### embedding\_termfreq embedding\_QSO **→ input → input** data/current\_sequences.csv data/current\_sequences.csv data/current\_words.csv output → postprocessing/QSO.csv postprocessing/termfreq.csv true\_similarity\_syntax **→ input** data/current\_sequences.csv output → similarity/true\_syntax.csv true\_similarity\_semantics **→ input** data/current\_sequences.csv output → similarity/true\_semantics\_BP.csv similarity/true\_semantics\_CC.csv similarity/true\_semantics\_MF.csv evaluation **→ input** similarity/QSO.csv similarity/QSO\_CCR.csv similarity/biophys.csv similarity/biophys\_CCR.csv similarity/biophys\_SIF.csv similarity/biophys\_SIF\_CCR.csv similarity/biophys\_TFIDF.csv similarity/biophys\_TFIDF\_CCR.csv similarity/random.csv similarity/random\_CCR.csv similarity/random\_SIF.csv similarity/random\_SIF\_CCR.csv similarity/random\_TFIDF.csv similarity/random\_TFIDF\_CCR.csv similarity/seq2vec.csv similarity/seq2vec\_CCR.csv similarity/seq2vec\_SIF.csv similarity/seq2vec\_SIF\_CCR.csv similarity/seq2vec\_TFIDF.csv similarity/seq2vec\_TFIDF\_CCR.csv similarity/termfreq.csv similarity/termfreq\_CCR.csv similarity/true\_semantics\_BP.csv similarity/true\_semantics\_CC.csv similarity/true\_semantics\_MF.csv similarity/true\_syntax.csv output → similarity/scores/QSO.txt similarity/scores/QSO\_CCR.txt similarity/scores/biophys.txt similarity/scores/biophys\_CCR.txt similarity/scores/biophys\_SIF.txt similarity/scores/biophys\_SIF\_CCR.txt similarity/scores/biophys\_TFIDF.txt similarity/scores/biophys\_TFIDF\_CCR.txt similarity/scores/random.txt similarity/scores/random\_CCR.txt similarity/scores/random\_SIF.txt similarity/scores/random\_SIF\_CCR.txt similarity/scores/random\_TFIDF.txt similarity/scores/random\_TFIDF\_CCR.txt similarity/scores/seq2vec.txt similarity/scores/seq2vec\_CCR.txt similarity/scores/seq2vec\_SIF.txt similarity/scores/seq2vec\_SIF\_CCR.txt similarity/scores/seq2vec\_TFIDF.txt similarity/scores/seq2vec\_TFIDF\_CCR.txt similarity/scores/termfreq.txt similarity/scores/termfreq\_CCR.txt downstream **∽ input** similarity/scores/QSO.txt similarity/scores/QSO\_CCR.txt similarity/scores/biophys.txt similarity/scores/biophys\_CCR.txt similarity/scores/biophys\_SIF.txt similarity/scores/biophys\_SIF\_CCR.txt similarity/scores/biophys\_TFIDF.txt similarity/scores/biophys\_TFIDF\_CCR.txt similarity/scores/random.txt similarity/scores/random\_CCR.txt similarity/scores/random\_SIF.txt similarity/scores/random\_SIF\_CCR.txt similarity/scores/random\_TFIDF.txt similarity/scores/random\_TFIDF\_CCR.txt similarity/scores/seq2vec.txt similarity/scores/seq2vec\_CCR.txt similarity/scores/seq2vec\_SIF.txt similarity/scores/seq2vec\_SIF\_CCR.txt similarity/scores/seq2vec\_TFIDF.txt similarity/scores/seq2vec\_TFIDF\_CCR.txt similarity/scores/termfreq.txt similarity/scores/termfreq\_CCR.txt output → scores\_w5\_d100.csv all **→ input** scores\_w5\_d100.csv similarity/scores/QSO.txt similarity/scores/QSO\_CCR.txt similarity/scores/biophys.txt similarity/scores/biophys\_CCR.txt similarity/scores/biophys\_SIF.txt similarity/scores/biophys\_SIF\_CCR.txt similarity/scores/biophys\_TFIDF.txt similarity/scores/biophys\_TFIDF\_CCR.txt similarity/scores/random.txt similarity/scores/random\_CCR.txt similarity/scores/random\_SIF.txt similarity/scores/random\_SIF\_CCR.txt similarity/scores/random\_TFIDF.txt similarity/scores/random\_TFIDF\_CCR.txt similarity/scores/seq2vec.txt similarity/scores/seq2vec\_CCR.txt similarity/scores/seq2vec\_SIF.txt similarity/scores/seq2vec\_SIF\_CCR.txt similarity/scores/seq2vec\_TFIDF.txt similarity/scores/seq2vec\_TFIDF\_CCR.txt similarity/scores/termfreq.txt

similarity/scores/termfreq\_CCR.txt

seq2vec\_weighting embedding\_seq2vec **→ input** data/TF\_IDF.csv data/current\_sequences.csv data/current\_words.csv data/ids\_hp\_w5.csv

sampling

data/proteome\_human.csv

data/current\_sequences.csv

data/current\_words.csv

output →

data/words\_hp.csv

**→ input** 

output →

**→ input** 

output →

data/current\_sequences.csv

data/current\_words.csv

data/weights\_w5\_d100.csv

postprocessing/seq2vec.csv

data/ids\_hp\_w5.csv

data/weights\_w5\_d100.csv output → postprocessing/seq2vec\_SIF.csv

postprocessing/seq2vec\_TFIDF.csv

# embedding\_biophys

**→ input** data/current\_sequences.csv output →

postprocessing/biophys.csv

biophys\_weighting **→ input** data/TF\_IDF.csv data/current\_sequences.csv data/current\_words.csv data/ids\_hp\_w5.csv data/weights\_w5\_d100.csv

postprocessing/random.csv output → postprocessing/biophys\_SIF.csv

postprocessing/biophys\_TFIDF.csv

embedding\_random **→ input** data/current\_sequences.csv output →

**→ input** data/TF\_IDF.csv data/current\_sequences.csv data/current\_words.csv data/ids\_hp\_w5.csv data/weights\_w5\_d100.csv output → postprocessing/random\_SIF.csv postprocessing/random\_TFIDF.csv

random\_weighting

## CCR

**→ input** postprocessing/QS0.csv postprocessing/biophys.csv postprocessing/biophys\_SIF.csv postprocessing/biophys\_TFIDF.csv postprocessing/random.csv postprocessing/random\_SIF.csv

postprocessing/random\_TFIDF.csv postprocessing/seq2vec.csv postprocessing/seq2vec\_SIF.csv postprocessing/seq2vec\_TFIDF.csv postprocessing/termfreq.csv

output → postprocessing/QSO\_CCR.csv postprocessing/biophys\_CCR.csv postprocessing/biophys\_SIF\_CCR.csv postprocessing/biophys\_TFIDF\_CCR.csv postprocessing/random\_CCR.csv postprocessing/random\_SIF\_CCR.csv postprocessing/random\_TFIDF\_CCR.csv postprocessing/seq2vec\_CCR.csv postprocessing/seq2vec\_SIF\_CCR.csv postprocessing/seq2vec\_TFIDF\_CCR.csv

postprocessing/termfreq\_CCR.csv

similarity **→ input** postprocessing/QSO.csv postprocessing/QS0\_CCR.csv postprocessing/biophys.csv postprocessing/biophys\_CCR.csv postprocessing/biophys\_SIF.csv postprocessing/biophys\_SIF\_CCR.csv postprocessing/biophys\_TFIDF.csv postprocessing/biophys\_TFIDF\_CCR.csv postprocessing/random.csv postprocessing/random\_CCR.csv postprocessing/random\_SIF.csv postprocessing/random\_SIF\_CCR.csv postprocessing/random\_TFIDF.csv postprocessing/random\_TFIDF\_CCR.csv postprocessing/seq2vec.csv postprocessing/seq2vec\_CCR.csv postprocessing/seq2vec\_SIF.csv postprocessing/seq2vec\_SIF\_CCR.csv postprocessing/seq2vec\_TFIDF.csv

postprocessing/termfreq\_CCR.csv output → similarity/QSO.csv similarity/QSO\_CCR.csv similarity/biophys.csv similarity/biophys\_CCR.csv

postprocessing/seq2vec\_TFIDF\_CCR.csv

postprocessing/termfreq.csv

similarity/biophys\_SIF.csv similarity/biophys\_SIF\_CCR.csv similarity/biophys\_TFIDF.csv similarity/biophys\_TFIDF\_CCR.csv similarity/random.csv similarity/random\_CCR.csv similarity/random\_SIF.csv similarity/random\_SIF\_CCR.csv similarity/random\_TFIDF.csv similarity/random\_TFIDF\_CCR.csv similarity/seq2vec.csv similarity/seq2vec\_CCR.csv similarity/seq2vec\_SIF.csv similarity/seq2vec\_SIF\_CCR.csv similarity/seq2vec\_TFIDF.csv similarity/seq2vec\_TFIDF\_CCR.csv similarity/termfreq.csv similarity/termfreq\_CCR.csv