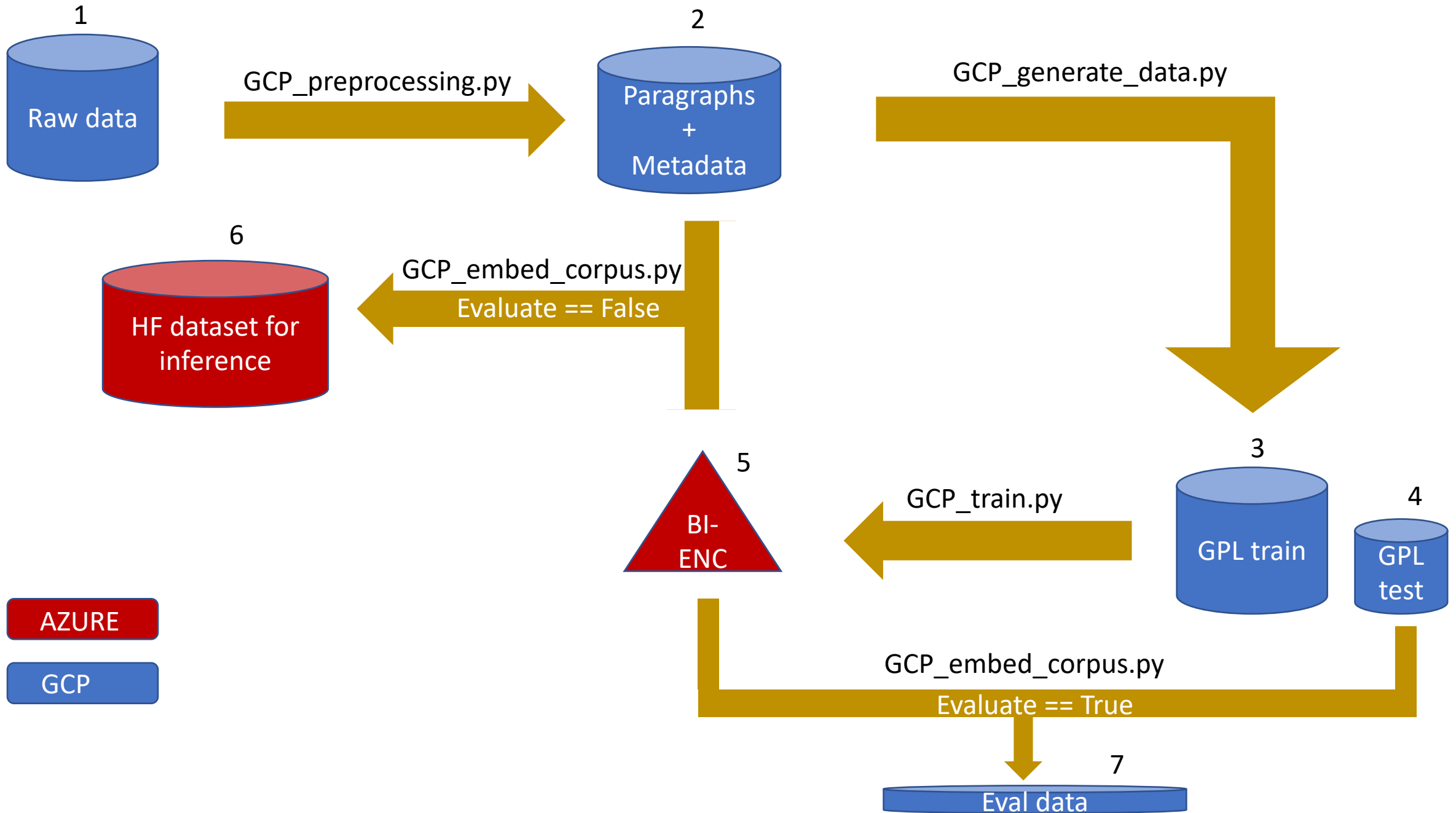




# Semantic Search Overdracht



# Stored where

1. **Raw data**  
GCP gsc-migrated-blobstore/raw/publicregister
2. **Paragraphs**  
semantic\_search\_blobstore/preprocessed\_corpus
3. **GPL train data**  
semantic\_search\_blobstore/gpl\_training\_data/corpus\_all\_length\_250/train
4. **GPL test data**  
semantic\_search\_blobstore/gpl\_training\_data/corpus\_all\_length\_250/test
5. **Bi-encoder**  
SemanticSearch\_Train/models
6. **HF inference dataset**  
SemanticSearch\_Train/workspaceblobstore/corpus\_ds\_with\_embedding
7. **Eval data**  
semantic\_search\_blobstore/corpus\_ds\_with\_embedding/evaluation/corpus\_all\_length\_250\_eval/multi-qa-mpnet-base-dot-v1\_TRAINED\_2022-10-21-14-33-13\_version\_45

\*All GCP stuff in Rebatch-sandbox project

\*corpus\_all\_length\_250 is the latest corpus (used for GSC demo)

\*multi-qa-mpnet-base-dot-v1\_TRAINED\_2022-10-21-14-33-13\_version\_45 it the latest model (used for GSC demo)

# Preprocessing

- Start from raw json files → extract info for every paragraph
  - `_id`: *string*, important as `unique_id` for GPL training later
  - `title`: *string*, title of document containing paragraph
  - `doc_id`: *string*, id of document containing paragraph (e.g. ST .... INIT)
  - `passage_id`: *integer*, paragraph number within document
  - `text`: *string*, paragraph text
  - `chunked`: *bool*, if True, paragraph was further split into pieces after initial document split on `'/n/n'` (not entirely relevant anymore I think)
- NOTES
  - To change max tokens (~words) per paragraph change `max_length` under `split_paragraph()`
  - Slow on GCP for large corpora: speed up by doing different parts in parallel and stitching back together with `merge_jsonfiles_gcp.py`
  - Lots of text cleaning + check for `bad_paragraph` because text in jsons extracted from pdf and often very low quality

# GPL generate data

Generate training data from preprocessed paragraphs

- Generate GPL data and training was 1 single script originally.
- Original script many, many parameters, most of default in training pipeline should be okay. Important parameters to specify yourself are in GPL/gpl/gpl\_args.yml.
- Original script only adapted by introducing custom functions
  - `prepend_title`: prepend title to paragraph text
  - `copy_corpus_and_filter`: use fraction of all available documents and keep 'good' ones
  - `postprocess_queries`: Select fraction of all 'best' synthetically generated questions
  - `upload_train_test_split`: upload train and test data to separate folders

# GPL generate data

- Structure of output folder important as train script will look for location of specific files
  - corpus.jsonl -> selection of entire preprocessed corpus used for generating queries
  - gpl-training-data.tsv -> triplets of query\_id, positive\_id, hard-negative\_id and cross-score (effective training lines)
  - hard-negatives.jsonl -> for every query\_id: positive\_id and list of all hard-negatives\_id
  - qrels
    - train.tsv (or test.tsv for test data) -> for every query\_id: positive\_id
  - qgen-queries.jsonl -> query\_id and query text

# GPL train

Start training from synthetic generated data

- Generate GPL data and training was 1 single script originally.
- Original script many, many parameters, most of default in training pipeline should be okay. Important parameters to specify yourself are in GPL/gpl/gpl\_args.yml.
- Make sure gpl\_score\_function corresponds to correct model!
  - e.g. 'dot' if using 'multi-qa-mpnet-base-dot'
- Original script only adapted by introducing custom functions
  - evaluation: compute precision, recall and mAP on test data
  - save\_callback: save model on azure with evaluation details
- Tricky to know how many training steps
  - $\text{steps} = \text{training\_lines} / \text{batch\_size}$
  - $\text{training\_lines} = \text{queries\_per\_paragraph} * \text{negatives\_per\_paragraph} * \text{number\_of\_paragraphs}$

# Embed corpus

Use trained bi-encoder to transform entire corpus into dataset with embeddings

- Evaluation == False: Transform entire corpus to dataset and upload to Azure
- Evaluation == True: Only embed data select for train+test,  
in the following step, use this smaller dataset for evaluation

## NOTES:

- If evaluation==True: checks automatically if dataset already exists, if so, only do evaluation step
- Possibility to prepend title to paragraph only for embedding
- Currently no FAISS index added as flat one is used