

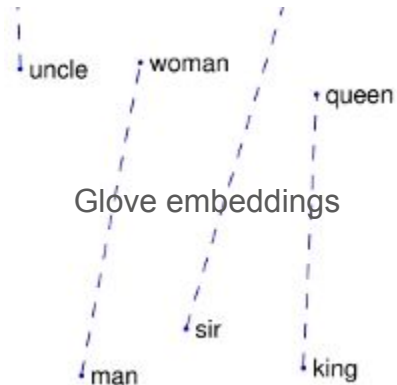
Word & System Analogies Retrieval

Akshara Prabhakar and Margarita Belova

<https://github.com/across-stars/pinecone/tree/master>

Methodology overview:

- Take datasets showing **word**, **relation** and **system analogies**
- Index word and relation embeddings in Pinecone
 - GloVe (glove42) – word embeddings (“kant” → vector)
 - ReBERT – word pair, or relation embeddings (“kant__philosopher” → vector)
- Use KNN to retrieve analogical words / pairs



Results:

- Compared knn analogical retrieval performance: *Glove* vs *ReBERT* embeddings
- Extended the method for complex system analogies retrieval

Word Analogy

$$w_1 : w_2 :: w_3 : w_4$$

For Glove embeddings check:

$$w_4 = w_3 + (w_2 - w_1) + \delta,$$

Algorithm:

- Retrieve w_1, w_2, w_3 from Pinecone db;
- Retrieve $\text{knn}(w_3 + w_2 - w_1)$
- $\text{label}(w_4)$ – ground truth
- Check: $\text{label}(w_4) \in \text{knn}(w_3 + w_2 - w_1)$

Glove database:

- 1_894_411 embeddings;
- dim = 300;
- Metric: cosine

source:

<https://nlp.stanford.edu/projects/glove/>
Common Crawl (42B tokens, 1.9M
vocab, uncased)

w/o non-ASCII symbols

```
base,target,top0,top1,top2,top3,top4,k_order,k_score
belgrade,serbia,belgrade,serbia,lebanon,macedonia,kosovo,1,0.678768814
manila,philippines,manila,philippines,philippine,singapore,cebu,1,0.738674104
paris,france,paris,france,french,belgium,pierre,1,0.692115188
lilongwe,malawi,malawi,lilongwe,gabon,botswana,zambia,0,0.703917921
```

Relation Analogy $w1_w2 :: w3_w4$

For Relbert embeddings check:

$w3_w4 \in \text{knn}(w1_w2)$

Algorithm:

- Retrieve $w1_w2$ from Pinecone db;
- Retrieve $\text{knn}(w1_w2)$
- $\text{label}(w3_w4)$ – ground truth
- Check: $\text{label}(w3_w4) \in \text{knn}(w1_w2)$

RelBERT database:

- 104_973 embeddings;
- dim = 1024;
- Metric: cosine

```
base,target,top0,top1,top2,top3,top4,k_order,k_score
hitler__dictator,strauss__composer,hitler__dictator,napoleon__emperor,charlatan__impostor,superman__emperor,truman__president,,
rousseau__writer,hegel__philosopher,rousseau__writer,tolstoi__novelist,andersen__writer,raphael__painter,wagner__composer,,
kant__philosopher,lincoln__president,kant__philosopher,mencius__philosopher,locke__philosopher,raphael__painter,tolstoi__novelist,,
stalin__dictator,hawking__physicist,stalin__dictator,newton__scientist,strauss__composer,stalin__napoleon,kepler__mathematician,,
```

System Analogy: definition

$$\langle concept_1 \rightarrow relation_1 \rightarrow concept_2 \rangle,$$
$$\{< u_1 \rightarrow r_1 \rightarrow u_2 >, < u_1 \rightarrow r_2 \rightarrow u_3 >, \dots\} :: \{< v_1 \rightarrow r_1 \rightarrow v_2 >, < v_1 \rightarrow r_3 \rightarrow v_3 >, \dots\}$$

Examples (from SCAN_dataset):

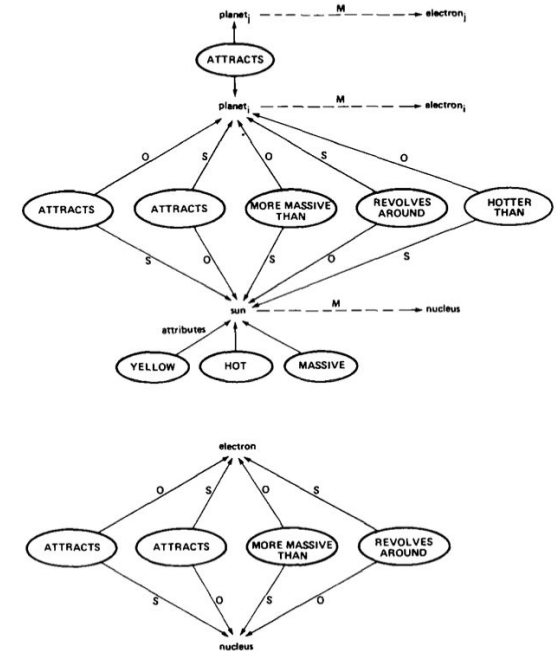
Solar system -- Atom:

Solar system:

- Sun
- Planet
- Mass
- Attracts
- Revolves
- Gravity

Atom:

- Nucleus
- Electron
- Charge
- Attracts
- Revolves
- Electromagnetism



Genter et. al. (1983)

System Analogy: task

$\langle concept_1 \rightarrow relation_1 \rightarrow concept_2 \rangle,$

$\{\langle u_1 \rightarrow r_1 \rightarrow u_2 \rangle, \langle u_1 \rightarrow r_2 \rightarrow u_3 \rangle, \dots\} :: \{\langle v_1 \rightarrow r_1 \rightarrow v_2 \rangle, \langle v_1 \rightarrow r_3 \rightarrow v_3 \rangle, \dots\}$

Solar system - Atom:

Solar system:

- Sun
- Planet
- Mass
- Attracts
- Revolves
- Gravity

Atom:

- Nucleus
- Electron
- Charge
- Attracts
- Revolves
- Electromagnetism

systems

entities

Task:

Given a **system** and its **entities**, retrieve the most analogous **system**

System analogy: methodology

1	target	source	targ_word	src_word
2	atom	solar system	nucleus	sun
3	atom	solar system	electron	planet
4	atom	solar system	charge	mass
5	atom	solar system	attracts	attracts
6	atom	solar system	revolves	revolves
7	atom	solar system	electromagnetism	gravity

GloVe + RelBERT for relations = System analogy

- Retrieve RelBERT for every pair source__source_word

Relbert(atom__nucleus), Relbert(atom__electron), Relbert(atom__charge)... → [w1, w2, w3...]

- Retrieve **labels** for knn(w1), knn(w2), knn(w3)... → <{ }, { }, { }...> – sets of labels pairs

E.g. knn(atom__charge) = {object__light, blossom__water, vital__energy...}

- Parse labels: pick the **most frequent label** across entities – this is the analogical system (*target*)

We have source (**s**) and target (**t**) labels that are words.

We can retrieve from Glove $r = t - s$ – relation and save it to our *library of relations*.

Evaluation Setup

Datasets

- BATS - 1799 example
- GOOGLE - 500 examples
- SAT - 337 examples
- SCAN (system analogy) - 45 examples

Metric

- Accuracy @ k

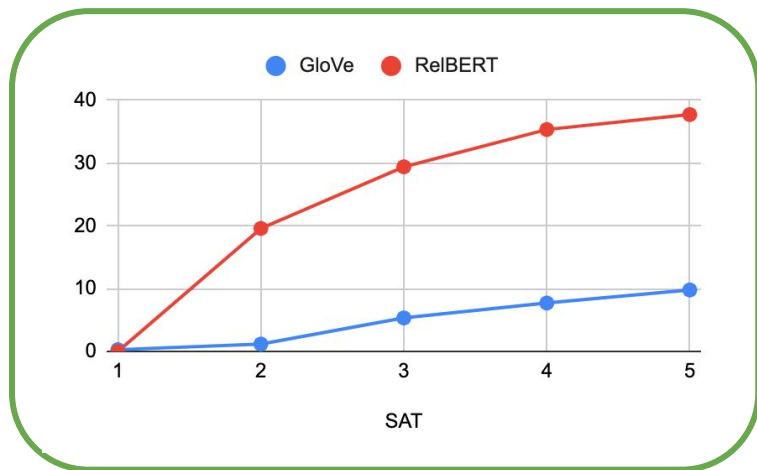
ReBERT

- Distilling relation embeddings from RoBERTa (roberta-large)
- Manual Prompting
- Triplet loss to finetune
 - x_a (anchor), x_p (positive), x_n (negative)

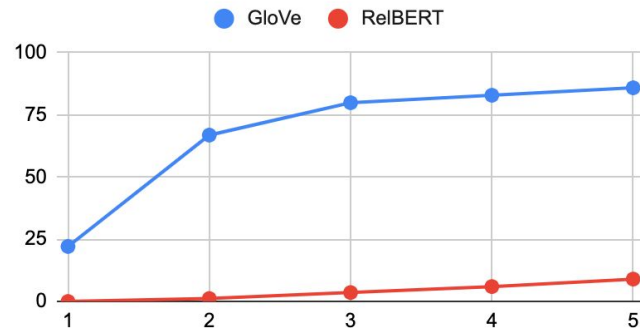
$$L_t = \max(0, \|x_a - x_p\| - \|x_a - x_n\| + \varepsilon)$$

1. Today, I finally discovered the relation between **[h]** and **[t]** : **[h]** is the <mask> of **[t]**
2. Today, I finally discovered the relation between **[h]** and **[t]** : **[t]** is **[h]**'s <mask>
3. Today, I finally discovered the relation between **[h]** and **[t]** : <mask>
4. I wasn't aware of this relationship, but I just read in the encyclopedia that **[h]** is the <mask> of **[t]**
5. I wasn't aware of this relationship, but I just read in the encyclopedia that **[t]** is **[h]**'s <mask>

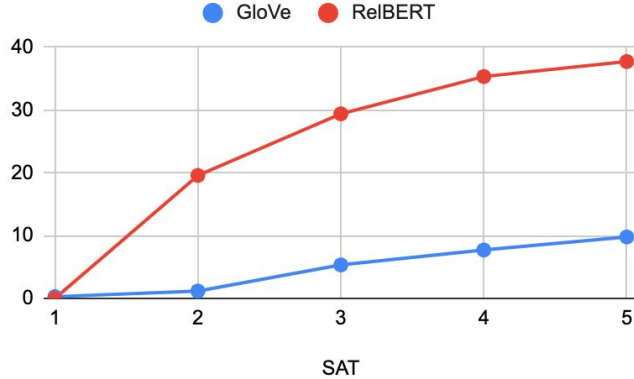
Word Analogy Results



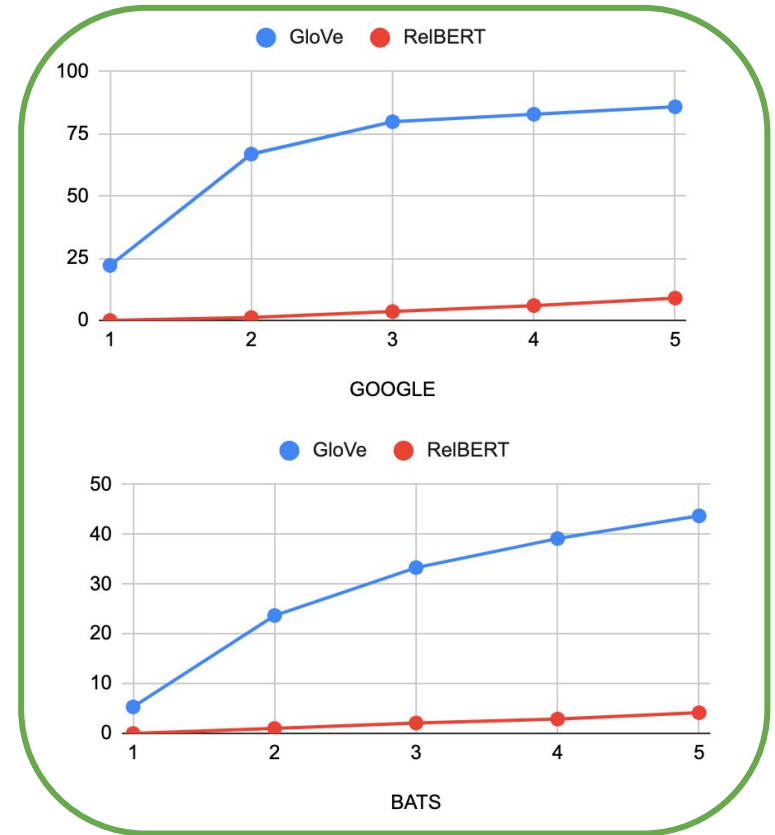
- Questions from real SAT exams



Word Analogy Results



- GOOGLE dataset has been shown to be biased towards word similarity



System Analogy Results

- Hard accuracy: 15.56%
 - Considering exact match only

Some predictions

`reasons_for_a_theory` \longleftrightarrow `grounds_for_a_building`

`gas_molecules` \longleftrightarrow `billiard_balls`

`respiration` \longleftrightarrow `combustion`

`bacterial_mutation` \longleftrightarrow `distortion`

NEW System Analogy: Library of relations (created early)

- 443 relations
- SVD of them to remove noise
- Take a random word \mathbf{w} , add a relation \mathbf{r} from library and retrieve $\text{knn}(\mathbf{w}+\mathbf{r})$ from Glove
- Hopefully, it gives some analogy...

NEW System Analogy: Library of relations (created early)

```
"pain": [  
  "headache",  
  "cause",  
  "discomfort",  
  "feel",  
  "symptoms",  
  "suffering",  
  "back",  
  "pain",  
  "pains"  
],
```

```
"professor": [  
  "university",  
  "dean",  
  "emeritus",  
  "lecturer",  
  "professor",  
  "prof."  
],
```

```
"music": [  
  "songs",  
  "music",  
  "dance",  
  "song"  
],
```

SVD of relation
matrix rank 10

SVD of relation
matrix rank 80

```
"pain": [  
  "anxiety",  
  "ache",  
  "low",  
  "pain",  
  "sun",  
  "hot",  
  "burn",  
  "suffering",  
  "headache",  
  "fire",  
  "surgery",  
  "symptoms",  
  "worse",  
  "high",  
  "agony",  
  "discomfort",  
  "pains",  
  "burning",  
  "fuel"  
],
```

```
"music": [  
  "sun",  
  "songs",  
  "tracks",  
  "song",  
  "sound",  
  "hot",  
  "jazz",  
  "musical",  
  "mp3",  
  "fire",  
  "records",  
  "music",  
  "listen",  
  "cd",  
  "high",  
  "dance",  
  "air",  
  "concert",  
  "fuel",  
  "artists",  
  "tune",  
  "well"  
],
```