

# Exploring Long-Term Prediction of Type 2 Diabetes Microvascular Complications

Elizabeth Remfry

Rafael Henkin

Michael R Barnes

*Queen Mary University of London, UK*

Aakanksha Naik

*Allen Institute for AI, USA*

E.A.REMFRY@QMUL.AC.UK

R.HENKIN@QMUL.AC.UK

M.R.BARNES@QMUL.AC.UK

AAKANKSHAN@ALLEN.AI.ORG

## Abstract

Electronic healthcare records (EHR) contain a huge wealth of data that can support the prediction of clinical outcomes. EHR data is often stored and analysed using clinical codes (ICD10, SNOMED), however these can differ across registries and healthcare providers. Integrating data across systems involves mapping between different clinical ontologies requiring domain expertise, and at times resulting in data loss. To overcome this, code-agnostic models have been proposed. We assess the effectiveness of a code-agnostic representation approach on the task of long-term microvascular complication prediction for individuals living with Type 2 Diabetes. Our method encodes individual EHRs as text using fine-tuned, pretrained clinical language models. Leveraging large-scale EHR data from the UK, we employ a multi-label approach to simultaneously predict the risk of microvascular complications across 1-, 5-, and 10-year windows. We demonstrate that a code-agnostic approach outperforms a code-based model and illustrate that performance is better with longer prediction windows but is biased to the first occurring complication. Overall, we highlight that context length is vitally important for model performance. This study highlights the possibility of including data from across different clinical ontologies and is a starting point for generalisable clinical models.

**Keywords:** Clinical language models, electronic healthcare records, multi-label classification, disease prediction, pretrained models, type 2 diabetes, time series

**Data and Code Availability** This study uses the Clinical Practice Research Datalink (CPRD), real-world anonymised patient data from primary care

across the UK and linked to other health related registries. CPRD AURUM includes routinely collected data on 19 million patients including demographics, diagnoses, symptoms, prescriptions, referrals, lifestyle factors and tests (Wolf et al., 2019). Data access is subject to approval from an Independent Scientific Advisory Committee (ISAC). Code is available [github.com/LizRem/diabetes-complications](https://github.com/LizRem/diabetes-complications)

**Institutional Review Board (IRB)** The application was reviewed by an (ISAC) and the data were used under license for the current study.

## 1. Introduction

Type 2 Diabetes (T2DM) is a long-term cardiometabolic condition associated with increased risk of microvascular complications; diabetic retinopathy, nephropathy and neuropathy. These complications can result in severe outcomes, such as vision loss, end stage renal disease and amputations, respectively (Brownrigg et al., 2016; Khanam et al., 2017). Approximately one-third of individuals living with T2DM develop at least one of these complications (Arnold et al., 2022), which in turn increases the risk of developing others (Deshpande et al., 2008). As various risk factors for microvascular complications are modifiable, timely identification of individuals at high risk of developing these diseases can help to inform treatment pathways and healthcare interventions (Khalil, 2017; Lu et al., 2023).

Recent research has demonstrated the utility of deep learning models for disease prediction tasks due to their ability to handle messy electronic healthcare record (EHR) data which is temporal, sparse and high-dimensional (Hassaine et al., 2020; Wornow et al., 2023). Deep learning approaches for such tasks

typically represent diseases as clinical codes, which requires mapping between heterogeneous clinical ontologies and manual curation or reduction of codes. Moreover, code-based representations also make such approaches less likely to generalize to unseen diseases and complications as well as across different health-care settings.

To address these caveats with code-based representations, our study explores a code-agnostic design taking inspiration from Munoz-Farre et al. (2022); Hur et al. (2022). This approach leverages existing clinical knowledge embedded in pre-trained language models and integrates a wider range of data from across different health registries. We combine this with a multi-label approach which enables us to construct shared representations of T2DM complications, which is beneficial as complications are closely related and often share various risk factors. We explore disease prediction over short-, mid- and long-term time windows.

## 2. Related Work

There are a plethora of pre-trained clinical language models, however, due to data privacy very few are publicly available and those that are, come with limitations due to the heterogeneity of code ontology used in the training data (Wornow et al., 2023).

To navigate this challenge of detaching models from the specific ontologies, research has started to utilise the natural language descriptions of the clinical codes. Munoz-Farre et al. (2022) utilised textual descriptors fed into an encoder only model pre-trained on clinical literature and then fine-tuned to predict various diseases. They reported improved performance compared to a model trained using traditional code embeddings. Hur et al. (2022) compared various model set ups; trained from scratch, continual pre-training and fine-tuning, on textual descriptors from MIMIC-III and eICU and found that BERT performed similarly to the models trained on clinical literature, even under different training approaches.

Our work builds on previous studies by including a broader range of clinical data at a granular level without aggregating codes in clinical hierarchies. We include all textual descriptors within the EHR, which includes diagnoses, prescriptions, symptoms, referrals and procedures. We particularly focus on the prediction of T2DM microvascular complications over different and longer time intervals.

## 3. Methods

### 3.1. Cohort

We analysed EHRs from CPRD AURUM and included all individuals  $\geq 18$ , permanently registered to any General Practice in London between 01/01/2010 and 01/01/2020, see Data and Code Availability for more details.

Our dataset included 133,784 patient records, with 44,820 experiencing at least one microvascular complication Table 1. A diagnosis of T2DM, retinopathy, neuropathy or nephropathy were identified using validated phenotype definitions and we used the first occurring diagnosis date (Eto, 2023). Patients with micro-vascular complications prior to a diagnosis of T2DM were excluded.

Study entry was defined as the first EHR event until the visit prior to the first recorded complication, or the last recorded event for those without complications. We evaluated 1-, 5- and 10- years risk prediction windows post first complication. Only patients with at least 3 unique events were included.

Table 1: Cohort Characteristics

Characteristic	
Number of patients	133,784
Total number of complications	
0	88,964
1	33,161
2	9,282
3	2,377
Number with each complication	
Retinopathy	31,396
Nephropathy	19,595
Neuropathy	7,865
Sex	
Male	72,012
Female	61,772
Age at first complication (SD)	63.06 (14.73)

### 3.2. EHR pre-processing

Every clinical code is associated with a textual descriptor, for example the ICD10 code *E11.9* is associated with *type 2 diabetes mellitus without complications*. For our text-based approach we take the textual descriptor for every event in a patient’s EHR (diagnosis, procedure, symptoms, prescription, etc.). All textual terms are then concatenated chronologically to generate text sentences for each patient. For our code-based approach, we take the clinical code

Table 2: Performance Comparison of Text- and Code-based Models

	Text-based		Code-based	
	Micro-F1	Micro-AUPRC	Micro-F1	Micro-AUPRC
<b>1 year</b>	0.45 (0.44-0.46)	0.44 (0.43-0.46)	0.43 (0.42- 0.44)	0.40 (0.38-0.41)
<b>5 year</b>	<b>0.50 (0.49-0.51)</b>	<b>0.51 (0.50-0.52)</b>	0.43 (0.42-0.44)	0.43 (0.41-0.44)
<b>10 year</b>	0.49 (0.48-0.50)	<b>0.50 (0.49-0.51)</b>	0.47 (0.46-0.49)	0.47 (0.45-0.48)

Note: Values in brackets represent 95% confidence intervals, bold indicates statistical significance

and concatenate them chronologically producing a sequence of codes Appendix A.

### 3.3. Model architecture

We utilised a pretrained clinical language model, GatorTron-base (Yang et al., 2022), to encode the tokenized EHR sequences. All sequences were truncated or padded to 512 tokens, the maximum length for GatorTron. We then fine-tuned the pretrained model, one for each risk prediction window. The models consisted of a fine-tuned encoder with a single linear output layer with 3 output nodes. We split our data 80/10/10 into training, test and validation using stratified sampling to ensure the imbalance remained the same. We used weighted cross entropy due to label imbalance and report on micro F1, micro recall and micro area under the precision recall curve (AUPRC). All results are presented calculated on the held out test set. For more information on pre-processing and architecture see Appendix A.

We assess the variation in model performance and calculate a 95% confidence interval (CI), by employing a bootstrap resampling technique. Using our test set of 13,314 patients we performed 1000 bootstrap iterations. Pairwise comparisons between models were conducted using a z-test approach, where the standard error of the difference was derived from the bootstrapped CIs. To account for multiple comparisons, we applied the Bonferroni correction adjusting our significance threshold.

## 4. Results

**Code-agnostic models outperform code-based models:** Models trained on textual descriptors performed significantly better than models trained on clinical codes although not at all time windows (Table 2). This suggests that there is utility in using the textual terms which may allow the model to take advantage of existing clinical knowledge.

**Models perform better over longer prediction timeframes:** the 5-year risk prediction window achieved a micro-AUPRC of 0.51 (Table 2). This is likely due to the number of additional labels providing a more balanced dataset, as the longer prediction windows increases the likelihood of observing a complication.

**The multi-label design is biased towards first-occurring T2DM complication:** Across all time frames, retinopathy is the highest achieving class (Table 3), this is likely due to being the most commonly occurring first condition (in 60.19% of cases) and the largest class. Nephropathy is the first complication in only 30.15% of cases and neuropathy 9.67%. As the model is only exposed to data up until the visit prior to the first complication and complications can occur at different timepoints across the life course this early data may not contain sufficient information for the model to make an accurate prediction about subsequent complications.

**Restrictions on context length affects performance:** we explored the average number of tokens in each individual’s EHR (median: 2272), which falls substantially over the capability of GatorTron at 512 maximum token length. This results in the truncation of 85.37% of sequences leading to data loss, see Appendix A for further exploration. In order to mimic clinicians behaviour, where they typically look at recent events first in an EHR, we mirror this by truncating from the left, removing the earliest data and preserving the most recent events Table 4.

Truncating from left led to improved performance across all prediction windows, suggesting that the EHR events recorded closer in time to a diagnosis of a complication are more important than events happening earlier. We also present the performance of pretrained models with longer context length (4096 token length) in Appendix A, and indicate that shorter context lengths negatively impact performance.

Table 3: F1 and Recall Scores for Microvascular Complications at 1, 5, and 10 years

Time	Nephropathy		Retinopathy		Neuropathy	
	F1	Recall	F1	Recall	F1	Recall
<b>1 year</b>	0.39	0.44	0.51	0.54	0.22	0.18
<b>5 year</b>	0.42	0.45	0.55	0.57	0.29	0.29
<b>10 year</b>	0.44	0.51	0.55	0.55	0.30	0.30

Table 4: F1 and Recall Scores for Microvascular Complications and Micro-F1/AUPRC at 1, 5, and 10 years for Models Truncated Left

Time	Nephropathy		Retinopathy		Neuropathy		Micro-F1	Micro-AUPRC
	F1	Recall	F1	Recall	F1	Recall		
<b>1 year</b>	0.53	0.54	0.70	0.75	0.34	0.35	0.61	0.64
<b>5 year</b>	0.53	0.59	0.73	0.73	0.38	0.35	0.62	0.66
<b>10 year</b>	0.57	0.61	0.74	0.75	0.39	0.40	0.64	0.69

## 5. Discussion and future work

We present a code-agnostic method for long-term microvascular complication prediction in Type 2 Diabetes that utilises textual descriptors associated with clinical codes, unifying data across different health registries and taking advantage of pretrained language models.

**Real-world assessment of reusability of pre-trained models:** Our study found that pre-trained models yielded relatively low performance on T2DM complication prediction over various time-frames despite being heralded as reusable and capable of saving time and resources. Other previous studies (Munoz-Farre et al., 2022; Hur et al., 2022) have also yielded varying performances depending on model design, indicating that we should more thoroughly investigate reusability of pretrained models for real-world prediction tasks. Some work has tried extract more utility out of pretrained models via continual pretraining. Munoz-Farre et al. (2022) conducted continual pretraining using a MLM task and then fine-tuned a pretrained model which demonstrated better performance with an average AUPRC 0.61 across 4 diseases. In the future, we plan to investigate continual pretraining with span-based MLM that masks out multiple tokens representing a medical concept or phrase (Joshi et al., 2020) as in our setting multiple tokens may represent a single concept (e.g. Type 2 Diabetes).

**Incorporating data beyond text:** In our approach we limited ourselves to text descriptions, however there are additional sources of data such as numerical test results that could improve performance. For instance, Hur et al. (2022) combined both textual descriptors and numerical embeddings in a fine-tuned BERT model and achieved a 0.59 AUPRC on a multi-label disease classification task. We leave investigation of this to future work.

**Addressing context length limitations:** As we include more records from across different registries, creating a richer picture of a patient, this limits model performance as models are unable to handle longer sequences and capture long term dependencies. We plan to assess models with longer context lengths (Beltagy et al., 2020), as well as hierarchical models to further improve this (Li et al., 2023).

**Intrinsic task difficulty:** Finally, some diseases may be clinically harder to predict using language models. Li et al. (2020) with a model trained on codes from scratch achieved AUPRC of 0.53 in a multi-label classification task across 301 diseases (classes). The performance varied, from 0.07 AUPRC for hearing loss, to 0.65 AUPRC for epilepsy although both diseases had roughly the same occurrence ratio of 0.02.

We believe that this work will prompt discussion around generalisable multi-purpose language models not tied to one specific healthcare setting or ontology and promote research comparing performance across different datasets.

## Acknowledgments

This work uses data provided by patients and collected by the NHS as part of their care and support. ER is funded by the Wellcome Trust Health Data in Practice (HDiP) Programme (218584/Z/19/Z). RH is funded by the AI for Multiple Long Term Conditions (AIM) Programme (NIHR203982). RM is supported by Barts Charity (MGU0504). This research utilised Queen Mary’s Apocrita HPC facility, supported by QMUL Research-IT.

## References

- Suzanne V Arnold, Kamlesh Khunti, Fengming Tang, Hungta Chen, Javier Cid-Ruzafa, Andrew Cooper, Peter Fenici, Marilia B Gomes, Niklas Hammar, Linong Ji, Gabriela Luporini Saraiva, Jesús Medina, Antonio Nicolucci, Larisa Ramirez, Wolfgang Rathmann, Marina V Shestakova, Iichiro Shimomura, Filip Surmont, Jiten Vora, Hirotaka Watada, and Mikhail Kosiborod. Incidence rates and predictors of microvascular and macrovascular complications in patients with type 2 diabetes: Results from the longitudinal global discover study. *American Heart Journal*, 243:232–239, January 2022. ISSN 0002-8703. doi: 10.1016/j.ahj.2021.10.181. URL <https://www.sciencedirect.com/science/article/pii/S0002870321004336>.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The Long-Document Transformer, December 2020. URL <http://arxiv.org/abs/2004.05150>. arXiv:2004.05150 [cs].
- Jack R. W. Brownrigg, Cian O. Hughes, David Burleigh, Alan Karthikesalingam, Benjamin O. Patterson, Peter J. Holt, Matthew M. Thompson, Simon de Lusignan, Kausik K. Ray, and Robert J. Hinchliffe. Microvascular disease and risk of cardiovascular events among individuals with type 2 diabetes: a population-level cohort study. *The Lancet Diabetes & Endocrinology*, 4(7):588–597, July 2016. ISSN 2213-8587, 2213-8595. doi: 10.1016/S2213-8587(16)30057-2. URL [https://www.thelancet.com/journals/landia/article/PIIS2213-8587\(16\)30057-2/abstract](https://www.thelancet.com/journals/landia/article/PIIS2213-8587(16)30057-2/abstract). Publisher: Elsevier.
- Anjali D Deshpande, Marcie Harris-Hayes, and Mario Schootman. Epidemiology of Diabetes and Diabetes-Related Complications. *Physical Therapy*, 88(11):1254–1264, November 2008. ISSN 0031-9023. doi: 10.2522/ptj.20080020. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3870323/>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- Fabiola Eto. MULTIPLY-Initiative, August 2023. URL <https://github.com/Fabiola-Eto/MULTIPLY-Initiative>. original-date: 2020-11-25T17:37:13Z.
- Abdelaali Hassaine, Gholamreza Salimi-Khorshidi, Dexter Canoy, and Kazem Rahimi. Untangling the complexity of multimorbidity with machine learning. *Mechanisms of Ageing and Development*, 190:111325, September 2020. ISSN 0047-6374. doi: 10.1016/j.mad.2020.111325. URL <https://www.sciencedirect.com/science/article/pii/S0047637420301214>.
- Kyunghoon Hur, Jiyoung Lee, Jungwoo Oh, Wesley Price, Younghak Kim, and Edward Choi. Unifying Heterogeneous Electronic Health Records Systems via Text-Based Code Embedding. In *Proceedings of the Conference on Health, Inference, and Learning*, pages 183–203. PMLR, April 2022. URL <https://proceedings.mlr.press/v174/hur22a.html>. ISSN: 2640-3498.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. SpanBERT: Improving Pre-training by Representing and Predicting Spans, January 2020. URL <http://arxiv.org/abs/1907.10529>. arXiv:1907.10529 [cs].
- H. Khalil. Diabetes microvascular complications—A clinical update. *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, 11: S133–S139, November 2017. ISSN 1871-4021. doi: 10.1016/j.dsx.2016.12.022. URL

- <https://www.sciencedirect.com/science/article/pii/S1871402116302648>.
- Parvin Akter Khanam, Sayama Hoque, Tanjima Begum, Samira Humaira Habib, and Zafar Ahmed Latif. Microvascular complications and their associated risk factors in type 2 diabetes mellitus. *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, 11:S577–S581, December 2017. ISSN 1871-4021. doi: 10.1016/j.dsx.2017.04.007. URL <https://www.sciencedirect.com/science/article/pii/S1871402117300747>.
- Thomas King, Simon Butcher, and Lukasz Zalewski. Apocrita - High Performance Computing Cluster for Queen Mary University of London. March 2017. doi: 10.5281/zenodo.438045. URL <https://zenodo.org/records/438045>. Publisher: Zenodo.
- Yikuan Li, Shishir Rao, José Roberto Ayala Solares, Abdelaali Hassaine, Rema Ramakrishnan, Dexter Canoy, Yajie Zhu, Kazem Rahimi, and Gholamreza Salimi-Khorshidi. BEHRT: Transformer for Electronic Health Records. *Scientific Reports*, 10(1):7155, April 2020. ISSN 2045-2322. doi: 10.1038/s41598-020-62922-y. URL <https://www.nature.com/articles/s41598-020-62922-y>. Number: 1 Publisher: Nature Publishing Group.
- Yikuan Li, Mohammad Mamouei, Gholamreza Salimi-Khorshidi, Shishir Rao, Abdelaali Hassaine, Dexter Canoy, Thomas Lukasiewicz, and Kazem Rahimi. Hi-BEHRT: Hierarchical Transformer-Based Model for Accurate Prediction of Clinical Events Using Multimodal Longitudinal Electronic Health Records. *IEEE journal of biomedical and health informatics*, 27(2):1106–1117, February 2023. ISSN 2168-2194. doi: 10.1109/JBHI.2022.3224727. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7615082/>.
- Yongxia Lu, Wei Wang, Jingyu Liu, Min Xie, Qiang Liu, and Sufang Li. Vascular complications of diabetes: A narrative review. *Medicine*, 102(40):e35285, October 2023. ISSN 0025-7974. doi: 10.1097/MD.00000000000035285. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10553000/>.
- Anna Munoz-Farre, Harry Rose, and Sera Aylin Cakiroglu. sEHR-CE: Language modelling of structured EHR data for efficient and generalizable patient cohort expansion, November 2022. URL <http://arxiv.org/abs/2211.17121>. arXiv:2211.17121 [cs, stat].
- Achim Wolf, Daniel Dedman, Jennifer Campbell, Helen Booth, Darren Lunn, Jennifer Chapman, and Puja Myles. Data resource profile: Clinical Practice Research Datalink (CPRD) Aurum. *International Journal of Epidemiology*, 48(6):1740–1740g, December 2019. ISSN 0300-5771, 1464-3685. doi: 10.1093/ije/dyz034. URL <https://academic.oup.com/ije/article/48/6/1740/5374844>.
- Michael Wornow, Yizhe Xu, Rahul Thapa, Birju Patel, Ethan Steinberg, Scott Fleming, Michael A. Pfeffer, Jason Fries, and Nigam H. Shah. The shaky foundations of large language models and foundation models for electronic health records. *npj Digital Medicine*, 6(1):1–10, July 2023. ISSN 2398-6352. doi: 10.1038/s41746-023-00879-8. URL <https://www.nature.com/articles/s41746-023-00879-8>. Publisher: Nature Publishing Group.
- Xi Yang, Aokun Chen, Nima PourNejatian, Hoo Chang Shin, Kaleb E. Smith, Christopher Parisien, Colin Compas, Cheryl Martin, Anthony B. Costa, Mona G. Flores, Ying Zhang, Tanja Magoc, Christopher A. Harle, Gloria Lipori, Duane A. Mitchell, William R. Hogan, Elizabeth A. Shenkman, Jiang Bian, and Yonghui Wu. A large language model for electronic health records. *npj Digital Medicine*, 5(1):1–9, December 2022. ISSN 2398-6352. doi: 10.1038/s41746-022-00742-2. URL <https://www.nature.com/articles/s41746-022-00742-2>. Publisher: Nature Publishing Group.

## Appendix A.

### A.1. Pre-processing

Patients were only included if they were eligible for data linkage to Hospital Episode Statistics (HES) and Office for National Statistics (ONS) registries. This ensured that only patients with primary and secondary care were included. Data was pre-processed to remove duplicate events (identical rows), impossible events (dates of events that occur before birth or after deregistration), events with missing dates, or missing clinical code (events without a textual descriptor). Due to data quality issues only records between 1985 and 2020 were included (Wolf et al., 2019).

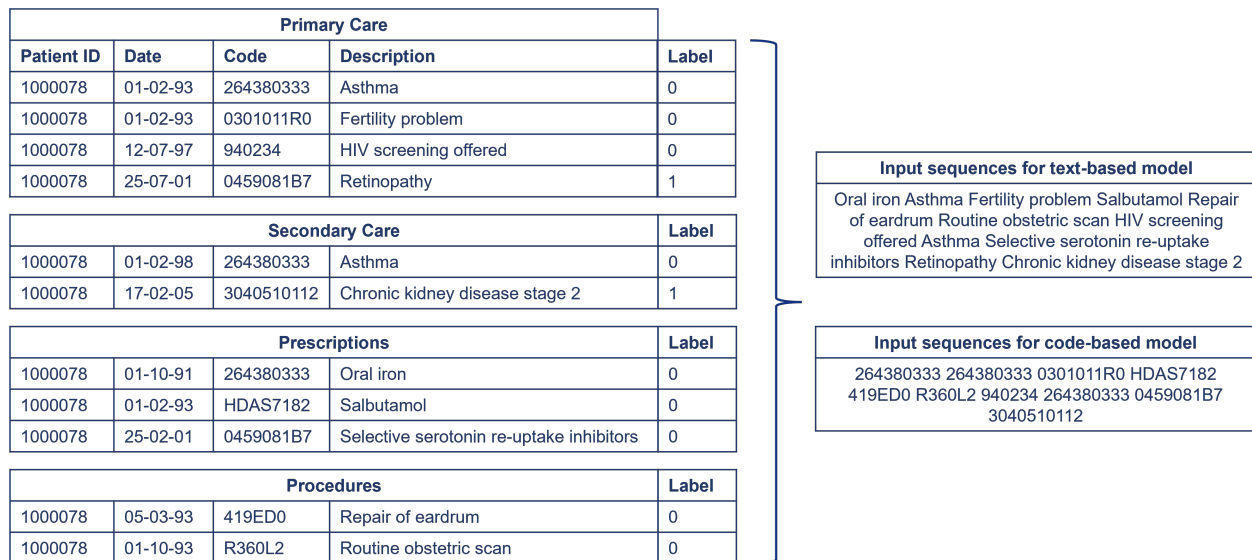


Figure 1: Input Format for Text and Code-based Approaches

For the code-based model we kept the clinical codes from each of the registries, whilst for the code-agnostic models we kept the textual descriptions (Figure 1). For data in primary care including diagnoses, symptoms, demographics etc, this follows Systematized Nomenclature of Medicine Clinical Terms (SNOMED CT), prescriptions within primary care follow the British National Formulary (BNF), within secondary care, diagnoses utilise the International Classification of Diseases, Tenth Revision (ICD10) and procedures use the OPCS Classification of Interventions and Procedures (OPCS 4).

### A.2. Model architecture

Gatortron-base is a smaller version of the original with 345M parameters. It was trained on scratch on 82B words of de-identified clinical notes, 6.1B words from PubMed, 2.5B words from WikiText and 0.5B words of de-identified clinical notes from MIMIC-III.

For all models, input was first tokenized and special token [CLS] added. The tokenized sequences, special tokens and positional embeddings were fed into the pretrained encoder-only model. The final hidden state of the [CLS] token was used as input to the fully connected layer. A sigmoid activation function was applied to logits to produce independent probabilities for each label (Figure 2).

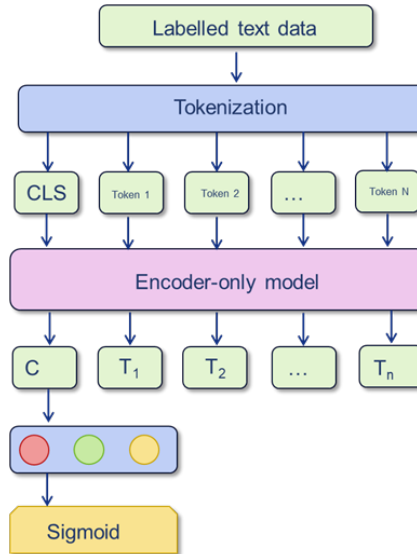


Figure 2: Multi-label Approach and Model Architecture

For each model we searched for a learning rate that gave the lowest F1 score (1e-3, 2e-5, 3e-5, 4e-5, 5e-5) and fine-tuned on the entire dataset for 48000 steps with early stopping. Losses were monitored for overfitting. Models were fine-tuned on an NVidia A100 GPU. This research utilised Queen Mary’s Apocrita HPC facility, supported by QMUL Research-IT (King et al., 2017).

### A.3. Comparison to other pretrained models

To assess the potential benefits gained from the existing knowledge encoded in the pre-trained clinical model, GatorTron-base, we also compare to an out of domain pre-trained model, BERT-base (Devlin et al., 2019) trained on Wikipedia and Google Books and additionally to Biomedical-longformer-base (Beltagy et al., 2020), a model trained on abstracts from PubMed and PubMed Central articles. The Biomedical-longformer is based on the Longformer architecture, which uses an attention mechanism that scales linearly enabling a max token length of 4096. These sequences are truncated from right, as the default.

We can see from Table 5 that all models perform better on text-based approaches, compared to code-based approaches. BERT performs similarly but marginally worse than GatorTron, indicating that large pre-trained models, even when not trained directly on clinical data still contain valuable knowledge.

Biomedical-longformer was significantly better when applying a text-based approach, over a code-based approach and outperformed all other models across the prediction tasks. This suggesting that there is additional information to be gained from capturing long term dependencies in the data. However, the improvement seen on Biomedical-longformer is at the expense of additional resources and time, taking on average  $\sim 20$  hours compared to GatorTron at  $\sim 3$  hours.

### A.4. Context length

The median number of tokens in each individual’s EHR is 2272 tokens, greater than both the maximum length of GatorTron at 512 tokens. We can see from Figure 3 that many EHRs are still truncated when using Biomedical-longformer as they fall over 4096 tokens.

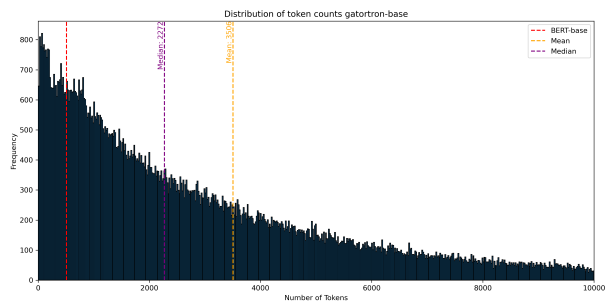


Figure 3: Distribution of Token Counts



Table 5: Performance Comparison of Text- and Code-based Models across Different Pre-trained Models and Prediction Windows

		Text-based		Code-based	
		Micro-F1	Micro-AUPRC	Micro-F1	Micro-AUPRC
<b>BERT-base</b>	<b>1 year</b>	0.44 (0.43, 0.45)	0.42 (0.40, 0.43)	0.42 (0.41, 0.43)	0.39 (0.37, 0.40)
	<b>5 year</b>	<b>0.47 (0.46, 0.48)</b>	<b>0.46 (0.44, 0.47)</b>	0.43 (0.43, 0.45)	0.42 (0.41, 0.43)
	<b>10 year</b>	<b>0.48 (0.47, 0.49)</b>	<b>0.49 (0.47, 0.50)</b>	0.46 (0.45, 0.47)	0.46 (0.44, 0.47)
<b>Biomedical-longformer-base</b>	<b>1 year</b>	<b>0.56 (0.55, 0.57)</b>	<b>0.57 (0.56, 0.59)</b>	0.53 (0.52, 0.54)	0.54 (0.52, 0.55)
	<b>5 year</b>	<b>0.60 (0.60, 0.61)</b>	<b>0.63 (0.62, 0.65)</b>	0.53 (0.52, 0.54)	0.55 (0.54, 0.57)
	<b>10 year</b>	<b>0.60 (0.59, 0.61)</b>	<b>0.64 (0.63, 0.66)</b>	0.56 (0.55, 0.57)	0.60 (0.58, 0.61)
<b>Gatortron-base</b>	<b>1 year</b>	0.45 (0.44, 0.46)	0.44 (0.43, 0.46)	0.43 (0.42, 0.44)	0.40 (0.39, 0.41)
	<b>5 year</b>	<b>0.50 (0.49, 0.51)</b>	<b>0.51 (0.50, 0.52)</b>	0.43 (0.42, 0.44)	0.43 (0.41, 0.44)
	<b>10 year</b>	0.49 (0.48, 0.50)	<b>0.50 (0.49, 0.51)</b>	0.47 (0.46, 0.49)	0.47 (0.45, 0.48)

Note: Values in parentheses represent 95% confidence intervals, bold indicates statistical significance