

한국어 임베딩

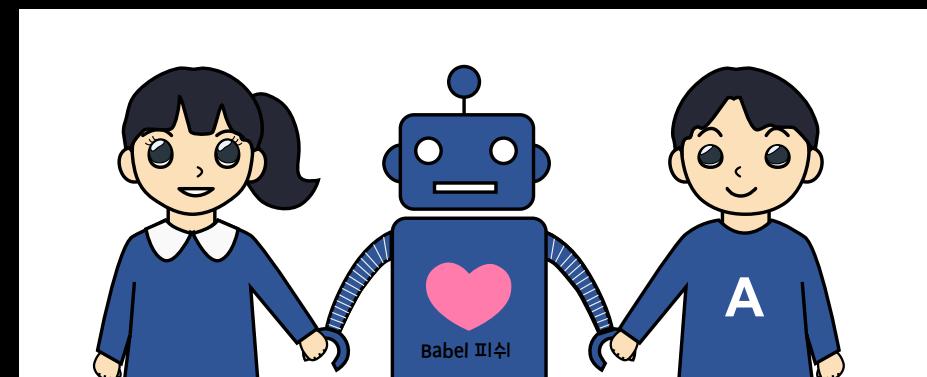
Korean Embedding

nlp4kor

<http://github.com/bage79/nlp4kor>



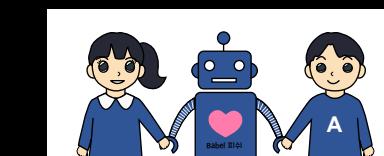
박혜웅



바벨피쉬

What is Word Embedding?

em·bed¹ | imbéd |
동사 타동사 (embed·ded, embed·ding)
1. 《보통 수동태》 〈물건을〉 (...에) 박아 넣다, 〈꽃 등을〉 (...에) 심다(*in ...*);



Word Embedding

TensorBoard PROJECTOR Points: 10000 | Dimension: 300 | Selected 101 points

DATA

1 tensor found
ko.wikipedia.org.top_n_10000

Label by spell Color by No color map

Edit by spell Tag selection as

Load Download Label

Sphereize data ?

Checkpoint: /Users/bage/tensorboard_log./ko.wikipedia.org.top_n_10000.ckpt

Metadata: ./ko.wikipedia.org.top_n_10000.tsv

T-SNE PCA CUSTOM

X Component #1 Y Component #2

Z Component #3

PCA is approximate. ?

Total variance described: 9.6%.

A | Points: 10000 | Dimension: 300 | Selected 101 points

언어
spell 언어
freq 0.0000579

Nearest points in the original space:

언어의	0.543
언어를	0.640
언어는	0.643
사용자	0.644
단어	0.650
인터페이스	0.653
언어로	0.656
언어이다.	0.657
문자	0.668
시각	0.682
영어	0.684
텍스트	0.690
학문	0.695
문화	0.698
프로그래밍	0.703

Show All Data Isolate 101 points Clear selection

Search 언어 by spell

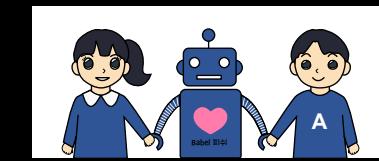
neighbors ? 10

distance COSINE EUCLIDEAN

BOOKMARKS (0) ?

Various Similarity

Semantic & Syntactic

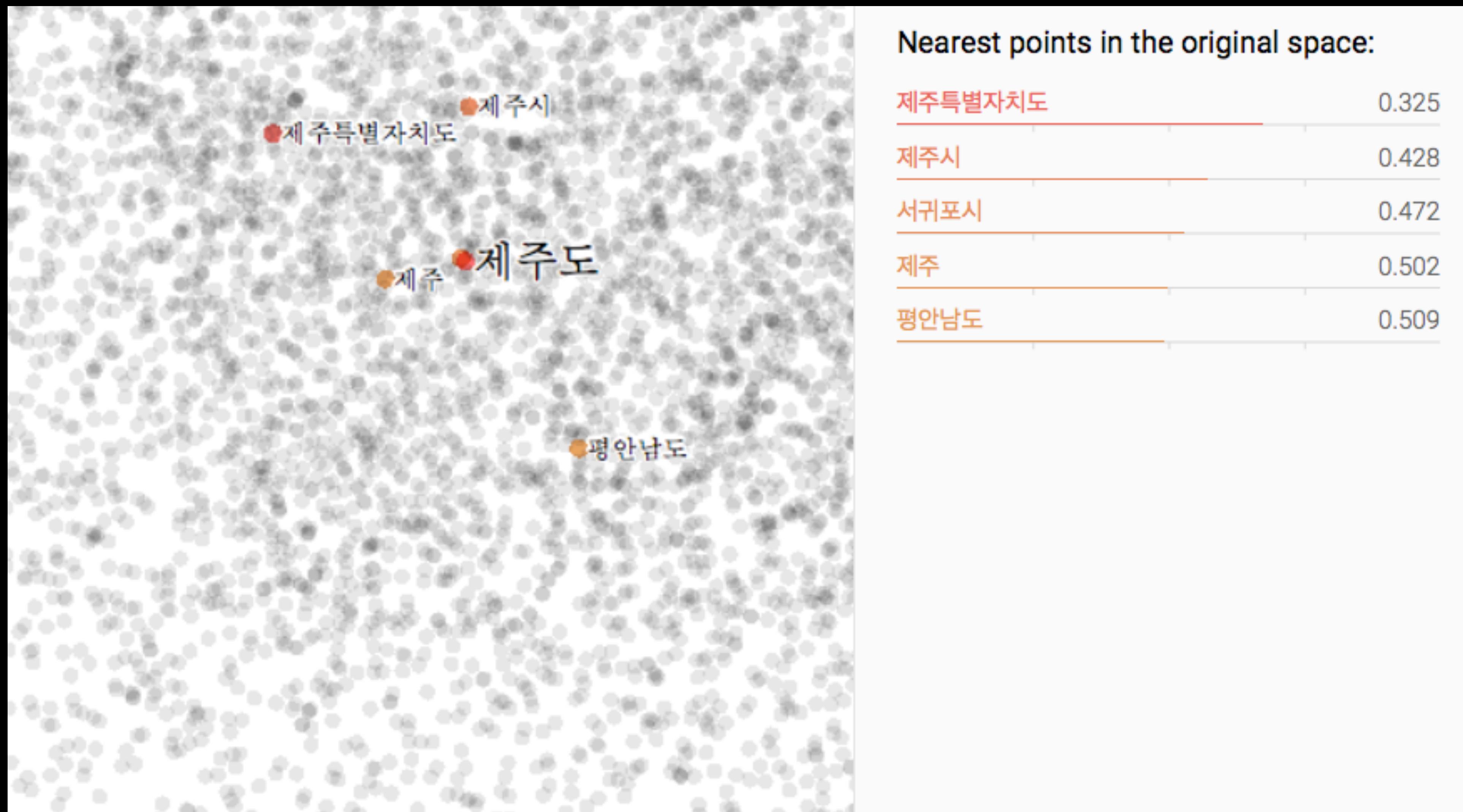


What are Similar two words?

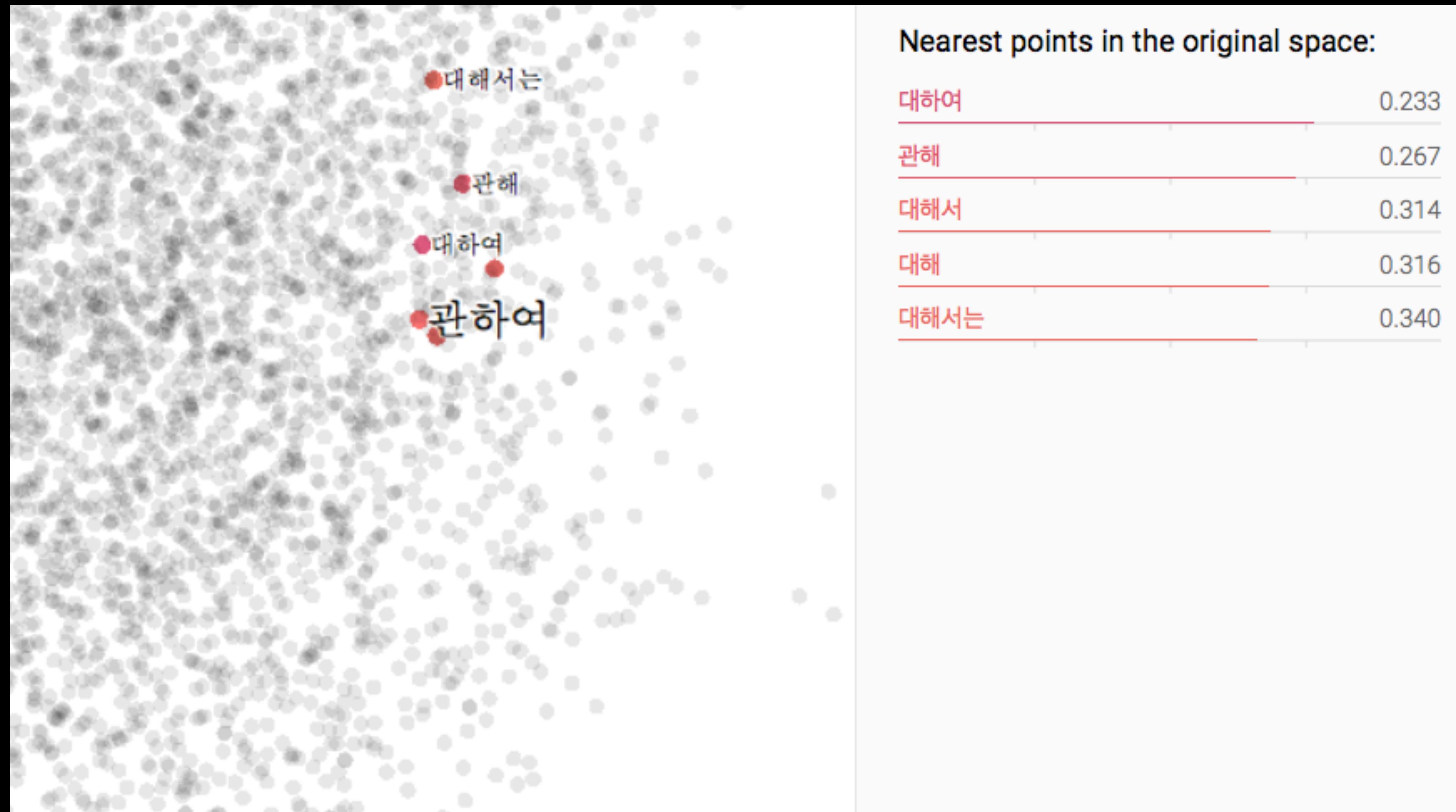


EBS 다큐프라임 <명사로 세상을 보는 서양인, 동사로 세상을 보는 동양인>

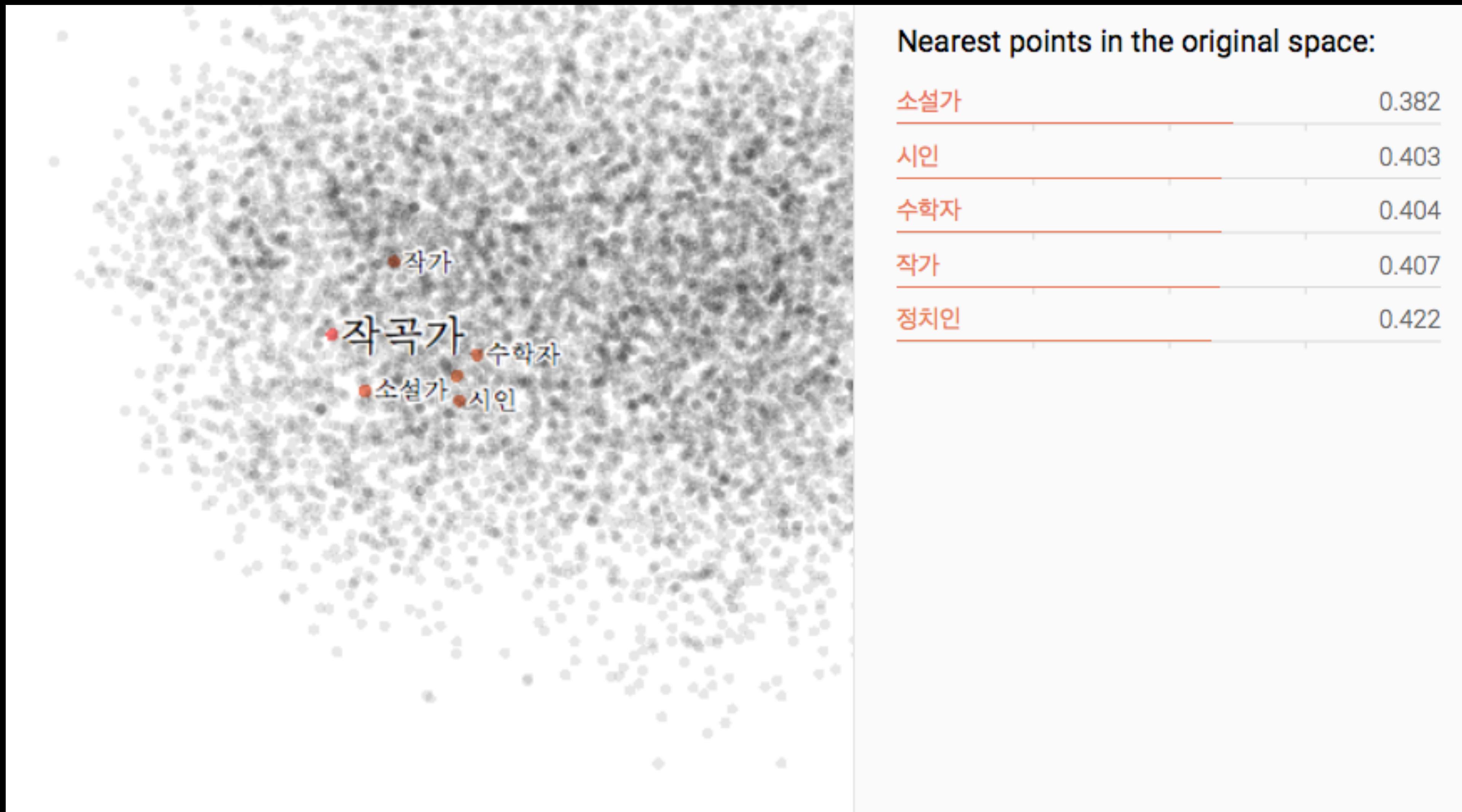
Synonym



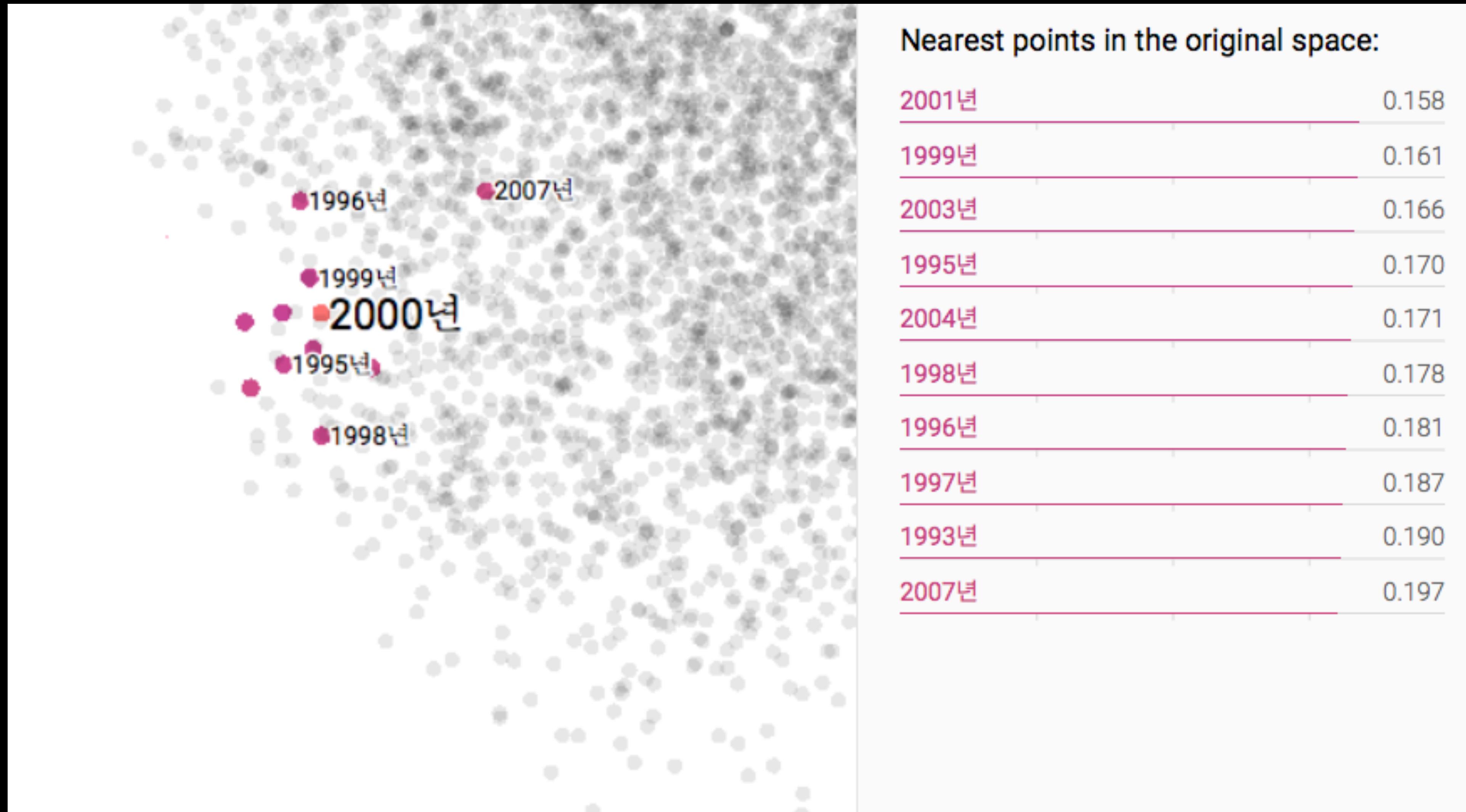
Synonym



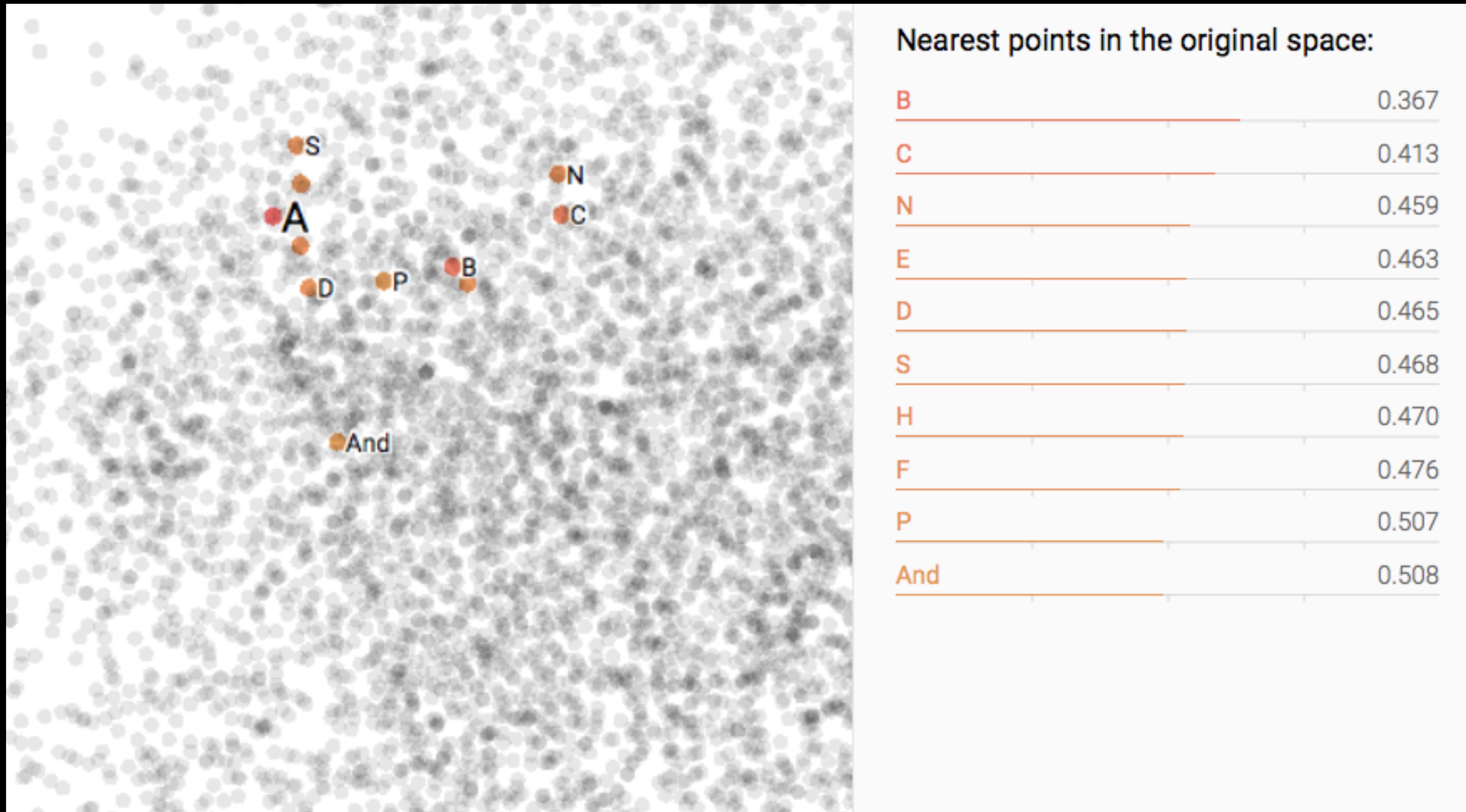
Similiar Types



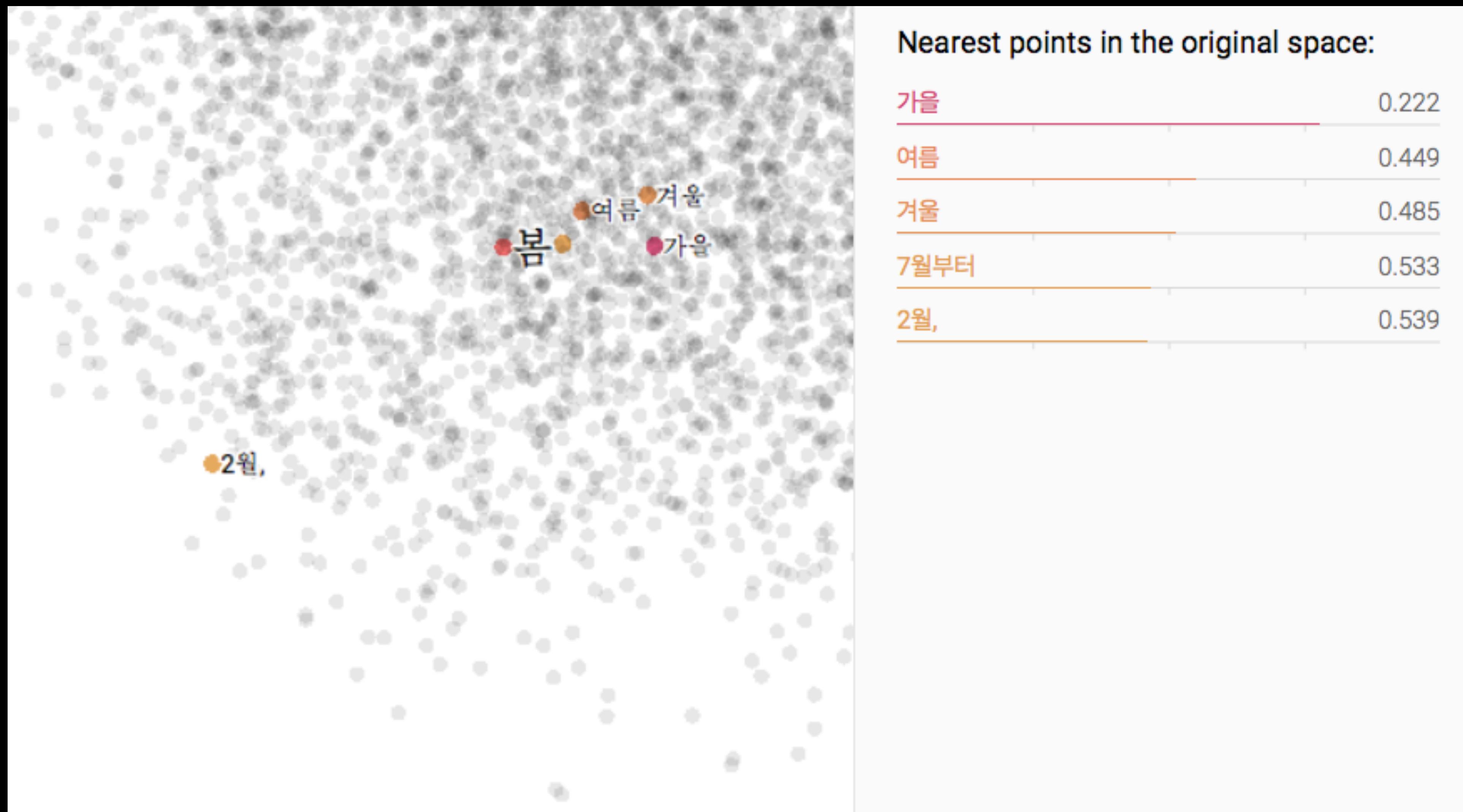
Similiar Types



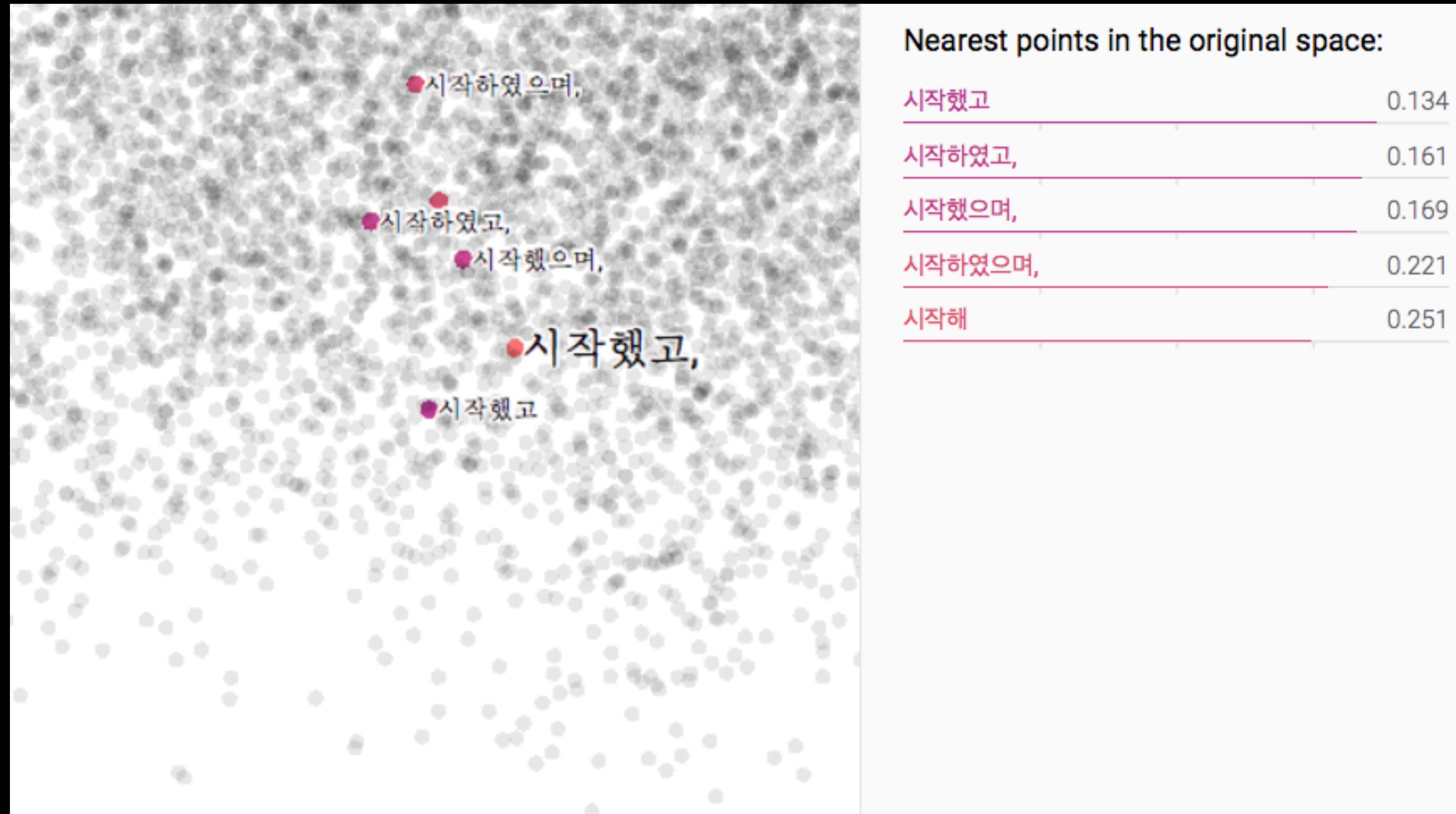
Similiar Types



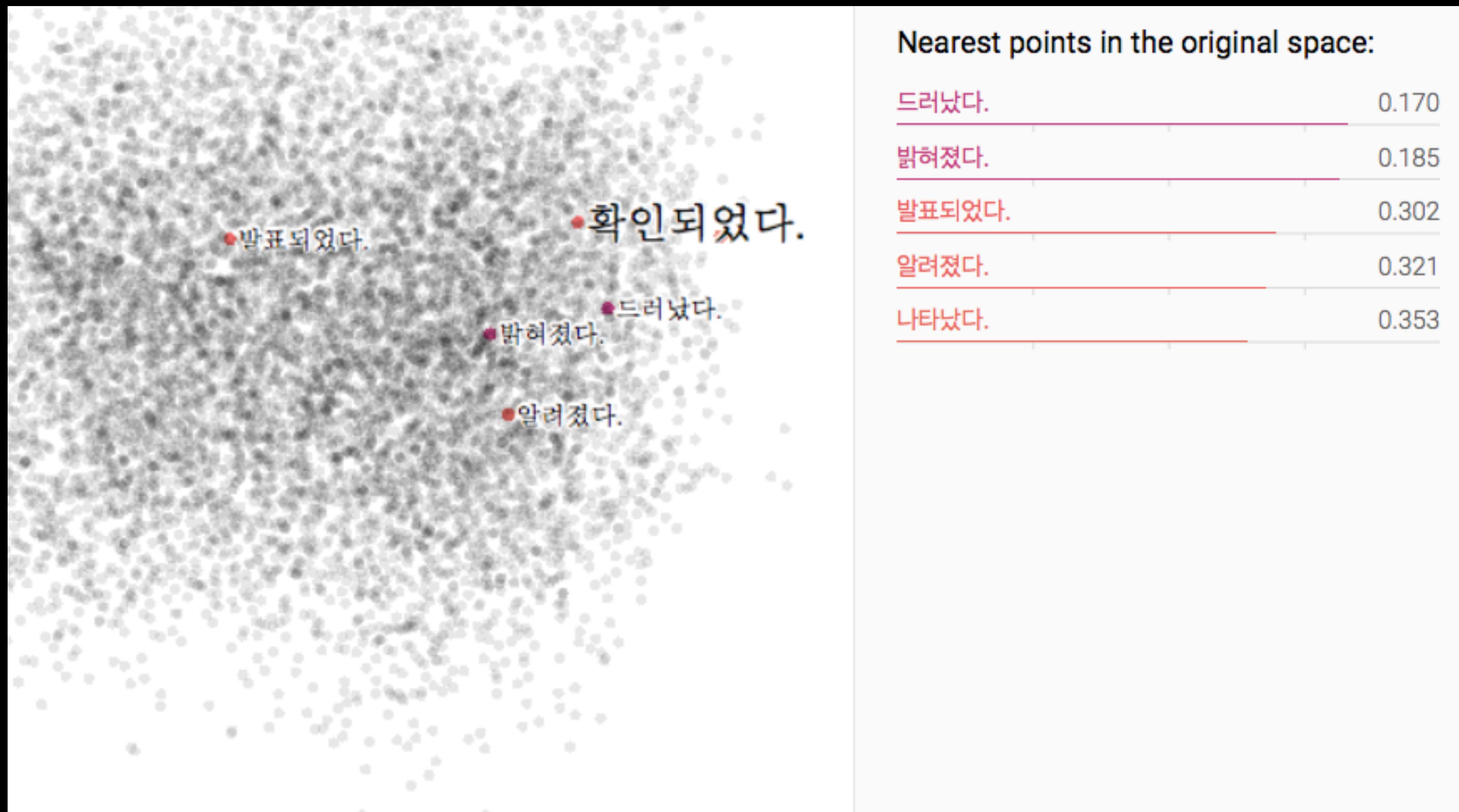
Similiar Types



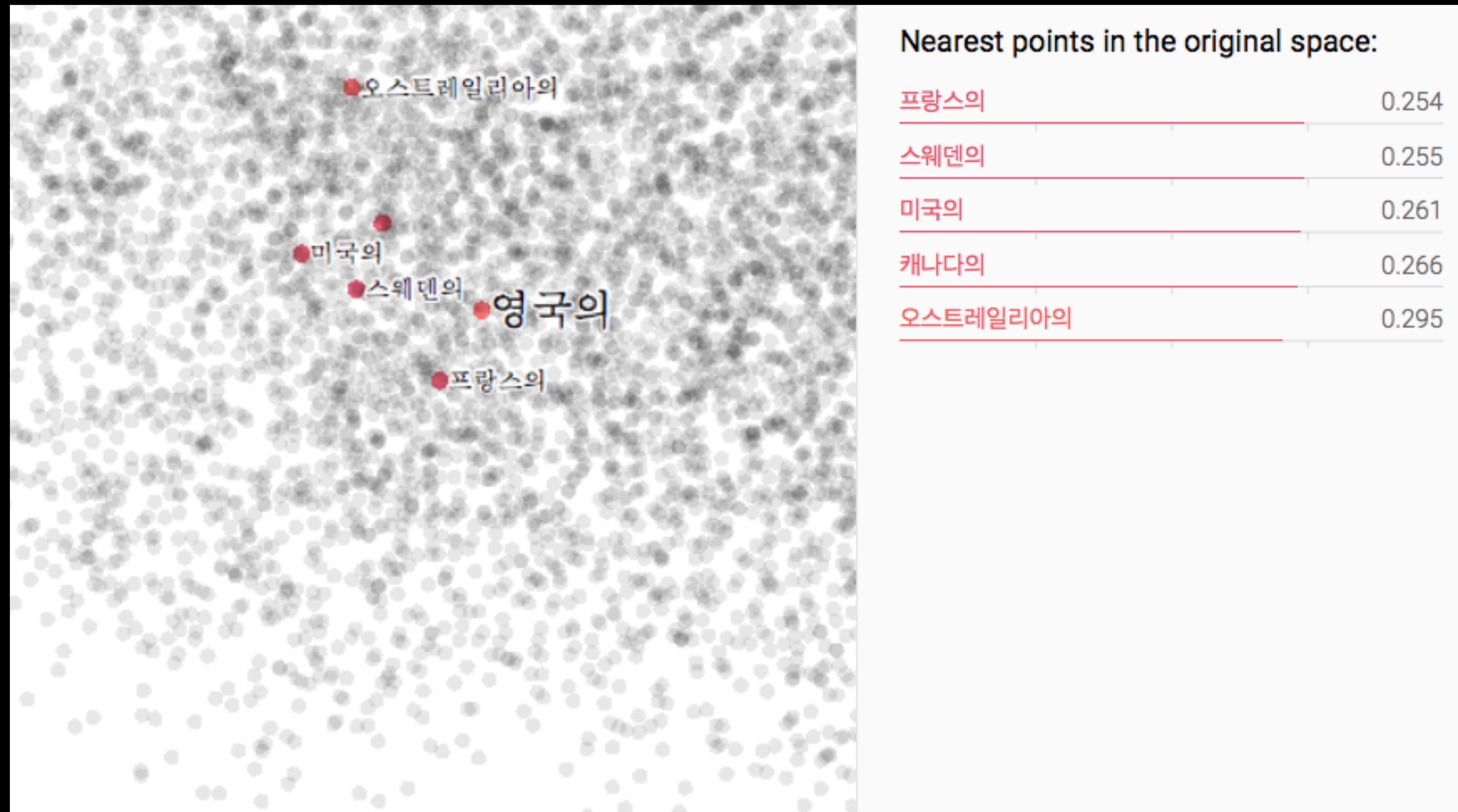
Derivative words



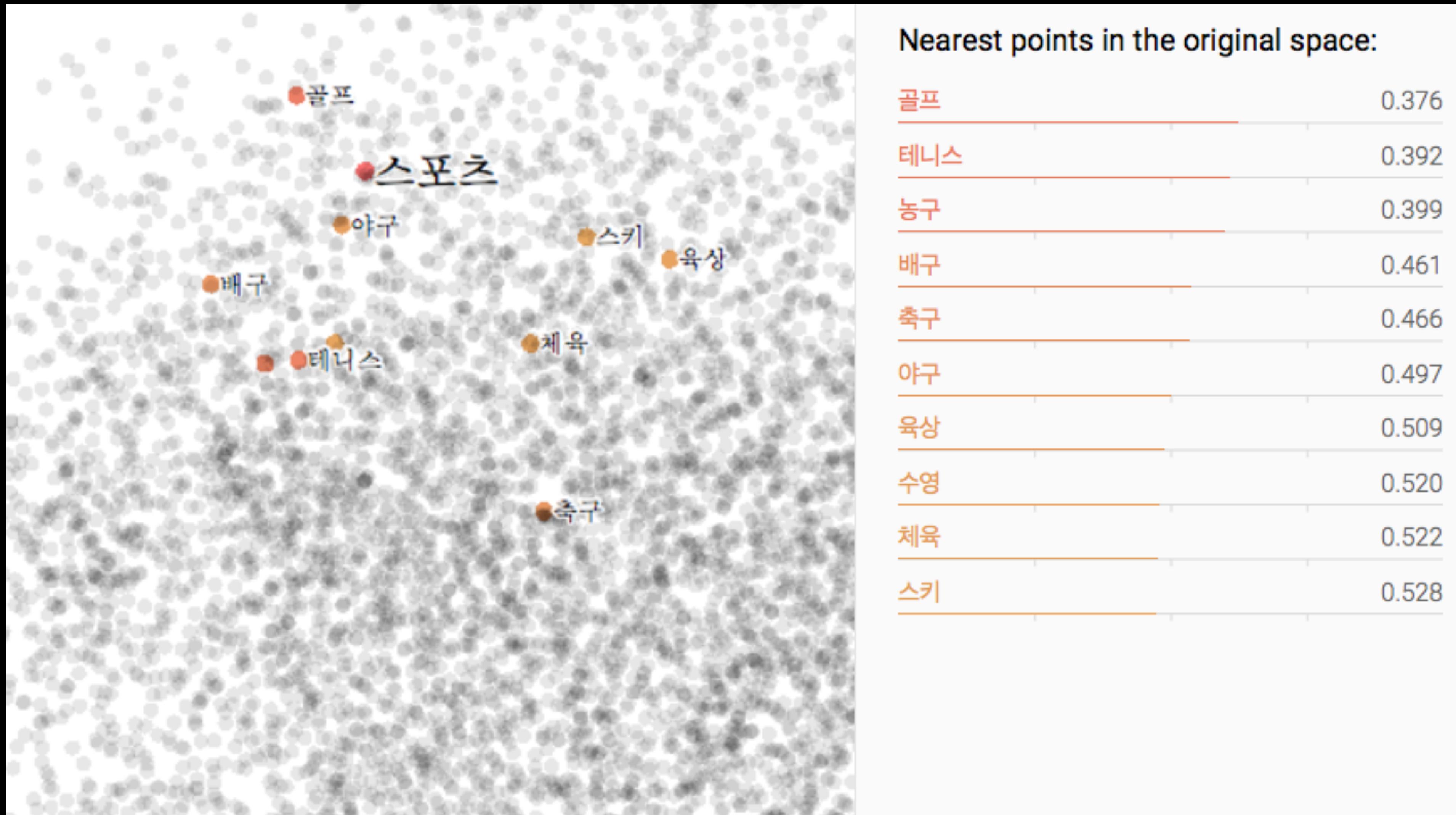
Synonym + Derivative words



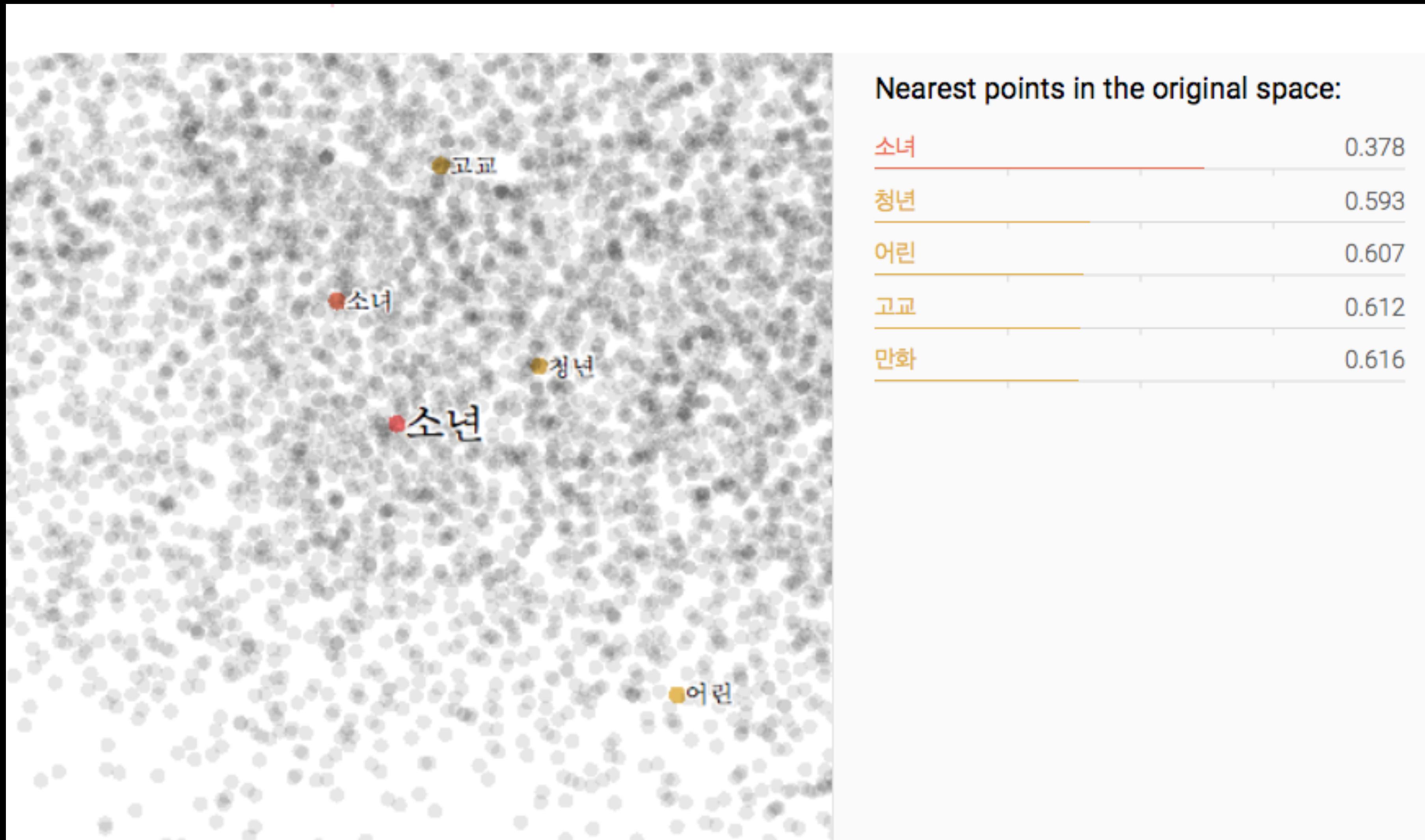
Similiar Types + Derivative words



a Part of

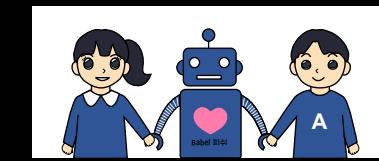


Antonym



Word Representaion in NLP

Natural Language Processing



Dictionary vs Embedding

남자² (男子) ★★★ 🔍 +

(명사)

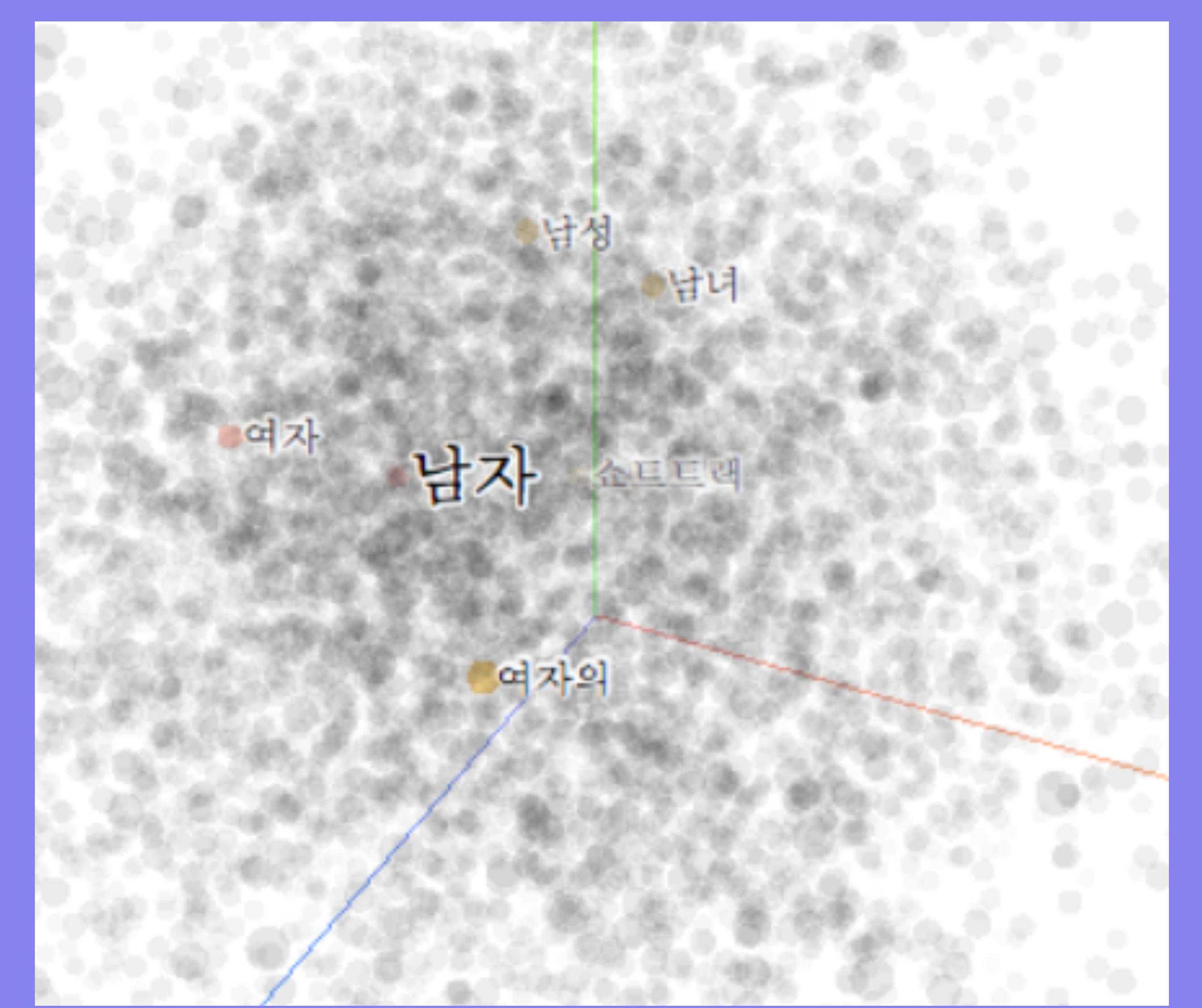
1. 남성으로 태어난 사람.
2. 사내다운 사내.
3. 한 여자의 남편이나 애인을 이르는 말.

여자² (女子) ★★★ 🔍 +

(명사)

1. 여성으로 태어난 사람.
2. 여자다운 여자.
3. 한 남자의 아내나 애인을 이르는 말.

유의어 : 부녀자, 부녀², 아낙



Dictionary vs Embedding

남자² (男子) ★★★ 🔍 +

(명사)

1. 남성으로 태어난 사람.
2. 사내다운 사내.
3. 한 여자의 남편이나 애인을 이르는 말.

늑대 [늑때] ★ 🔍 +

(명사)

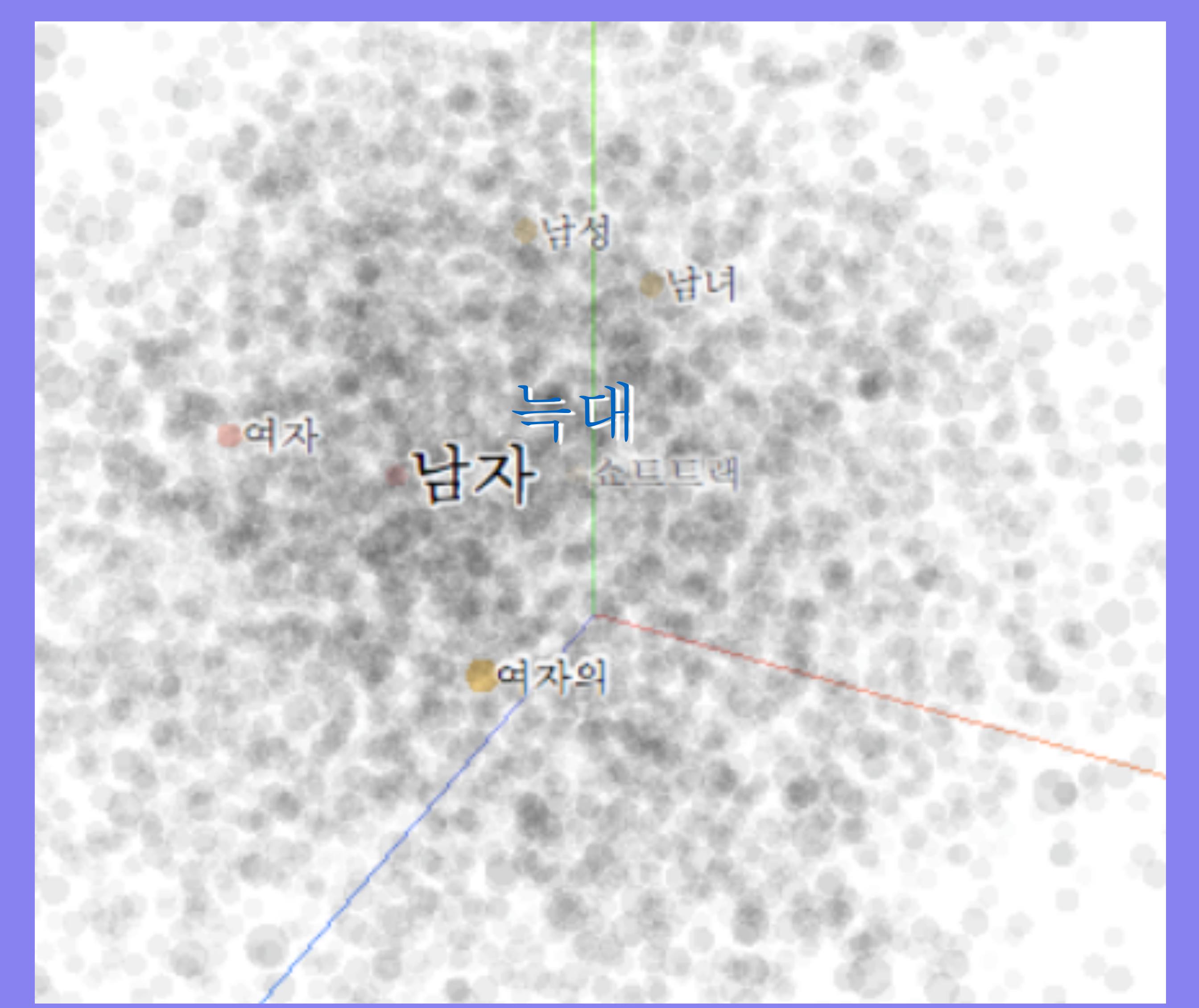
1. <동물>갯과의 포유류. 몸의 길이는 120cm, 꼬리는 35cm,
2. 여자에게 음흉한 마음을 품은 남자를 비유적으로 이르는 말.

여자² (女子) ★★★ 🔍 +

(명사)

1. 여성으로 태어난 사람.
2. 여자다운 여자.
3. 한 남자의 아내나 애인을 이르는 말.

유의어 : 부녀자, 부녀², 아낙



Dictionary based Representation

남자² (男子) ★★★ 🔍 +

[명사]

- 남성으로 태어난 사람.
- 사내다운 사내.
- 한 여자의 남편이나 애인을 이르는 말.

늑대 [늑때] ★ 🔍 +

[명사]

- 〈동물〉갓과의 포유류. 몸의 길이는 120cm, 꼬리는 35cm, 등 높이는 100cm.
- 여자에게 음흉한 마음을 품은 남자를 비유적으로 이르는 말.

여자² (女子) ★★★ 🔍 +

[명사]

- 여성으로 태어난 사람.
- 여자다운 여자.
- 한 남자의 아내나 애인을 이르는 말.

유의어 : 부녀자, 부녀², 아낙

	남성	사람	사내	여자	남편	애인	여성	남자	아내	애인	포유류	몸	음흉	마음
--	----	----	----	----	----	----	----	----	----	----	-----	---	----	----

“남자”	0	0	0	0	0	0								
“늑대”								0			0	0	0	0
“여자”		0		0			0	0	0	0				

Character based Representation

	남	여	자	늑	대
“남자”	0		0		
“늑대”				0	0
“여자”		0	0		

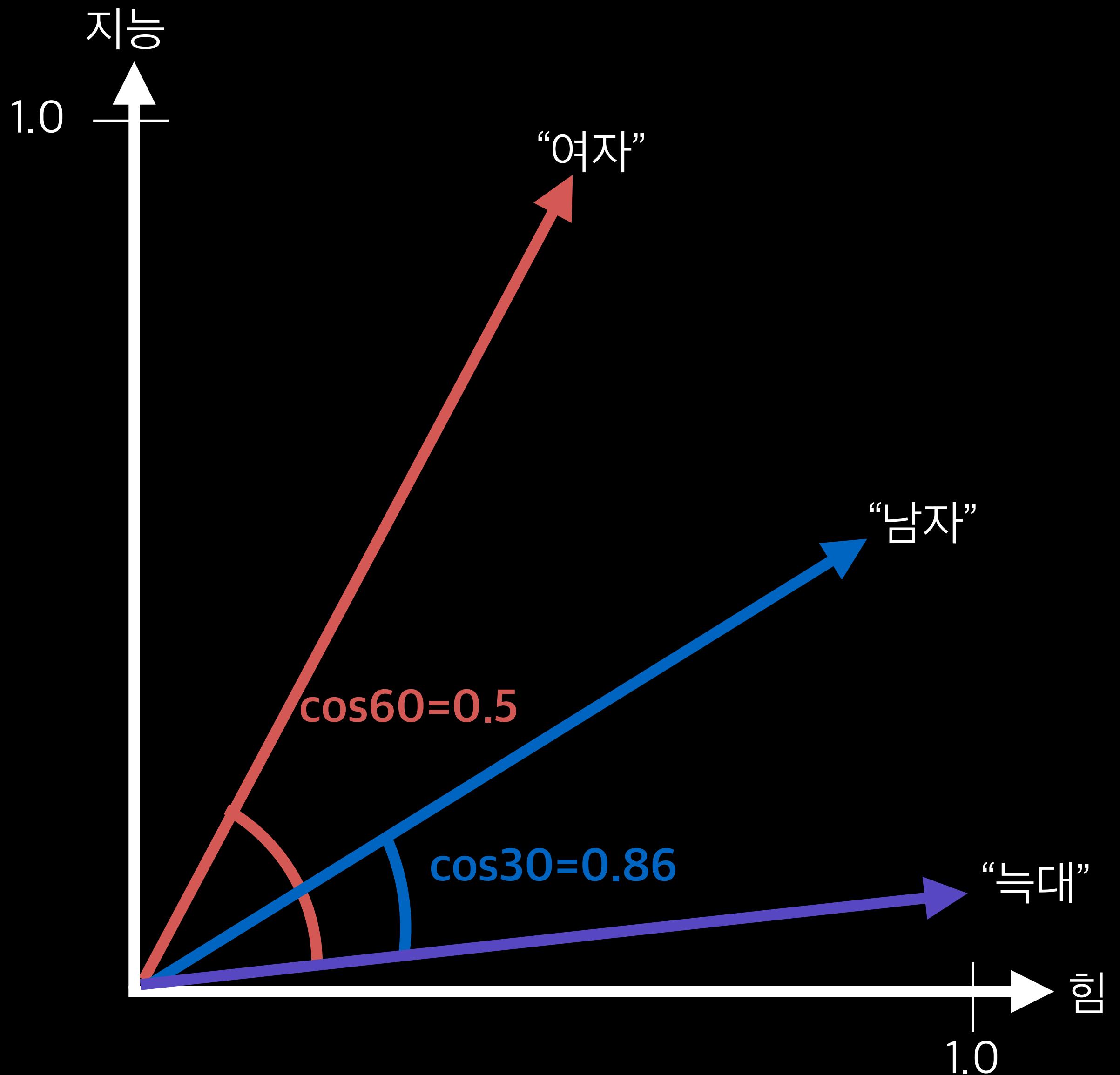
Vector based Representation=Embedding

	힘	지능	덩치	고음
“남자”	0.9	0.5	0.8	0.2
“늑대”	1.0	0.1	0.9	0.7
“여자”	0.4	0.9	0.5	1.0

Word Similarity

	힘	지능
“남자”	0.9	0.5
“늑대”	1.0	0.1
“여자”	0.4	0.9

$$\text{sim}(\text{“남자”, “늑대”}) > \text{sim}(\text{“여자”, “늑대”})$$

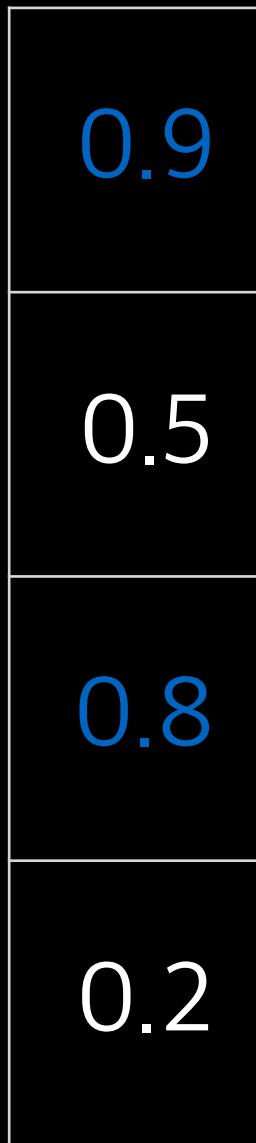


남자 ≈ 늑대

Embedding

Text

남자



Find Opposite word
(ML / DL)

Embedding

0.4

0.9

0.5

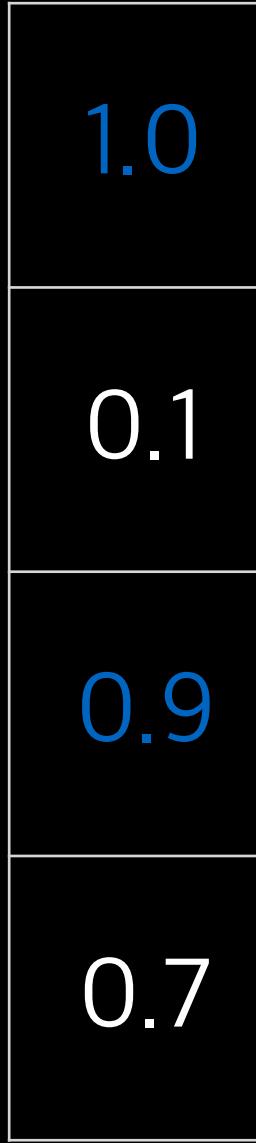
1.0

Text

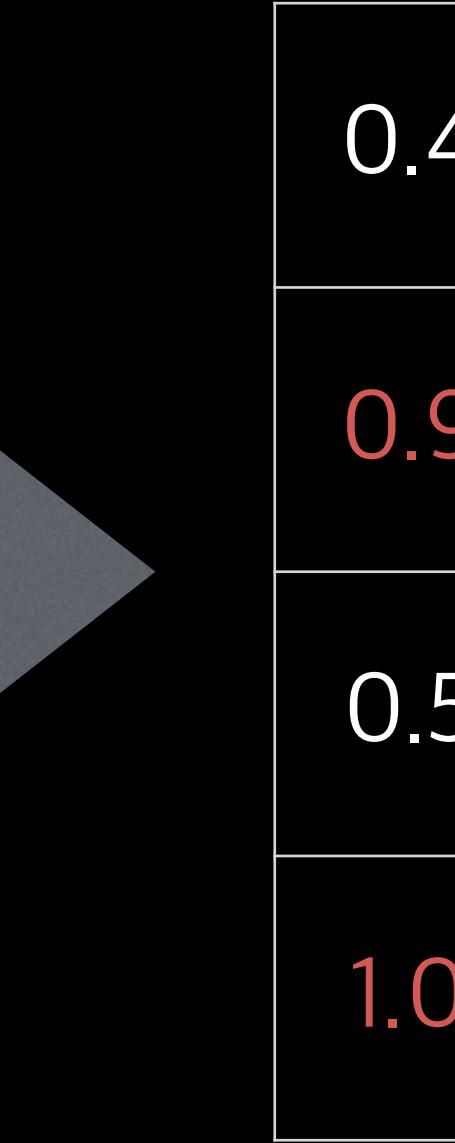
여자

Text

늑대



Find Opposite word
(ML / DL)

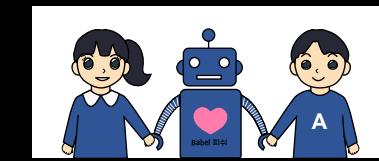


Text

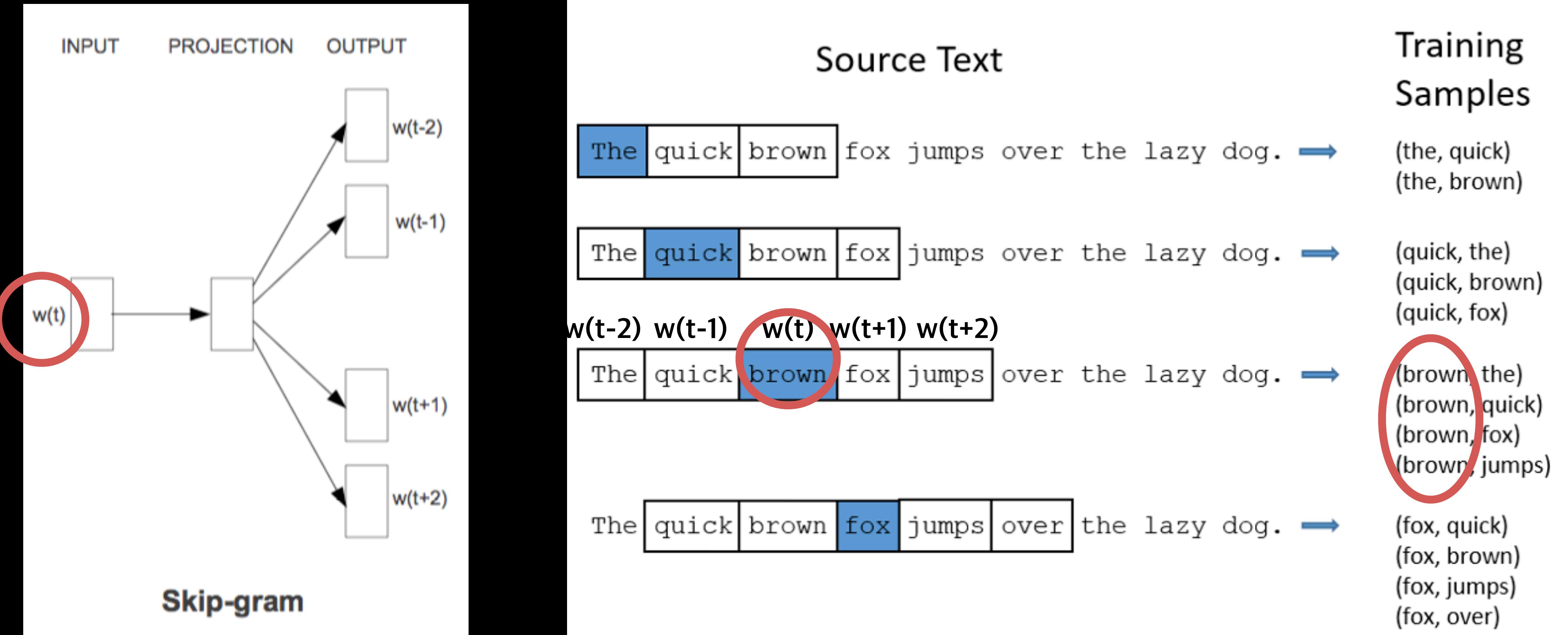
여자

Language Models

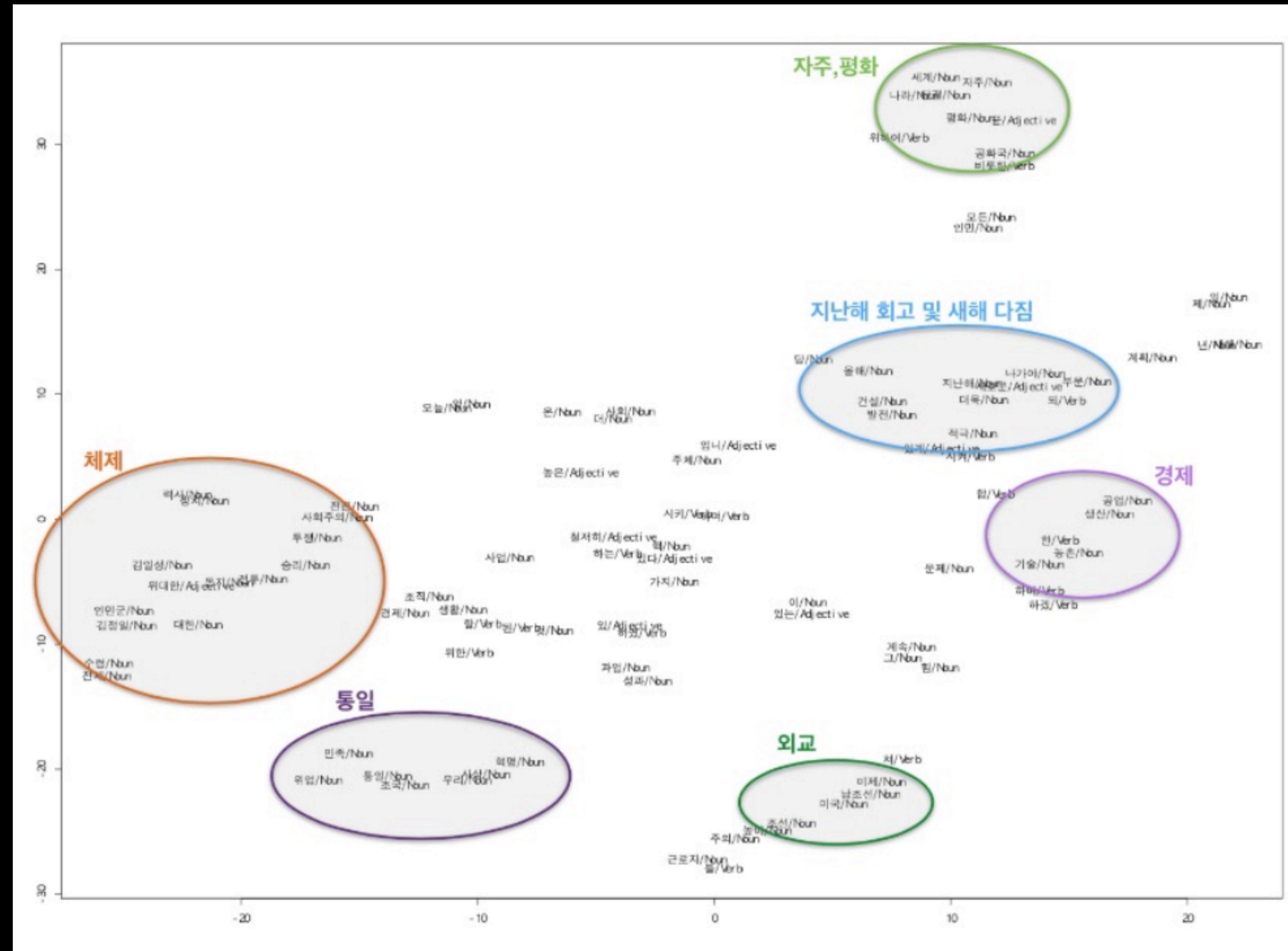
N-Gram, Word2vec, Fasttext, Glove...



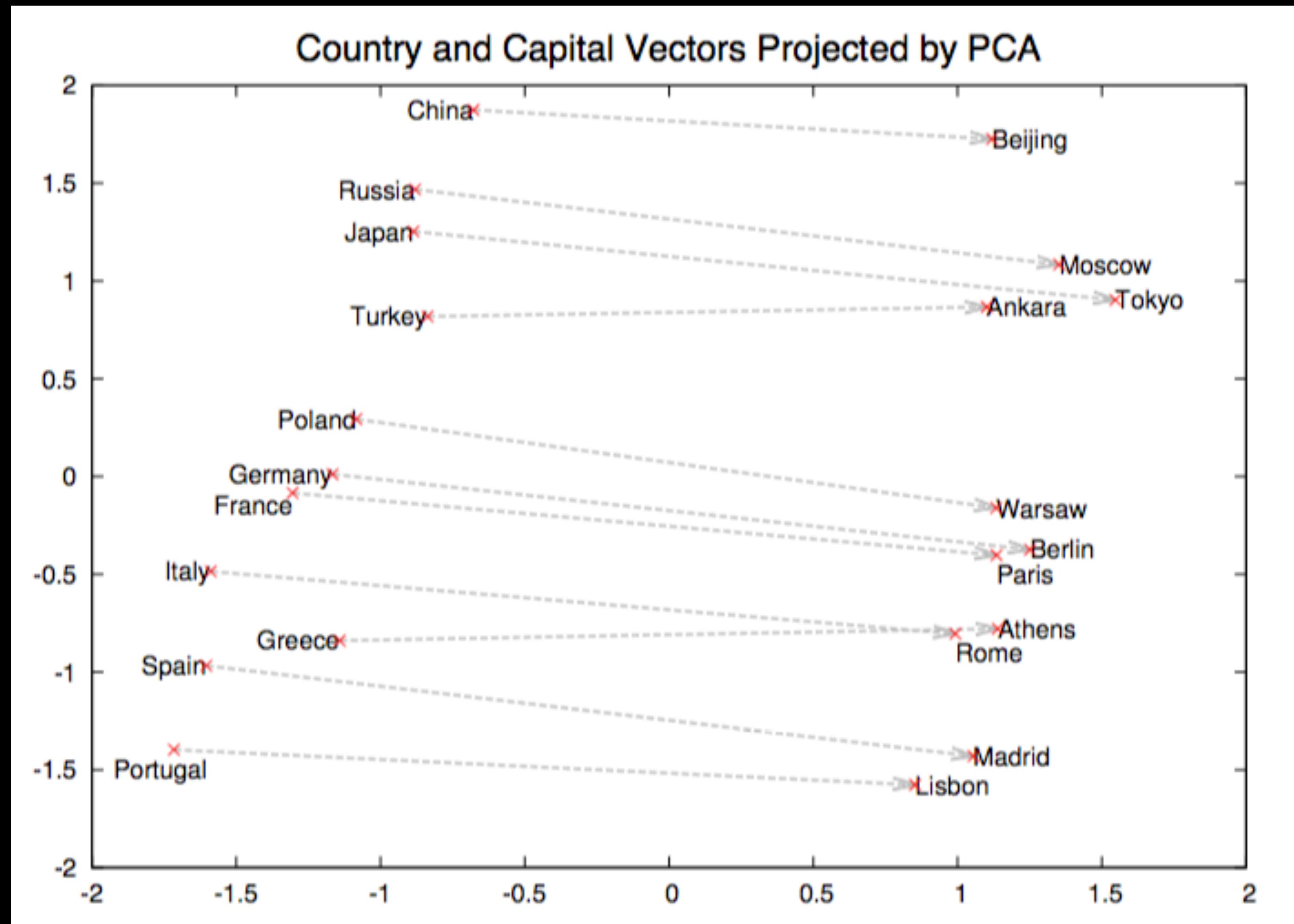
Word2vec Skip-Gram Model



Word Similarity in word2vec

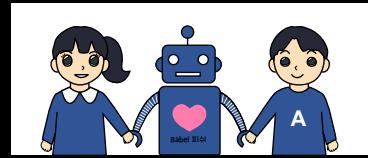


Word Relationship in word2vec



Create your own word embedding

with open-sources



Prepare Corpus

- Download Korean Wikipedia dump file
 - <https://dumps.wikimedia.org/kowiki/20180220/>
- Parse dump file(mediawiki format) to text file
 - <https://pypi.python.org/pypi/mediawiki-parser/>

mediawiki-parser 0.4.1

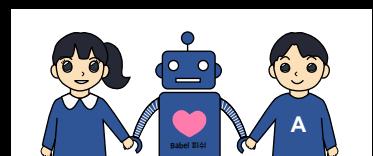
A parser for the MediaWiki syntax, based on Pijnu.

build passing

[Download mediawiki-parser-0.4.1.tar.gz](#)

Presentation

This is a parser for MediaWiki's (MW) syntax. Its goal is to transform wikitext into an abstract syntax tree (AST) and then render this AST into various formats such as plain text and HTML.



Word2vec open source

- <https://github.com/theeluwin/pytorch-sgns>
 - Fast training Word2vec Skip-gram with Pytorch
 - Fully class based source codes



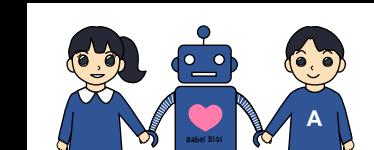
PyTorch SGNS

Word2Vec's **SkipGramNegativeSampling** in Python.

Yet another but quite general [negative sampling loss](#) implemented in [PyTorch](#).

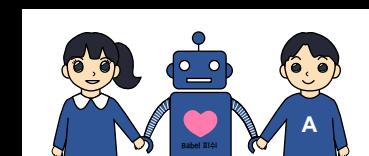
It can be used with ANY embedding scheme! Pretty fast, I bet.

```
vocab_size = 20000
word2vec = Word2Vec(vocab_size=vocab_size, embedding_size=300)
sgns = SGNS(embedding=word2vec, vocab_size=vocab_size, n_negs=20)
optim = Adam(sgns.parameters())
for batch, (iword, owords) in enumerate(dataloader):
    loss = sgns(iword, owords)
    optim.zero_grad()
    loss.backward()
    optim.step()
```



Korean Embedding on Github

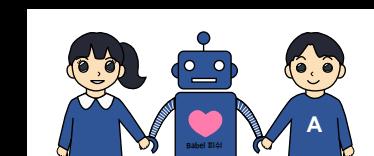
- <https://github.com/bage79/word2vec4kor>
- File format for Tensorboard Embedding(Projector)
- Corpus: <https://ko.wikipedia.org>
 - Total sentences: about 3,115,431
 - Total Unique words: 10,000
 - Tokenized: white-space
 - Embedding Dimension: 300
 - Skip-Gram + Negative Sampling + Subsampling



Q & A

[바벨피쉬 모임] Word2vec Embedding 실습

https://www.facebook.com/events/529492287451286/?notif_t=event_description_mention¬if_id=1519395371494159



Tensorboard Demo

