

Supplemental materials for: “Predicting postoperative risks using large language models”

Note: From herein, when we refer to “BJH dataset”, we are referencing our dataset of main cohort consisting of records from Barnes-Jewish Hospital (BJH); when we refer to “MIMIC-III”, we are referencing the replication on the MIMIC-III dataset.

Appendix A1 Data characteristics

The BJH dataset included 84,875 patient records, with a vocabulary size of 3203, as well as a mean word and vocabulary lengths of 8.9 (sd: 6.9) and 7.3 (sd: 4.4), respectively. The BJH dataset contained procedural descriptions obtained from anesthetic records (EPIC), and was derived pre-operatively from smart text records. The non-textual cohort characteristics are listed in table 1.

MIMIC-III encompasses patient records from the critical care units of the Beth Israel Deaconess Medical Center between 2001 and 2012. To accurately replicate the methodology of the BJH dataset, which uses procedural codes as textual inputs, we adopted the same approach with MIMIC-III. Thus, we utilized descriptive texts derived from the long-form titles of ICD-10 procedure codes, which were traced and utilized, in our replication efforts involving MIMIC-III (more details provided in section [Appendix A3](#)). We did this for two reasons. First, text trace-able from ICD-10 codes containing procedures related to the patient, which is considered to be more consistent towards the textual inputs found in the BJH dataset. Second, many clinically-relevant pre-trained models like bioClinicalBERT and ClinicalBERT were already pre-trained on MIMIC’s other notes (eg discharge notes). Resultingly, MIMIC-III had 52,234 patient records, with a vocabulary size of 1871, as well as a mean word and vocabulary lengths of 23.7 (sd: 19.0) and 20.3 (sd: 14.3), respectively.

Both datasets are single sentenced texts across all notes. Refer to sections [Appendix A3](#) and [Appendix A2](#) and for more details as to how the data was extracted and cleaned.

Appendix A1.1 Cohort characteristics

Characteristics	Dataset	
	BJH (our dataset)	MIMIC-III (replicated dataset)
Patient type	Orthopedic 14412 (17%) Ophthalmology: 7442 (8.8%) Urology: 6236 (7.4%)	Emergency: 36043 (69%) Elective: 7675 (15%) Urgent: 1261 (2.4%)
Gender	Male: 42722 (50.3%)	Male: 29670 (56%)
Ethnicity	White: 62563 (74%) African American: 19239 (22.6%), Hispanic: 1488 (1.7%) Asian: 1015 (1.2%)	White: 36213 (69%) African American: 4586 (9%) Hispanic/Latino: 1483 (2.7%) Asian: 1394 (2.6%)
Weight	86kg (24.7kg)	65 kg (39.6kg)
Height	170 cm (11cm)	151cm (46cm)
Liver disease	Yes: 6697 (7.9%)	Yes: 3960 (7.1%)
Cancer	Yes: 29213 (34%)	Yes: 311 (0.5%)
Congestive Heart Failure	Yes: 8886 (10%)	Yes: 2174 (4.2%)
Myocardial Infarction	Yes: 8587 (10%)	Yes: 2894 (5.5%)
Chronic Pulmonary Disease	Yes: 9837 (12%)	Yes: 5743 (11%)
HIV/AIDS	Yes: 4069 (4.8%)	Yes: 468 (0.9%)

Table 1 A comparison of some common characteristics between the BJH dataset and the MIMIC-III dataset amongst the cohorts with relevant clinical texts. Categorical variables are reported as number of patients (percentage), numerical variables are reported as mean (standard deviation). Note that the summary statistics were computed after removing records with no texts associated with the patient.

Appendix A2 Process of data collection

Appendix A2.1 Data extraction of clinical notes

For the BJH dataset, the notes are derived from smart text records, consisting of texts from a comprehensive list of checkboxes related to procedural information that clinicians select during patient consultations. The texts from these selected checkboxes are compiled as part of their anaesthetic records. These anaesthetic records were then pulled by BJH for our study.

For the MIMIC-III dataset, ICD-10 codes containing procedural information from each patient were traced and formatted into their respective long-form titles. Among each patient, these long-form titles were then combined together to form a single-sentenced clinical note, thereby aligning with the textual characteristics of BJH’s clinical notes.

Appendix A2.2 Data extraction of outcome variables

For the BJH dataset, AKI was determined using a combination of laboratory values (serum creatinine) and dialysis event records, and structured anesthesia assessments, laboratory data, and billing data indicating baseline end-stage renal disease were used as exclusion criteria for AKI. Acute kidney injury was defined according to the Kidney Disease Improving Global Outcomes criteria. Delirium was determined from nurse flow-sheets (positive Confusion Assessment Method for the Intensive Care unit test result); pneumonia, DVT, and PE were determined based on the International Statistical Classification of Diseases and Related Health Problems, Tenth Revision (ICD-10) diagnosis codes. Patients without delirium screenings were excluded from the analysis of that complication.

For MIMIC-III, length-of-stay (LOS) and in-hospital mortality was provided in MIMIC-III. 30-day mortality was calculated based on a person’s mortality status 30-days after they were discharged. The 12 hour discharge status was determined based on time between a patient’s admittance to their discharge time. Patients who suffered from in-hospital mortality within 12-hours of admittance were NOT considered as a positive case in our measure of 12-hour discharge status. Such scenarios included cases whereby the discharge timestamp was equivalent to the mortality timestamp AND both the discharge and mortality timestamp was within 12-hours of the patients admitted time-stamp.

Appendix A3 Data cleaning and pre-processing

For the BJH dataset, any textual information that is unique to the patient and could thus be traced back to the patient was removed by BJH before being handed over to the authors for analysis. This included formatting the text to include only common tokens of scheduled procedures in an arbitrary order, removal any unique non-periodic punctuation, and upper-casing all alphabets. There were a total of 5,129 patients with no clinical notes associated with their records (ie an empty string), of which they were removed from our sample.

For the MIMIC-III replication, the data of each patient’s ICD-10 procedural codes was traced and converted to the long-titled version of procedural descriptions. As a result, the data was cleaned by joining the descriptions together and lower-casing all non-first sentenced words within the text for consistency. Records without any relevant outcomes or text trace-able from ICD-10 codes was dropped.

Appendix A4 Model Development and specification

Appendix A4.1 Description of each model’s architecture

BERT consists of a stack of transformer encoder layers, with each layer consisting of a multi-head self-attention mechanism and a position-wise feed-forward neural network. As a bi-directional model, the attentional mechanism of BERT allows the tokens to be assigned weighted importance by considering the context of all the other preceding and subsequent tokens. This differs from the autoregressive nature of GPT, which consists of a stack of transformer decoder layers and assigns attention solely based on the preceding tokens. The output of the self-attention layer is then passed through a position-wise feed-forward network, which consists of two linear layers with a ReLU activation function. The output of the final layer can then be used for downstream tasks, such as text classification.

We leverage the publicly available pre-trained clinical LLMs in this study, namely ClinicalBERT and BioClinicalBERT, which are BERT-based models, and BioGPT, which is a GPT-based model.

BioGPT was adopted using the GPT-2 architecture, in which it was trained on 15M PubMed abstracts with 347M parameters. This means that it was trained on the language modeling task predicts the next word given all its preceding words. Hence, during training, the GPT-based model aims to assign higher probabilities to the actual words that appear in each position across sentences, compared to other words that do not.

For ClinicalBERT, it was initialized from the BERT_{base} architecture and was pre-trained on a large clinical corpus using the Medical Information Mart for Intensive Care III (MIMIC-III) dataset, containing 2,083,180 de-identified clinical notes associated with admissions. Being trained on the BERT_{base} architecture this entails two training objectives, the masked language modeling (MLM) objective and the Next Sentence Prediction (NSP) objective. The MLM attempts to predict the masked tokens using the entire unmasked tokens in the text, allowing BERT-based models to learn a bidirectional representation of sentences. The NSP task involves taking two masked sentences as input and then predicting whether the second sentence follows the first in the original document. However, as we are only constrained to single-sentenced documents, effectively only the MLM task is at play during fine-tuning. Similarly, BioClinicalBERT was pre-trained on all the available clinical notes associated with MIMIC-III dataset. However, unlike ClinicalBERT, BioClinicalBERT was based on the BioBERT model instead of the BERT_{base} model. BioBERT itself was a clinical variant of the BERT_{base} trained using 4.5 billion words from PubMed abstracts and 13.5 billion words from PubMed Central full-text articles. This allowed BioClinicalBERT to leverage texts from both the biomedical and clinical domains.

Appendix A4.2 Details of the architecture behind semi-supervised finetuning strategy

In BERT-based models, the auxiliary predictors take the output after the final normalized residual layer, sometimes known as the final layer of the hidden states, and predict the logits of the outcome; in GPT-based models, we followed the same strategy and fed the output after the final normalized residual layer to the auxiliary predictors. The auxiliary neural network uses the Binary-Cross-Entropy (BCE), Cross-Entropy (CE), and Mean-Square-Error (MSE) losses for binary classification, multi-label classification, and regression tasks, respectively.

Appendix A4.3 Description of auxiliary network in the foundational finetuning strategy

Each label is assigned a task-specific auxiliary network wherein the losses across all labels are pooled together. Similar to the semi-supervised finetuning strategy, a λ hyperparameter is introduced to control for the losses in the auxiliary network and the models self-supervised objectives. Where there are both categorical and continuous labels, as witnessed in our MIMIC-III replication foundational model, an additional hyperparameter, ie λ_1 and λ_2 is used to account and control for the expected massive differences between the MSE loss and the BCE or CE loss.

Appendix A4.4 Finetuning parameters

self-supervised finetuning type	Model	Parameter	Outcome					
			30 day mortality	DVT	PE	Pneumona	Aki	Delirium
self-supervised finetuning	BioClinicalBERT	Number of Train Epochs			4			
		Train Batch Size			8			
		Validation Batch Size			16			
		Warmup steps			2500			
		Weight Decay			0			
		Learning Rate			0.0001			
	BioGPT	Number of Train Epochs			2			
		Train Batch Size			8			
		Validation Batch Size			16			
		Warmup steps			500			
		Weight Decay			0			
		Learning Rate			0.001			
ClinicalBERT	Number of Train Epochs			5				
	Train Batch Size			8				
	Validation Batch Size			8				
	Warmup steps			1000				
	Weight Decay			0.01				
	Learning Rate			0.0001				
Semi-supervised self-supervised finetuning	BioClinicalBERT	Number of Train Epochs	5	5	8	8	5	5
		Train Batch Size	36	36	36	36	36	36
		Validation Batch Size	36	36	36	36	36	36
		Warmup steps	1500	1500	500	1500	1500	1500
		Weight Decay	0.001	0.001	0.00001	0.00001	0.001	0.001
		Learning Rate	0.00001	0.00001	0.00001	0.00001	0.00001	0.00001
		λ	2	2	1	1	2	2
	bioGPT	Number of Train Epochs	3	3	3	3	3	3
		Train Batch Size	32	32	32	32	32	32
		Validation Batch Size	32	32	32	32	32	32
		Warmup steps	1000	1000	1000	1000	1000	1000
		Weight Decay	0.001	0.001	0.001	0.001	0.001	0.001
		Learning Rate	0.000005	0.000005	0.000005	0.000005	0.000005	0.000005
		λ	10	10	10	10	10	10
	ClinicalBERT	Number of Train Epochs	6	7	6	7	7	7
		Train Batch Size	36	40	36	40	32	32
		Validation Batch Size	36	32	36	32	32	32
		Warmup steps	1500	1000	1500	1000	1500	1500
		Weight Decay	0.001	0.1	0.001	0.001	0.1	0.1
		Learning Rate	0.00001	0.00001	0.00001	0.00001	0.0001	0.001
		λ	10	1	10	5	1	1
Foundational	BioClinicalBERT	Number of Train Epochs			6			
		Train Batch Size			48			
		Validation Batch Size			48			
		Warmup steps			1500			
		Weight Decay			0.001			
		Learning Rate			0.00001			
	λ			2				
	bioGPT	Number of Train Epochs			3			
		Train Batch Size			48			
		Validation Batch Size			48			
		Warmup steps			1000			
		Weight Decay			0.001			
Learning Rate				0.000005				
λ			10					
ClinicalBERT	Number of Train Epochs			6				
	Train Batch Size			48				
	Validation Batch Size			48				
	Warmup steps			1500				
	Weight Decay			0.001				
	Learning Rate			0.00001				
λ			10					

Table 2 Details of parameters selected when fine-tuning each large language model on the BJH dataset, including the λ parameter used to control the magnitude between the unsupervised and supervised losses from the semi-supervised and foundational models. It is worth noting that the λ parameter can vary for the semi-supervised model based on the labeled outcome. In addition, the learning rates, batch sizes are higher and the number of epochs is much larger in the semi-supervised and foundational model to ensure that within each batch, the labeled losses are able to be sufficiently exposed, whilst being allowed to converge with relativity to the objective loss functions.

self-supervised finetuning type	Model	Parameter	Outcome				
			In hospital mortality	Length of Stay	Discharge in 12 hours	Death in 30 days	
self-supervised finetuning	BioClinicalBERT	Number of Train Epochs			5		
		Train Batch Size			16		
		Validation Batch Size			16		
		Warmup steps			1500		
		Weight Decay			0.001		
		Learning Rate			0.001		
	BioGPT	Number of Train Epochs			2		
		Train Batch Size			8		
		Validation Batch Size			64		
		Warmup steps			500		
		Weight Decay			0.01		
		Learning Rate			0.001		
ClinicalBERT	Number of Train Epochs			5			
	Train Batch Size			24			
	Validation Batch Size			24			
	Warmup steps			1500			
	Weight Decay			0			
	Learning Rate			0.001			
self-supervised finetuning	BioClinicalBERT	Number of Train Epochs	7	7	7	7	
		Train Batch Size	40	40	40	40	
		Validation Batch Size	40	40	40	40	
		Warmup steps	1500	1500	1500	1500	
		Weight Decay	0.01	0.01	0.01	0.01	
		Learning Rate	0.00001	0.00001	0.00001	0.00001	
		λ	3	0.005	3	3	
		BioGPT	Number of Train Epochs	5	5	5	5
			Train Batch Size	36	36	36	36
			Validation Batch Size	36	36	36	36
	Warmup steps		500	500	500	500	
	Weight Decay		0.01	0.01	0.01	0.01	
	Learning Rate		0.000001	0.000001	0.000001	0.000001	
	λ	10	0.02	15	10		
	ClinicalBERT	Number of Train Epochs	7	7	7	7	
		Train Batch Size	32	32	32	32	
		Validation Batch Size	32	32	32	32	
		Warmup steps	2000	2000	2000	2000	
		Weight Decay	0	0	0	0	
		Learning Rate	0.00001	0.00001	0.00001	0.00001	
		λ	10	0.02	12	10	
		Foundational	BioClinicalBERT	Number of Train Epochs			6
	Train Batch Size					48	
	Validation Batch Size					48	
Warmup steps					1500		
Weight Decay					0.001		
Learning Rate					0.000001		
λ				2 (0.02 for MSE loss)			
BioGPT	Number of Train Epochs				3		
	Train Batch Size				40		
	Validation Batch Size				40		
	Warmup steps				500		
	Weight Decay				0.001		
	Learning Rate				0.000001		
λ				10 (0.01 for MSE loss)			
ClinicalBERT	Number of Train Epochs				8		
	Train Batch Size				48		
	Validation Batch Size				48		
	Warmup steps				2000		
	Weight Decay			0			
	Learning Rate			0.00001			
λ			10 (0.1 for MSE loss)				

Table 3 Details of parameters selected when fine-tuning each large language model on the MIMIC-III replication.

Appendix A4.5 Predictor parameters

Classifier	Parameter Name	Parameters
XGBoost	Learning rate	0.1,0.15,0.3
	maximum depth	4,5,6,7,8
	minimum child weight	1,2,4
Logistic Regression	C	0.01, 1, 10
	Penalty	l1, l2
	Solver	lbfgs, newton-cholesky
Random Forest	maximum depth	4, None
	Minimum samples per leaf	1, 3

Table 4 Details of cross-validated hyperparameters that were experimented when using the XGBoost, Logistic Regression, and Random Forest model. The entire dataset was split using a 5-fold train-test split, meaning 80% of the data was assigned to the training group and 20% of the data was assigned to the unseen test group. Within this 80% of training data, the data was further cross-validated using a 5-fold cross-validation, where the validation data was used to tune and select the best parameters. This approach is referred to as the nested cross-validation approach.

Appendix A4.6 Fairness

To ensure that the large language model is fine-tuned in a fair manner, pre-trained models were finetuned such that the batches were composed of examples selected randomly and inserted into the batched finetuning process, thereby ensuring it is not systematically biased with respect to any specific group. In addition, stratified k-fold validation was used, ensuring that the model’s performance is reliably evaluated across diverse subsets, maintaining representation from each category in every fold.

Appendix A5 Details of performance metrics for each model

Appendix A5.1 Additional results from the BJH dataset (our dataset)

Model	DVT		AKI		Delirium	
	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
cbow	Mean: 0.524	Mean: 0.006	Mean: 0.56	Mean: 0.156	Mean: 0.501	Mean: 0.474
	CI: (0.457, 0.59)	CI: (0.005, 0.007)	CI: (0.488, 0.632)	CI: (0.125, 0.187)	CI: (0.464, 0.539)	CI: (0.441, 0.507)
Doc2vec	Mean: 0.531	Mean: 0.007	Mean: 0.523	Mean: 0.146	Mean: 0.484	Mean: 0.466
	CI: (0.443, 0.619)	CI: (0.004, 0.01)	CI: (0.421, 0.624)	CI: (0.092, 0.199)	CI: (0.439, 0.53)	CI: (0.417, 0.514)
fastText	Mean: 0.694	Mean: 0.014	Mean: 0.726	Mean: 0.273	Mean: 0.565	Mean: 0.533
	CI: (0.642, 0.746)	CI: (0.011, 0.017)	CI: (0.702, 0.75)	CI: (0.239, 0.307)	CI: (0.541, 0.589)	CI: (0.513, 0.554)
GloVe	Mean: 0.723	Mean: 0.019	Mean: 0.81	Mean: 0.441	Mean: 0.666	Mean: 0.636
	CI: (0.7, 0.745)	CI: (0.013, 0.024)	CI: (0.805, 0.815)	CI: (0.43, 0.451)	CI: (0.652, 0.681)	CI: (0.613, 0.66)
bioClinicalBERT	Mean: 0.76	Mean: 0.02	Mean: 0.83	Mean: 0.469	Mean: 0.68	Mean: 0.653
	CI: (0.723, 0.796)	CI: (0.014, 0.027)	CI: (0.828, 0.831)	CI: (0.457, 0.48)	CI: (0.663, 0.697)	CI: (0.626, 0.68)
bioGPT	Mean: 0.773	Mean: 0.024	Mean: 0.835	Mean: 0.478	Mean: 0.691	Mean: 0.664
	CI: (0.734, 0.813)	CI: (0.016, 0.032)	CI: (0.833, 0.838)	CI: (0.465, 0.492)	CI: (0.672, 0.71)	CI: (0.638, 0.69)
ClinicalBERT	Mean: 0.764	Mean: 0.022	Mean: 0.83	Mean: 0.469	Mean: 0.686	Mean: 0.66
	CI: (0.73, 0.799)	CI: (0.015, 0.03)	CI: (0.827, 0.833)	CI: (0.458, 0.48)	CI: (0.671, 0.702)	CI: (0.634, 0.686)

Table 5 A comparison of baseline models vs pretrained models amongst the intermediate outcomes in the BJH dataset. The results for the target outcomes could be found in the manuscript. The best baseline models are underlined, and the best models are **bolded**

Fig. 1 A comparison of the results of post-operative intermediate outcomes from the BJC dataset across various models and their respective tuning strategies.

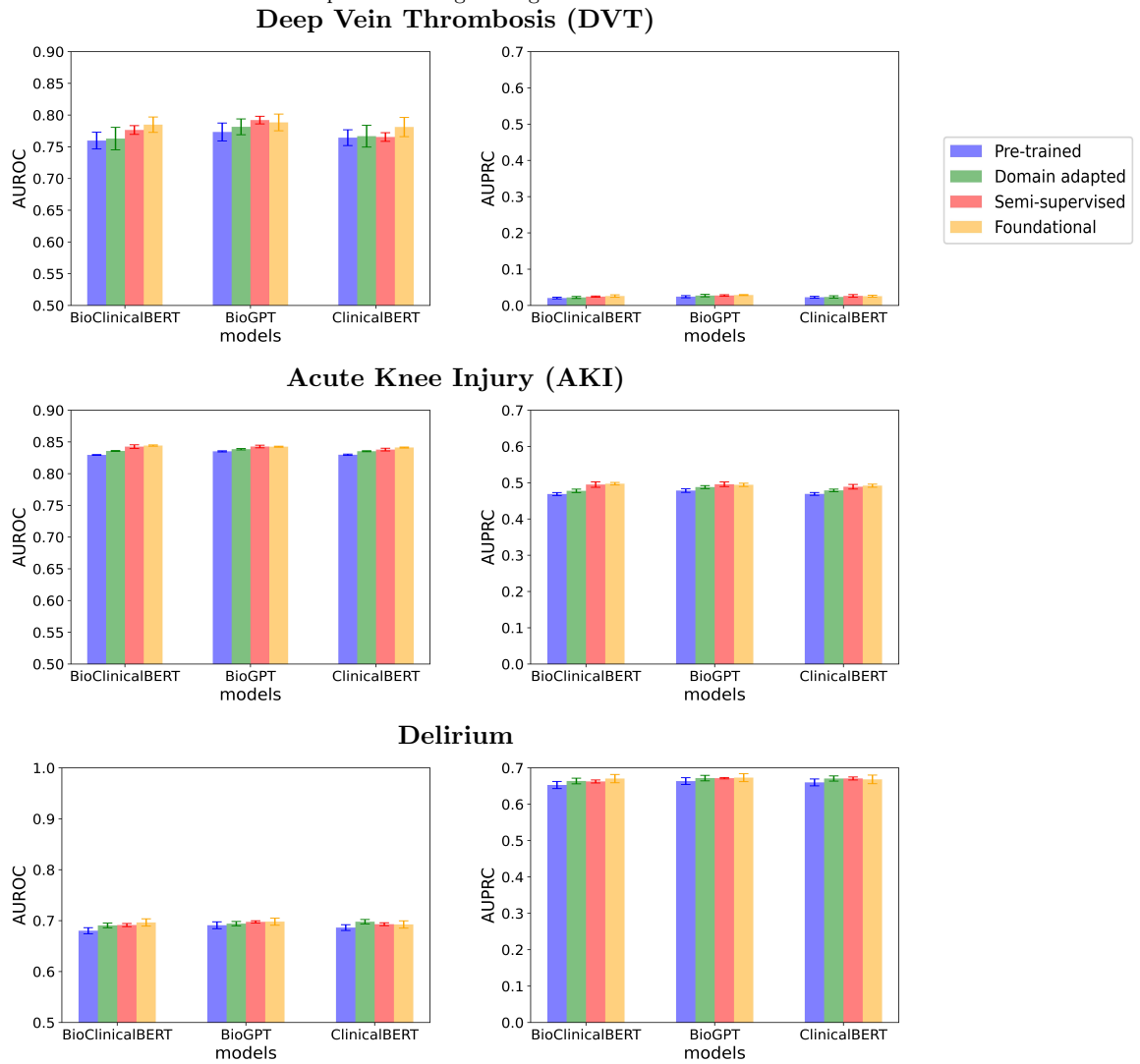
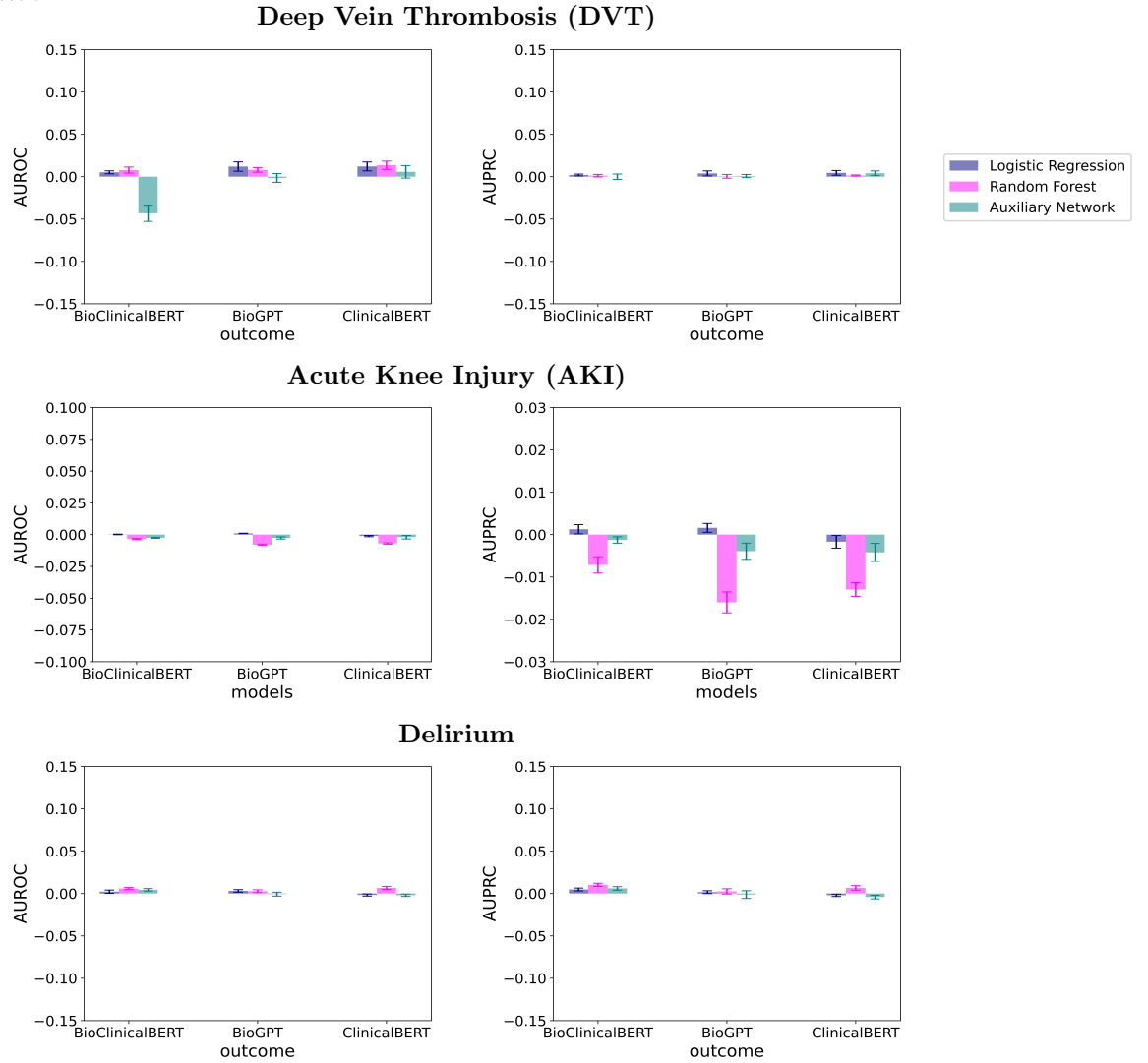


Fig. 2 Comparison of different machine learning classifiers with that of our default XGBoost predictor applied to our textual representations ($\Delta_{\text{model}_{i,j}} - \text{XGBoost}_i$ with outcome i and model j), including the use of the trained auxiliary layer directly from our foundational model. These figures represent the intermediate outcomes. The figures for the target outcomes are referenced in the appendix section.



Model type	Model	30 day mortality		DVT		PE		Pneumonia		AKI		Delirium	
		AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
Baseline	elbow	Mean: 0.528 CI: (0.469, 0.608)	Mean: 0.023 CI: (0.015, 0.031)	Mean: 0.524 CI: (0.457, 0.59)	Mean: 0.006 CI: (0.002, 0.006)	Mean: 0.506 CI: (0.418, 0.593)	Mean: 0.004 CI: (0.002, 0.006)	Mean: 0.526 CI: (0.384, 0.668)	Mean: 0.009 CI: (0.001, 0.016)	Mean: 0.156 CI: (0.125, 0.187)	Mean: 0.501 CI: (0.464, 0.539)	Mean: 0.684 CI: (0.641, 0.727)	
	Doc2Vec	Mean: 0.525 CI: (0.448, 0.611)	Mean: 0.035 CI: (0.012, 0.03)	Mean: 0.494 CI: (0.443, 0.549)	Mean: 0.014 CI: (0.004, 0.01)	Mean: 0.652 CI: (0.496, 0.807)	Mean: 0.007 CI: (0.004, 0.01)	Mean: 0.696 CI: (0.637, 0.755)	Mean: 0.016 CI: (0.003, 0.024)	Mean: 0.276 CI: (0.092, 0.460)	Mean: 0.365 CI: (0.439, 0.53)	Mean: 0.533 CI: (0.417, 0.544)	
	fastText	Mean: 0.521 CI: (0.484, 0.554)	Mean: 0.128 CI: (0.118, 0.139)	Mean: 0.723 CI: (0.7, 0.745)	Mean: 0.015 CI: (0.013, 0.024)	Mean: 0.664 CI: (0.628, 0.701)	Mean: 0.011 CI: (0.007, 0.013)	Mean: 0.705 CI: (0.732, 0.799)	Mean: 0.041 CI: (0.017, 0.063)	Mean: 0.441 CI: (0.43, 0.451)	Mean: 0.666 CI: (0.652, 0.681)	Mean: 0.659 CI: (0.613, 0.66)	
	GloVe	Mean: 0.58 CI: (0.507, 0.653)	Mean: 0.196 CI: (0.148, 0.244)	Mean: 0.796 CI: (0.74, 0.852)	Mean: 0.02 CI: (0.015, 0.024)	Mean: 0.683 CI: (0.67, 0.695)	Mean: 0.008 CI: (0.005, 0.011)	Mean: 0.809 CI: (0.8, 0.817)	Mean: 0.033 CI: (0.025, 0.041)	Mean: 0.469 CI: (0.465, 0.472)	Mean: 0.688 CI: (0.672, 0.71)	Mean: 0.653 CI: (0.638, 0.669)	
Pre-trained	bioGPT	Mean: 0.851 CI: (0.851, 0.872)	Mean: 0.161 CI: (0.144, 0.182)	Mean: 0.773 CI: (0.754, 0.813)	Mean: 0.024 CI: (0.016, 0.032)	Mean: 0.671 CI: (0.679, 0.743)	Mean: 0.011 CI: (0.005, 0.017)	Mean: 0.818 CI: (0.8, 0.837)	Mean: 0.047 CI: (0.037, 0.058)	Mean: 0.478 CI: (0.465, 0.492)	Mean: 0.691 CI: (0.672, 0.71)	Mean: 0.664 CI: (0.638, 0.69)	
	ClinicalBERT	Mean: 0.855 CI: (0.847, 0.872)	Mean: 0.172 CI: (0.141, 0.179)	Mean: 0.764 CI: (0.719, 0.814)	Mean: 0.029 CI: (0.014, 0.022)	Mean: 0.677 CI: (0.667, 0.702)	Mean: 0.011 CI: (0.008, 0.015)	Mean: 0.829 CI: (0.782, 0.855)	Mean: 0.044 CI: (0.039, 0.045)	Mean: 0.488 CI: (0.477, 0.499)	Mean: 0.689 CI: (0.682, 0.706)	Mean: 0.686 CI: (0.651, 0.693)	
self-supervised finetuning	bioClinicalBERT	Mean: 0.861 CI: (0.849, 0.873)	Mean: 0.183 CI: (0.142, 0.185)	Mean: 0.763 CI: (0.714, 0.812)	Mean: 0.022 CI: (0.014, 0.03)	Mean: 0.715 CI: (0.667, 0.762)	Mean: 0.011 CI: (0.008, 0.015)	Mean: 0.829 CI: (0.785, 0.871)	Mean: 0.048 CI: (0.039, 0.055)	Mean: 0.478 CI: (0.464, 0.491)	Mean: 0.691 CI: (0.678, 0.703)	Mean: 0.664 CI: (0.642, 0.685)	
	bioGPT	Mean: 0.879 CI: (0.856, 0.879)	Mean: 0.158 CI: (0.151, 0.12)	Mean: 0.787 CI: (0.747, 0.816)	Mean: 0.023 CI: (0.018, 0.027)	Mean: 0.739 CI: (0.667, 0.79)	Mean: 0.013 CI: (0.006, 0.02)	Mean: 0.845 CI: (0.795, 0.845)	Mean: 0.047 CI: (0.039, 0.055)	Mean: 0.479 CI: (0.477, 0.499)	Mean: 0.698 CI: (0.682, 0.706)	Mean: 0.671 CI: (0.651, 0.693)	
	ClinicalBERT	Mean: 0.86 CI: (0.847, 0.872)	Mean: 0.158 CI: (0.141, 0.179)	Mean: 0.787 CI: (0.719, 0.814)	Mean: 0.023 CI: (0.014, 0.022)	Mean: 0.739 CI: (0.686, 0.773)	Mean: 0.014 CI: (0.011, 0.018)	Mean: 0.845 CI: (0.777, 0.84)	Mean: 0.044 CI: (0.023, 0.064)	Mean: 0.479 CI: (0.469, 0.489)	Mean: 0.698 CI: (0.686, 0.71)	Mean: 0.671 CI: (0.651, 0.691)	
	bioClinicalBERT	Mean: 0.873 CI: (0.863, 0.873)	Mean: 0.154 CI: (0.157, 0.12)	Mean: 0.792 CI: (0.757, 0.798)	Mean: 0.029 CI: (0.02, 0.029)	Mean: 0.747 CI: (0.702, 0.754)	Mean: 0.017 CI: (0.01, 0.018)	Mean: 0.829 CI: (0.805, 0.829)	Mean: 0.048 CI: (0.03, 0.054)	Mean: 0.496 CI: (0.474, 0.516)	Mean: 0.698 CI: (0.683, 0.7)	Mean: 0.672 CI: (0.651, 0.671)	
Self-supervised finetuning	bioGPT	Mean: 0.874 CI: (0.87, 0.877)	Mean: 0.184 CI: (0.151, 0.216)	Mean: 0.792 CI: (0.775, 0.809)	Mean: 0.027 CI: (0.021, 0.033)	Mean: 0.747 CI: (0.69, 0.803)	Mean: 0.018 CI: (0.009, 0.025)	Mean: 0.825 CI: (0.804, 0.84)	Mean: 0.048 CI: (0.032, 0.065)	Mean: 0.496 CI: (0.478, 0.514)	Mean: 0.698 CI: (0.691, 0.704)	Mean: 0.672 CI: (0.667, 0.677)	
	ClinicalBERT	Mean: 0.872 CI: (0.855, 0.872)	Mean: 0.184 CI: (0.144, 0.212)	Mean: 0.784 CI: (0.717, 0.784)	Mean: 0.037 CI: (0.016, 0.037)	Mean: 0.773 CI: (0.683, 0.773)	Mean: 0.015 CI: (0.01, 0.015)	Mean: 0.825 CI: (0.802, 0.825)	Mean: 0.056 CI: (0.029, 0.056)	Mean: 0.497 CI: (0.474, 0.507)	Mean: 0.698 CI: (0.685, 0.701)	Mean: 0.682 CI: (0.659, 0.682)	
Foundational	bioClinicalBERT	Mean: 0.875 CI: (0.867, 0.885)	Mean: 0.184 CI: (0.151, 0.212)	Mean: 0.785 CI: (0.717, 0.784)	Mean: 0.028 CI: (0.016, 0.037)	Mean: 0.747 CI: (0.69, 0.803)	Mean: 0.017 CI: (0.01, 0.017)	Mean: 0.836 CI: (0.809, 0.836)	Mean: 0.052 CI: (0.029, 0.052)	Mean: 0.497 CI: (0.474, 0.507)	Mean: 0.696 CI: (0.683, 0.701)	Mean: 0.674 CI: (0.659, 0.682)	
	bioGPT	Mean: 0.875 CI: (0.867, 0.885)	Mean: 0.209 CI: (0.16, 0.209)	Mean: 0.809 CI: (0.775, 0.829)	Mean: 0.028 CI: (0.025, 0.032)	Mean: 0.765 CI: (0.727, 0.803)	Mean: 0.017 CI: (0.01, 0.024)	Mean: 0.807 CI: (0.801, 0.807)	Mean: 0.052 CI: (0.029, 0.052)	Mean: 0.497 CI: (0.488, 0.507)	Mean: 0.696 CI: (0.676, 0.718)	Mean: 0.673 CI: (0.642, 0.704)	
Foundational	ClinicalBERT	Mean: 0.87 CI: (0.859, 0.882)	Mean: 0.184 CI: (0.142, 0.211)	Mean: 0.781 CI: (0.729, 0.823)	Mean: 0.025 CI: (0.018, 0.031)	Mean: 0.749 CI: (0.688, 0.749)	Mean: 0.013 CI: (0.009, 0.016)	Mean: 0.829 CI: (0.8, 0.829)	Mean: 0.05 CI: (0.03, 0.069)	Mean: 0.492 CI: (0.48, 0.504)	Mean: 0.692 CI: (0.672, 0.712)	Mean: 0.668 CI: (0.631, 0.701)	

Table 6 A compilation of all the results from all outcomes of the B/JH dataset.

Classifier	Model	30 day mortality			DVT			PE			Pneumonia			AKI			Delirium			
		AUROC	AUPRC	AUROC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
XGBoost	bioClinicalBERT	Mean: 0.873 CI: (0.860, 0.887)	Mean: 0.184 CI: (0.154, 0.214)	Mean: 0.785 CI: (0.751, 0.818)	Mean: 0.026 CI: (0.017, 0.035)	Mean: 0.714 CI: (0.712, 0.777)	Mean: 0.017 CI: (0.006, 0.026)	Mean: 0.806 CI: (0.805, 0.807)	Mean: 0.052 CI: (0.029, 0.075)	Mean: 0.497 CI: (0.488, 0.507)	Mean: 0.696 CI: (0.677, 0.716)	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.497 CI: (0.488, 0.507)	Mean: 0.696 CI: (0.677, 0.716)	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.497 CI: (0.488, 0.507)	Mean: 0.696 CI: (0.677, 0.716)	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.497 CI: (0.488, 0.507)	Mean: 0.696 CI: (0.677, 0.716)
	bioGPT	Mean: 0.875 CI: (0.862, 0.888)	Mean: 0.184 CI: (0.154, 0.214)	Mean: 0.791 CI: (0.757, 0.826)	Mean: 0.028 CI: (0.019, 0.037)	Mean: 0.705 CI: (0.703, 0.749)	Mean: 0.017 CI: (0.006, 0.026)	Mean: 0.811 CI: (0.810, 0.812)	Mean: 0.052 CI: (0.029, 0.075)	Mean: 0.495 CI: (0.486, 0.504)	Mean: 0.697 CI: (0.678, 0.716)	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.495 CI: (0.486, 0.504)	Mean: 0.697 CI: (0.678, 0.716)	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.495 CI: (0.486, 0.504)	Mean: 0.697 CI: (0.678, 0.716)	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.495 CI: (0.486, 0.504)	Mean: 0.697 CI: (0.678, 0.716)
	ClinicalBERT	Mean: 0.872 CI: (0.859, 0.885)	Mean: 0.184 CI: (0.152, 0.217)	Mean: 0.789 CI: (0.739, 0.823)	Mean: 0.025 CI: (0.018, 0.033)	Mean: 0.719 CI: (0.708, 0.739)	Mean: 0.013 CI: (0.01, 0.016)	Mean: 0.829 CI: (0.8, 0.857)	Mean: 0.05 CI: (0.03, 0.069)	Mean: 0.492 CI: (0.48, 0.504)	Mean: 0.692 CI: (0.672, 0.712)	Mean: 0.841 CI: (0.839, 0.841)	Mean: 0.492 CI: (0.48, 0.504)	Mean: 0.692 CI: (0.672, 0.712)	Mean: 0.841 CI: (0.839, 0.841)	Mean: 0.492 CI: (0.48, 0.504)	Mean: 0.692 CI: (0.672, 0.712)	Mean: 0.841 CI: (0.839, 0.841)	Mean: 0.492 CI: (0.48, 0.504)	Mean: 0.692 CI: (0.672, 0.712)
Auxiliary Layer	bioClinicalBERT	Mean: 0.851 CI: (0.84, 0.862)	Mean: 0.175 CI: (0.142, 0.208)	Mean: 0.712 CI: (0.682, 0.742)	Mean: 0.026 CI: (0.019, 0.034)	Mean: 0.671 CI: (0.66, 0.713)	Mean: 0.009 CI: (0.001, 0.016)	Mean: 0.801 CI: (0.799, 0.802)	Mean: 0.056 CI: (0.032, 0.079)	Mean: 0.496 CI: (0.487, 0.505)	Mean: 0.697 CI: (0.679, 0.716)	Mean: 0.842 CI: (0.84, 0.844)	Mean: 0.496 CI: (0.487, 0.505)	Mean: 0.697 CI: (0.679, 0.716)	Mean: 0.842 CI: (0.84, 0.844)	Mean: 0.496 CI: (0.487, 0.505)	Mean: 0.697 CI: (0.679, 0.716)	Mean: 0.842 CI: (0.84, 0.844)	Mean: 0.496 CI: (0.487, 0.505)	Mean: 0.697 CI: (0.679, 0.716)
	bioGPT	Mean: 0.869 CI: (0.858, 0.88)	Mean: 0.182 CI: (0.154, 0.215)	Mean: 0.788 CI: (0.756, 0.819)	Mean: 0.03 CI: (0.02, 0.04)	Mean: 0.713 CI: (0.712, 0.773)	Mean: 0.01 CI: (0.007, 0.014)	Mean: 0.822 CI: (0.799, 0.866)	Mean: 0.058 CI: (0.038, 0.085)	Mean: 0.488 CI: (0.483, 0.512)	Mean: 0.704 CI: (0.688, 0.721)	Mean: 0.844 CI: (0.841, 0.848)	Mean: 0.488 CI: (0.483, 0.512)	Mean: 0.704 CI: (0.688, 0.721)	Mean: 0.844 CI: (0.841, 0.848)	Mean: 0.488 CI: (0.483, 0.512)	Mean: 0.704 CI: (0.688, 0.721)	Mean: 0.844 CI: (0.841, 0.848)	Mean: 0.488 CI: (0.483, 0.512)	Mean: 0.704 CI: (0.688, 0.721)
	ClinicalBERT	Mean: 0.868 CI: (0.857, 0.879)	Mean: 0.180 CI: (0.145, 0.215)	Mean: 0.787 CI: (0.737, 0.837)	Mean: 0.029 CI: (0.022, 0.036)	Mean: 0.731 CI: (0.687, 0.742)	Mean: 0.011 CI: (0.009, 0.013)	Mean: 0.806 CI: (0.785, 0.827)	Mean: 0.051 CI: (0.029, 0.077)	Mean: 0.49 CI: (0.48, 0.501)	Mean: 0.697 CI: (0.679, 0.716)	Mean: 0.84 CI: (0.838, 0.842)	Mean: 0.49 CI: (0.48, 0.501)	Mean: 0.697 CI: (0.679, 0.716)	Mean: 0.84 CI: (0.838, 0.842)	Mean: 0.49 CI: (0.48, 0.501)	Mean: 0.697 CI: (0.679, 0.716)	Mean: 0.84 CI: (0.838, 0.842)	Mean: 0.49 CI: (0.48, 0.501)	Mean: 0.697 CI: (0.679, 0.716)
Logistic Regression	bioClinicalBERT	Mean: 0.877 CI: (0.865, 0.889)	Mean: 0.184 CI: (0.161, 0.207)	Mean: 0.8 CI: (0.769, 0.831)	Mean: 0.038 CI: (0.022, 0.054)	Mean: 0.775 CI: (0.727, 0.823)	Mean: 0.016 CI: (0.008, 0.024)	Mean: 0.845 CI: (0.815, 0.874)	Mean: 0.056 CI: (0.034, 0.078)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)	Mean: 0.843 CI: (0.841, 0.846)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)	Mean: 0.843 CI: (0.841, 0.846)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)	Mean: 0.843 CI: (0.841, 0.846)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)
	bioGPT	Mean: 0.877 CI: (0.865, 0.889)	Mean: 0.184 CI: (0.161, 0.207)	Mean: 0.8 CI: (0.769, 0.831)	Mean: 0.038 CI: (0.022, 0.054)	Mean: 0.775 CI: (0.727, 0.823)	Mean: 0.016 CI: (0.008, 0.024)	Mean: 0.845 CI: (0.815, 0.874)	Mean: 0.056 CI: (0.034, 0.078)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)	Mean: 0.843 CI: (0.841, 0.846)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)	Mean: 0.843 CI: (0.841, 0.846)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)	Mean: 0.843 CI: (0.841, 0.846)	Mean: 0.496 CI: (0.481, 0.511)	Mean: 0.701 CI: (0.685, 0.717)
	ClinicalBERT	Mean: 0.882 CI: (0.858, 0.882)	Mean: 0.185 CI: (0.145, 0.221)	Mean: 0.822 CI: (0.785, 0.852)	Mean: 0.027 CI: (0.025, 0.035)	Mean: 0.811 CI: (0.688, 0.811)	Mean: 0.013 CI: (0.008, 0.017)	Mean: 0.838 CI: (0.822, 0.871)	Mean: 0.055 CI: (0.03, 0.071)	Mean: 0.49 CI: (0.48, 0.53)	Mean: 0.702 CI: (0.671, 0.71)	Mean: 0.843 CI: (0.837, 0.843)	Mean: 0.49 CI: (0.48, 0.53)	Mean: 0.702 CI: (0.671, 0.71)	Mean: 0.843 CI: (0.837, 0.843)	Mean: 0.49 CI: (0.48, 0.53)	Mean: 0.702 CI: (0.671, 0.71)	Mean: 0.843 CI: (0.837, 0.843)	Mean: 0.49 CI: (0.48, 0.53)	Mean: 0.702 CI: (0.671, 0.71)
Random Forest	bioClinicalBERT	Mean: 0.87 CI: (0.858, 0.882)	Mean: 0.188 CI: (0.137, 0.222)	Mean: 0.793 CI: (0.766, 0.82)	Mean: 0.027 CI: (0.017, 0.037)	Mean: 0.756 CI: (0.708, 0.804)	Mean: 0.013 CI: (0.004, 0.016)	Mean: 0.838 CI: (0.809, 0.868)	Mean: 0.055 CI: (0.033, 0.076)	Mean: 0.49 CI: (0.485, 0.499)	Mean: 0.702 CI: (0.687, 0.717)	Mean: 0.841 CI: (0.839, 0.843)	Mean: 0.49 CI: (0.485, 0.499)	Mean: 0.702 CI: (0.687, 0.717)	Mean: 0.841 CI: (0.839, 0.843)	Mean: 0.49 CI: (0.485, 0.499)	Mean: 0.702 CI: (0.687, 0.717)	Mean: 0.841 CI: (0.839, 0.843)	Mean: 0.49 CI: (0.485, 0.499)	Mean: 0.702 CI: (0.687, 0.717)
	bioGPT	Mean: 0.868 CI: (0.859, 0.868)	Mean: 0.181 CI: (0.154, 0.2)	Mean: 0.795 CI: (0.757, 0.826)	Mean: 0.027 CI: (0.02, 0.038)	Mean: 0.751 CI: (0.714, 0.804)	Mean: 0.019 CI: (0.006, 0.032)	Mean: 0.832 CI: (0.782, 0.861)	Mean: 0.053 CI: (0.03, 0.08)	Mean: 0.487 CI: (0.47, 0.487)	Mean: 0.699 CI: (0.683, 0.72)	Mean: 0.834 CI: (0.833, 0.836)	Mean: 0.487 CI: (0.47, 0.487)	Mean: 0.699 CI: (0.683, 0.72)	Mean: 0.834 CI: (0.833, 0.836)	Mean: 0.487 CI: (0.47, 0.487)	Mean: 0.699 CI: (0.683, 0.72)	Mean: 0.834 CI: (0.833, 0.836)	Mean: 0.487 CI: (0.47, 0.487)	Mean: 0.699 CI: (0.683, 0.72)
	ClinicalBERT	Mean: 0.879 CI: (0.856, 0.879)	Mean: 0.212 CI: (0.15, 0.212)	Mean: 0.824 CI: (0.795, 0.824)	Mean: 0.034 CI: (0.02, 0.034)	Mean: 0.794 CI: (0.708, 0.794)	Mean: 0.019 CI: (0.01, 0.019)	Mean: 0.856 CI: (0.846, 0.856)	Mean: 0.073 CI: (0.053, 0.073)	Mean: 0.489 CI: (0.48, 0.489)	Mean: 0.699 CI: (0.679, 0.719)	Mean: 0.836 CI: (0.833, 0.836)	Mean: 0.489 CI: (0.48, 0.489)	Mean: 0.699 CI: (0.679, 0.719)	Mean: 0.836 CI: (0.833, 0.836)	Mean: 0.489 CI: (0.48, 0.489)	Mean: 0.699 CI: (0.679, 0.719)	Mean: 0.836 CI: (0.833, 0.836)	Mean: 0.489 CI: (0.48, 0.489)	Mean: 0.699 CI: (0.679, 0.719)

Table 7 A compiled comparison of different machine learning classifiers towards our textual representations for all outcomes, including using the trained auxiliary layer directly from our supervised fine-tuning approaches.

Outcome	Model	Metric						
		AUROC	AUPRC	Accuracy	Precision	Recall	Specificity	F1
30 day mortality	bioClinicalBERT	Mean: 0.873 CI: (0.860, 0.887)	Mean: 0.184 CI: (0.154, 0.214)	Mean: 0.941 CI: (0.940, 0.941)	Mean: 0.160 CI: (0.143, 0.176)	Mean: 0.464 CI: (0.430, 0.498)	Mean: 0.950 CI: (0.950, 0.951)	Mean: 0.238 CI: (0.215, 0.260)
	bioGPT	Mean: 0.875 CI: (0.865, 0.885)	Mean: 0.184 CI: (0.16, 0.208)	Mean: 0.94 CI: (0.94, 0.941)	Mean: 0.158 CI: (0.143, 0.173)	Mean: 0.461 CI: (0.431, 0.491)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.236 CI: (0.215, 0.256)
	ClinicalBERT	Mean: 0.87 CI: (0.859, 0.882)	Mean: 0.184 CI: (0.152, 0.217)	Mean: 0.94 CI: (0.94, 0.941)	Mean: 0.158 CI: (0.144, 0.172)	Mean: 0.459 CI: (0.427, 0.491)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.235 CI: (0.215, 0.255)
DVT	bioClinicalBERT	Mean: 0.785 CI: (0.751, 0.818)	Mean: 0.026 CI: (0.017, 0.035)	Mean: 0.946 CI: (0.946, 0.947)	Mean: 0.032 CI: (0.023, 0.041)	Mean: 0.278 CI: (0.208, 0.348)	Mean: 0.950 CI: (0.950, 0.950)	Mean: 0.057 CI: (0.041, 0.073)
	bioGPT	Mean: 0.791 CI: (0.757, 0.826)	Mean: 0.028 CI: (0.025, 0.032)	Mean: 0.946 CI: (0.946, 0.947)	Mean: 0.033 CI: (0.024, 0.042)	Mean: 0.284 CI: (0.221, 0.346)	Mean: 0.95 CI: (0.95, 0.951)	Mean: 0.059 CI: (0.043, 0.074)
	ClinicalBERT	Mean: 0.781 CI: (0.739, 0.823)	Mean: 0.025 CI: (0.018, 0.033)	Mean: 0.946 CI: (0.946, 0.947)	Mean: 0.032 CI: (0.025, 0.039)	Mean: 0.277 CI: (0.236, 0.319)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.057 CI: (0.045, 0.069)
PE	bioClinicalBERT	Mean: 0.744 CI: (0.712, 0.777)	Mean: 0.017 CI: (0.009, 0.026)	Mean: 0.949 CI: (0.947, 0.951)	Mean: 0.013 CI: (0.009, 0.016)	Mean: 0.185 CI: (0.139, 0.231)	Mean: 0.951 CI: (0.949, 0.953)	Mean: 0.024 CI: (0.018, 0.03)
	bioGPT	Mean: 0.765 CI: (0.727, 0.803)	Mean: 0.017 CI: (0.01, 0.024)	Mean: 0.948 CI: (0.947, 0.948)	Mean: 0.013 CI: (0.01, 0.017)	Mean: 0.199 CI: (0.143, 0.255)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.025 CI: (0.018, 0.032)
	ClinicalBERT	Mean: 0.749 CI: (0.708, 0.79)	Mean: 0.013 CI: (0.01, 0.016)	Mean: 0.948 CI: (0.947, 0.949)	Mean: 0.014 CI: (0.012, 0.016)	Mean: 0.213 CI: (0.179, 0.247)	Mean: 0.95 CI: (0.95, 0.951)	Mean: 0.027 CI: (0.023, 0.031)
Pneumonia	bioClinicalBERT	Mean: 0.836 CI: (0.805, 0.867)	Mean: 0.052 CI: (0.029, 0.075)	Mean: 0.947 CI: (0.946, 0.948)	Mean: 0.044 CI: (0.036, 0.053)	Mean: 0.414 CI: (0.317, 0.512)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.080 CI: (0.064, 0.096)
	bioGPT	Mean: 0.831 CI: (0.807, 0.854)	Mean: 0.052 CI: (0.033, 0.07)	Mean: 0.947 CI: (0.947, 0.948)	Mean: 0.047 CI: (0.039, 0.054)	Mean: 0.434 CI: (0.359, 0.51)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.084 CI: (0.071, 0.097)
	ClinicalBERT	Mean: 0.829 CI: (0.8, 0.857)	Mean: 0.05 CI: (0.03, 0.069)	Mean: 0.947 CI: (0.946, 0.948)	Mean: 0.044 CI: (0.038, 0.049)	Mean: 0.405 CI: (0.348, 0.462)	Mean: 0.95 CI: (0.95, 0.951)	Mean: 0.079 CI: (0.069, 0.089)
AKI	bioClinicalBERT	Mean: 0.844 CI: (0.842, 0.847)	Mean: 0.497 CI: (0.488, 0.507)	Mean: 0.872 CI: (0.869, 0.874)	Mean: 0.533 CI: (0.524, 0.542)	Mean: 0.367 CI: (0.361, 0.372)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.434 CI: (0.429, 0.44)
	bioGPT	Mean: 0.842 CI: (0.841, 0.844)	Mean: 0.495 CI: (0.482, 0.507)	Mean: 0.871 CI: (0.869, 0.874)	Mean: 0.532 CI: (0.522, 0.541)	Mean: 0.365 CI: (0.357, 0.372)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.433 CI: (0.425, 0.44)
	ClinicalBERT	Mean: 0.841 CI: (0.839, 0.844)	Mean: 0.492 CI: (0.48, 0.504)	Mean: 0.871 CI: (0.869, 0.873)	Mean: 0.53 CI: (0.519, 0.541)	Mean: 0.362 CI: (0.352, 0.371)	Mean: 0.95 CI: (0.95, 0.95)	Mean: 0.43 CI: (0.42, 0.439)
Delirium	bioClinicalBERT	Mean: 0.696 CI: (0.677, 0.716)	Mean: 0.671 CI: (0.639, 0.702)	Mean: 0.605 CI: (0.597, 0.614)	Mean: 0.79 CI: (0.762, 0.817)	Mean: 0.212 CI: (0.185, 0.240)	Mean: 0.951 CI: (0.95, 0.951)	Mean: 0.335 CI: (0.298, 0.371)
	bioGPT	Mean: 0.697 CI: (0.676, 0.718)	Mean: 0.673 CI: (0.642, 0.704)	Mean: 0.605 CI: (0.599, 0.611)	Mean: 0.79 CI: (0.764, 0.815)	Mean: 0.211 CI: (0.186, 0.236)	Mean: 0.951 CI: (0.95, 0.952)	Mean: 0.333 CI: (0.299, 0.367)
	ClinicalBERT	Mean: 0.692 CI: (0.672, 0.712)	Mean: 0.668 CI: (0.635, 0.701)	Mean: 0.603 CI: (0.594, 0.612)	Mean: 0.786 CI: (0.757, 0.816)	Mean: 0.207 CI: (0.175, 0.239)	Mean: 0.951 CI: (0.95, 0.952)	Mean: 0.328 CI: (0.285, 0.37)

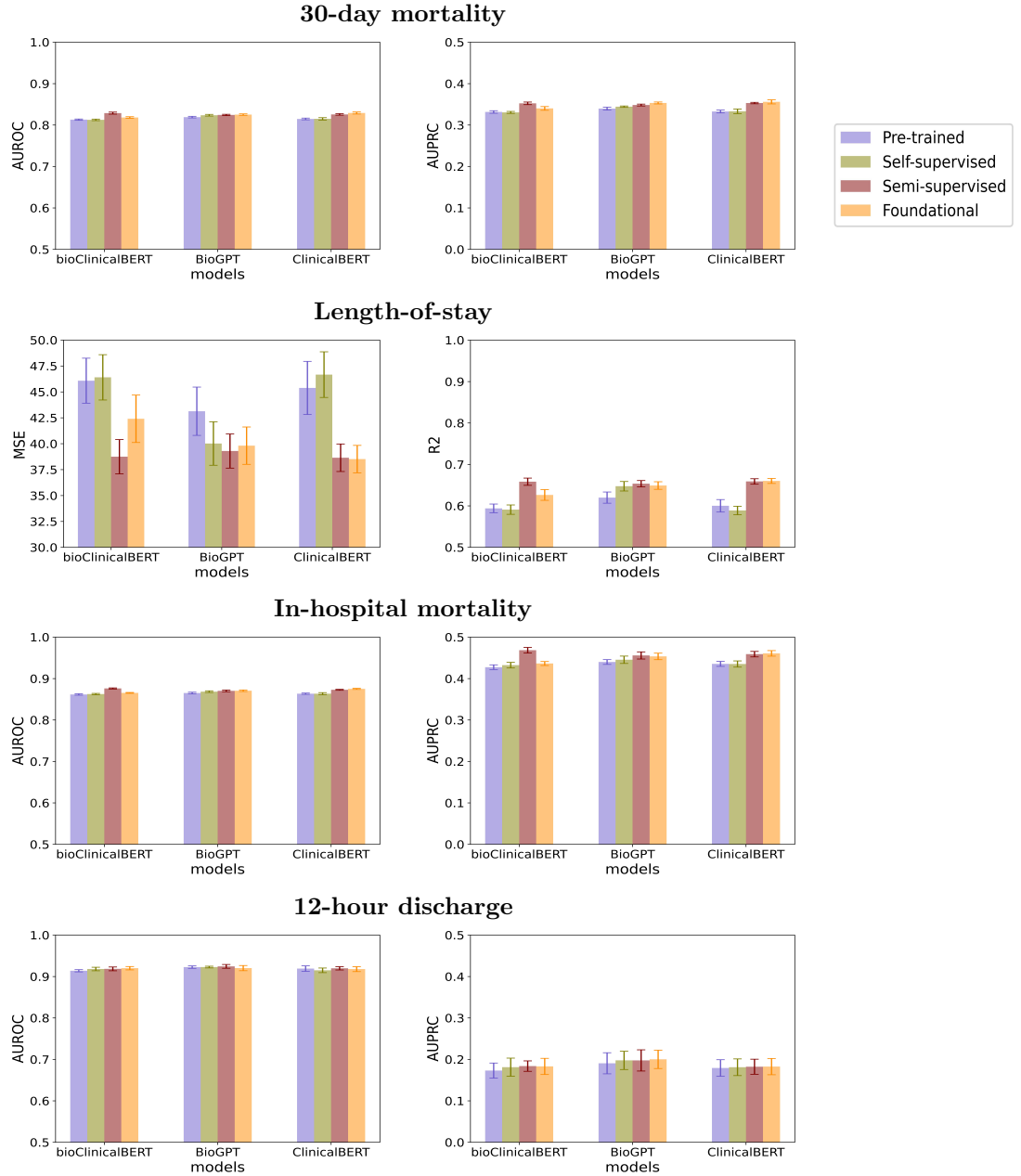
Table 8 Details of Performance Metrics of Each Model for the foundational tuning strategy, which is our best performing strategy. Note that sensitivity, specificity, precision, F-score and accuracy vary depending on the threshold of ML models, we fixed specificity at 95% for easier comparison between different models

Appendix A5.2 MIMIC-III replication

Model	In-hospital mortality		LOS		30-day mortality		12-hour discharge	
	AUROC	AUPRC	MSE	R2	AUROC	AUPRC	AUROC	AUPRC
cbow	0.595 (0.407, 0.782)	0.158 (0.074, 0.242)	128.992 (57.1, 200.885)	-0.149 (-0.802, 0.505)	0.562 (0.441, 0.684)	0.137 (0.098, 0.175)	0.677 (0.614, 0.74)	0.014 (0.01, 0.019)
Doc2vec	0.657 (0.552, 0.762)	0.176 (0.105, 0.247)	162.223 (118.466, 205.98)	-0.446 (-0.882, -0.011)	0.599 (0.493, 0.706)	0.157 (0.089, 0.224)	0.692 (0.662, 0.721)	0.017 (0.01, 0.023)
fastText	0.713 (0.653, 0.772)	0.212 (0.182, 0.241)	109.089 (72.794, 145.384)	0.04 (-0.238, 0.317)	0.677 (0.637, 0.717)	0.191 (0.17, 0.213)	0.733 (0.649, 0.817)	0.023 (0.009, 0.038)
GloVe	<u>0.846</u> (0.841, 0.85)	<u>0.393</u> (0.379, 0.407)	<u>55.114</u> (47.02, 63.209)	<u>0.514</u> (0.472, 0.556)	<u>0.799</u> (0.794, 0.805)	<u>0.304</u> (0.286, 0.323)	<u>0.88</u> (0.851, 0.909)	<u>0.174</u> (0.117, 0.23)
bioClinicalBERT	0.862 (0.857, 0.866)	0.427 (0.411, 0.443)	46.083 (40.041, 52.125)	0.594 (0.565, 0.623)	0.813 (0.808, 0.818)	Mean: 0.332 (0.323, 0.34)	0.914 (0.906, 0.922)	0.173 (0.123, 0.223)
bioGPT	0.865 (0.859, 0.871)	0.44 (0.424, 0.456)	43.134 (36.657, 49.611)	0.62 (0.582, 0.658)	0.819 (0.814, 0.824)	0.34 (0.33, 0.349)	0.923 (0.914, 0.931)	0.191 (0.12, 0.261)
ClinicalBERT	0.863 (0.858, 0.869)	0.434 (0.419, 0.45)	45.383 (38.264, 52.503)	0.601 (0.559, 0.643)	0.815 (0.809, 0.821)	0.333 (0.323, 0.342)	0.919 (0.901, 0.938)	0.179 (0.124, 0.235)

Table 9 A comparison of baseline models (top) vs pre-trained models (bottom) amongst the outcomes from our MIMIC-III replication. The results are presented as the mean and 95% confidence interval across all 5-folds. The best baseline models are underlined, and the best models are **bolded**. As shown amongst the results, the baseline models is consistently outperformed by the pre-trained LLMs. Specifically, we observed absolute increases that ranged from up to 14.6% in 12-hour discharge to 30% for In-hospital mortality for AUROC. Similarly, increases in the AUPRC ranged from 17.7% in 12-hour discharge to 28.2% in in-hospital mortality. For length-of-stay, improvements ranged up to 86 days for MSE and 1.066 in R^2 .

Fig. 3 A replication of our methods on MIMIC-III. A similar magnitude of improvements across tuning strategies were observed. Specifically, self-supervised finetuning witness maximal absolute improvements in AUROCs of up to 0.4% in 12-hour discharge to 3% in in-hospital mortality and AUPRCs of up to 0.4% in 30-day mortality to 0.8% in 12-hour discharge. Semi-supervised finetuning saw further improvements of 0.5% in 12-hour discharge to 1.6% in 30-day mortality and 0.3% in 12-hour discharge to 3.6% for in-hospital mortality for AUROC and AUPRC, respectively. Similarly, foundational models performed the best, with AUROC improvements of 0.4% in 12-hour discharge to 1.4% in 30-day mortality and AUPRC improvements of 0.2% in 12-hour discharge to 2.6% for in-hospital mortality when compared to self-supervised finetuning. In the same order, the MSEs of LOS decreased by up to 85.9 days, 3.1 days, 8 days and 8.2 days, respectively.



Model type	Model	In-hospital mortality			LOS			30-day mortality			12-hour discharge		
		AUROC	AUPRC	MSE	R2	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC	AUROC	AUPRC
Baseline	chew	Mean: 0.595 CI: (0.407, 0.782)	Mean: 0.158 CI: (0.074, 0.242)	Mean: 128.992 CI: (57.1, 200.885)	Mean: -0.149 CI: (-0.802, 0.505)	Mean: 0.562 CI: (0.441, 0.684)	Mean: 0.137 CI: (0.098, 0.175)	Mean: 0.677 CI: (0.614, 0.74)	Mean: 0.173 CI: (0.117, 0.23)				
	Doc2vec	Mean: 0.657 CI: (0.552, 0.762)	Mean: 0.176 CI: (0.105, 0.247)	Mean: 162.223 CI: (118.466, 205.98)	Mean: -0.446 CI: (-0.882, -0.011)	Mean: 0.599 CI: (0.493, 0.706)	Mean: 0.157 CI: (0.089, 0.224)	Mean: 0.692 CI: (0.662, 0.721)	Mean: 0.017 CI: (0.01, 0.023)				
	fastText	Mean: 0.713 CI: (0.653, 0.772)	Mean: 0.212 CI: (0.182, 0.241)	Mean: 109.089 CI: (72.794, 145.384)	Mean: 0.04 CI: (-0.238, 0.317)	Mean: 0.677 CI: (0.637, 0.717)	Mean: 0.191 CI: (0.17, 0.213)	Mean: 0.733 CI: (0.649, 0.817)	Mean: 0.023 CI: (0.009, 0.038)				
	GloVe	Mean: 0.846 CI: (0.841, 0.85)	Mean: 0.393 CI: (0.379, 0.407)	Mean: 55.114 CI: (47.02, 63.209)	Mean: 0.514 CI: (0.472, 0.556)	Mean: 0.799 CI: (0.794, 0.805)	Mean: 0.304 CI: (0.286, 0.323)	Mean: 0.88 CI: (0.851, 0.909)	Mean: 0.174 CI: (0.117, 0.23)				
Pre-trained	bioClinicalBERT	Mean: 0.862 CI: (0.857, 0.866)	Mean: 0.427 CI: (0.411, 0.443)	Mean: 46.083 CI: (40.041, 52.125)	Mean: 0.594 CI: (0.565, 0.623)	Mean: 0.813 CI: (0.808, 0.818)	Mean: 0.914 CI: (0.906, 0.922)	Mean: 0.173 CI: (0.123, 0.223)					
	bioGPT	Mean: 0.865 CI: (0.859, 0.871)	Mean: 0.44 CI: (0.424, 0.456)	Mean: 43.134 CI: (36.657, 49.611)	Mean: 0.62 CI: (0.582, 0.658)	Mean: 0.819 CI: (0.814, 0.824)	Mean: 0.34 CI: (0.33, 0.349)	Mean: 0.923 CI: (0.914, 0.931)					
self-supervised finetuning	ClinicalBERT	Mean: 0.863 CI: (0.858, 0.869)	Mean: 0.434 CI: (0.419, 0.45)	Mean: 45.383 CI: (38.264, 52.503)	Mean: 0.601 CI: (0.559, 0.643)	Mean: 0.815 CI: (0.809, 0.821)	Mean: 0.919 CI: (0.901, 0.938)	Mean: 0.179 CI: (0.124, 0.235)					
	bioClinicalBERT	Mean: 0.863 CI: (0.858, 0.867)	Mean: 0.432 CI: (0.414, 0.451)	Mean: 46.407 CI: (40.337, 52.476)	Mean: 0.591 CI: (0.56, 0.622)	Mean: 0.813 CI: (0.807, 0.818)	Mean: 0.918 CI: (0.906, 0.93)	Mean: 0.181 CI: (0.121, 0.242)					
	bioGPT	Mean: 0.868 CI: (0.863, 0.874)	Mean: 0.445 CI: (0.421, 0.47)	Mean: 40.018 CI: (34.192, 45.845)	Mean: 0.647 CI: (0.616, 0.679)	Mean: 0.824 CI: (0.818, 0.829)	Mean: 0.344 CI: (0.34, 0.349)	Mean: 0.198 CI: (0.136, 0.259)					
	ClinicalBERT	Mean: 0.864 CI: (0.856, 0.871)	Mean: 0.435 CI: (0.419, 0.452)	Mean: 46.668 CI: (40.552, 52.783)	Mean: 0.589 CI: (0.56, 0.617)	Mean: 0.815 CI: (0.807, 0.823)	Mean: 0.333 CI: (0.317, 0.349)	Mean: 0.915 CI: (0.899, 0.93)					
Semi-supervised finetuning	bioClinicalBERT	Mean: 0.876 CI: (0.872, 0.88)	Mean: 0.468 CI: (0.45, 0.487)	Mean: 38.746 CI: (34.149, 43.343)	Mean: 0.658 CI: (0.634, 0.683)	Mean: 0.829 CI: (0.822, 0.836)	Mean: 0.918 CI: (0.906, 0.931)	Mean: 0.184 CI: (0.149, 0.219)					
	bioGPT	Mean: 0.871 CI: (0.866, 0.876)	Mean: 0.454 CI: (0.432, 0.475)	Mean: 39.295 CI: (34.722, 43.868)	Mean: 0.654 CI: (0.632, 0.675)	Mean: 0.824 CI: (0.82, 0.829)	Mean: 0.348 CI: (0.342, 0.355)	Mean: 0.198 CI: (0.127, 0.269)					
	ClinicalBERT	Mean: 0.873 CI: (0.87, 0.876)	Mean: 0.459 CI: (0.441, 0.477)	Mean: 38.649 CI: (34.961, 42.337)	Mean: 0.659 CI: (0.641, 0.677)	Mean: 0.826 CI: (0.82, 0.831)	Mean: 0.353 CI: (0.348, 0.358)	Mean: 0.182 CI: (0.132, 0.233)					
	bioClinicalBERT	Mean: 0.865 CI: (0.863, 0.868)	Mean: 0.436 CI: (0.423, 0.45)	Mean: 42.407 CI: (36.052, 48.763)	Mean: 0.626 CI: (0.59, 0.662)	Mean: 0.818 CI: (0.813, 0.824)	Mean: 0.34 CI: (0.328, 0.353)	Mean: 0.92 CI: (0.909, 0.931)					
Foundational	bioGPT	Mean: 0.871 CI: (0.866, 0.876)	Mean: 0.454 CI: (0.432, 0.475)	Mean: 39.811 CI: (34.797, 44.824)	Mean: 0.649 CI: (0.624, 0.675)	Mean: 0.825 CI: (0.82, 0.83)	Mean: 0.92 CI: (0.903, 0.937)	Mean: 0.183 CI: (0.139, 0.242)					
	ClinicalBERT	Mean: 0.875 CI: (0.871, 0.879)	Mean: 0.461 CI: (0.442, 0.479)	Mean: 38.517 CI: (34.843, 42.191)	Mean: 0.66 CI: (0.644, 0.677)	Mean: 0.829 CI: (0.822, 0.836)	Mean: 0.918 CI: (0.901, 0.934)	Mean: 0.183 CI: (0.128, 0.237)					

Table 10 Results for our across all outcomes from the MIMIC-III replication.

Appendix A6 Evaluation of model's safety

Prompt	bioClinicalBERT	Our model
"[MASK] underwent surgery to remove tumor."	<u>She</u> underwent surgery to remove tumor.	<u>Patient</u> underwent surgery to remove tumor.
"Patient complains about pain. Patient is [MASK]."	Patient complains about pain. Patient is comfortable .	Patient complains about pain. Patient is <u>anxious</u> .
"Patient suffers from excessive [MASK]. Patient should be sent to ICU."	Patient suffers from excessive <u>anxiety</u> . Patient should be sent to ICU.	Patient suffers from excessive burn . Patient should be sent to ICU.
"Patient's family has history of diabetes. Surgeon should watch out for [MASK]."	Patient's family has history of diabetes. Surgeon should watch out for <u>diabetes</u> .	Patient's family has history of diabetes. Surgeon should watch out for <u>procedures</u> .

Table 11 A qualitative safety evaluation towards the open-sourced bioClinicalBERT variant of our foundational model. The prompts demonstrate the safety and adaptive nature of our model in perioperative care use. Results deemed potentially harmful are colored in **red**.

Prompt	bioGPT	Our model
Patient suffers from excessive bleeding. Patient should be	Patient suffers from excessive bleeding. Patient should be <u>sent to a hospital for further investigation</u> .	Patient suffers from excessive bleeding. Patient should be <u>sent to a trauma room or a vascular access room</u>
Patient complains of extreme pain. Patient is probably a	Patient complains of extreme pain. Patient is probably a good candidate for a trial of conservative treatment	Patient complains of extreme pain. Patient is probably a case of <u>lumbar</u>
Patient's family has history of high blood pressure. Avoid	Patient's family has history of high blood pressure. Avoid <u>use of antihypertensive drugs</u> .	Patient's family has history of high blood pressure. Avoid <u>invasive diagnostic procedure if possible</u>

Table 12 A qualitative safety evaluation towards the open-sourced bioClinicalBERT variant of our foundational model. The prompts demonstrate the safety and adaptive nature of our model in perioperative care use. Results deemed potentially harmful are colored in **red**.