UNIVERSITY OF BRISTOL

Liu, Y., Elsworth, B., & Gaunt, T. R. (2023). Using language models and ontology topology to perform semantic mapping of traits between biomedical datasets. *Bioinformatics*, *39*(4), [btad169]. https://doi.org/10.1093/bioinformatics/btad169

OXFORD

## Data and text mining

# Using language models and ontology topology to perform semantic mapping of traits between biomedical datasets

Yi Liu ⓘ [1,‡], Benjamin L. Elsworth[2,‡], Tom R. Gaunt ⓘ [1]*

[1]MRC Integrative Epidemiology Unit, Bristol Medical School, University of Bristol, Bristol, United Kingdom
[2]Our Future Health, Manchester, United Kingdom

*Corresponding author. MRC Integrative Epidemiology Unit, Bristol Medical School, University of Bristol, Oakfield House, Clifton, Bristol, BS8 2BN, United Kingdom. E-mail: Tom.Gaunt@bristol.ac.uk
‡Equal contribution.

### Abstract

**Motivation:** Human traits are typically represented in both the biomedical literature and large population studies as descriptive text strings. Whilst a number of ontologies exist, none of these perfectly represent the entire human phenome and exposome. Mapping trait names across large datasets is therefore time-consuming and challenging. Recent developments in language modelling have created new methods for semantic representation of words and phrases, and these methods offer new opportunities to map human trait names in the form of words and short phrases, both to ontologies and to each other. Here, we present a comparison between a range of established and more recent language modelling approaches for the task of mapping trait names from UK Biobank to the Experimental Factor Ontology (EFO), and also explore how they compare to each other in direct trait-to-trait mapping.

**Results:** In our analyses of 1191 traits from UK Biobank with manual EFO mappings, the BioSentVec model performed best at predicting these, matching 40.3% of the manual mappings correctly. The BlueBERT-EFO model (fine-tuned on EFO) performed nearly as well (38.8% of traits matching the manual mapping). In contrast, Levenshtein edit distance only mapped 22% of traits correctly. Pairwise mapping of traits to each other demonstrated that many of the models can accurately group similar traits based on their semantic similarity.

**Availability and implementation:** Our code is available at https://github.com/MRCIEU/vectology.

## 1 Introduction

Population health and medical research are increasingly reliant on large population studies, such as UK Biobank (https://www.ukbiobank.ac.uk), The Million Women Study (http://www.millionwomenstudy.org), Our Future Health (https://ourfuturehealth.org.uk), The Million Veterans Program (https://www.research.va.gov/mvp), China Kadoorie Biobank (https://www.ckbiobank.org), and others to discover new predictive biomarkers and interventions. Such studies measure many thousands of phenotypic variables. Systematic analyses, such as phenome-wide association studies (Jones et al. 2005, Denny et al. 2010, Millard et al. 2019), can describe relationships between thousands of variables, producing large datasets. However, many variables are inconsistently named across studies, and can prove difficult to map to each other or an existing ontology, such as the Experimental Factor Ontology (EFO) (Malone et al. 2010), Human Phenotype Ontology (Robinson et al. 2008) or the

Disease Ontology (Kibbe et al. 2015). In parallel, the biomedical literature contains a wealth of data on human diseases, traits, and risk factors described using free text (with some mappings to Medical Subject Headings). Systematically integrating knowledge across these different datasets and domains would enable us to triangulate the evidence (Lawlor et al. 2016) for different risk factor/disease combinations, but at the moment this is hindered by the inconsistencies in trait nomenclature.

The complexity of variable names is illustrated by UK Biobank, an internationally important population study that has collected a wealth of data on half a million people. Examples of text labels for variables in UK Biobank include easily recognizable traits, such as "systolic blood pressure", and disease names, such as "coronary heart disease". However, the study also includes more complex variables, including those derived from questionnaire data, including "able to walk or cycle unaided for 10 minutes" and "cough on most days". An array of other variables also exists, including

International Statistical Classification of Diseases and Related Health Problems 10th Revision (ICD10) codes, such as "anaemia due to enzyme disorders" (D55) and "syncope and collapse" (R55), the former mapping directly to the EFO (EFO: 0009529), but the latter not. Direct mapping by text matching to ontology terms is therefore not realistic, and whilst manual mapping to ontologies is sometimes appropriate, this is time-consuming, especially if mapping to multiple different ontologies (which cover different domains of the human phenome and exposome).

Given this, there are four potential solutions to link two datasets based on their lists of trait (variable) names:

1. Manual mapping to an ontology to find shared terms between datasets.
2. Using automated tools to map each variable to an ontology to find shared terms between datasets.
3. Direct mapping of variables using a generalizable text embedding model to identify semantically similar terms.
4. Direct mapping of variables using a bespoke model trained on the particular datasets to identify semantically similar terms.

Each of these options has different strengths and weaknesses. Option 1 can only really be used in cases where the numbers of variables is low, or the requirement of human assigned ontological terms is essential. Option 2 relies on existing tools, such as OnToma (https://github.com/opentargets/OnToma), Zooma (https://www.ebi.ac.uk/spot/zooma/) or MetaMap Lite (Demner-Fushman et al. 2017) for common ontologies, such as EFO and UMLS (https://www.nlm.nih.gov/research/umls/index.html). These rule-based tools can work well, but the mapping to ontology may identify a more generic parent term in the ontology losing valuable information in the process. Options 3 and 4 may offer benefits in mapping variables between datasets by avoiding the intermediate step of an ontology term (Fig. 1).

The development of methods based on text embeddings, such as word2vec (Mikolov et al. 2013), sent2vec (Pagliardini et al. 2018), and doc2vec (Le and Mikolov 2014), have opened up the potential to map terms based on semantic similarity. These methods have been applied to data from the biomedical domain e.g. BioWordVec (Zhang et al. 2019) and BioSentVec (Chen et al. 2019) and have been applied to real world problems (Duong et al. 2017, Jaeger et al. 2018, Allot et al. 2019, Blagec et al. 2019, Karadeniz and Özgür 2019, Tshitoyan et al. 2019).

Further development to shallow/non-contextual text embeddings gives rise to contextualized methods, such as the transformer model architecture (Vaswani et al. 2017) and its implementation in language modelling [e.g. BERT (Devlin et al. 2019)] applied in a range of contexts [e.g. GLUE (Wang et al. 2018), BLUE (Peng et al. 2019), and BLURB (Gu et al. 2021)]. These models can be finetuned to tackle specific problems with great effect (Duong et al. 2019, Koroleva et al. 2019, Liu et al. 2019, Fabian et al. 2020). However, despite their merits, transformer models are slower and more resource intensive compared to the word2vec architecture.

Here, we apply a range of text embedding methods and BERT language models (including one trained on EFO) to the problem of mapping biomedical variables (from UK Biobank) to an ontology (EFO) and compare their performance, strengths, and weaknesses. We also illustrate the use of these models on a direct trait-to-trait mapping problem.
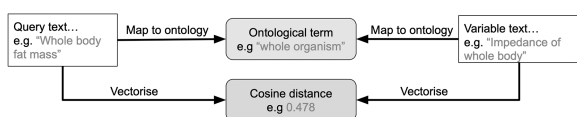


**Figure 1** Example of potential benefits of using text embeddings to connect two biomedical strings compared to using a shared ontology

## 2 System and methods

### 2.1 The EFO dataset
The **EFO** contains parts of several biological ontologies as well as variables from many large scale databases. Whilst many other ontologies exist, this particular ontology is widely used for human traits and is well documented, so was considered a good choice for this evaluation. A version of the EFO dataset was downloaded from the EBI RDF platform (https://www.ebi.ac.uk/rdf/services/sparql) in March 2021 containing 25 390 terms. This was used for all subsequent analysis and is available in the supplementary material (Supplementary Files S1 and S2).

### 2.2 Ontology distance metric
To understand the relative distance between any two EFO terms and enable us to measure how well a trait was mapped, we used the nxontology Python library (https://github.com/related-sciences/nxontology). By creating a parent–child network of EFO terms, we could compute a similarity measure between any pair of EFO terms and use this to create a measure of how close two terms are within the EFO hierarchy. For this analysis, we used the Batet (parameter "batet" in the library) measure (Batet et al. 2011) as this was developed using biomedical taxonomy data and produced good correlation results to manual biomedical concept comparisons. The measure ([0, 1]) is a ratio calculated from the shared and non-shared information between a pair of concepts, where the lower the score the less shared ancestry between the two ontology concepts have. From here on we will refer to this metric as the **EFO-Batet score**.

To create a nxontology instance, we provided the parent/child EFO edge data to the 'NXOntology' class (Supplementary Code block S1).

### 2.3 Trait-to-trait mapping distance score
The different models use different approaches for measuring distance between text terms (Supplementary Table S2). For simplicity, we refer to these metrics (edit distance, cosine similarity, and semantic distance) as "trait similarity score" throughout.

### 2.4 Mapping methods
We used a range of existing string comparison language models representing different approaches to language representation and different pre-training datasets to enable us to evaluate the impact of these differences on mapping performance.

#### 2.4.1 String comparison methods
**Levenshtein edit distance ratio** (Levenshtein 1966) was used to understand how well a basic string comparison performs. Using the implementation from the python-Levenshtein library (https://pypi.org/project/python-Levenshtein/), we obtained a measure of similarity between two strings.

**Zooma** is an established tool to map text to ontologies using a combination of curated mapping to existing datasets and standard text matching (the exact method is undocumented). For this analysis, we utilized the Zooma API setting the "required" parameter set to "None" and "ontologies" parameter set to "efo" (Supplementary Code block S2) to avoid circularity.

#### 2.4.2 Text embedding methods
**BioSentVec** is an established model created using sent2vec (Pagliardini et al. 2018), pre-trained on over 28 million titles and abstracts from PubMed (https://pubmed.ncbi.nlm.nih.gov/) and 2 million clinical notes from MIMIC III (Johnson et al. 2016). The BioSentVec (Chen et al. 2019) model was downloaded from the project GitHub repository (https://github.com/ncbi-nlp/BioSentVec) and installed following the examples in the tutorial (https://github.com/ncbi-nlp/BioSentVec/blob/master/BioSentVec_tutorial.ipynb) (Supplementary Code block S3).

**Google Universal Sentence Encoder v4** (**GUSE**) is a generalized text embedding model trained and optimized for sentence level tasks (Cer et al. 2018). The model was trained on Wikipedia and other generalized texts with no focus on biomedical information. The model was downloaded from the project home page (https://tfhub.dev/google/universal-sentence-encoder/4) and implemented as described in the documented example (Supplementary Code block S4).

**spaCy** is a natural language processing platform, which provides various tools, methods, and pipelines, one of which is word embeddings (Montani et al. 2021). The 'en_core_web_lg' model was downloaded and installed as described in the documentation (https://spacy.io/usage/linguistic-features#vectors-similarity) (Supplementary Code block S5).

**ScispaCy** is built on spaCy and provides models for processing biomedical, scientific or clinical text (Neumann et al. 2019). The 'en_core_sci_lg' model was downloaded and installed as described in the documentation (https://allenai.github.io/scispacy). The model is accessed via the same spaCy methods as above.

**BlueBERT** (Peng et al. 2019) (NCBI_BERT_pubmed_mimic_uncased_L-12_H-768_A-12) and **BioBERT** (Lee et al. 2020) (biobert_v1.1_pubmed) are biomedical language model implementations based on the original BERT pre-trained weights, with further language model training with biomedical corpora to improve language understanding tasks in the biomedical domain. For transformer models, the vector representation of the entity is computed as the average of the hidden state tensor of the $N - 1$ layer as a fixed representation of the tokenized sequence (i.e. the default strategy in bert-as-service, https://github.com/hanxiao/bert-as-service). These models were obtained from their respective model repositories, then served via the bert-as-service API (see Supplementary Code block S6 for example usage and code repository for detailed set up).

### 2.4.3 Bespoke ontology classifier
In addition to established language models, we also explored the effect of tailoring a transformer model to the EFO using transfer learning.

**BlueBERT-EFO** was developed by finetuning BlueBERT with an ontology entity alignment training process designed as a sequence classification task (for details see Supplementary Text S1). To create a similarity matrix of the entities, for each pair of terms the model produces a score representing the inferred ontology distance, where the lower number of steps between two entities as predicted by the model, the closer they are represented in an ontology graph. The model can be used for inference using the Huggingface Transformers (Wolf et al. 2019) package (see Supplementary Code block S7 for example usage and code repository for detailed set up).

Supplementary Table S2 shows a summary of the models and methods.

Whilst we only evaluated the benefits of finetuning on one base model (BlueBERT), we recognize that a number of other models could also be finetuned.

### 2.5 Mapping to ontology (EFO)
To assess how the models perform when mapping biomedical variables to an existing ontology, we utilized the **EBI UK Biobank EFO dataset** (Pendlington et al. 2019). This is a list of around 1500 UK Biobank variables that have been manually mapped to EFO terms. In addition, each mapping has been assigned a mapping type (Exact, Broad, Narrow, and Other). The original dataset was modified in the following ways: first, any query that had been assigned multiple EFO terms was dropped. Second, exact matches were excluded as uninformative (i.e. the query term is identical to the EFO label). Third, due to our use of an EFO hierarchy distance method (EFO-Batet), we only included those rows containing an EFO term present in our parent/child EFO dataset. Fourth, all EFO and variable terms were lower-cased. Lastly, duplicates were removed. These filtering steps created a dataset with 1191 entries (Supplementary File S3). Supplementary Table S1 displays the numbers of each by mapping type and a brief description of each mapping type as described in the original dataset.

Using this dataset, we applied the models described above to conduct pairwise comparison between the UK Biobank variables and the EFO terms to measure their semantic similarity and ontology distance. Specifically, a UK Biobank variable $A$ is associated with a manually mapped EFO term $a$ in the source dataset, then for an EFO term $b$, we calculated the similarity score between $A$ and $b$ as well as the EFO-Batet distance score between $a$ and $b$. Therefore, for the variable of interest $A$, the results dataset gives us a measure of how close the top ranking (by a specific similarity score metric) EFO term predictions $b_0 \ldots b_N$ are to the variable's equivalent EFO representation $a$ in the ontology space (by the EFO-Batet score).

### 2.6 Direct trait-to-trait mapping
In some scenarios mapping trait names between two datasets directly (without using an ontology) might be preferable. To compare how the different methods perform when predicting the similarity between two biomedical variables, we again used the EBI UK Biobank EFO dataset. This time, we limited the entries to those labelled as "Exact" on the assumption that these would provide a better dataset for assessing pairwise distances, both semantically and using the same ontology based method. Additional filtering steps were taken to create a dataset with one query per predicted EFO term, resulting in 530 entries (Supplementary File S4). For the purposes of visualization, we then manually selected a subset of 43 traits that represented a broad spectrum of variables, covering measurements, questionnaire data, and disease (Supplementary File S5). For each of the pre-trained models, pairwise cosine distances were generated for each query text. For Levenshtein, the similarity ratio was calculated as before. For BlueBERT-EFO, we generate the inferred ontology distance for each pair of terms. Whilst we were not mapping trait terms to an ontology, we also compared how close these pairs of traits are in the EFO for comparison using the EFO-Batet score for each pair of terms.

## 3 Implementation

### 3.1 Comparison to other approaches for automated mapping to ontology

#### 3.1.1 Top ranking results
We first explored how well the top prediction of each method compared to the manual annotation (Fig. 2). For results that exactly agree with the manual annotation (Fig. 2A), the best performing methods were BioSentVec (40.3%), BlueBERT-EFO (38.8%), Zooma (37.2%), and ScispaCy (36.5%), the results of which were notably higher than those of the methods included in the analysis. Pairwise proportions $Z$-test results (Supplementary Table S4) between each of the mapped proportions confirm that there is a notable difference between results of the best performing group and the those of the other methods, but the differences are minimal within the group (largest difference is between BioSentVec and ScispaCy, $P$-value $= .058$).

Whilst none of the methods exceeded 40.3% exact mapping, it is important to consider three key points: (i) some of the manual predictions are likely to be incorrect; (ii) the methods and models used here to automate this approach are quick and easy to use, and would scale to a task size that would be impractical for manual annotation; (iii) even the most sophisticated natural language processing models will struggle to predict the same result as a human, particularly in cases where the query string contains two un-linked entities, or even a negated term, e.g. "enduring personality changes not attributable to brain damage and disease".

In some situations (e.g. a recommender of similar concepts), an exact match may not be required, and if the top prediction from a model is sufficiently close to the manual annotation, this may be a suitable result. We then examine how well the top predictions from a method align with the manual annotation in terms of their EFO-Batet score distance to the manual EFO terms. Fig. 2B shows the aggregate results for the subset (see Supplementary Fig. S9 for full results) of methods over different range of EFO-Batet score
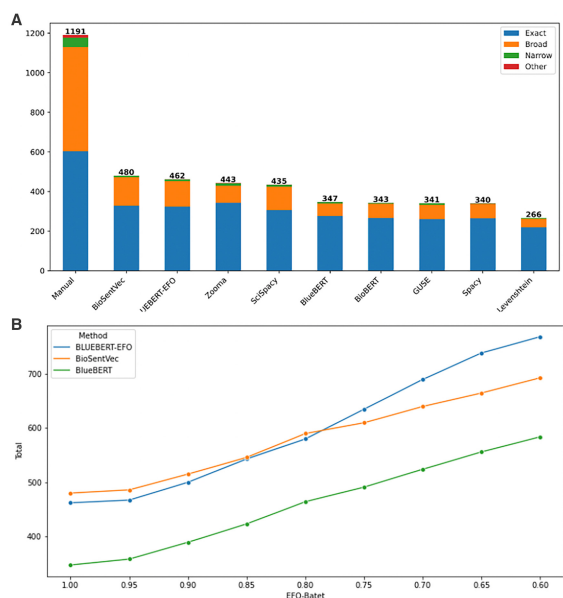
**Figure 2** Distribution of top matching predictions. (A) Number of top matching predictions by MAPPING_TYPE. The Total bar contains all manual mappings, subdivided into Exact, Broad (parent term), Narrow (child term), and Other. Each other bar represents the number of traits exactly matched by the named method to the manual mapping for that trait, with the same subdivisions. (B) Total number of top matching predictions that are equal or above an EFO-Batet threshold, i.e. if a method produces greater number of matched predictions with a threshold closer to 1, greater number of predictions exhibit close ontology relationship to the manual mapping results. Points at EFO-Batet thresholds 1.0, 0.9, 0.8, and 0.7 are equivalent to the Total values for each method in Fig. 2A, Supplementary Figs S1–S3. Full results for all methods can be found in Supplementary Fig. S9

threshold for top predictions to be included, from total number of top predictions that are strictly identical to manual annotation (threshold = 1, i.e. Fig. 2B), to those that are sufficiently close to the manual annotation in the ontology space (e.g. threshold ≥0.9), and then to results with a greater ontology distance tolerance (e.g. threshold ≥0.6). Supplementary Figs S1–S3 show the detailed distributions for thresholds of 0.9, 0.8, and 0.7, respectively.

For inexact mapping results, BlueBERT-EFO and BioSentVec retrieved similar number of concepts that are close (e.g. under an EFO-Batet threshold of 0.9 or 0.8) to manual annotation, where notably greater number of predictions by BlueBERT-EFO have more ontology similarity to their manual annotation counterparts then the rest of the methods. In other words, BlueBERT-EFO as a finetuned model on BlueBERT with EFO structural information is able to enhance the performance of the foundational BlueBERT to be on par with BioSentVec, and able to incorporate EFO knowledge on candidate retrieval (we note that BioSentVec has not been finetuned, and is not directly comparable to BlueBERT-EFO).

### 3.1.2 Overall results for top N predictions

With methods that produce a distance or score, there may still be significant value in a set of top predictions (which we would expect to be enriched for related terms, and potentially contain the correct mapping term). We then investigated the distribution of EFO-Batet scores for both the top prediction (Fig. 3A) and the top 10 predictions (weighted average EFO-Batet scores, Fig. 3B), and the aggregate results of generalized top ranges, to determine which models prioritize the most relevant set of traits. As shown in Fig. 3A, for top predictions BlueBERT-EFO is able to retrieve higher number of candidates that have high ontology relevance to the manual annotation (greater mass in the upper tail) and lower number of candidates that have low relevance (lower mass in the lower tail), which is also confirmed by the pairwise Kolmogorov–Smirnov two sample tests
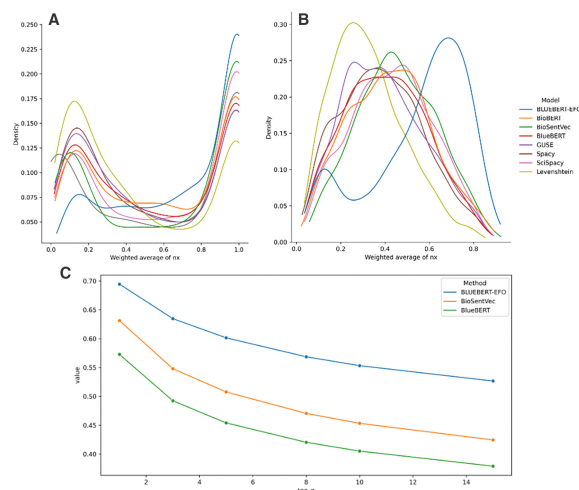


**Figure 3** Distribution of predicted EFO-Batet scores by method. (A) Distribution of EFO-Batet score for the highest-ranking (top 1) match for each query term; (B) distribution of weighted average EFO-Batet score for the top 10 matches for each query term. (C) Averaged sum of the top *N* weighted averaged EFO-Batet score of the predicted EFO candidates for a query term, for subset methods of BlueBERT-EFO, BlueBERT, and BioSentVec (full results are available in Supplementary Fig. S10)

(Supplementary Table S4) on the statistical difference of its distribution to those of other methods ($P - value \leq 3.3 \times 10^{-9}$).

We then extended the analysis to consider a set of top results. Fig. 3B shows the distribution of the weighted average EFO-Batet scores for the top 10 EFO predictions for each method (see Supplementary Fig. S4 for violin plot and Supplementary Table S3 for descriptive statistics on the same data). For top 10 predicted EFO terms, we computed the EFO-Batet score vis-a-vis the manual annotation counterpart, then averaged with the ranking weights (i.e. top prediction getting a weighting of 10, 9, and so on) to show the aggregate ontology relevance of the retrieved candidates. Fig. 3C shows the averaged sum of the weighted average scores for each top *N* level to provide an overall measure on the general ontology relevance of the candidate retrieval for a subset of methods (see Supplementary Fig. S10 for full results). The results suggest that BlueBERT-EFO will generally return a set of traits that are more closely associated with the correct part of the EFO ontology compared to other methods, and corroborates with earlier analysis findings that the finetuning of the BlueBERT language model with EFO structure information will notably improve EFO candidate retrieval.

We also investigated on the performance of a hybrid method (BioSentVec-X-BlueBERT-EFO) where BioSentVec is applied in the first stage to select the top X (e.g. 30) candidates, then BlueBERT-EFO is applied in the second stage to select the top *N* (e.g. 5) candidates, with the aim to improve inference efficiency as transformer models are more computationally expensive than simpler model architectures, such as BioSentVec. Supplementary Figs S5–S7 show the weighted average score distribution for top 1, 5, and 10 matches, and Supplementary Fig. S8 shows the averaged sum of weighted average scores for generalized top *N* levels. These results suggest that top matching results produced by the second stage BlueBERT-EFO in the hybrid methods retain the overall behaviour of BlueBERT-EFO, and is robust to the first stage filtering via BioSentVec.

To try and understand why certain traits are challenging to map, and why others are not, we extracted the top UK Biobank queries, which were most and least variable in EFO-Batet score between methods. Details of this can be found in Supplementary Text S2.

### 3.2 Comparison to other approaches for trait-to-trait mapping

Our final set of analyses explores the differences in direct trait-to-trait mapping of the different models. For each model, we estimated
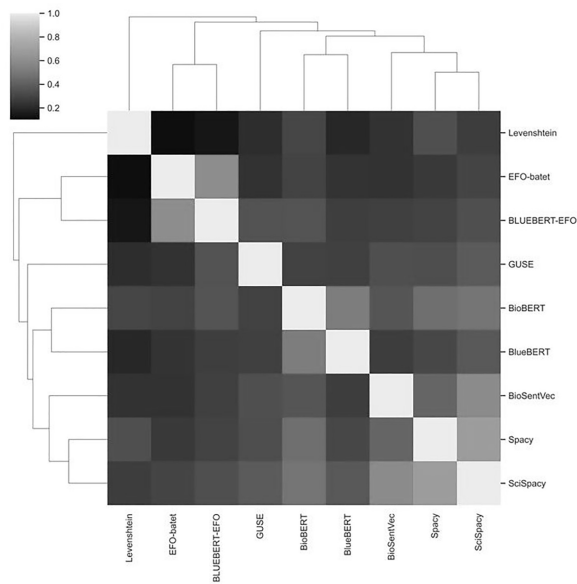
**Figure 4** Pairwise plot of spearman correlations between methods based on a matrix of cosine similarity (or equivalent) scores for all pairwise combines of traits (excluding self)

trait similarity scores between each trait ($n = 530$, see Section 2) and all others (excluding itself). Fig. 4 shows the results of a Spearman rank correlation analysis comparing the matrices of these pairwise trait-mapping scores between each pair of models. The results broadly indicate three clusters of models. One contains the EFO-Batet (manual mapping) and BlueBERT-EFO scores, suggesting again that the BlueBERT-EFO model, as expected, is predicting distances most similar to that which we find in the EFO hierarchy. A second group contains the other BERT models (BioBERT and BlueBERT) highlighting the similarity between those two transformer models. A third group contains the spaCy, ScispaCy, and BioSentVec models, which may represent their shared underlying methodology (i.e. variations of word2vec). Whilst this analysis cannot tell us which method performs "best" at trait-to-trait mapping, it highlights that these models do perform differently at this task, which should be taken into account in the development of future automated trait-to-trait mapping methods.

Finally, we present visualizations of the trait similarity scores for all pairwise trait-to-trait mappings for a selected set of 43 traits (representing a mixture of disease, continuous traits, and medications) to illustrate how these models perform at this task. Supplementary Fig. S11 is provided as a reference and shows a clustered dendrogram of EFO-Batet scores for the distance between traits in the EFO hierarchy. The clusters represent the relationships between EFO terms as determined by the EFO hierarchy and Batet scores. We observe a sharp separation between measurement-based quantitative traits and disease traits. This reflects the structure of the EFO, with quantitative traits falling into the "information entity" and disease traits into the "material property" top-level branches of EFO.

Using the same 43 traits, we then produced a matrix of trait-to-trait distance scores for each model, but this time based on cosine distances (or equivalent—see Section 2). These matrices were compared to each other using the Mantel test in scikit-bio (http://scikit-bio.org/), a method to compute correlation distances between matrices (Supplementary Fig. S12). Here, we see a similar pattern, with the BlueBERT-EFO and EFO-Batet (i.e. position in the EFO hierarchy) scores clustered together. This similarity is obvious in the BlueBERT-EFO clustermap (Supplementary Fig. S13) where there are some clear differences, but the major distinction between quantitative traits and disease is present, with almost exactly the same traits clustering into the same two groups. This likely reflects the finetuning of this model to EFO.

## 4 Discussion

A number of approaches exist for text matching and semantic representation of text. We set out to investigate the use of these approaches for the automated mapping of human trait names to ontologies (using the specific example of EFO) and explore how they perform at direct trait-to-trait mapping.

### 4.1 Comparison of approaches for automated mapping to ontology

Our analyses illustrate that using text embeddings to map biomedical variables to EFO has a fairly high error rate, but is at least comparable to existing approaches (e.g. Zooma). Given the ease of use and scalability of some of the models, we recommend this approach when tackling problems that involve many thousands of variables and manual annotation is not feasible. When attempting an exact match (i.e. top match) BioSentVec (Chen et al. 2019) appears to perform best in terms of speed, precision, and accuracy. However, if it is more important that the top $N$ predictions are close to the truth, then BlueBERT-EFO consistently out-performed all other models. The increase in performance of the BlueBERT model when finetuned to the EFO suggests that finetuning of other models to EFO could yield similar gains, and not that there is something inherently superior about the BlueBERT model.

It is important to note that several of the models had similar performance at finding a top match, with the group including BioSentVec, BlueBERT-EFO, Zooma, and ScispaCy (Neumann et al. 2019) showing little statistical evidence of a difference. In contrast to the other models, the standard Zooma tool also brings the benefit of continually updated manually curated mappings.

Embedding methods appear to perform well when the query string describes a single event or entity, e.g. "whooping cough/pertussis". They perform poorly when the query string describes multiple entities, e.g. "hiv disease resulting in malignant neoplasms". This is perhaps not surprising, as the embedding of this phrase is unlikely to be close to either HIV or cancer terms. Addressing such traits therefore remains a complex challenge, i.e. properly identifying mentioned concepts via named entity recognition (NER) and then incorporating pre-trained concept embeddings from the knowledge base to the document embeddings (Park et al. 2019, Chen et al. 2020). In other words, a complex processing system, which includes major components of NER, document level embeddings, and concept embeddings, is required to approach mapping of complex traits in a generalized and robust manner, though we are keen to explore this aspect in future research.

We compared our models to a manually mapped set of trait names, but it is important to recognize this may itself contain errors. Supplementary File S7 lists examples where no models predicted an EFO term with an EFO-Batet score >0.95. Here, e.g. the query term "malignant neoplasm of colon" was manually mapped to "colon carcinoma". However, six of the models predicted the EFO term "malignant colon neoplasm", which has an EFO-Batet score of 0.86 and is therefore a better fit (it is possible these differences reflect changes in the EFO since the initial mapping rather than a mapping error).

### 4.2 Comparison of approaches for trait-to-trait mapping

Mapping traits directly between two datasets has potential value, but in the absence of a benchmark it is hard to validate. We therefore focussed on variables that had been mapped to a single EFO term, and then refined that further for closer inspection. The use of clustering methods enabled us to manually inspect groups of traits and describe events that agree with standard biomedical knowledge. Our analyses show that by including topological information from a well-established ontology like the EFO, the BlueBERT-EFO model can create sensible pairwise distances between variables, without actually mapping to ontology.

When focussing on a specific set of traits, we see that whilst the finetuning of BlueBERT-EFO has produced a model which reflects major patterns in the EFO hierarchy, there are some differences.

One example is the loss of the "angina", "worrier/anxious feeling" cluster (present in EFO, Supplementary Fig. S11), with "angina" joining the larger disease cluster next to "atrial fibrillation and flutter" and "worrier/anxious feeling" moving next to "neuroticism score" (Supplementary Fig. S13). The manual EFO term assigned to "angina" was "EFO_0003913" (angina pectoris, http://www.ebi.ac.uk/efo/EFO_0003913), which can be found within the "material phenotype" EFO group as it is listed as a "Phenotype abnormality" and not a disease. Even though the BlueBERT-EFO model has been finetuned on the EFO hierarchy, the biomedical literature underpinning the model has created distances placing "angina" with other diseases rather than measurements. This highlights the subtle balance of information contained within this model.

Interestingly, the BlueBERT-EFO model fails to group together the neurological illnesses ("parkinson's disease", "alzheimer's disease", and "secondary parkinsonism"). Looking at the other models, several also fail to do this, often grouping traits with the word "disease" together (Supplementary Figs S14–S20). However, BioSentVec, BlueBERT, and BioBERT appear to group these appropriately. This highlights one of the key challenges that the developers of these models face: how to distinguish between informative words and ignore the generic (e.g. "disease"). This point is again present in the BioBERT cluster map (Supplementary Fig. S19), with "weight" an outlier to all other traits, suggesting this term was not sufficiently similar to anthropometric traits.

It is worth noting, that the alternative methods to using language models for this type of distance analysis appear to perform less well (e.g. Levenshtein edit distance, Supplementary Fig. S14). Other established methods, such as Zooma, are just not possible to use when comparing data in this way.

At the moment there is no practical alternative automated approach to trait-to-trait mapping, so our results using language models are promising. However, they are far from perfect with many cases of traits not grouping together as we might expect, and the models often focussing on generic words, such as disease over and above other more defining terms. This approach therefore requires further development before it can be of practical use.

### 4.3 Use cases of these models
The models are imperfect but are still successful in mapping 40% of trait names in the dataset we used. One obvious use case would be a semi-automated mapping tool which would provide a suggestion for the user to approve or edit. As highlighted above, many simple trait names map well, and it is the more complex traits (e.g. combinations of entities) that would need manual intervention.

Another scenario in which an imperfect one-to-many mapping tool like those presented here may be useful is in a "trait name recommender". One example of this is our OpenGWAS (Elsworth et al. 2020) recommender, which provides recommended trait matches from amongst thousands of GWAS datasets to enable a user to see other relevant GWAS traits they may be interested in. The OpenGWAS recommender uses a combination of ScispaCy and BlueBERT-EFO to search for the top matching GWAS traits in the semantic embedding vector space and optionally predict the ontology relationships between the query term and the match candidates (Liu et al. 2021).

In a follow-up study (Liu and Gaunt 2022), we applied ScispaCy and BlueBERT-EFO as an ontology mapper in a hybrid architecture, where a first stage model is used to efficiently filter EFO ontology candidates associated with the query ULMS terms, and in the second stage BlueBERT-EFO is then used to predict the ranking of the top $N$ results (similar to results in Supplementary Figs S5–S8 where BioSentVec was applied as the first stage model). The retrieval results have shown to be sensible for the systematic analysis on medRxiv submission abstracts, without sacrificing inference performance due to the computationally expensive nature of transformer models whilst retaining relevancy in candidate retrieval.

## 5 Conclusions
We have shown that current text matching and embedding approaches offer some promise in the task of mapping traits to ontologies and to each other. However, the mapping is imperfect and unsuitable for fully automated mapping. Models trained on the biomedical literature perform better than more generalized models, and finetuning to the EFO ontology may further improve the performance of a specific model. Some trait names present in population health datasets, such as UK Biobank, are complex and their embeddings are unlikely to be very representative; future work should focus on how to handle such trait names.

## Supplementary data
Supplementary data is available at *Bioinformatics* online.

## Funding

## References
Allot A, Chen Q, Kim S *et al.* LitSense: making sense of biomedical literature at sentence level. *Nucleic Acids Res* 2019;**47**:W594–9.

Batet M, Sánchez D, Valls A *et al.* An ontology-based measure to compute semantic similarity in biomedicine. *J Biomed Inform* 2011;**44**:118–25.

Blagec K, Xu H, Agibetov A *et al.* Neural sentence embedding models for semantic similarity estimation in the biomedical domain. *BMC Bioinformatics* 2019;**20**:178.

Cer D, Yang Y, Kong S *et al.* Universal sentence encoder. *arXiv* 2018. https://doi.org/10.48550/arXiv.1803.11175

Chen Q, Peng Y, Lu Z. BioSentVec: creating sentence embeddings for biomedical texts. In: *2019 IEEE International Conference on Healthcare Informatics (ICHI)*, pp. 1–5. Xi'an, China: IEEE, 2019.

Chen Q, Lee K, Yan S *et al.* BioConceptVec: creating and evaluating literature-based biomedical concept embeddings on a large scale. *PLoS Comput Biol* 2020;**16**:e1007617.

Demner-Fushman D, Rogers WJ, Aronson AR. MetaMap Lite: an evaluation of a new Java implementation of MetaMap. *J Am Med Inform Assoc* 2017;**24**:841–4.

Denny JC, Ritchie MD, Basford MA *et al.* PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene-disease associations. *Bioinformatics* 2010;**26**:1205–10.

Devlin J, Chang MW, Lee Kenton, *et al.* BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of the 2019 Conference of the North*. 4171–86. Minneapolis, Minnesota: Association for Computational Linguistics, 2019.

Duong D, Ahmad WU, Eskin E *et al.* Word and sentence embedding tools to measure semantic similarity of Gene Ontology terms by their definitions. *bioRxiv* 2017. https://doi.org/10.1101/103648.

Duong D, Uppunda A, Ju C *et al.* Evaluating representations for Gene Ontology terms. *bioRxiv* 2019. https://doi.org/10.1101/765644.

Elsworth B, Lyon M, Alexander T *et al.* The MRC IEU OpenGWAS data infrastructure. *bioRxiv* 2020. https://doi.org/10.1101/2020.08.10.244293.

Fabian B, Edlich T, Gaspar H *et al.* Molecular representation learning with language models and domain-relevant auxiliary tasks. *bioRxiv* 2020. https://doi.org/10.48550/arXiv.2011.13230.

Gu Y, Tinn R, Cheng H *et al.* Domain-specific language model pretraining for biomedical natural language processing. *arXiv* 2021. https://doi.org/10.48550/arXiv.2007.15779.

Jaeger S, Fulle S, Turk S *et al.* Mol2vec: unsupervised machine learning approach with chemical intuition. *J Chem Inf Model* 2018;**58**:27–35.

Johnson AEW, Pollard TJ, Shen L *et al.* MIMIC-III, a freely accessible critical care database. *Sci Data* 2016;**3**:160035.

Jones R, Pembrey M, Golding J *et al.* The search for genotype/phenotype associations and the phenome scan. *Paediatr Perinat Epidemiol* 2005;**19**: 264–75.

Karadeniz İ, Özgür A. Linking entities through an ontology using word embeddings and syntactic re-ranking. *BMC Bioinformatics* 2019;**20**:156.

Kibbe WA, Arze C, Felix V *et al.* Disease Ontology 2015 update: an expanded and updated database of human diseases for linking biomedical knowledge through disease data. *Nucleic Acids Res* 2015;**43**:D1071–8.

Koroleva A, Kamath S, Paroubek P *et al.* Measuring semantic similarity of clinical trial outcomes using deep pre-trained language representations. *J Biomed Inform* 2019;**100**:100058.

Lawlor DA, Tilling K, Davey Smith G *et al.* Triangulation in aetiological epidemiology. *Int J Epidemiol* 2016;**45**:1866–86.

Le QV, Mikolov T. Distributed representations of sentences and documents. *arXiv* 2014. https://doi.org/10.48550/arXiv.1405.4053.

Lee J, Yoon W, Kim S *et al.* BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 2020;**36**:1234–40.

Levenshtein VI. Binary codes capable of correcting deletions, insertions, and reversals. *Sov Phys Dokl* 1966;**10**:707–10.

Liu H, Perl Y, Geller J. Transfer learning from BERT to support insertion of new concepts into SNOMED CT. *AMIA Annu Symp Proc* 2019;**2019**:1129–38.

Liu Y, Elsworth B, Erola P *et al.* EpiGraphDB: a database and data mining platform for health data science. *Bioinformatics* 2021;**37**:1304–11.

Liu Y, Gaunt TR. Triangulating evidence in health sciences with Annotated Semantic Queries. *medRxiv* 2022. https://doi.org/10.1101/2022.04.12.22273803.

Malone J, Holloway E, Adamusiak T *et al.* Modeling sample variables with an Experimental Factor Ontology. *Bioinformatics* 2010;**26**:1112–8.

Mikolov T *et al.* Efficient estimation of word representations in vector space. *arXiv* 2013. https://doi.org/10.48550/arXiv.1301.3781.

Millard LAC, Davies NM, Tilling K *et al.* Searching for the causal effects of body mass index in over 300 000 participants in UK Biobank, using Mendelian randomization. *PLoS Genet* 2019;**15**:e1007951.

Montani I, Honnibal M, Honnibal M *et al.* explosion/spaCy: v3.1.0: new pipelines for Catalan & Danish, SpanCategorizer for arbitrary overlapping spans, use predicted annotations during training, bug fixes & more. 2021. https://doi.org/10.5281/zenodo.5079800.

Neumann M, King D, Beltagy I *et al.* ScispaCy: fast and robust models for biomedical natural language processing. In: *Proceedings of the 18th BioNLP Workshop and Shared Task*, pp. 319–27. Florence, Italy: Association for Computational Linguistics, 2019.

Pagliardini M, Gupta P, Jaggi M. Unsupervised learning of sentence embeddings using compositional n-gram features. In: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Vol. 1 (Long Papers)*, pp. 528–40. New Orleans, Louisiana: Association for Computational Linguistics. 2018.

Park J, Kim K, Hwang W *et al.* Concept embedding to measure semantic relatedness for biomedical information ontologies. *J Biomed Inform* 2019;**94**: 103182.

Pendlington ZM, Roncaglia P, Mountjoy E *et al.* Mapping UK Biobank to the experimental factor ontology. ISMB/ECCB 2019.

Peng Y, Yan S, Lu Z. Transfer learning in biomedical natural language processing: an evaluation of BERT and ELMo on ten benchmarking datasets. *arXiv* 2019. https://doi.org/10.48550/arXiv.1906.05474.

Robinson PN, Köhler S, Bauer S *et al.* The Human Phenotype Ontology: a tool for annotating and analyzing human hereditary disease. *Am J Hum Genet* 2008;**83**:610–5.

Tshitoyan V, Dagdelen J, Weston L *et al.* Unsupervised word embeddings capture latent knowledge from materials science literature. *Nature* 2019;**571**: 95–8.

Vaswani A, Shazeer N, Parmar N *et al.* Attention is all you need. *arXiv* 2017. https://doi.org/10.48550/arXiv.1706.03762.

Wang A, Singh A, Michael J *et al.* GLUE: a multi-task benchmark and analysis platform for natural language understanding. In: *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 353–5. Brussels, Belgium: Association for Computational Linguistics, 2018.

Wolf T, Debut L, Sanh V *et al.* HuggingFace's transformers: state-of-the-art natural language processing. *arXiv* 2019. https://doi.org/10.48550/arXiv.1910.03771.

Zhang Y, Chen Q, Yang Z *et al.* BioWordVec, improving biomedical word embeddings with subword information and MeSH. *Sci Data* 2019;**6**:52.