# TraitHunter

Mapping and extraction of biomedical traits via text embeddings

IEU Programme 3 talk
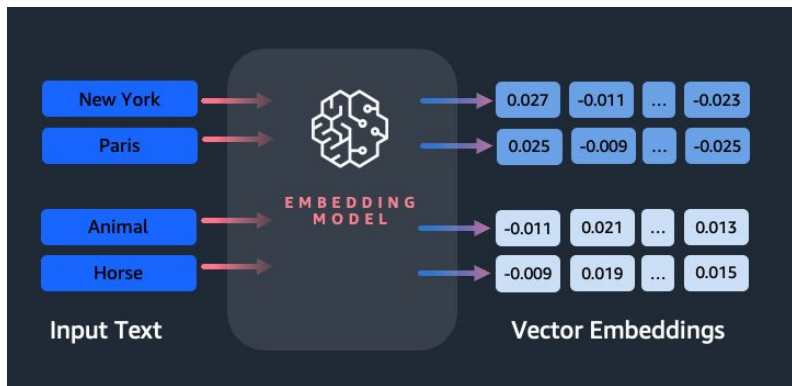
Yi Liu
3rd September 2024

# Today's talk

The main aim for today's talk is to showcase a web service for trait mapping.

This is really quick and dirty hack, however feedbacks on various aspects are deeply appreciated.

Outline

- Concepts
- Background
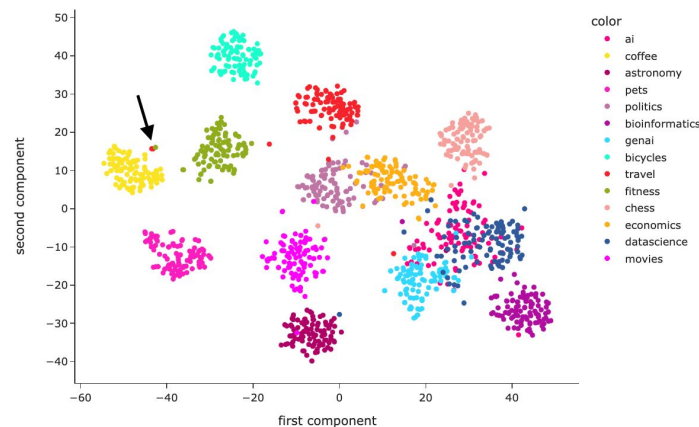- TraitHunter
  - Concepts
  - Live demo
- Next steps

# Concept: Text embeddings



**Input Text** → **EMBEDDING MODEL** → **Vector Embeddings**

| New York | → | 0.027 | -0.011 | ... | -0.023 |
| Paris | → | 0.025 | -0.009 | ... | -0.025 |
| Animal | → | -0.011 | 0.021 | ... | 0.013 |
| Horse | → | -0.009 | 0.019 | ... | 0.015 |

Convert text into their (semantic) vector representations via an encoder model

Simple(r) use case: mapping text

t-SNE embeddings



color
- ai
- coffee
- astronomy
- pets
- politics
- bioinformatics
- genai
- bicycles
- travel
- fitness
- chess
- economics
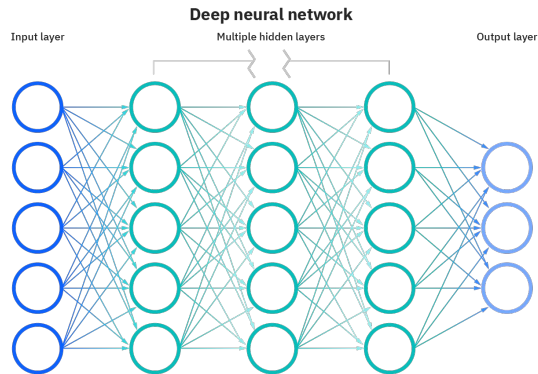- datascience
- movies

Retrieval Augmented Generation

**Retrieval augmentation**



Who is the president of the US? → **Generator** (Language Model) → Joe Biden

**Retriever**

Knowledge can be expanded & updated (new domain, news, etc.)

Interpretability (reference to source)

Memory

WIKIPEDIA The Free Encyclopedia

Joe Biden is the 46th and current president of the United States, assumed office on January 20, 2021.

Retrieved document

# Concept: Large Language Models (LLM)

**Deep neural network**

Input layer — Multiple hidden layers — Output layer

**Transformer vs LLaMA**

**Transformer**
("Attention is all you need")

Output Probabilities
Softmax
Linear
Add & Norm
Feed Forward
Add & Norm
Multi-Head Attention
Add & Norm
Masked Multi-Head Attention
Positional Encoding
Input Embedding
Inputs
Output Embedding
Outputs (shifted right)
Positional Encoding
Nx

**LLaMA**
Softmax
Linear
RMS Norm
Feed Forward SwiGLU
RMS Norm
Self-Attention (Grouped Multi-Query Attention) with KV Cache
Q K V
RMS Norm
Embeddings
Input
Nx
Rotary Positional Encodings

Large Language Model

Pretraining
Large data → Base Model → 

Fine-tuning
Base Model → Small data → Fine-tuned Model

**Data**
- Text
- Images
- Speech
- Structured Data
- 3D Signals

→ Training → Foundation Model → Adaptation →

**Tasks**
- Question Answering
- Sentiment Analysis
- Information Extraction
- Image Captioning
- Object Recognition
- Instruction Following

# Background: related works

In Liu, et al. (2023) we approached the trait mapping problem with a two stage approach:

- Identify a sub set of related traits with naive text embeddings, from the broad set
- Map the traits using an LLM task model finetuned on ontologies

This approach is put into use in the ASQ platform (Liu & Gaunt, 2024) for querying EpiGraphDB evidence.

Limitations:

- Locked to a specific ontology (EFO)

- Finetuned task model not widely adaptable

- Just using naive embeddings might be preferable with the latest-and-greats of LLMs (in ASQ next gen)

# Background: related works in LLM based service

**Vectology**

- Provides naive embedding and mapping of text with
  - BioSentVec
  - BERT
- Direct precursor to TraitHunter
- Some form of vectology is still internally accessible

**epigraphdb-neural**

- Pre-computes embedding of EpiGraphDB entities in a elasticsearch vector store
  - ScispaCy
- Powers the search functionalities in EpiGraphDB

# Why yet another LLM based recommender service

Models

- LLaMA3
  ('nuff said!)



Tooling

- Huggingface and other tools (ollama) are now much more mature
- Elasticsearch only recently supports 4095 dim dense vector embeddings

Research

- We collaborated with a few research teams on the problem of trait mapping for pre-select traits of interest

- So we probably should have a formal service and short paper up so collaborators can cite it

- In preparation for a few research projects

# TraitHunter

TraitHunter is

- A publicly accessible IEU web service for
    - **Mapping traits across biomedical dictionaries**
    - Linking traits with biomedical text and literature
- Embedding and extraction of traits
    - Will replace vectology and epigraphdb-neural
    - Semi-publicly accessible with rate limit

- Web app: https://traithunter.epigraphdb.org
- API: https://traithunter-api.epigraphdb.org
- code: https://github.com/MRCIEU/traithunter

# TraitHunter: concepts

**Dictionary**

A vocabulary set of trait entities

Regular examples:

- UK Biobank variables
- Trait names of OpenGWAS studies
- EpiGraphDB terms

Ontologies

- Hierarchical from general to specific concepts

**Trait entity**

- ID
- Label
    - Concept singleton
    - Complex trait

- (optional) Description
- (optional) Synonyms
- (optional) Parents, children, and siblings

**Full text document**
Abstracts, text descriptions, etc.

## Glycogen storage disease due to aldolase A deficiency MONDO:0012747

ORPHA:57

Glycogen storage disease due to aldolase A deficiency is an extremely rare glycogen storage disease characterized by hemolytic anemia with or without myopathy or intellectual deficit. Myopathy can be severe enough to result in fatal rhabdomyolysis in some patients. A family with episodic rhabdomyolysis (triggered by fever) without hemolytic anemia has recently been reported.

⬇ Export Associations ⚙ Report Entry Issue

---

## breast cancer `IMPORTED`

http://purl.obolibrary.org/obo/MONDO_0007254 📋 Copy

A primary or metastatic malignant neoplasm involving the breast. The vast majority of cases are carcinomas arising from the breast parenchyma or the nipple. Malignant breast neoplasms occur more frequently in females than in males. ⓘ

**Defined by** MONDO

**Also appears in** CPONT GENEPIO OBA

**Synonym** BC | breast cancer ⓘ | breast tumor ⓘ | breast tumour ⓘ | cancer of breast ⓘ | malignant breast neoplasm ⓘ | malignant breast tumor ⓘ | malignant breast tumour ⓘ

---

# 114480

## BREAST CANCER

*Alternative titles; symbols*

BREAST CANCER, FAMILIAL

Other entities represented in this entry:

BREAST CANCER, FAMILIAL MALE, INCLUDED

▼ **Description**

Breast cancer (referring to mammary carcinoma, not mammary sarcoma) is histopathologically and almost certainly etiologically and genetically heterogeneous. Important genetic factors have been indicated by familial occurrence and bilateral involvement.

▼ **Clinical Features**

Cady (1970) described a family in which 3 sisters had bilateral breast cancer. Together with reports in the literature, this suggested to him the existence of families with a particular tendency to early-onset, bilateral breast cancer. The genetic basis might, of course, be multifactorial. ⊕

Anderson (1974) concluded that the sisters of women with breast cancer whose mothers also had breast cancer have a risk 47 to 51 times than in control women; a revised estimate was 39 times (Anderson, 1976). The disease in these women usually developed before menopause, was often bilateral, and seemed to be associated with ovarian function. About 30% of daughters with early-onset, bilateral breast cancer inherited the susceptibility. The risk of breast cancer to women with affected relatives is higher when the diagnosis is made at an early age and when the disease is bilateral. Ottman et al. (1983) provided tables that give the cumulative risk of breast cancer to mothers and sisters at various ages. The highest risk group is sisters of premenstrual probands with bilateral disease. Among the sisters of women with breast cancer, Anderson and Badzioch (1985) found the highest lifetime risks when the proband had bilateral disease, an affected mother (25 +/- 7.2%), or an affected sister (28 +/- 11%). The risks were reduced to 18 +/- 3.3% and 14 +/- 2.6%, respectively, with unilateral disease. An early example of familial breast cancer was provided by

Diagram of TraitHunter web service

# TraitHunter: Demo

# Next steps, in infra development

Assessment of embedding quality with benchmarks

- How well does an embedding model objectively perform for the task of mapping biomedical traits

Curation of full text into the vector store

- OMIM descriptions
- Medline abstracts

Named entity recognition

- Identification and extraction of a trait entity from full text docs.

Molecular symbols

- Special treatment for gene names, etc.

# Next steps, in research

The TraitHunter paper

- Describing the service and methods

Your input is needed on what this platform can help with your research.

**Clustering analysis**

Investigation on the different clusters of biomedical relationships for novel insights

- Semantics
- Genetic associations

**RAG**

RAG on knowledge graphs to investigate about biomedical assertions (i.e. ASQ-next-gen).

# Acknowledgements