

CS 598 DLH Project Proposal

Amit Jangid, Constantin Kappel, Daniel Sanchez

March 24th 2024

ajangid2@illinois.edu, normank2@illinois.edu, daniel43@illinois.edu

1 Citation to original Paper

We picked a publication by Hur et al. [4] titled “*Unifying Heterogeneous Electronic Health Records Systems via Text-Based Code Embedding*”.

2 General Problem

While EHR (Electronic Health Records) provide an attractive data source to do research on exposures and disease, there are many heterogeneous medical code formats used by different healthcare providers. Thus, clinical studies using EHRs are difficult to scale up due to data incompatibilities. The technical hurdle is that previous systems learned hidden representations (latent embeddings) for each code system which projected the same medical concept with different encodings into different, incompatible semantic spaces. The paper [4] learns from the unstructured textual descriptions instead of the medcodes themselves, a strategy which the authors termed *description embedding* (*DescEmb*). Using *DescEmb* the authors were able to pool differently structured EHRs, namely *MIMIC-III* and *eICU*, into one larger dataset and achieved higher accuracy.

3 Scientific Approach

1. First, we will try to replicate the results of [4] as a baseline.
2. We plan on using `pyhealth` applicable. Especially for handling EHR data this should be very helpful.
3. Thirdly, we intend to test a hypothesis described in more detail in section 4. In short, we would like to replace the BERT¹-derived encoders by one which was pretrained from scratch on biomedical data.

As we outlined in section 7 we will not do comprehensive pre-training from scratch, but work with pre-trained architectures, which we may fine-tune for our experiments.

4 Hypotheses to be tested

The authors of [4] conducted experiments with several BERT-based architectures, such as BioBERT (author?) [7], ClinicalBERT (author?) [3] and BlueBERT, which were partially trained or fine-tuned on medical literature. In contrast, Gu et al. [2] found that pre-training on generic NLP corpora,

¹Bidirectional Encoder Representations from Transformers (author?) [1]

such as derived from Wikipedia, actually perform worse than ones which have been initialized with random values and only pre-trained on PubMed. They called their model PubMedBERT and were able to demonstrate superior performance to the models employed by [4].

We would like to test if PubMedBERT is capable of matching or even superceding the performance of the other models utilized, which were not fully pre-trained on pure corpi from the biomedical domain.

5 Ablations planned

A true ablation study with PubMedBERT and e.g. ClinicalBERT or BioBERT would mean we would have to pre-train all these architectures from scratch on both, generic NLP corpi as well as PubMed. This will not be computationally feasible. Instead, we will cite the work of [2] where applicable and rather try to replicate their results with different, already pre-trained architectures.

6 Data access

The two datasets we plan on utilizing, namely MIMIC-III (**author?**) [5, 6] and eICU (**author?**) [8], are both publicly available. We will need to take into account some regulations and license conditions concerning data safety and privacy.

7 Feasibility of the computation

As far as computational cost is concerned let’s consider the size of the data and the size of the model architecture we need to train.

7.1 Model architecture

According to (**author?**) [4] p. 4, bottom right paragraph, the largest architecture the authors trained was BERT-base, which has 110 M (10^6) parameters. Other models the authors experimented with include BioBERT and ClinicalBERT, which are based on BERT-base, too. Let’s do a quick estimate² of how much GPU memory is needed to load and train such a model:

	Bytes per parameter
Model parameters	4
Adam optimizer	8
Gradients	4
Activations and temp memory	8
TOTAL	24

Thus, a model with 110 M parameters would consume $110 \cdot 10^6 \times 24\text{bytes} = 2.46\text{GB}$ to just load the model. This managable on modern-day GPUs. One of the authors has a GTX 3080 with 12 GB of VRAM available for full training and a GTX 1070 with 8 GB is available to do some code testing on small batches. As mentioned in sections 3 and 4 we are not planning to do pretraining. While Hur et al. (**author?**) [4] did not explicetely provide information on training time, the authors of BioBERT mentioned that pre-training took 10 days on a V100 GPU. Given the iterative nature of the doing model training and the amount of time we can spend, the cost for pre-training from scratch would be prohibitive for us. So, we will only do inferencing with available pre-trained architectures and fine-tuning where applicable.

²Assuming 32-bit single precision computation

7.2 Datasets

We will use two datasets : MIMIC-III and eICU. According to ... the MIMIC-III dataset is a large, de-identified, and publicly-available collection of medical records. It consists of 112,000 clinical reports records with an average length of 709.3 tokens. The dataset includes 1,159 top-level ICD-9 codes, and each report is assigned to 7.6 codes on average. This makes for about 78056000 tokens in MIMIC-III.

(?? expand the calculation ??)

The eICU Collaborative Research Database holds data associated with over 200,000 patient stays, providing a large sample size for research studies.

(?? expand this ??)

According to **(author?)** [1] the size of the token embeddings is 768, the number of transformer blocks is 12 and the number of self-attention heads is 12.

(?? use for building an argument over computation time, also using ??)

8 Existing Code

References

- [1] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv:1810.04805 [cs]*, May 2019. arXiv: 1810.04805.
- [2] Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23, January 2022. arXiv:2007.15779 [cs].
- [3] Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. ClinicalBERT: Modeling Clinical Notes and Predicting Hospital Readmission, November 2020. arXiv:1904.05342 [cs].
- [4] Kyunghoon Hur, Jiyoung Lee, Jungwoo Oh, Wesley Price, Young-Hak Kim, and Edward Choi. Unifying Heterogeneous Electronic Health Records Systems via Text-Based Code Embedding, March 2022. arXiv:2108.03625 [cs].
- [5] Alistair Johnson, Tom Pollard, and Roger Mark. MIMIC-III Clinical Database, 2015.
- [6] Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. MIMIC-III, a freely accessible critical care database. *Scientific Data*, 3(1):160035, May 2016. Publisher: Nature Publishing Group.
- [7] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, page btz682, September 2019. arXiv:1901.08746 [cs].
- [8] Tom J. Pollard, Alistair E. W. Johnson, Jesse D. Raffa, Leo A. Celi, Roger G. Mark, and Omar Badawi. The eICU Collaborative Research Database, a freely available multi-center database for critical care research. *Scientific Data*, 5(1):180178, September 2018.