

# Search Enhancements using NLP Techniques

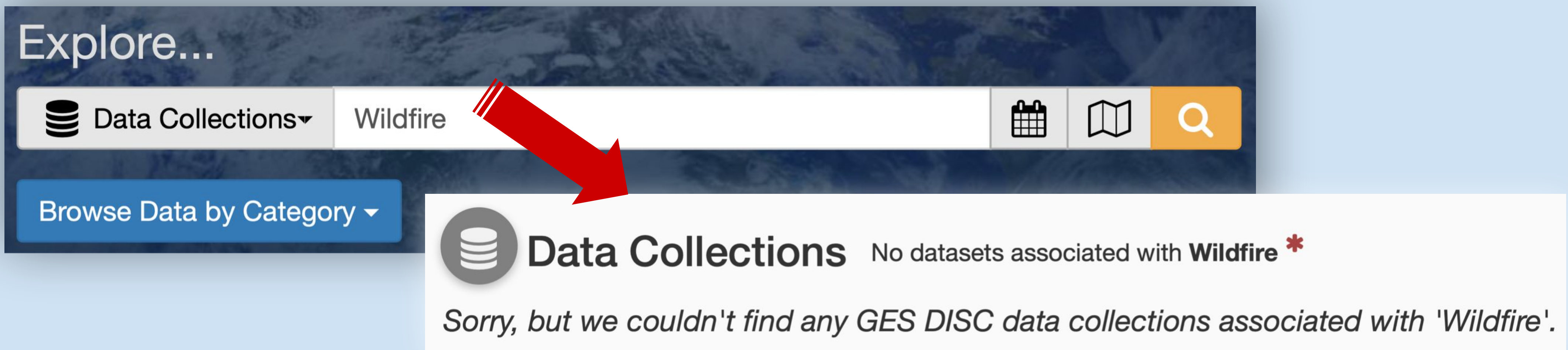


**GES DISC**  
SUMMER 2021

Armin Mehrabian<sup>1,2</sup>, Irina Gerasimov<sup>1,2</sup>, Mohammad Khayyat<sup>1,2</sup>, Jennifer Wei<sup>1</sup>, Binita KC<sup>1,2</sup>, and Hegde Mahabaleshwara<sup>1</sup>  
 1 NASA Goddard Space Flight Center  
 2 ADNET Systems INC.  
**Corresponding Author:** Armin Mehrabian (armin.mehrabian@nasa.gov)

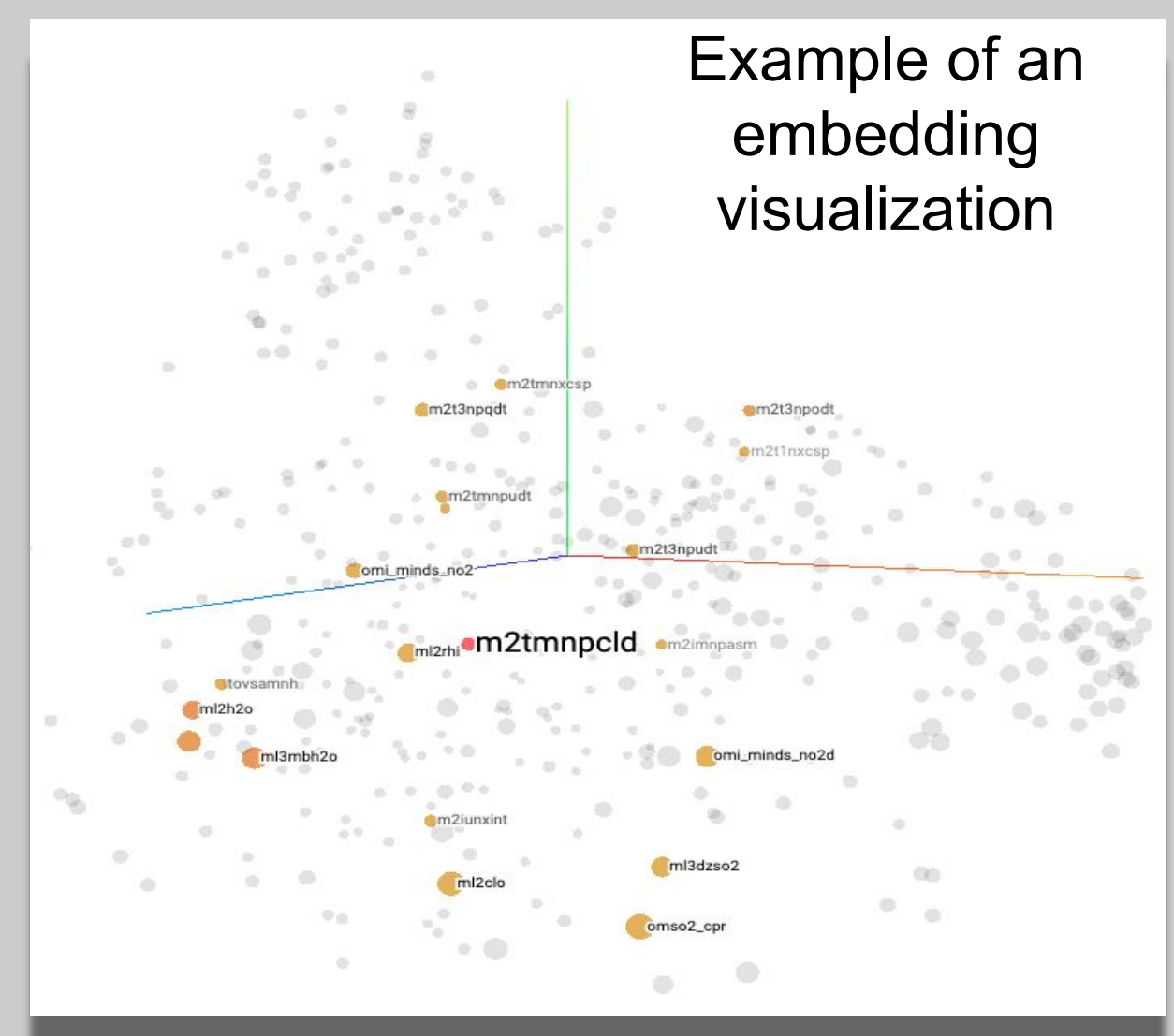
## PAIN POINT

- Traditional search engines rely on hard-matching of query string and some form of document metadata
- The above approach fails to find semantic similarities between the search query and metadata that are not exact matches

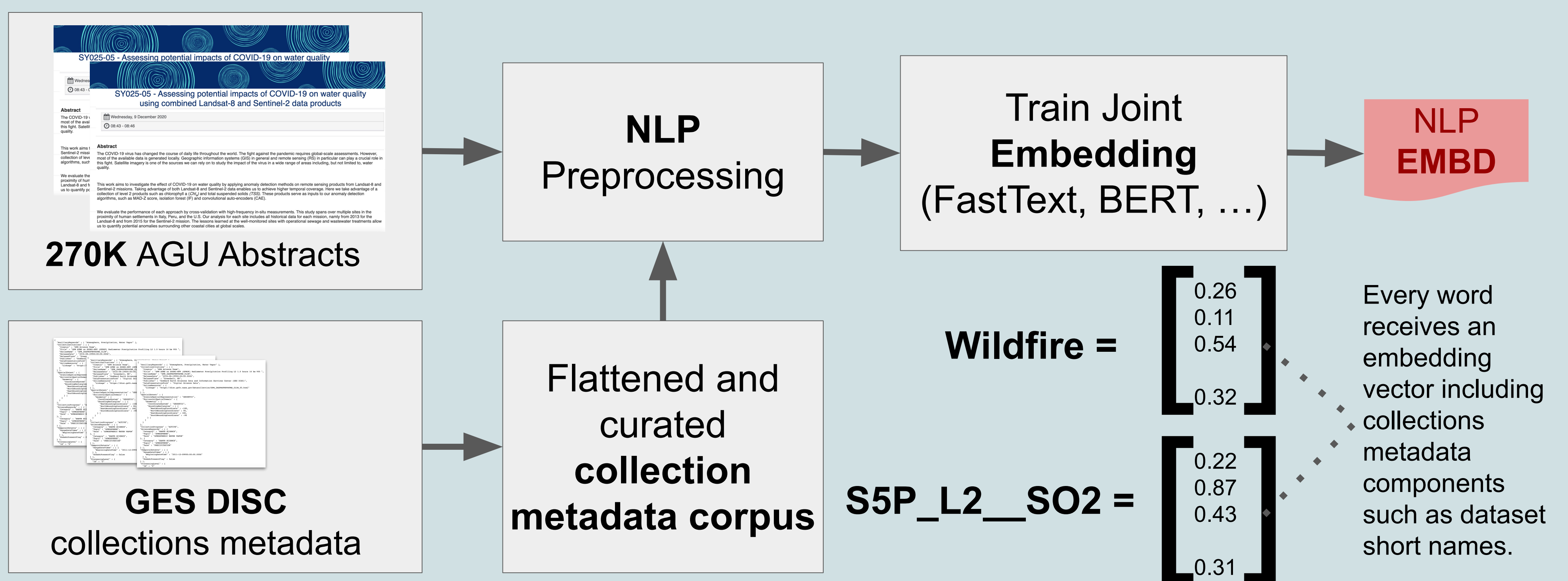


## EMBEDDINGS

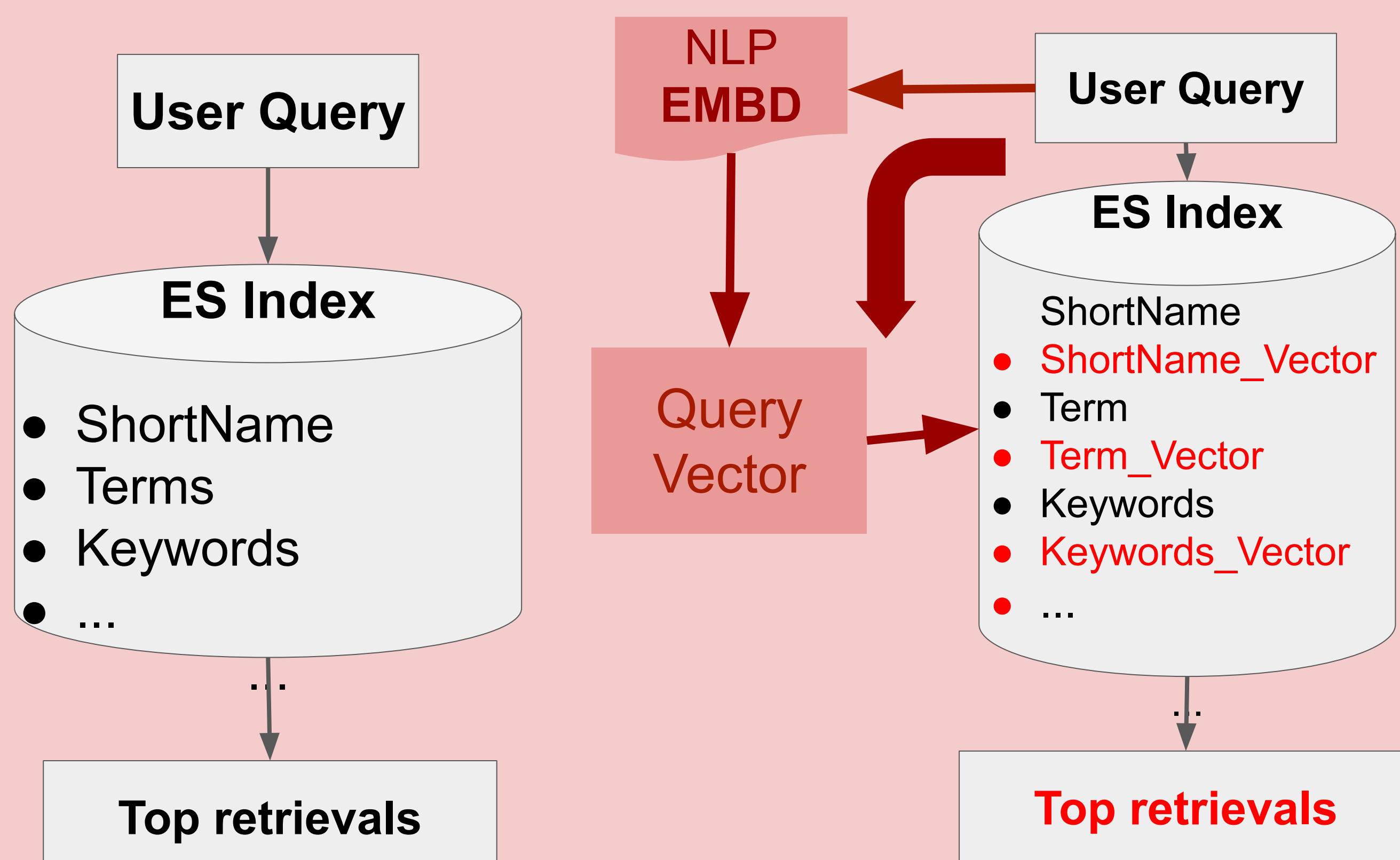
- In order to perform mathematical operations in NLP one needs to convert language components such as words, sentences, ... into numerical values, namely **vectors** or **embeddings**
- These vectors are high-dimensional (100~300 dimensions)
- Embeddings capture semantics and context. I.E. words **“precipitation”** and **“rainfall”** should have a shorter distance to one another relative to a word such as **“wildfire”** in an embedding space
- **Here** we train a joint embedding from a corpus of relevant scientific text along with a curated text corpus from **GES DISC** collections metadata



## TRAINING



## ELASTICSEARCH



- Elasticsearch offers **“dense\_vector”** field types to store embedding vectors
- For a curated set of fields in the collections metadata we ingest the fields accompanied by their embedding vectors
- During **“query time”**, the user query string is converted into its vector representation
- A similarity score is calculated based on the query vector and the vectors of ingested fields
- This allow us to have a search return based on similarity score even in the absence of exact matches