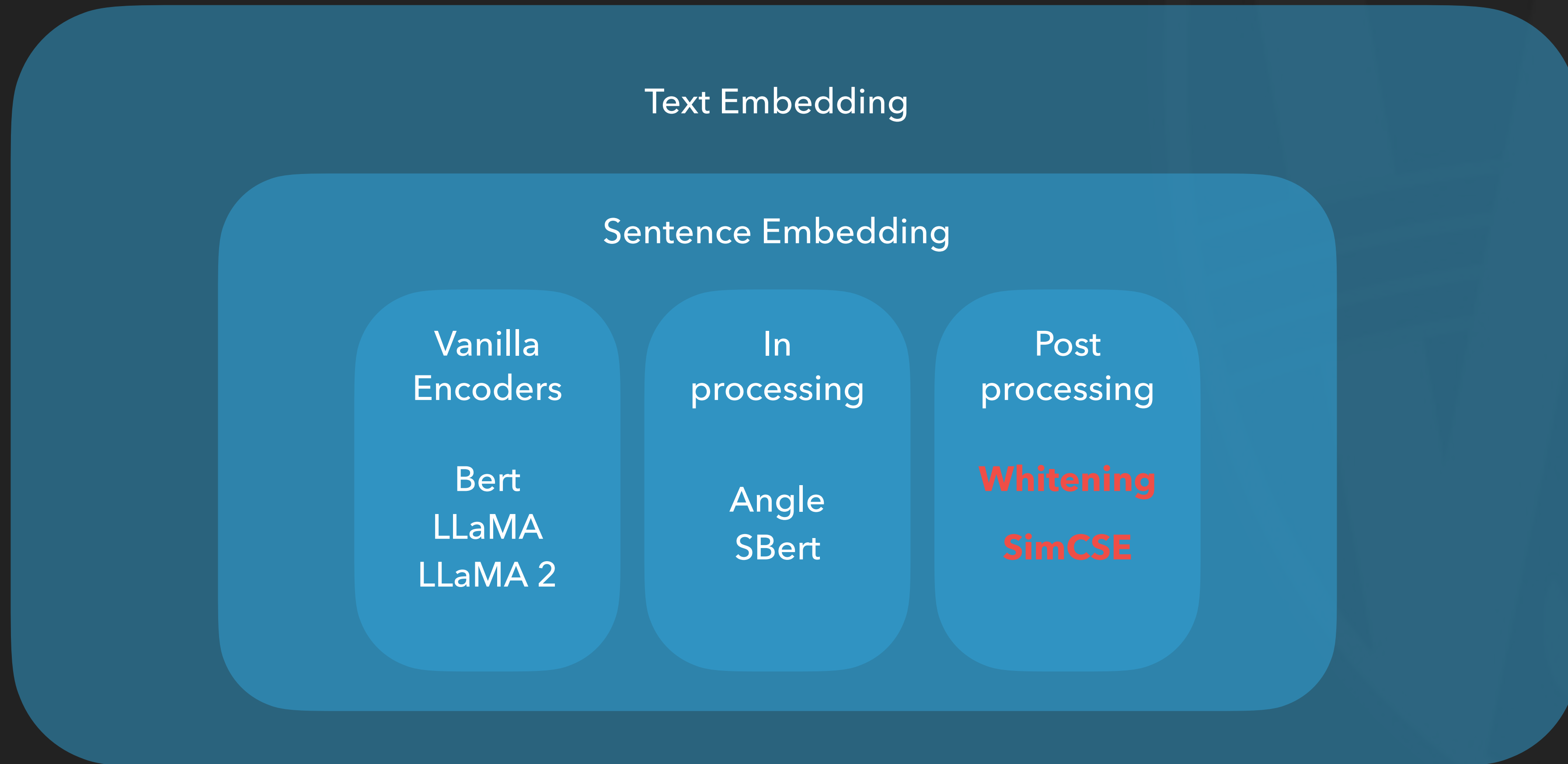FEB 2, 2024

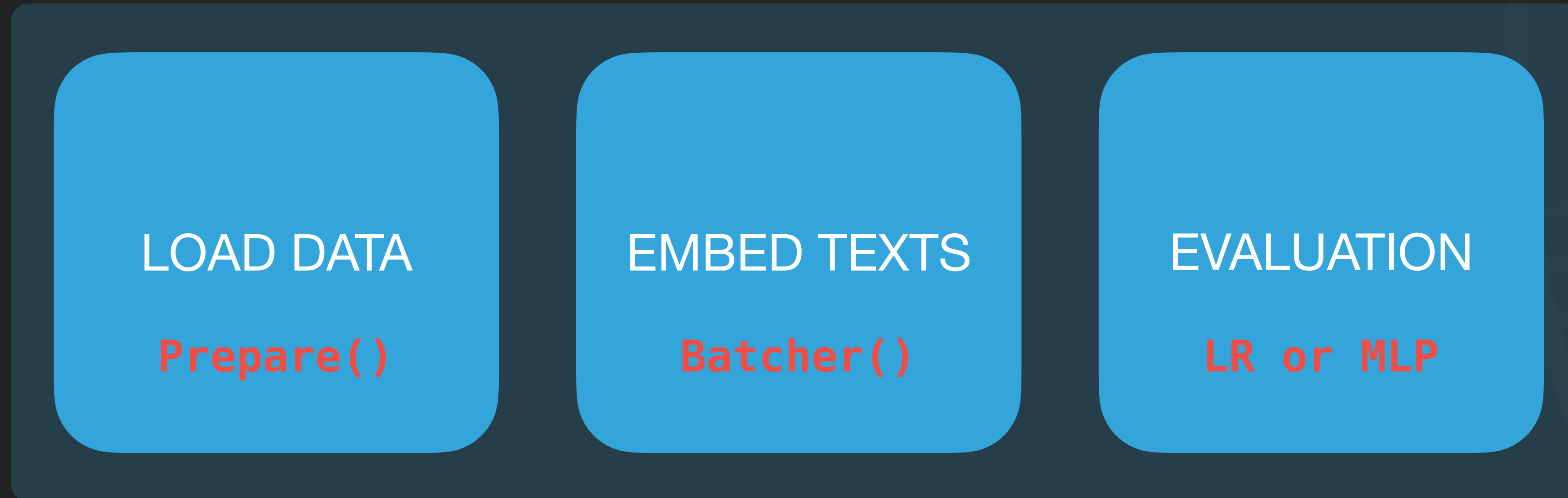# EFFECT OF WHITENING ON TEXT CLASSIFICATION

# GOAL

# WHITENING

▸ In which scenarios does the process of whitening demonstrate effectiveness?

  ▸ Which Data?

  ▸ Which Tasks?

  ▸ Which Encoders?

# WHITENING

▸ Which Data?

  ▸ MR, CR, STS, …

▸ Which Tasks?

  ▸ Sentence Similarity, Text Classification, …

▸ Which Encoders?

  ▸ LLaMA, Bert, ChatGPT, …

# SENTEVAL

Pipeline of SentEval:

| LOAD DATA | EMBED TEXTS | EVALUATION |
|:---:|:---:|:---:|
| Prepare() | Batcher() | LR or MLP |

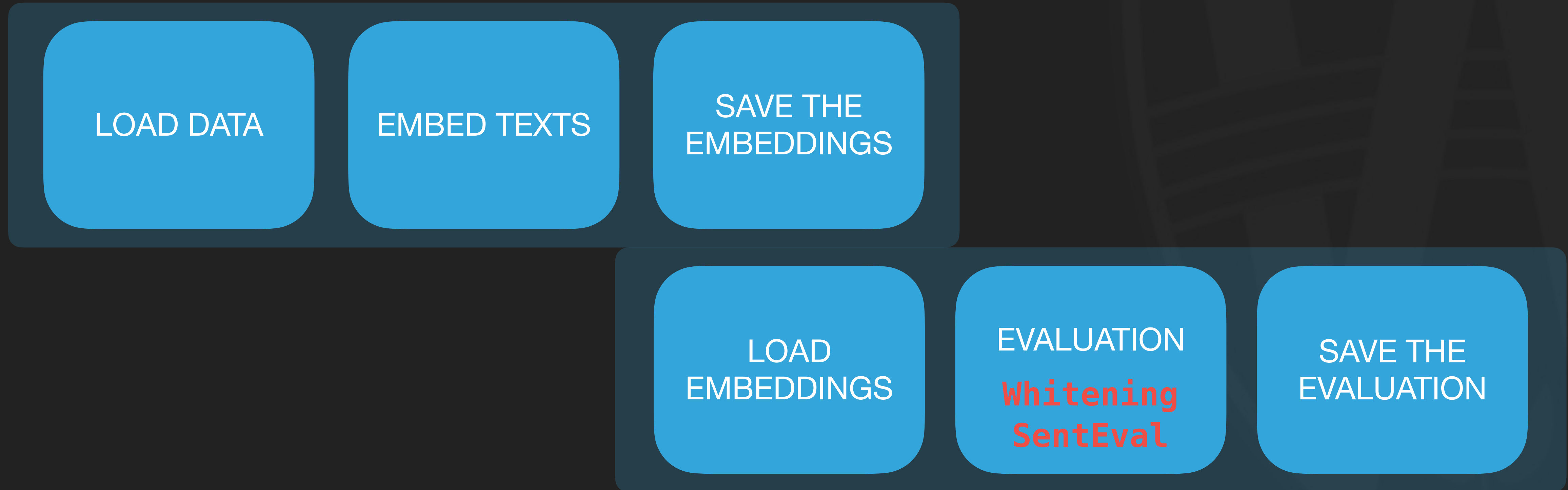```
PARAMS = {
'TASK_PATH':PATH_TO_DATA,
'USEPYTORCH': TRUE,
'KFOLD': 10
}
```

```
PARAMS['CLASSIFIER'] = {
    'NHID': 0,
    'OPTIM': 'ADAM',
    'BATCH_SIZE': 64,
    'TENACITY': 5,
    'EPOCH_SIZE': 4
}
```

You are going to lose the embeddings in this pipeline. This could be costly for ChatGPT, LLaMA, …
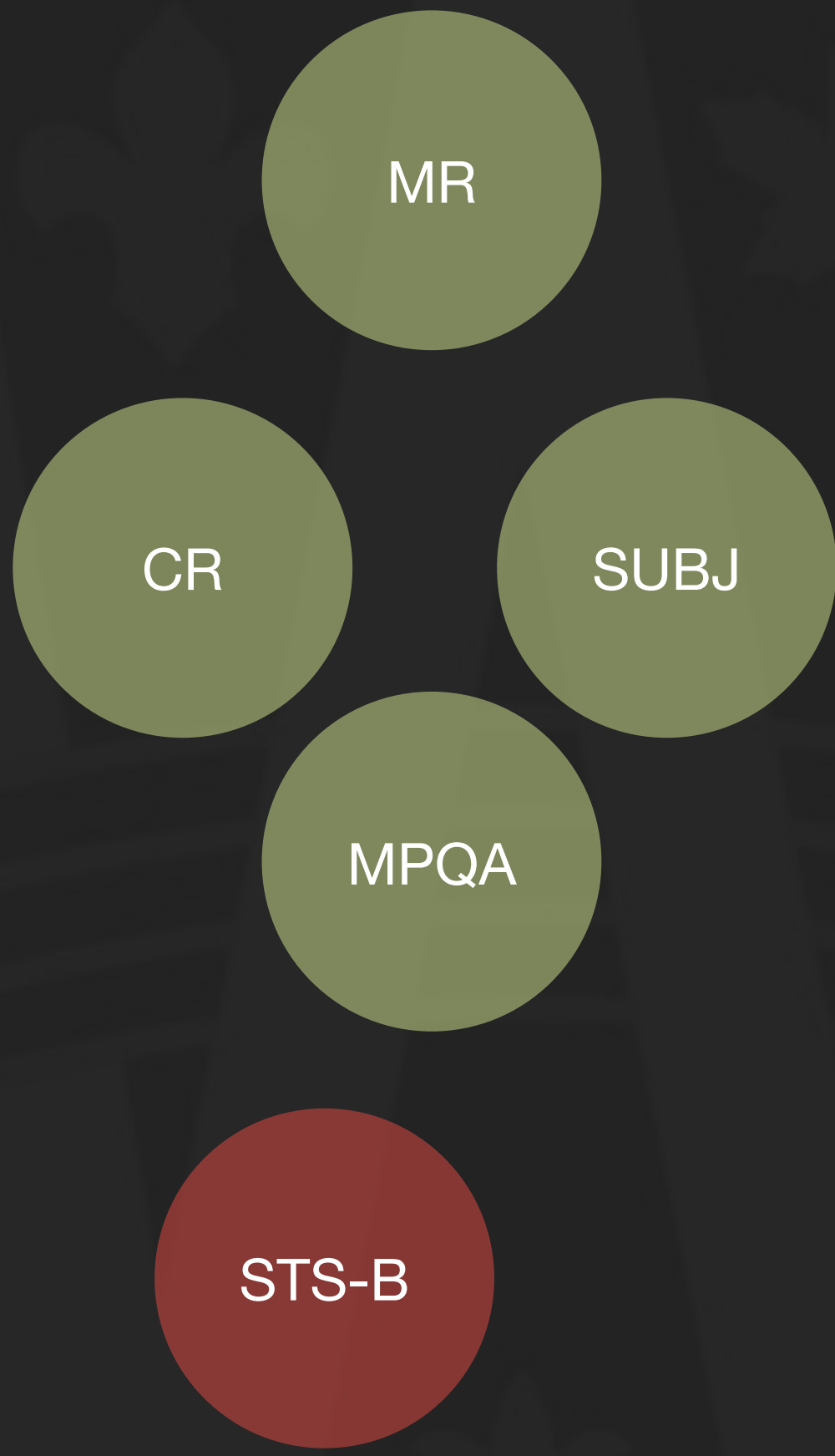
# OUR PIPELINE

1) Embedding Generation

| | | |
|---|---|---|
| LOAD DATA | EMBED TEXTS | SAVE THE EMBEDDINGS |

| | | |
|---|---|---|
| LOAD EMBEDDINGS | EVALUATION  Whitening SentEval | SAVE THE EVALUATION |

2) Embedding Evaluation

BERT · SIMCSE · S-BERT · ANGLE BERT · ANGLE LLAMA · LLAMA · LLAMA2 · CHATGPT

llama_embedding
- __pycache__
- .ipynb_checkpoints
- data
- embeddings
- env
- llama_converted
- llama2_converted
- results
- whitening                                    M
- .gitignore
- angle_base.py
- bert_base.py
- bert_finetune.py
- chatGPT.py
- clean_embeddings.py
- data_investigate.py
- data.py
- emb_util.py
- eval_cls.py
- eval_sts.py
- llama_base.py
- llama_eval.py
- llama_finetune.py
- main.py
- plotting.ipynb
- preprocessing.py
- README.md
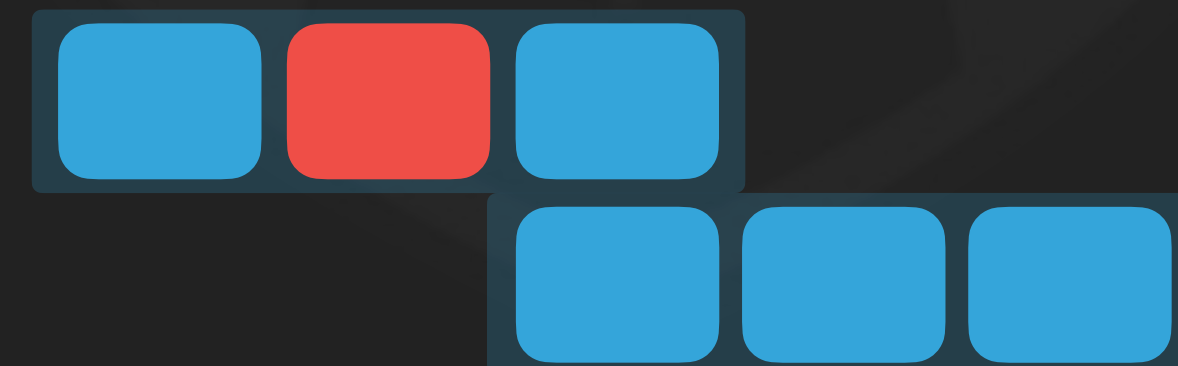- requirements.txt
- sbert_base.py
- simcse_base.py
- whitening.py

‣ Load Models
‣ Manage loading weights into multiple GPUs
‣ Generate Embeddings
‣ Store Embeddings

‣ Loop for generating Embeddings over datasets and encoders

```python
'''
options for models_llama:
    "llama-7B", "llama-13B", "llama-30B", "llama-65B", "llama2-7B", "llama2-13B", "llama2-70B"
'''
models_llama = ["llama-7B", "llama2-7B"]
'''
options for models_angle:
    angle-bert   : fine tuned bert on nli dataset
    angle-llama  : fine tuned llama2 with lora technique on nli dataset
'''
models_angle = ["angle-bert", "angle-llama"]
'''
options for models_chatGPT:
    text-embedding-3-small  : 62.3% in MTEB, 62,500 pages per dollor
    text-embedding-3-large  : 64.6% in MTEB, 9,615 pages per dollor
    text-embedding-ada-002  : 61.0& in MTEB, 12,500 pages per dollor
'''
models_chatGPT = ["text-embedding-3-small"]
'''
options for datasets:
    built-in train/test split:
        "yelpp", "imdb", "agnews", "yelpf", "trec"
    no built-in train/test split:
        "mr", "cr", "subj", "mpqa"
    similarity tasks:
        "sts1", "sts2"
    !!! if dataset has a predifined split you need to uncomment the code section for splited data IN
'''
datasets = ["mr", "cr", "subj", "mpqa"]

models = models_bert + models_llama + models_angle + models_simcse + models_chatGPT

for model in models:
    if(model in models_sbert):
        SBert_Embeddings(model, datasets)
    elif(model in models_bert):
        Bert_Embeddings(model, datasets)
    elif(model in models_simcse):
        SimCSE_Embeddings(model, datasets)
    elif(model in models_angle):
        Angle_Embeddings(model, datasets)
    elif(model in models_llama):
        Llama_Embeddings(model, datasets)
    elif(model in models_chatGPT):
        ChatGPT_Embeddings(model, datasets)
```

File tree:
```
∨ llama_embedding
    > __pycache__
    > .ipynb_checkpoints
    > data
    > embeddings
    > env
    > llama_converted
    > llama2_converted
    > results
    > whitening
    ◇ .gitignore
    ≈ angle_base.py
    ≈ bert_base.py
    ≈ bert_finetune.py
    ≈ chatGPT.py
    ≈ clean_embeddings.py
    ≈ data_investigate.py
    ≈ data.py
    ≈ emb_util.py
    ≈ eval_cls.py
    ≈ eval_sts.py
    ≈ llama_base.py  ⟵
    ≈ llama_eval.py
    ≈ llama_finetune.py
    ≈ main.py
    ◼ plotting.ipynb
    ≈ preprocessing.py
    ⓘ README.md
    ≡ requirements.txt
    ≈ sbert_base.py
    ≈ simcse_base.py
    ≈ whitening.py
```

```python
36      # Set device to auto to utilize GPU
37      device = "auto"  # balanced_low_0, auto, balanced, sequential
38
39      if self.model_name == "llama-310B":
40          print("loading llama 30B takes much longer time due to GPU management issues.")
41          self.model = LlamaForCausalLM.from_pretrained(
42              PATH_TO_CONVERTED_WEIGHTS,
43              device_map=device,
44              max_memory={0: "12GiB", 1: "12GiB", 2: "12GiB", 3: "12GiB"},
45              offload_folder="offload"
46          )
47      else:
48          self.model = LlamaForCausalLM.from_pretrained(
49              PATH_TO_CONVERTED_WEIGHTS,
50              device_map=device,
51              output_hidden_states=True
52          )
53
54      self.tokenizer = LlamaTokenizer.from_pretrained(PATH_TO_CONVERTED_WEIGHTS)
55
56      # unknow tokens. we want this to be different from the eos token
57      self.tokenizer.pad_token_id = (0)
58      self.tokenizer.padding_side = "left"
```

```python
85          tokens = self.tokenizer(
86              data_row['text'],
87              padding=True,
88              truncation=True,
89              return_tensors='pt',
90              max_length=64,
91              return_attention_mask = True
92          )
93          with torch.no_grad():
94              output = self.model(**tokens, return_dict=True)
95              hidden_states = output.hidden_states
96              if self.strategy == 'layer' and isinstance(self.pooling, int):
97                  embedding = (hidden_states[self.pooling]).mean(dim=1)
98              elif self.strategy == 'range' and isinstance(self.pooling, int):
99                  embedding = np.array(hidden_states[-self.pooling:]).mean(axis=0)
100                 embedding = np.array(embedding).mean(axis=1)
101             elif self.strategy == 'pair' and isinstance(self.pooling, tuple):
102                 embedding = (hidden_states[self.pooling[0]] + hidden_states[self.pooling[1]]).mean(dim=1)
103             else:
104                 raise Exception("unknown pooling")
105
106         embeddings.append(embedding[0])
```
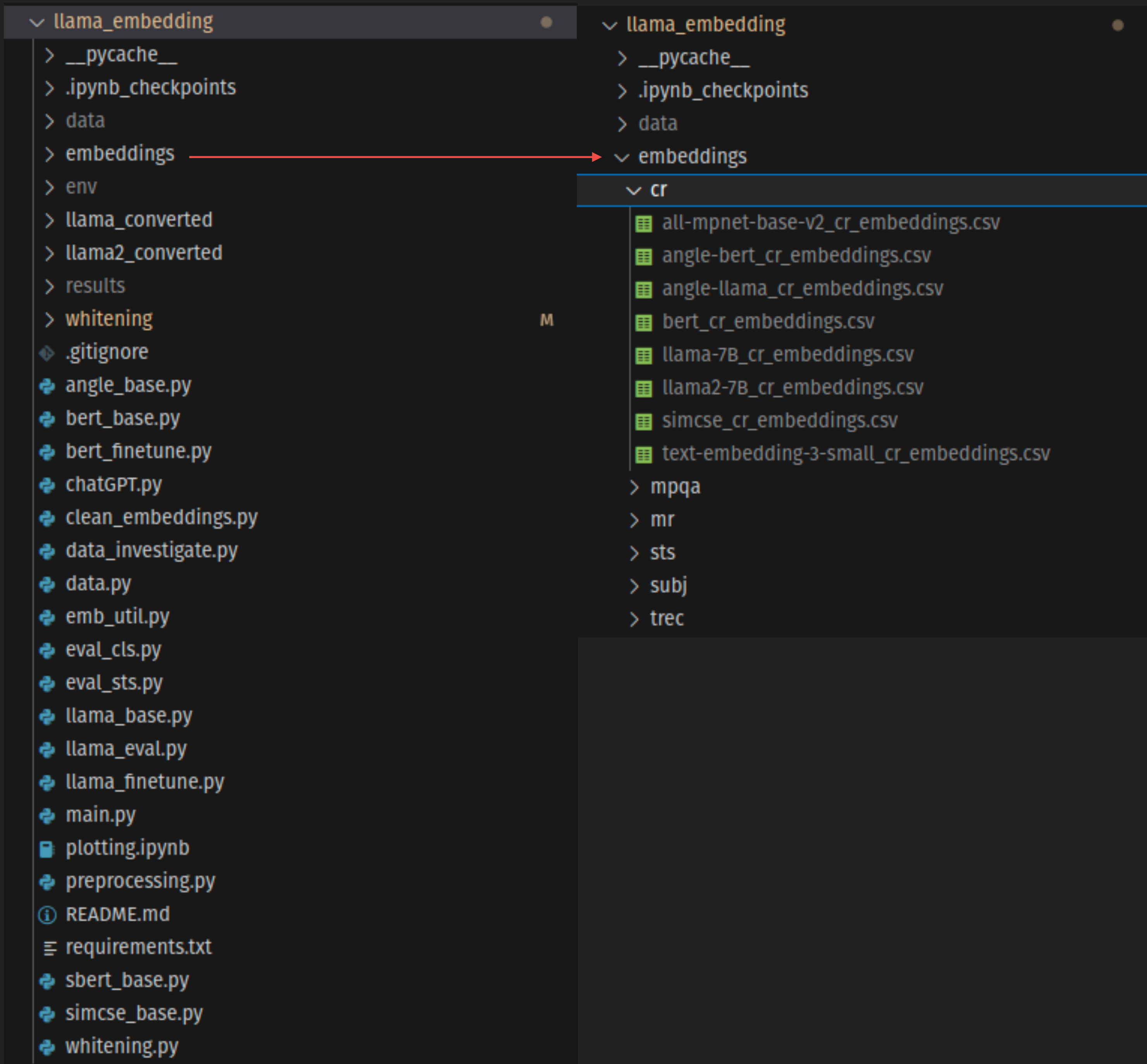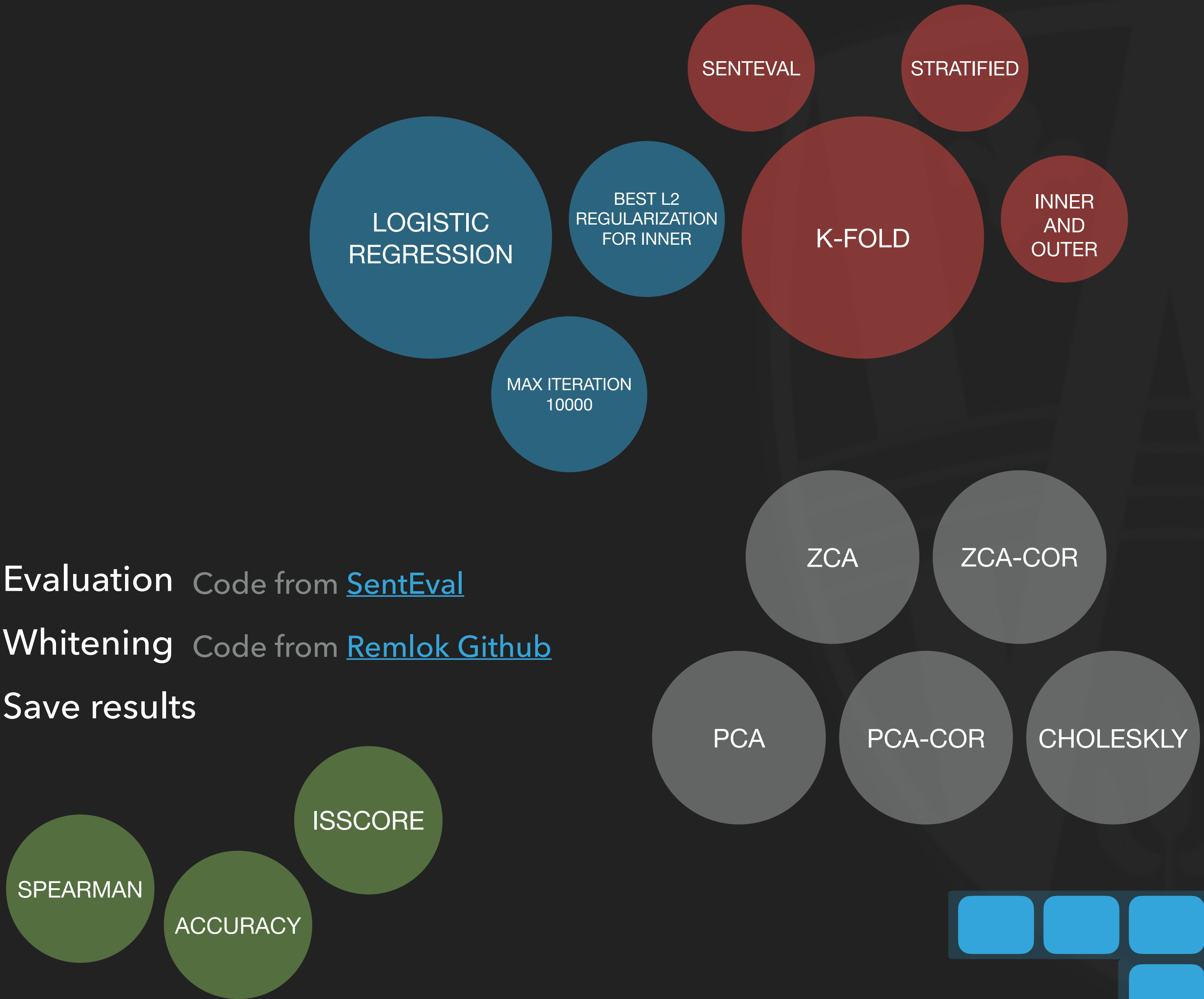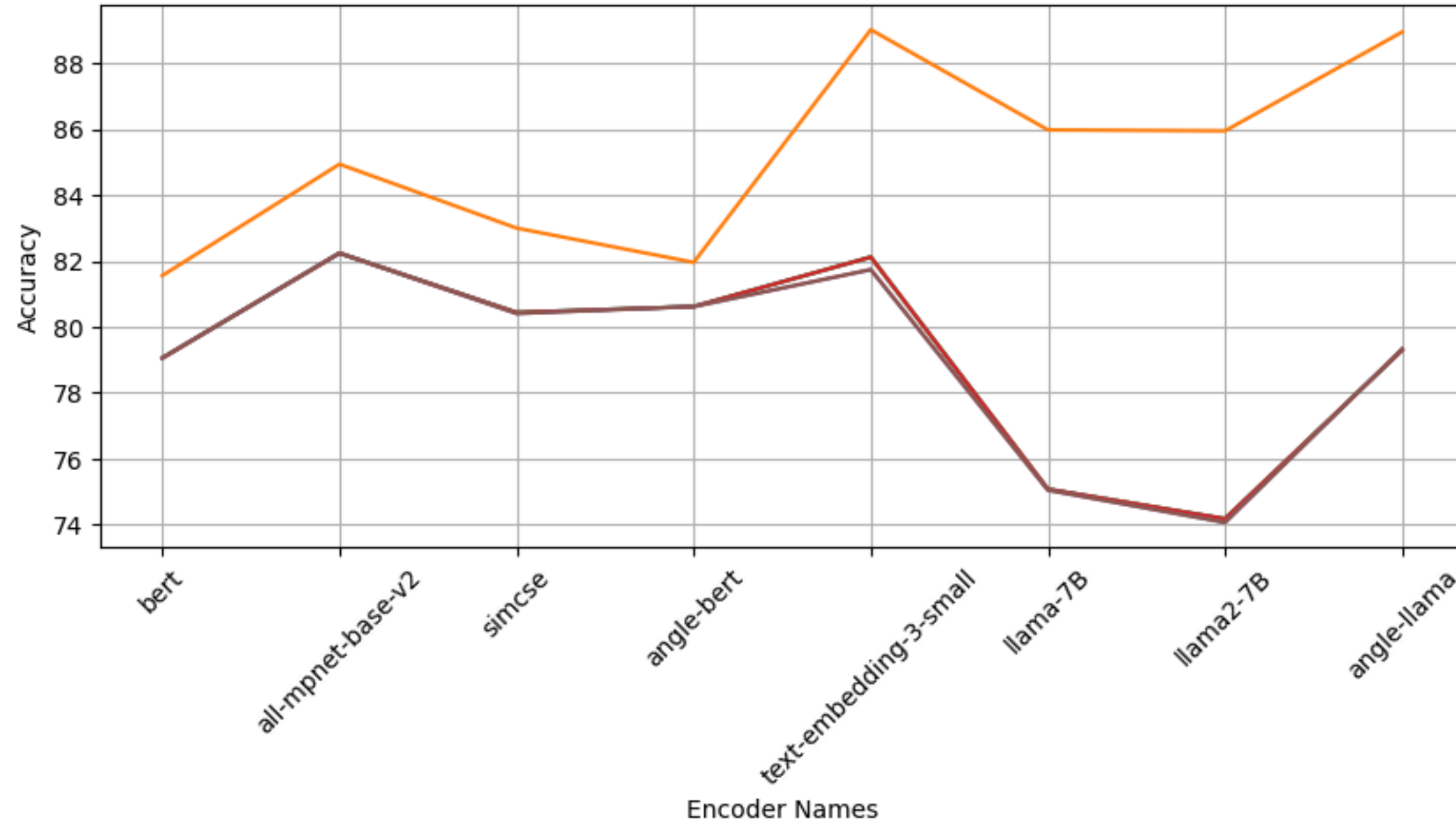
Storing embeddings of different datasets and models

CLASSIFICATION RESULTS

IMPROVEMENT / DIMINSHMENT PLOTS

PCA BEFORE AND AFTER WHITENING

pca of llama-7B-pair-31-and-32

pca of whitened llama-7B-pair-31-and-32
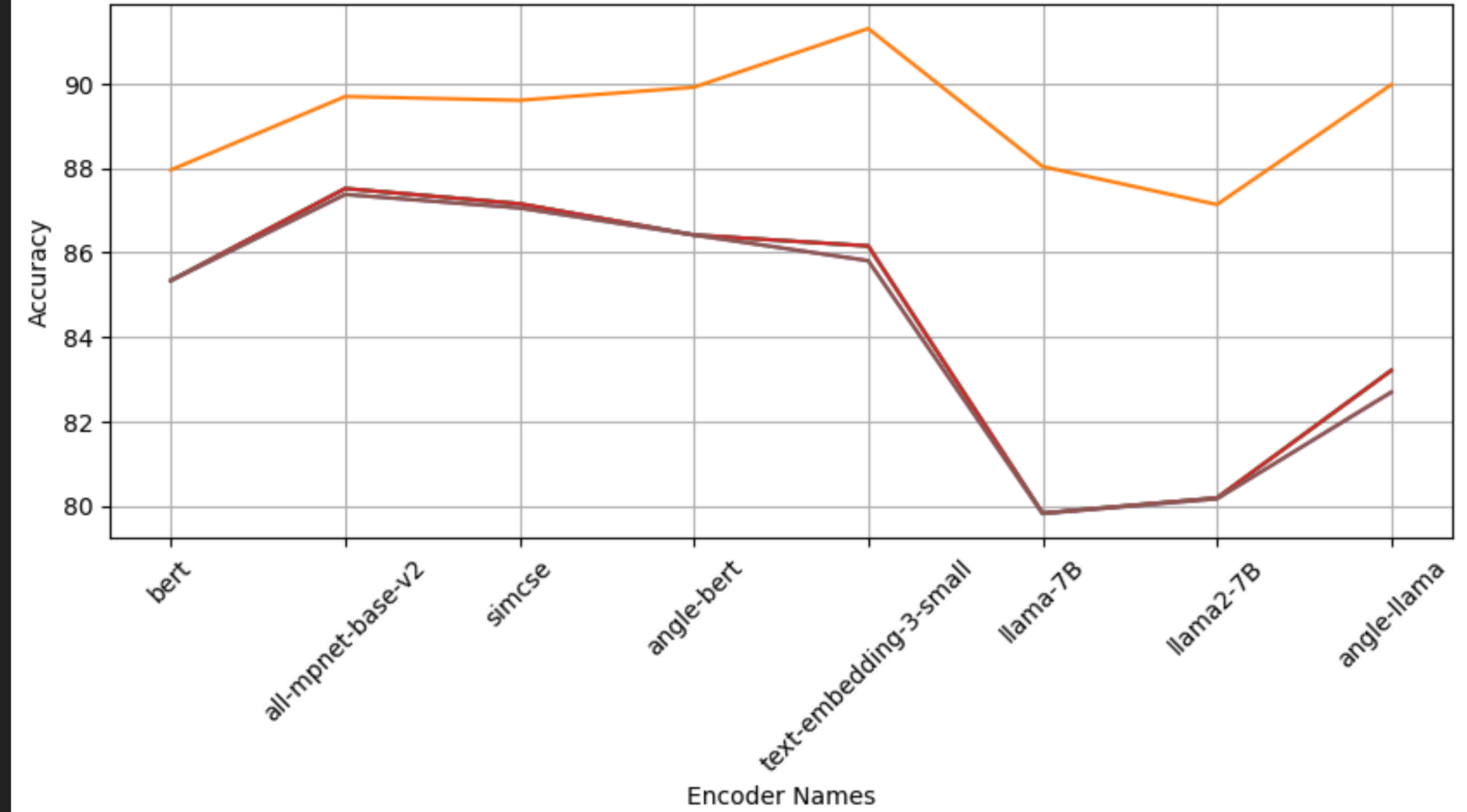
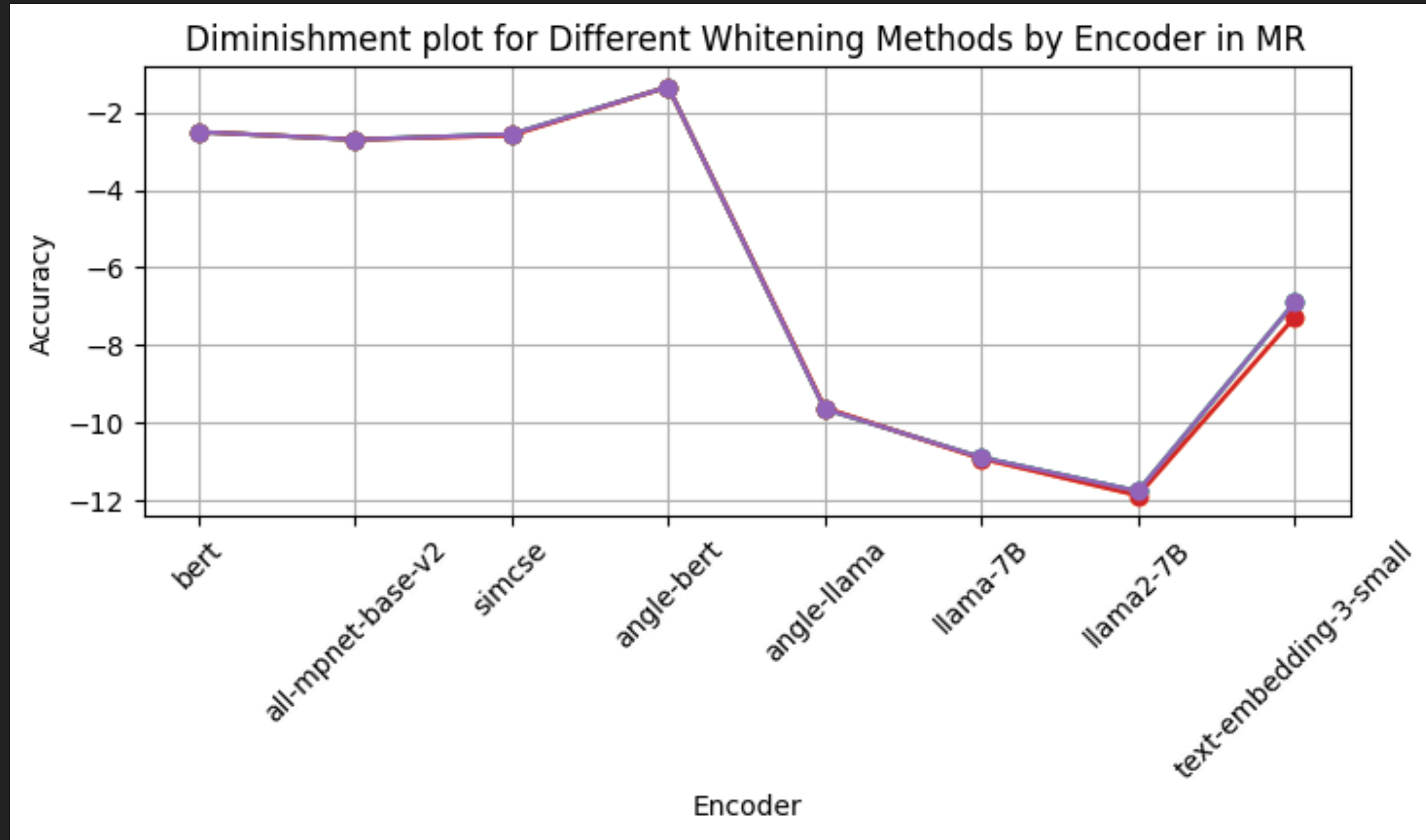pca of ChatGPT

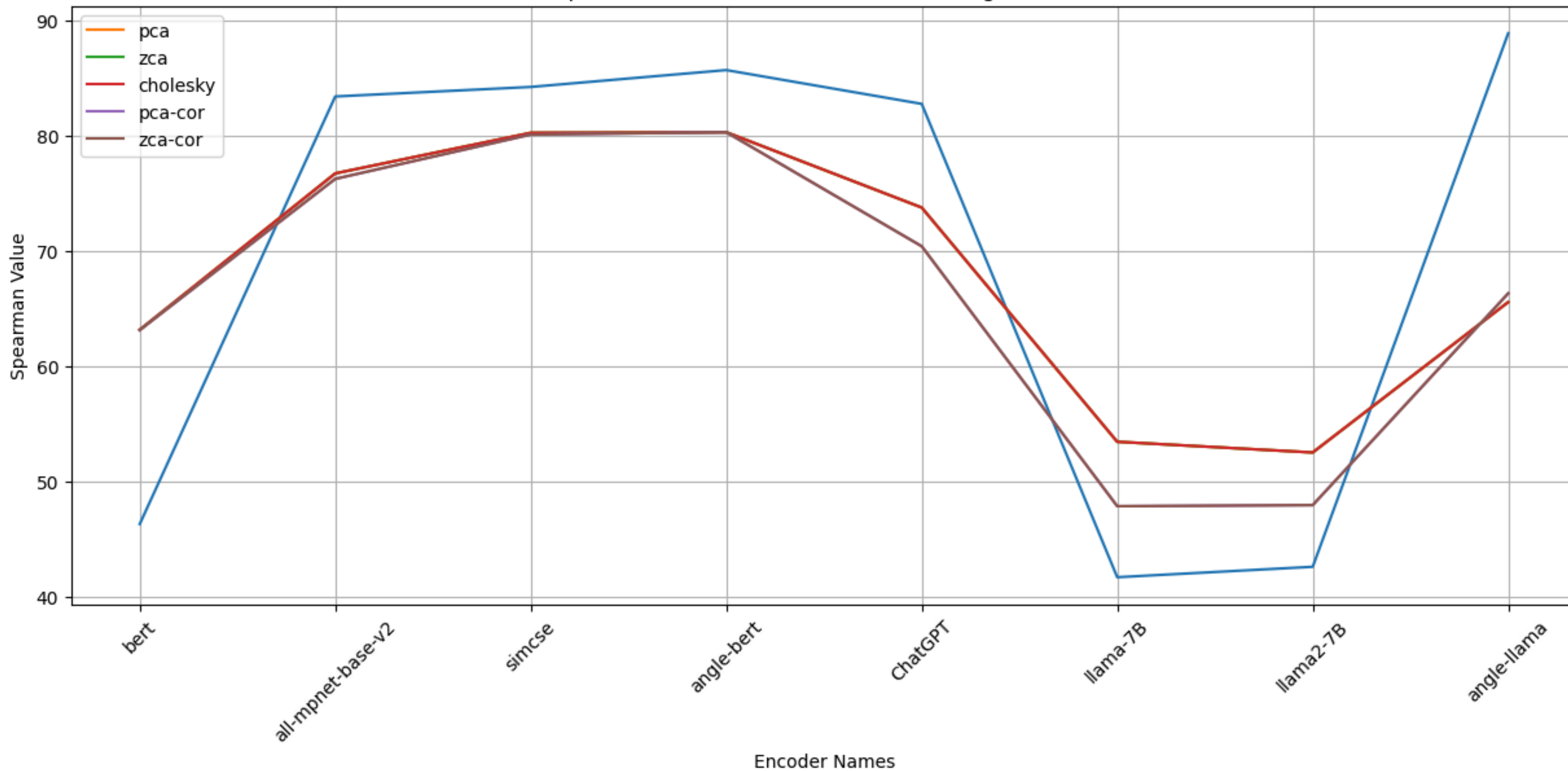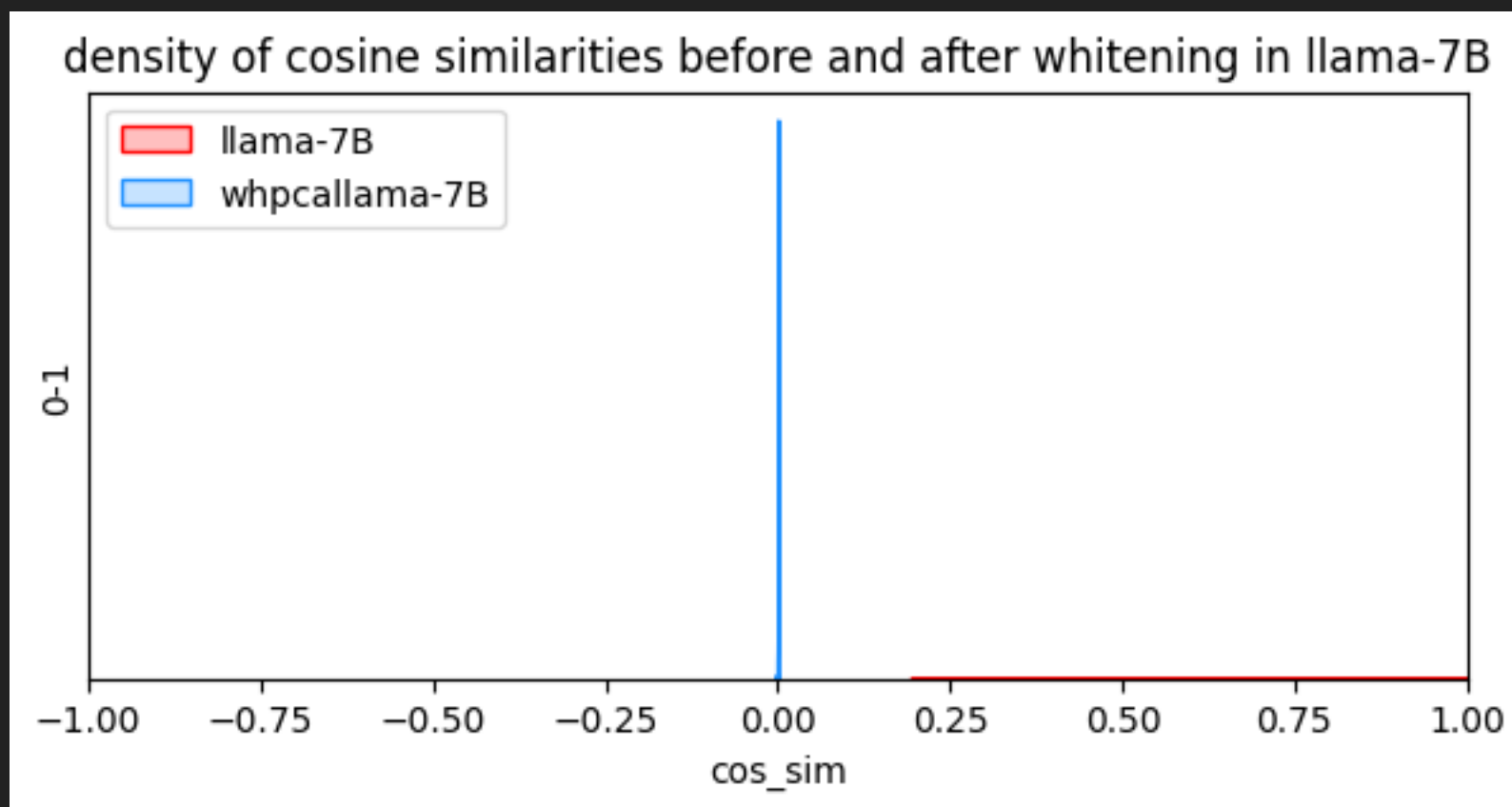pca of whitened ChatGPT
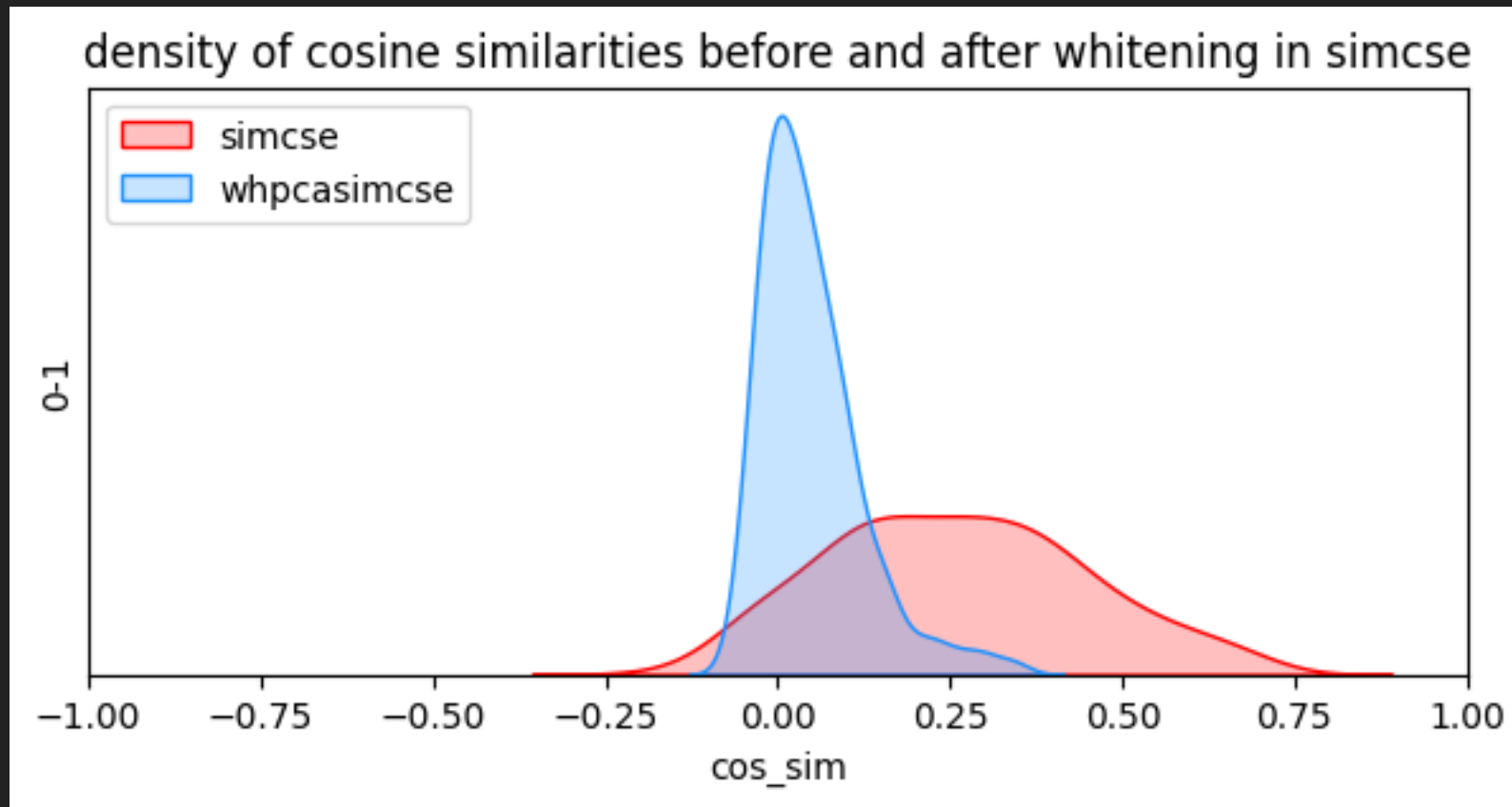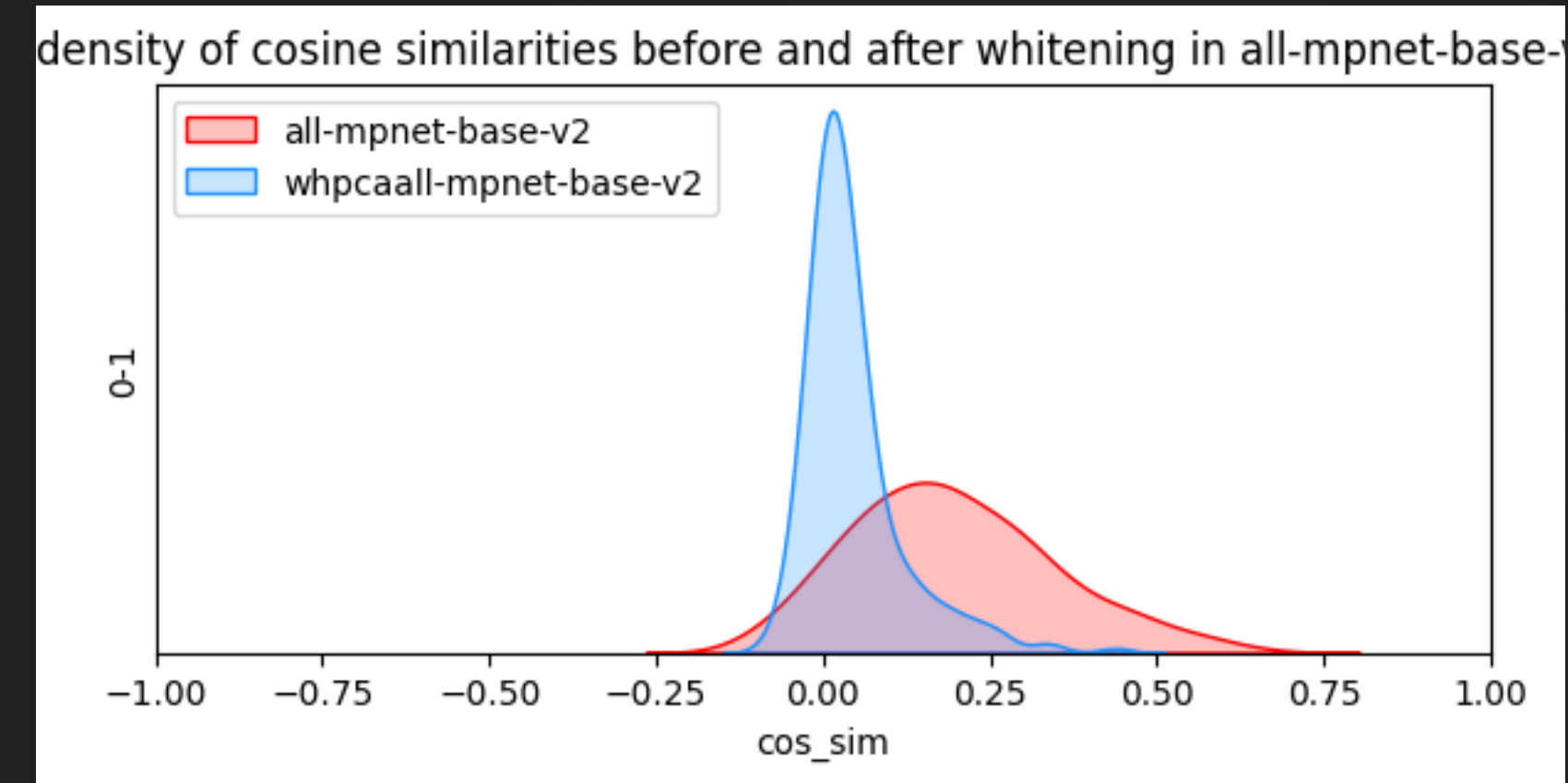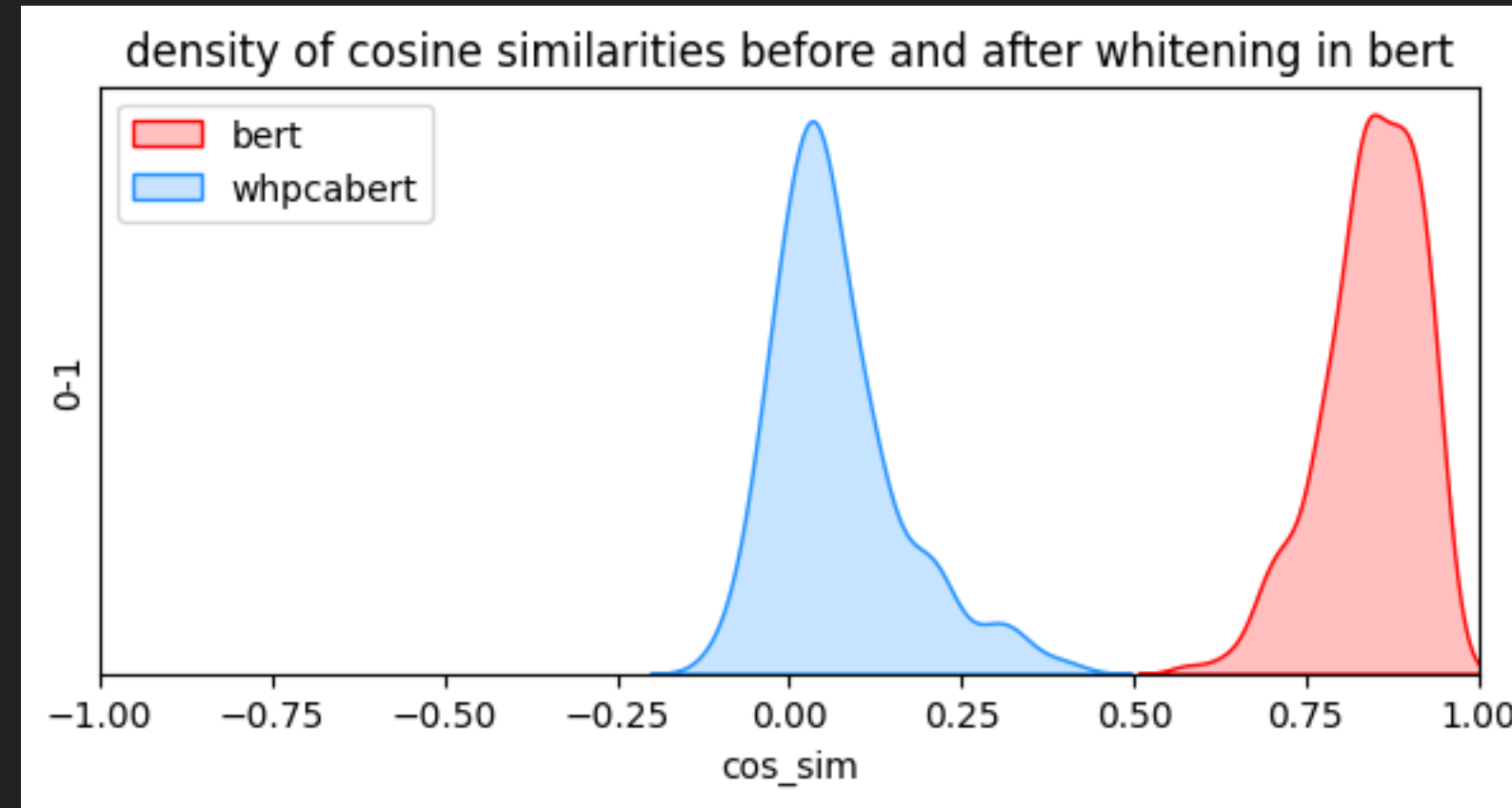
Spearman Values for Different Whitening Methods

COSINE SIMILARITY DENSITY BEFORE AND AFTER PCA WHITENING

# ANISTROPY PLOTS