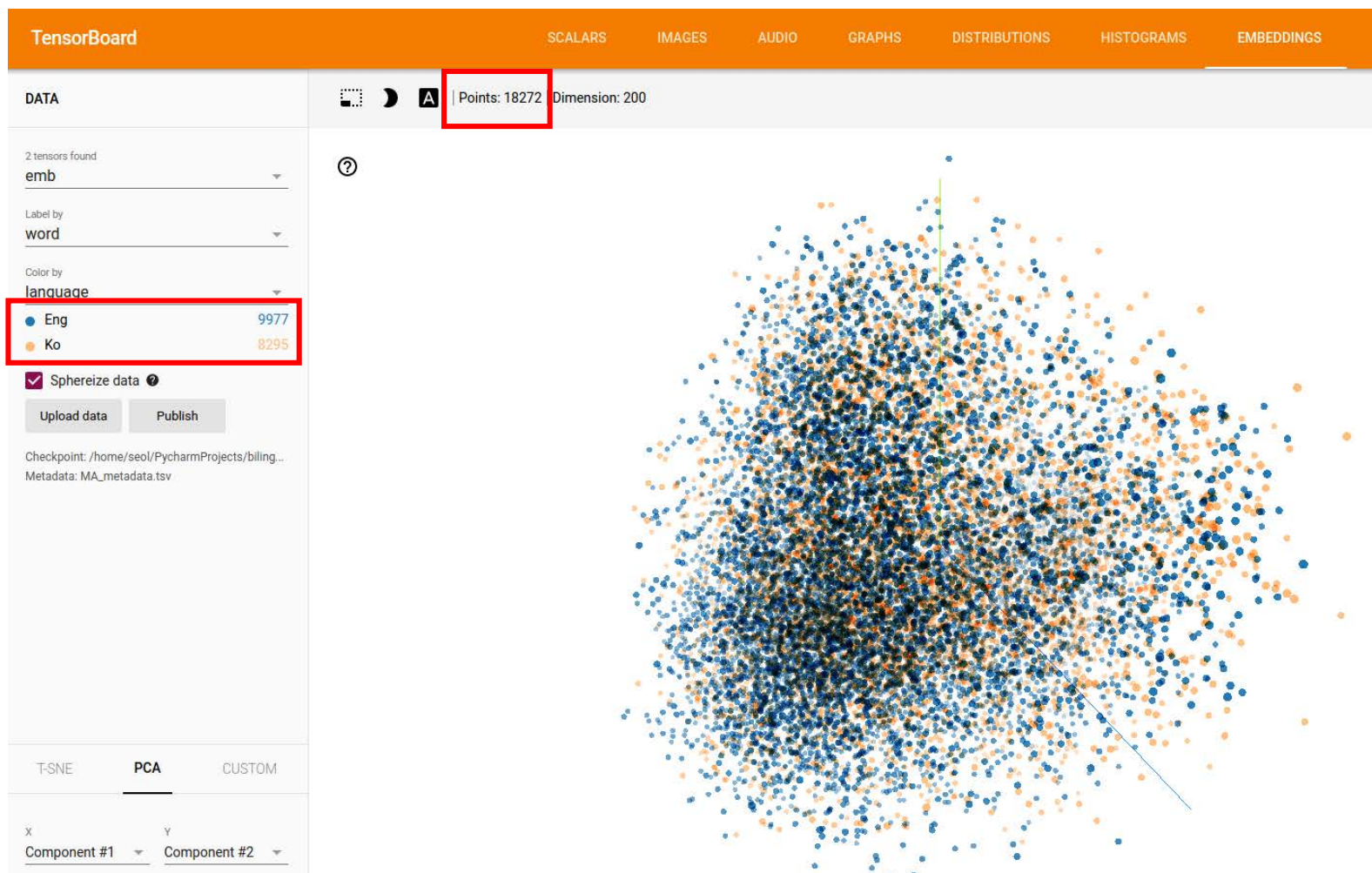
30



어절 단위의 결과

Visualization of seed lexicons_형태소 단위

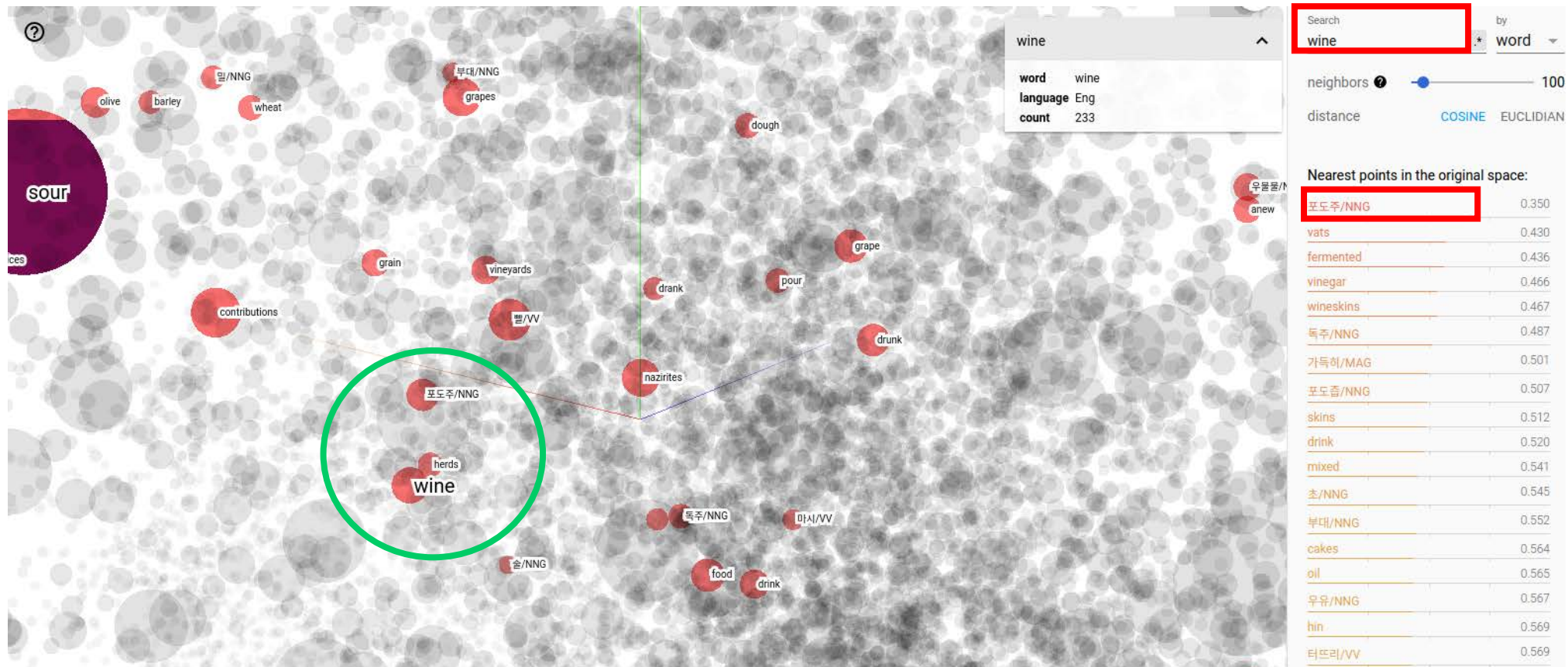
31



Visualization of seed lexicons _ 형태소 단위

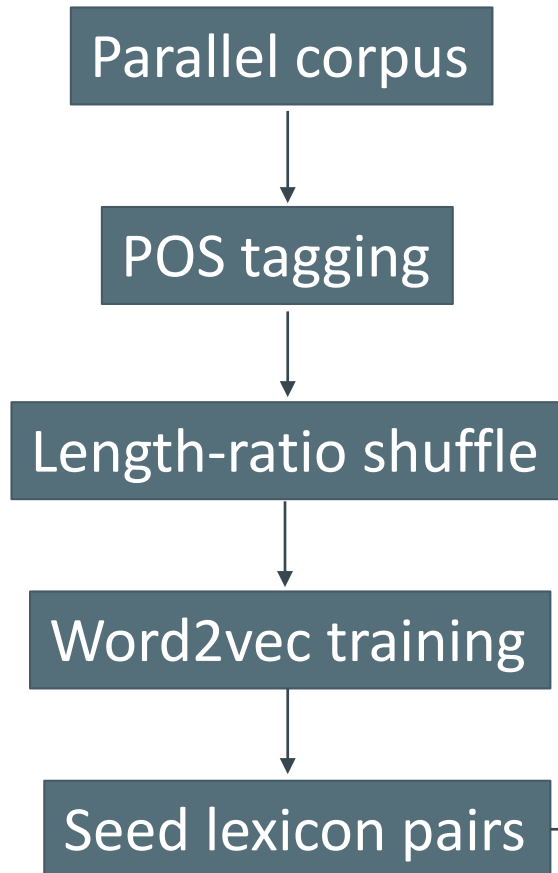
32

“Wine” 검색 => “포도주”가 가장 유사한 결과로 나옴



Mapping function

33



L2-regularized least-squares error objective

$$\min_{\mathbf{W} \in \mathbb{R}^{d_S \times d_T}} \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{W}\|_F^2$$

\mathbf{W} : 2차원 행렬
Mapping matrix

\mathbf{X} : Source 언어의 word 벡터 \mathbf{Y} : Target 언어의 word 벡터

\mathbf{X}, \mathbf{Y} : Seed lexicon

Mapping function

Seed lexicon demo URL

34

<http://blplab.iptime.org:6006>

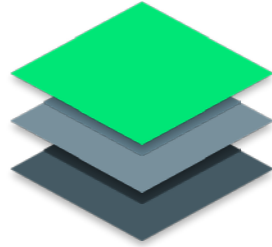
References

35

Vulic, Ivan, and Anna Korhonen. "On the role of seed lexicons in learning bilingual word embeddings." ACL, 2016.

Vulić, Ivan, and Marie-Francine Moens. "Bilingual distributed word representations from document-aligned comparable data." Journal of Artificial Intelligence Research 55 (2016): 953-994.

Zou, Will Y., et al. "Bilingual Word Embeddings for Phrase-Based Machine Translation." EMNLP. 2013.



감사합니다.

이설화

E-mail : whiteldark@korea.ac.kr