

# Labrador: Exploring the limits of masked language modeling for laboratory data

David R. Bellamy<sup>\*1,2</sup>

Bhawesh Kumar<sup>2</sup>

Cindy Wang<sup>3</sup>

Andrew Beam<sup>1,2</sup>

<sup>1</sup> Harvard Epidemiology Department

<sup>2</sup> Harvard Biostatistics Department

<sup>3</sup> Harvard Statistics Department

BELLAMYRD@GMAIL.COM

BHAWESH\_KUMAR@HSPH.HARVARD.EDU

CINDYWANG@COLLEGE.HARVARD.EDU

ANDREW\_BEAM@HMS.HARVARD.EDU

## Abstract

In this work we introduce LABRADOR, a pre-trained Transformer model for laboratory data. LABRADOR and BERT were pre-trained on a corpus of 100 million lab test results from electronic health records (EHRs) and evaluated on various downstream outcome prediction tasks. Both models demonstrate mastery of the pre-training task but neither consistently outperform XGBoost on downstream supervised tasks. Our ablation studies reveal that transfer learning shows limited effectiveness for BERT and achieves marginal success with LABRADOR. We explore the reasons for the failure of transfer learning and suggest that the data generating process underlying each patient cannot be characterized sufficiently using labs alone, among other factors. We encourage future work to focus on joint modeling of multiple EHR data categories and to include tree-based baselines in their evaluations.

**Keywords:** Transformer, EHR, lab data

**Data and Code Availability** The MIMIC-IV database can be accessed via PhysioNet after completing a short online training. The lab data that was used to pre-train LABRADOR and BERT come from the labevents table within MIMIC-IV and our pre-processing code can be found in our [codebase](#). The COVID-19, cancer diagnosis and alcohol consumption fine-tuning datasets are public and therefore we provide our pre-processed versions of these datasets directly in our codebase.

The codebase also contains all of the code necessary to replicate the pre-training data pre-processing, model pre-training, fine-tuning and other evaluations.

We cannot directly share the sepsis dataset due to MIMIC-IV’s data sharing policy, however with access to MIMIC-IV the steps described in Section F.3.2 recreate our evaluation dataset. Finally, the weights for the pre-trained LABRADOR and BERT models can be downloaded from our HuggingFace repository [here](#).

## 1. Introduction

In recent years, self-supervised pre-training of masked language models (MLMs) (see Appendix A for background) has demonstrated remarkable success across a wide range of machine learning problems and has led to significant downstream improvements across diverse tasks in natural language processing (Liu et al., 2019; Devlin et al., 2019; Raffel et al., 2020). There is considerable excitement surrounding the potential of large pre-trained MLMs to achieve similar success in medical applications. For instance, existing applications of MLMs in medicine have already yielded promising results in tasks related to medical text understanding (Lee et al., 2020; Alsentzer et al., 2019; Huang et al., 2019; Yang et al., 2019; Beltagy et al., 2019). Despite the success of self-supervised models in certain areas of biomedicine (Palepu and Beam, 2022; Beam et al., 2023; Singhal et al., 2023; Moor et al., 2023; Tiu et al., 2022; Ingraham et al., 2022; Watson et al., 2023; Shay et al., 2023), there has been no previous work on pre-trained models for laboratory measurements. Laboratory data is abundant, routinely collected, less biased compared to other types of data in electronic health records (EHRs) like billing codes (Beam et al., 2021), and directly mea-

---

\* Corresponding author.

sure a patient’s physiological state, offering a valuable opportunity for creating a medical foundation model.

However, there is a large body of evidence showing that deep learning is consistently outperformed on so-called “tabular” data prediction tasks by traditional machine learning techniques like random forests, XGBoost, and even simple regression models (Bellamy et al., 2020; Finlayson et al., 2023; Sharma, 2013). The reasons for this are only partially understood, but previous work (Grinsztajn et al., 2022) has suggested that this phenomenon may be caused by a rotational invariance in deep learning models that is harmful for tabular data. More broadly, the success of deep learning is thought to be largely due to inductive biases that can be leveraged for images, text, and graphs. These inductive biases are absent or only weakly present in tabular data. Conversely, tree-based methods are scale invariant and robust to uninformative features.

In this work, we introduce LABRADOR, a novel continuous Transformer architecture, which models permutation-invariant data consisting of (integer, float) tuples (Figure 1). We pre-trained LABRADOR alongside a standard implementation of BERT on 100 million lab test results from over 260,000 patients using an MLM objective. We evaluated both models on several downstream outcome prediction tasks and validated the success of pre-training with a set of intrinsic evaluations.

We discuss the design and implementation of LABRADOR and BERT, their training process, and their evaluation in Sections 3 and 4, respectively. We conclude in Section 5 with the implications of our findings, shedding light on the limitations and areas for improvement in the application of pre-trained MLMs to laboratory data and more generally, the creation of foundation models for numeric EHR data.

## 2. Related Work

In previous work, the Transformer architecture (Vaswani et al., 2017) has been adapted to handle continuous inputs alongside discrete inputs (i.e. tokens). For example, Gorishniy et al. (2021) introduce the Feature-Tokenizer (FT) Transformer, which is a similar architecture to LABRADOR. They benchmark FT Transformer against gradient boosted decision trees and Random Forests on several supervised learning tasks (without pre-training) and conclude that there is no universally superior solution. Later work by these authors (Gorishniy et al., 2022)

introduced an additional variant of continuous Transformer that uses a piecewise linear encoding system for continuous values. These authors argue that the design of embeddings for numerical features is an important consideration in adapting Transformer models to continuous data, a notion that we support strongly. Our work may be viewed as an extension of the FT Transformer to the setting of MLM-based pre-training.

Rossi et al. (2019) learned embeddings of laboratory data using Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) by encoding each lab test result as the concatenation of an octal indicator of abnormality and the test’s LOINC code<sup>1</sup>. The ablations showed improved mortality prediction when the octal indicator was included in the embeddings despite this being only a discrete representation of the underlying continuous lab test result.

Hegselmann et al. (2023) showed that pre-trained large language models (LLMs) can directly perform prediction tasks involving numeric features by serializing the features and a description of the task into a natural language string – an approach the authors called TabLLM. They found that this approach could outperform strong tree-based baselines in the very-few-shot setting (0 to 8 labeled examples). However, Dinh et al. (2022) could not reproduce this result with GPT-3 and instead found that a fine-tuned GPT-3 model performed worse than logistic regression for up to 250 training examples. Hegselmann et al. also found that the downstream performance is highly task dependent and that performance on medical prediction tasks was notably worse than others in a more general knowledge domain. However, regardless of the domain, XGBoost outperforms TabLLM on average across all tasks.

TabTransformer (Huang et al., 2020) is another architecture that can model categorical and continuous numeric features, however only categorical features pass through the attention layer. Numeric features are passed directly to the prediction head where they are concatenated to the contextualized embeddings for the categorical features. TabTransformer is distinct from our architecture, which handles data that consist of (integer, float) samples and performs attention jointly over these two types of data. To the best of our knowledge, the recently introduced xVal encoding scheme (Golkar et al., 2023) is the only existing work that adapts the Transformer architecture

---

1. <https://loinc.org/get-started/what-loinc-is/>

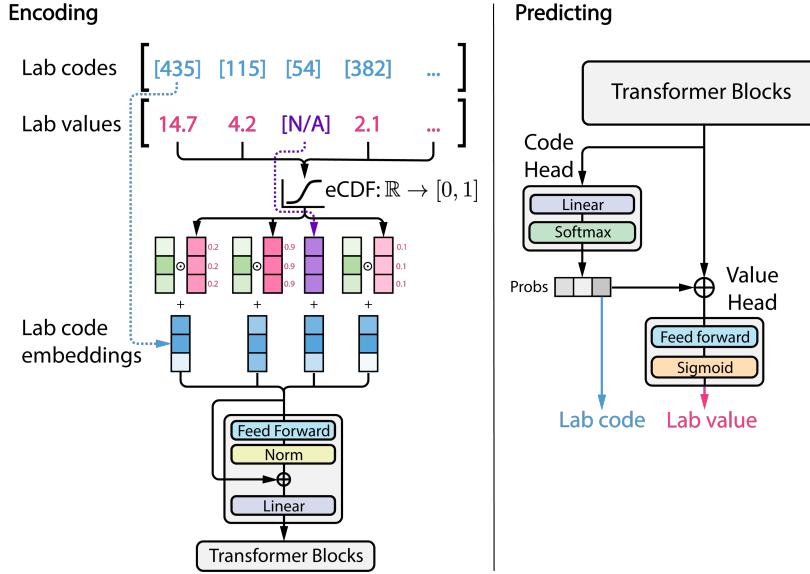


Figure 1: LABRADeR model architecture.

to jointly model these two types of data during pre-training. Their multi-task MLM pre-training objective and model architecture are similar to LABRADeR so we designed Figure 1 to facilitate their comparison. However, the authors evaluate xVal on a distinct set of tasks (learning arithmetic, temperature forecasting and predicting planetary orbits) and do not include a baseline comparison with gradient boosted decision trees.

Recent work by Zhu et al. (2024) evaluated 14 language models alongside conventional clinical predictive models and found that even with well-designed prompts, the best LLMs’ zero-shot performance does not consistently beat traditional machine learning models when there are 100 or more examples to train on.

### 3. Methods

In this section, we describe the LABRADeR and BERT model architectures and pre-training. For further experimental details, see Appendix F.

BERT is primarily distinguished from other Transformers by its learning objective rather than its model architecture. BERT’s architecture is classified as encoder-only and consists of an embedding layer, N transformer blocks with the canonical layers (Multi-HeadAttention, Dropout, LayerNorm, Feedforward, Dropout, LayerNorm), and a prediction head.

The primary contribution of the original work was the approach of masked language modeling (MLM). The key advantage of MLM is that it does not rely on the input data having a natural sequential ordering, unlike autoregressive language modeling. The lab data in our study does not have a natural sequential ordering so we had to use an MLM-based approach. We considered ways to add positional embeddings to the lab data such that they would have a sequential ordering (even if only artificially), but we concluded that it is impossible to find a natural ordering because batches of lab tests are requested at the same moment in time by physicians. So there is no way to order labs sequentially besides alphabetically or some other arbitrary order. We did not think that adding such arbitrary ordering would be helpful.

Besides encoder-only transformer architectures, GPT-style models have decoder-only architectures and older seq2seq transformer models (e.g. t5) have encoder-decoder architectures. Both decoder-only and encoder-decoder architectures are useful for generating sequences like text. In other words, these architectures were designed for seq2seq tasks. But since lab data do not have a sequential ordering, we cannot pose it as a seq2seq task so these transformer architectures are not applicable to our setting. This is why we chose a BERT-like (i.e. encoder-only) MLM-based approach in this work.

### 3.1. Labrador architecture

LABRADOR consists of two embedding modules (categorical and continuous), N transformer blocks, and two prediction heads (categorical and continuous). The categorical embedding layer uses the standard integer lookup method, whereas the continuous embedding layer involves a position-wise linear projection. A pseudo-code implementation is in Section B of the Appendix.

Each Transformer block consists of the canonical layer ordering: MultiHeadAttention, Dropout, LayerNorm, Feedforward, Dropout, LayerNorm. However, we use a distinct `key_dim` size for the multi-head attention layers. We set this hyperparameter equal to LABRADOR’s embedding dimension, whereas it is commonly set to be equal to `d_model // num_heads` where `d_model` is the model’s embedding size and `num_heads` is the number of attention heads.

The categorical prediction head consists of a ReLU-activated position-wise dense layer followed by a softmax layer over the 529 token vocabulary. A pseudo-code implementation of the continuous prediction head is also in Section B of the Appendix. The output of this layer is simply a float on the interval [0, 1], since all input values are standardized to this interval using the empirical cumulative distribution function (see Section 3.3). In our evaluations, we use 10 transformer blocks, an embedding dimension of 1024, 4 attention heads, and a feedforward dimension of 1024, resulting in 196,642,645 parameters.

### 3.2. BERT architecture

BERT is implemented using HuggingFace’s `TFBertForMaskedLM`. To eliminate positional embeddings, we pass the zero vector in place of the `position_ids` key.

The model configuration (HuggingFace’s `BERTConfig`) that is used in our experiments has a vocabulary size of 4251, a hidden size of 1024, 10 hidden layers, 4 attention heads and an intermediate (i.e. feedforward) size of 1024 with a ReLU activation.

All other parameters use their default `BERTConfig` value. Most notably, the dropout probabilities are set to 0.1 throughout the model, which is consistent with LABRADOR and common practice.

This model shares the same feedforward layer size (`intermediate_size`), number of attention heads and transformer blocks (`num_attention_heads` and `num_hidden_layers`), dropout rate, and embedding

size (`hidden_size`) with LABRADOR. Therefore, this BERT model has the same backbone as LABRADOR and differs only in its embedding module, prediction head and the `key_dim` in the multi-head attention layers.

This BERT instance has 68,522,139 parameters, which is less than LABRADOR due to their difference in `key_dim` size. To ensure that this difference in parameter count does not impact our conclusions, we performed all experiments with a parameter-matched BERT that uses the same `key_dim` as LABRADOR (Appendix E). The results from these experiments are indistinguishable from the results in Section 4.

### 3.3. Preprocessing of pre-training lab data

We pre-processed the lab data obtained from MIMIC (version IV) before MLM pre-training. First, we removed all lab codes (i.e. “itemid”) from the labevents data table that had 500 or fewer occurrences to ensure that there was a sufficient number of samples for BERT and LABRADOR to learn about each lab code. This results in 529 unique lab codes. We split patients into training, validation and test sets. 70% ( $\approx 229,000$ ) of the patients were selected uniformly at random for the training split, whereas 10% and 20% formed the validation and test splits, respectively.

For each lab code, the empirical cumulative distribution function (eCDF) is computed using the training split. Normally, mapping values of a random variable to their probabilities under an eCDF requires the entire dataset. However, in practice, there are far fewer unique lab values in the dataset. Therefore, we can represent each eCDF using only the unique lab values for each lab code as well as the probabilities that they map to on their corresponding eCDF. This compression of the eCDF is lossless, which means that it perfectly represents the eCDF obtained from the full training split. These compressed eCDFs are used in all downstream pre-training and evaluation steps and can be found in the `data/` directory of our codebase.

To tokenize the lab data for LABRADOR, lab values are mapped to the interval [0, 1] using their eCDFs whereas lab codes are mapped to their integer frequency ranking in the training split, indexed at 1. For example, hematocrit, the most frequently ordered lab test, receives the integer token 1. We also include a mask token (the integer 530) and a special null token (the integer 531). The null token is used when

a lab test is recorded but it has no associated lab value. Approximately 10% of lab tests in MIMIC-IV have no associated lab value. Although we used frequency ranking to assign integer tokens to lab codes, this ordering has no consequence so even an arbitrary ordering of the lab codes could be used.

For BERT, we construct an integer-only token vocabulary by assigning a unique integer to each decile of the eCDF for each lab code, sorted in descending order of test frequency. For example, tokens 1 through 10 represent the 10 deciles of the hematocrit lab value distribution because it is the lab test with the highest frequency ranking. We also include an eleventh token that represents a missing value for hematocrit. The second most frequent test was creatinine and therefore tokens 12 through 21 represent the 10 deciles of its lab value distribution. Once again, token 22 stands for a missing value of creatinine. This produces a vocabulary size of 4250 tokens plus one additional token that represents the `<MASK>` indicator. Although there are 529 unique lab codes, there are not  $529 \times 11 = 5819$  tokens in BERT’s vocabulary because 157 of the lab codes have no numeric values and, instead, are interpreted as binary presence/absence indicators. For example, the toxicology screens for opiates, cocaine, amphetamines and barbiturates are entered as positive/negative test results.

Finally, the labs for each patient are divided into order sets. A bag is defined as all labs ordered at the same time for a specific patient in the MIMIC database, also called an order set. Bags with fewer than 3 labs are removed from the training, validation and test splits as we consider these bags to be of insufficient size for learning embeddings. Order sets are converted to an appropriate input structure for LABRADeR and BERT separately. For LABRADeR, each bag is structured as a dictionary with keys `categorical_input` (integer) and `continuous_input` (float), whereas BERT’s input dictionary has keys `input_ids` (all integers) and `position_ids` (the zero vector). Bags of different lengths are padded to the maximum size in the current batch. In order to minimize step time during pre-training, we performed random masking of one element per bag in advance and sharded these bags across approximately 220 TensorFlow TFRecord files.

### 3.4. Transformer pre-training

Both LABRADeR and BERT were pre-trained using a single 40GB A100 GPU for 500,000 steps (80 hours) and 1.5M steps (237 hours), respectively. Our stopping criterion for each model was to observe that the validation loss was sufficiently converged. We also wanted to ensure that BERT’s degree of convergence was greater than or equal to LABRADeR’s in order to avoid falsely representing the capabilities of BERT in downstream evaluations. Both models were pre-trained with a dropout rate of 0.1, learning rate of 1e-5, a batch size of 256, and an embedding size of 1024.

We used Adam optimization for stochastic gradient descent. Model checkpoints were saved every 14,000 training steps. BERT optimized a standard categorical cross-entropy loss for measuring its accuracy at predicting masked tokens. LABRADeR optimized a simple multi-task loss that was the sum of its categorical cross-entropy for lab code prediction and mean-squared error (MSE) for lab value prediction.

## 4. Results

In the following sections, we assess LABRADeR and BERT on intrinsic evaluations to understand the extent to which each model accomplishes the pre-training task. Following that, we fine-tune both models for downstream outcome prediction on four datasets: COVID-19 diagnosis, cancer diagnosis, and the prediction of sepsis-related mortality and alcohol consumption level.

### 4.1. Intrinsic evaluations

#### 4.1.1. ASSESSMENT OF PRE-TRAINING LOSS

BERT and LABRADeR were pre-trained until each converged in their pre-training loss (Figure 6 and 7). After  $1.5 \times 10^6$  steps, BERT reaches a validation perplexity of 1.3 tokens, whereas after  $5 \times 10^5$  steps LABRADeR reaches a validation perplexity of 1.02 tokens and MSE of 0.013, respectively. Recall that BERT and LABRADeR have vocabulary sizes of 4251 and 531 tokens, respectively. This shows that both models clearly master their categorical prediction tasks.

Also, LABRADeR’s final validation MSE of 0.013 suggests that its average error is  $\pm\sqrt{0.013} \approx \pm 0.114$ , which corresponds approximately to quintile resolution in the lab value distribution because

LABRADOR’s lab values are uniform on the interval  $[0, 1]$  (described in Section 3.3).

#### 4.1.2. EMBEDDING SPACE VISUALIZATION

Despite the fact that BERT and LABRADOR have comparable pre-training perplexity, the structure of their embedding space differs greatly. We performed dimensionality reduction on BERT and LABRADOR’s embeddings for the labs in the test split using the UMAP algorithm (McInnes et al., 2018). Section F.1 describes the details of the experiment. Figure 2 presents the two-dimensional structure of the 70 most frequent lab tests, which are colored according to the panel of labs that they are most often ordered with. Appendix I defines each lab panel.

LABRADOR has a well-separated cluster of embeddings for each lab code, whereas BERT has far less separation. Qualitatively, we do not observe any clinical meaning to the relative distances between each pair of LABRADOR’s embedding clusters. In contrast, the relative positioning of the large embedding islands in BERT’s UMAP plot may possess some basic clinical meaning. For instance, urinalysis tests (pink) are directly adjacent to urine toxicology screening tests (light blue). Similarly, the CBC tests (dark blue) form a continuation with the CBC with differential tests (purple). We note that the two most common panels of lab tests (CBC and BMP) do not appear to be well-separated in BERT’s embedding space, whereas they are completely separated in LABRADOR’s embedding space.

Panel B of Figure 2 visualizes the embedding space for the four most frequently ordered lab tests and colors these embeddings by their lab value. While LABRADOR learns a smooth gradient representing the quantitative value for each lab test, BERT does not. There is some directionality in the BERT embedding space that correlates with lab value, but it is notably less monotonic when compared to LABRADOR.

It is important to note that while low-dimensional representations of embeddings can be helpful in providing high-level, qualitative intuitions for how a model behaves, they can be misleading and unfaithful to the original, high-dimensional space (Chari and Pachter, 2023). We provide these visualizations to provide some insight into the different ways LABRADOR and BERT learn to represent the input space, but caution should be used when attempting to make fine-grained, quantitative conclusions.

#### 4.1.3. IMPUTATION

We also evaluated LABRADOR and BERT for their ability to impute missing lab values in the test split of their pre-training data. To do so, we mask a randomly selected lab value from each bag of labs in the test split and run each model’s forward pass to obtain its prediction for the masked value (see Section F.2 for more details). In panel A of Figure 3, we see that the predictions from both pre-trained models achieve a Pearson correlation  $r^2 > 0.8$ . As an ablation, we also evaluated LABRADOR and BERT with randomly initialized parameters to assess the effect of pre-training on the imputations. We see that the strong correlation between the imputed and true lab values stems entirely from pre-training and is not an artifact of the model architectures or the data.

In panel B of Figure 3, we present individual scatter plots for the four best lab tests as measured by Pearson correlation. For LABRADOR, there is a clear relationship between the strength of the correlation and the frequency of the lab test in the training data. Interestingly, we found that this relationship was less pronounced for BERT, as there were some frequent lab tests where the model performed relatively poorly and a small number of rare tests where the model performed well. In panel C, we show the four worst lab tests as measured by Pearson correlation. It should be noted that although the majority of MIMIC-IV laboratory data consists of common test results, such as the tests from the panels shown in Figure 2, there is a long tail of rare tests in the database. LABRADOR and BERT may be less capable of capturing the statistical dependencies in this long tail.

## 4.2. EXTRINSIC EVALUATIONS

In this section, we evaluate LABRADOR and BERT for their ability to improve predictive performance in real-world, downstream use cases. We also conduct an ablation study to understand which parts of each model contribute most to its performance.

#### 4.2.1. FINE-TUNING EXPERIMENTS AND OUTCOME PREDICTION

We compare LABRADOR and BERT’s performance to a set of baseline models in four distinct outcome scenarios where predictions were made on the basis of laboratory measurements: ICU mortality due to sepsis, cancer diagnosis, COVID-19 diagnosis, and the prediction of alcohol consumption. Baseline methods

are defined in Section F.5 and pre-processing for each evaluation dataset is described in Section F.3.

For each of these experiments, we add a prediction head to LABRADOR or BERT consisting of multiple fully-connected layers followed by extensive hyperparameter tuning for this prediction head on each outcome dataset and with each base model. Section F.4 describes this fine-tuning procedure in detail. Table 1 shows the performance of LABRADOR and BERT compared to simple logistic (or linear) regression, Random Forest and XGBoost.

XGBoost narrowly outperforms all other methods on sepsis mortality prediction, cancer diagnosis and COVID-19 diagnosis. LABRADOR performs best on alcohol consumption prediction, which is a regression problem (as opposed to the other 3 classification tasks) and is also the smallest outcome dataset. Additionally, LABRADOR outperforms BERT across all four evaluations, suggesting there is value in a model that explicitly accounts for the continuous nature of the data. For a precise definition of each evaluation task see Section F.3.

#### 4.2.2. ABLATION STUDY

It is important to differentiate the inductive bias of the model architecture from the effect of pre-training, since previous studies have shown that many architectures (including Transformers) can perform surprisingly well even before they are trained on any data (Rives et al., 2021). In order to assess this, we perform an ablation study where the parameters of BERT and LABRADOR are randomly initialized and undergo the same hyperparameter tuning as their pre-trained counterparts. The difference between a pre-trained model’s performance and its ablation performance isolates the effect of pre-training on these downstream outcome prediction tasks. In Table 2, we observe that the pre-trained LABRADOR outperforms its ablation in 3 of 4 downstream tasks with the exception of cancer diagnosis. In contrast, the pre-trained BERT underperforms its ablation in 3 of 4 downstream tasks with the exception of COVID-19 diagnosis. Table 2 demonstrates that LABRADOR’s transfer learning is strictly superior to BERT’s.

## 5. Discussion

Through extensive experiments on downstream tasks, we demonstrated that LABRADOR consistently outperforms the standard BERT architecture on both

intrinsic and extrinsic evaluations. Our results show the value of an architecture specialized for laboratory data compared to the traditional approach of discretizing numeric values for tokenization. A key contribution of our work is the creation of an effective self-supervised objective for pre-training on the challenging distribution of laboratory data in EHRs.

Both LABRADOR and BERT were successfully optimized as masked language models, learning useful representations of laboratory tests. This was evidenced by their ability to accurately impute missing values in the test set. However, our fine-tuning experiments reveal that transfer learning from pre-training on laboratory data has limited success. On outcome prediction tasks, neither LABRADOR nor BERT consistently surpass simple baseline methods. Our ablation study (Section 4.2.2) indicates that this stems from insufficient representation learning during pre-training, rather than issues with the inductive biases or fine-tuning.

Next, we posit several reasons why pre-training on laboratory data may fail to produce significant downstream gains, and provide a summary of the available support for each.

**Lab data lack correlation structure.** If lab values were statistically independent then the pre-training task would be impossible. Likewise, if there is a trivial correlation structure (e.g. some pairs of labs are perfectly correlated) then we would not expect an MLM objective to outperform simple models. However, we believe that there is a meaningful and complex correlation structure amongst the lab measurements, as evidenced in Figure 8.

**MLM masking rate is too low.** LABRADOR and BERT may satisfy the pre-training objective without learning enough about labs and patients to be adept at outcome prediction. During our pre-training, the average mask rate is just 5% per bag. Unfortunately, increasing the mask rate may be impossible because each bag of labs is permutation invariant. As a result, if multiple elements are masked within a bag, LABRADOR and BERT’s predictions for each mask token will be identical. There is no way for the model to distinguish one mask token from the next in a permutation invariant setting.

Positional embeddings are typically used to break this symmetry, but we are not aware of a variable to encode using via positional embeddings. For example, the index in the input cannot be used since the relative ordering of multiple mask tokens is arbitrary. The timestamp when a lab was measured could

Table 1: Comparison of predictive performance among LABRADeR, BERT, and baseline models across four outcome prediction settings.

	Sepsis mortality ( $\downarrow$ )	Cancer diagnosis ( $\downarrow$ )	COVID-19 diagnosis ( $\downarrow$ )	Alcohol consumption ( $\downarrow$ )
LR <sup>1</sup>	0.410 (0.398, 0.424)	1.177 (1.100, 1.272)	0.462 (0.434, 0.529)	9.73 (4.83, 13.93)
RF	<b>0.387</b> (0.372, 0.401)	0.942 (0.902, 0.987)	0.431 (0.404, 0.487)	9.65 (4.62, 13.35)
XGB	<b>0.387</b> (0.376, 0.399)	<b>0.918</b> (0.847, 0.989)	<b>0.419</b> (0.371, 0.504)	10.84 (5.54, 15.39)
BERT	0.406 (0.368, 0.431)	1.01 (0.859, 1.127)	0.441 (0.391, 0.513)	8.34 (5.80, 12.71)
LABRADeR	0.400 (0.371, 0.420)	0.978 (0.889, 1.163)	0.425 (0.369, 0.487)	<b>7.11</b> (5.33, 9.15)

Note: We report the metrics as mean (min, max) across 5 random replicates. The best model for each evaluation dataset is in bold. We use cross-entropy as evaluation metrics for all but alcohol consumption, where we use MSE. Lower values indicate better predictive performance in all cases.

<sup>1</sup> LR = logistic regression for all evaluations but alcohol consumption, where LR = linear regression.

Table 2: Ablation study evaluating the influence of pre-training on the predictive performance of BERT and LABRADeR in four outcome prediction settings.

	BERT		LABRADeR	
	Pre-trained	Random weights (%) <sup>1</sup>	Pre-trained	Random weights (%) <sup>1</sup>
Sepsis mortality	0.406 (0.368, 0.431)	+0.5%	<b>0.400</b> (0.371, 0.420)	-2%
Cancer diagnosis	1.01 (0.859, 1.127)	+8.8%	0.978 (0.889, 1.163)	+5.1%
COVID-19 diagnosis	<b>0.441</b> (0.391, 0.513)	-3.2%	<b>0.425</b> (0.369, 0.487)	-3.8%
Alcohol consumption	8.34 (5.80, 12.71)	+11.5%	<b>7.11</b> (5.33, 9.15)	-5.9%

<sup>1</sup> A positive percentage corresponds to an improvement in the evaluation metric relative to the corresponding pre-trained model.

be used but labs are ordered in sets at the exact same moment in time in a patient’s medical record. Even if a bag of labs includes tests that span multiple order sets, there would still be a symmetry issue within each order set. We believe that this is a general challenge for any MLM of electronic health record data, since measurements are often made in order sets.

**Pre-training and downstream evaluations are a poor match.** The set of features used for cancer diagnosis is the least similar to the pre-training vocabulary and is also the evaluation that both LABRADeR and BERT perform the worst on. However, sepsis mortality prediction uses lab data from the same ICU data distribution as in pre-training and yet it is not the evaluation where pre-training delivers the most value, which would be COVID-19 diagnosis. Therefore, we do not believe that a mismatch between pre-training and downstream data distribu-

tions is the primary factor explaining the failure in transfer learning.

**Lab data only partially capture patient state.** Lab data, unlike text, do not provide a comprehensive view of a patient’s condition, necessitating the incorporation of additional patient data categories. Thus, it could be that even simple models are close to the Bayes error rate given the information captured by lab measurements. We propose that future work should jointly learn from multiple EHR data categories, like diagnostic codes, procedure codes, and medication prescriptions. This could lead to a multi-modal model predicting next tokens across structured EHR data, medical images, and clinical text.

**Insufficient data scale.** The capabilities of frontier LLM’s have emerged as a function of the scale of the pre-training data. We pre-trained LABRADeR and BERT on approximately 100 million lab tests but this only corresponds to 4.5 million input sam-

ples. In comparison, GPT-3 (Brown et al., 2020) was trained on a corpus of 300 billion tokens where each input sequence contained 2048 tokens or about 145 million input sequences. This is nearly two orders of magnitude larger than our pre-training dataset in terms of input sequences and over 3 orders of magnitude larger in token count. We believe that this is at least 1 order of magnitude larger than the largest available lab dataset today. Therefore, we encourage others to contribute to the harmonization of datasets across many sources, which requires large organized coding efforts to unify the data into a compatible vocabulary. This also requires open access to such data sources, as MIMIC has provided. Importantly, if we take into account the Chinchilla scaling laws from Hoffmann et al. (2022), it is plausible that scaling model size will not yield significant benefits in downstream performance given the relative scarcity of lab data compared to text.

## 6. Conclusion

In this work, we presented LABRADOR, a novel Transformer architecture that can be pre-trained on (integer, float) data such that the float values are treated in a continuous manner. We showed that LABRADOR outperforms BERT across all downstream evaluations and appears to learn monotonic axes corresponding to lab value in its embedding space. However, despite the success seen in pre-training, both models fail to outperform XGBoost on four outcome prediction tasks. Our ablation studies reveal that transfer learning largely fails in this setting despite being more successful for LABRADOR than for BERT. We hypothesize that this failure is largely due to a lack of pre-training data scale and an insufficient characterization of the data generating process underlying each patient using labs in isolation from other clinical data types. As a result, we urge future work to focus on harmonizing disparate datasets and jointly modeling several categories of EHR data.

A persistent challenge in medical machine learning for structured data is the lack of universally accepted evaluations. Model predictions are difficult to verify, since there is no human baseline for performance. As a result, we have a poor sense of the Bayes error rate for these types of prediction tasks. We argue that a prerequisite to progressing the self-supervised learning efforts on EHR data is the establishment of universally accepted evaluations that will provide a clear

picture of the field’s progress over time in comparison to tree-based methods (Bellamy et al., 2020).

## References

- Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. Publicly available clinical BERT embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*, pages 72–78, Minneapolis, Minnesota, USA, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-1909. URL <https://aclanthology.org/W19-1909>.
- Kristyn Beam, Puneet Sharma, Bhawesh Kumar, Cindy Wang, Dara Brodsky, Camilia R Martin, and Andrew Beam. Performance of a large language model on practice questions for the neonatal board examination. *JAMA Pediatr.*, July 2023.
- Kristyn S Beam, Matthew Lee, Keith Hirst, Andrew Beam, and Richard B Parad. Specificity of international classification of diseases codes for bronchopulmonary dysplasia: an investigation using electronic health record data and a large insurance database. *J. Perinatol.*, 41(4):764–771, April 2021.
- David Bellamy, Leo Celi, and Andrew L Beam. Evaluating progress on machine learning for longitudinal electronic healthcare data. *arXiv preprint arXiv:2010.01149*, 2020.
- Iz Beltagy, Kyle Lo, and Arman Cohan. SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1371. URL <https://aclanthology.org/D19-1371>.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Federico Cabitza, Andrea Campagner, Davide Ferrari, Chiara Di Resta, Daniele Ceriotti,

- Eleonora Sabetta, Alessandra Colombini, Elena De Vecchi, Giuseppe Banfi, Massimo Locatelli, and Anna Carobene. Development, evaluation, and validation of machine learning models for covid-19 detection based on routine blood tests. <https://zenodo.org/records/4081318#.X4RWqd7TIU>, 2020. Accessed: 2023-01-07.
- Federico Cabitza, Andrea Campagner, Davide Ferrari, Chiara Di Resta, Daniele Ceriotti, Eleonora Sabetta, Alessandra Colombini, Elena De Vecchi, Giuseppe Banfi, Massimo Locatelli, et al. Development, evaluation, and validation of machine learning models for covid-19 detection based on routine blood tests. *Clinical Chemistry and Laboratory Medicine (CCLM)*, 59(2):421–431, 2021.
- Tara Chari and Lior Pachter. The specious art of single-cell genomics. *PLOS Computational Biology*, 19(8):e1011288, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *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>.
- Tuan Dinh, Yuchen Zeng, Ruisu Zhang, Ziqian Lin, Michael Gira, Shashank Rajput, Jy-yong Sohn, Dimitris Papailiopoulos, and Kangwook Lee. Lift: Language-interfaced fine-tuning for non-language machine learning tasks. *Advances in Neural Information Processing Systems*, 35:11763–11784, 2022.
- Samuel G Finlayson, Andrew L Beam, and Maarten van Smeden. Machine learning and statistics in clinical research articles—moving past the false dichotomy. *JAMA pediatrics*, 177(5):448–450, 2023.
- Siavash Golkar, Mariel Pettee, Michael Eickenberg, Alberto Bietti, Miles Cranmer, Geraud Krawezik, Francois Lanusse, Michael McCabe, Ruben Ohana, Liam Parker, et al. xval: A continuous number encoding for large language models. *arXiv preprint arXiv:2310.02989*, 2023.
- Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning models for tabular data. *Advances in Neural Information Processing Systems*, 34:18932–18943, 2021.
- Yury Gorishniy, Ivan Rubachev, and Artem Babenko. On embeddings for numerical features in tabular deep learning. *Advances in Neural Information Processing Systems*, 35:24991–25004, 2022.
- Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? *Advances in Neural Information Processing Systems*, 35:507–520, 2022.
- Stefan Hergelmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaoyi Jiang, and David Sontag. Tabllm: Few-shot classification of tabular data with large language models. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 5549–5581. PMLR, 25–27 Apr 2023. URL <https://proceedings.mlr.press/v206/hegelmann23a.html>.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. An empirical analysis of compute-optimal large language model training. *Advances in Neural Information Processing Systems*, 35:30016–30030, 2022.
- Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. Clinicalbert: Modeling clinical notes and predicting hospital readmission. *arXiv preprint arXiv:1904.05342*, 2019.
- Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data modeling using contextual embeddings. *arXiv preprint arXiv:2012.06678*, 2020.
- John Ingraham, Max Baranov, Zak Costello, Vincent Frappier, Ahmed Ismail, Shan Tie, Wujie Wang, Vincent Xue, Fritz Obermeyer, Andrew Beam, and Others. Illuminating protein space with a programmable generative model. *bioRxiv*, pages 2022–2012, 2022.
- Alistair Johnson, Lucas Bulgarelli, Tom Pollard, Steven Horng, Leo Anthony Celi, and Roger Mark. Mimic-iv. *PhysioNet*.

Available online at: <https://physionet.org/content/mimiciv/1.0/> (accessed August 23, 2021), 2020.

Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. Mimic-iv, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1, 2023.

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*, 36(4):1234–1240, February 2020.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, April 2023.

Anil Palepu and Andrew L Beam. Tier: Text-image entropy regularization for clip-style models. *arXiv preprint arXiv:2212.06710*, 2022.

Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. URL <https://aclanthology.org/D14-1162/>.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.

Alexander Rives, Joshua Meier, Tom Serdu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proc. Natl. Acad. Sci. U. S. A.*, 118(15), April 2021.

Lorenzo A Rossi, Chad Shawber, Janet Munu, and Finly Zachariah. Evaluation of embeddings of laboratory test codes for patients at a cancer center. *arXiv preprint arXiv:1907.09600*, 2019.

Dhruv Sharma. *Elements of Optimal Predictive Modeling Success in Data Science: An Analysis of Survey Data for the 'Give Me Some Credit' Competition Hosted on Kaggle*. SSRN, 2013.

Denys Shay, Bhawesh Kumar, David Bellamy, Anil Palepu, Mark Dershitz, Jens M Walz, Maximilian S Schaefer, and Andrew Beam. Assessment of ChatGPT success with specialty medical knowledge using anaesthesiology board examination practice questions. *Br. J. Anaesth.*, May 2023.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.

Ekin Tiu, Ellie Talius, Pujan Patel, Curtis P Langlotz, Andrew Y Ng, and Pranav Rajpurkar. Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. *Nat Biomed Eng*, 6(12):1399–1406, December 2022.

I-Jung Tsai, Wen-Chi Shen, Chia-Ling Lee, Horng-Dar Wang, and Ching-Yu Lin. Machine learning in prediction of bladder cancer on clinical laboratory data. *Diagnostics*, 12(1):203, 2022.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you

need. In I Guyon, U V Luxburg, S Bengio, H Wallach, R Fergus, S Vishwanathan, and R Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017.

Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, Basile I M Wicky, Nikita Hanikel, Samuel J Pellock, Alexis Courbet, William Sheffler, Jue Wang, Preetham Venkatesh, Isaac Sappington, Susana Vázquez Torres, Anna Lauko, Valentin De Bortoli, Emile Mathieu, Sergey Ovchinnikov, Regina Barzilay, Tommi S Jaakkola, Frank DiMaio, Minkyung Baek, and David Baker. De novo design of protein structure and function with RFdiffusion. *Nature*, July 2023.

An Yang, Quan Wang, Jing Liu, Kai Liu, Yajuan Lyu, Hua Wu, Qiaoqiao She, and Sujian Li. Enhancing pre-trained language representations with rich knowledge for machine reading comprehension. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2346–2357. Association for Computational Linguistics, Jul 2019. URL <https://aclanthology.org/P19-1226>.

Yinghao Zhu, Junyi Gao, Zixiang Wang, Weibin Liao, Xiaochen Zheng, Lifang Liang, Yasha Wang, Chengwei Pan, Ewen M Harrison, and Liantao Ma. Is larger always better? evaluating and prompting large language models for non-generative medical tasks. *arXiv preprint arXiv:2407.18525*, 2024.

## Appendix A. Background: transformers and masked language models

### A.1. Transformers

Transformers are a class of deep learning models introduced by Vaswani et al. (2017) that have significantly impacted various areas of machine learning, including natural language processing, computer vision, and reinforcement learning. Transformers are based on the principle of self-attention, which allows them to efficiently capture long-range dependencies in input data by computing weighted combinations

of all input elements instead of relying on the fixed-size sliding windows used in traditional convolutional and recurrent neural networks.

The core building block of a Transformer is the attention mechanism, which computes the similarity between all pairs of input elements and replaces each input element with a weighted average of all other inputs, weighted by their similarity. Transformers are built by stacking multiple layers of self-attention mechanisms, enabling the model to learn complex patterns and representations in the input data.

### A.2. Masked language models

Masked language models (MLMs) are a popular pre-training technique for natural language processing tasks. The key idea behind MLMs is to train a model to predict missing words in a given text, with some portion of the input words being masked out. By learning to predict the masked words, the model is forced to capture the underlying structure and semantics of the language, resulting in a more robust and generalizable representation.

MLMs are usually pre-trained on large-scale text corpora and fine-tuned on specific downstream tasks, such as text classification, sentiment analysis, or named entity recognition. This two-step process, consisting of pre-training and fine-tuning, is known as transfer learning. Pre-training on large corpora allows the model to learn general language features, while fine-tuning adapts the model to the nuances of the specific task at hand.

The success of MLMs is primarily attributed to their ability to learn contextualized word representations, which capture both syntactic and semantic information from the input text. This is in contrast to traditional word embedding techniques, such as Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014), which learn static word representations that are context-independent.

In this work, we leverage the power of Transformers and the masked language model pre-training technique to explore their applicability to laboratory data in medicine. We aim to assess whether the successes of Transformers and MLMs in other domains can be replicated in the context of medical data, particularly in outcome prediction tasks based on laboratory data.

## Appendix B. Pseudo-code Implementations

The continuous embedding layer is implemented as follows:

```
def ContinuousEmbeddingLayer(
    lab_values, token_embeddings):
    x = Dense_pw(lab_values,
                 out_size=d_model,
                 activation='linear')
    x = x + token_embeddings
    x = Dense_pw(x, out_size=d_model,
                 activation='relu')
    x = LayerNorm(x)
    return x
```

Where `Dense_pw` represents a position-wise dense layer and `token_embeddings` is the output of the categorical embedding layer.

The continuous prediction head is implemented as follows:

```
def ContinuousPredictionHead(
    final_layer_embeddings, probs):
    x = concat([final_layer_embeddings,
                probs])
    x = Dense_pw(x,
                 out_size=embed_dim + probs_dim,
                 activation='relu')
    x = Dense_pw(x, out_size=1,
                 activation='sigmoid')
    return x
```

Where `probs` is the categorical head's probability distribution over the lab code vocabulary. This supplies the continuous prediction head with each token's predicted probability when attempting to impute the masked lab value.

## Appendix C. Embedding Space Visualization

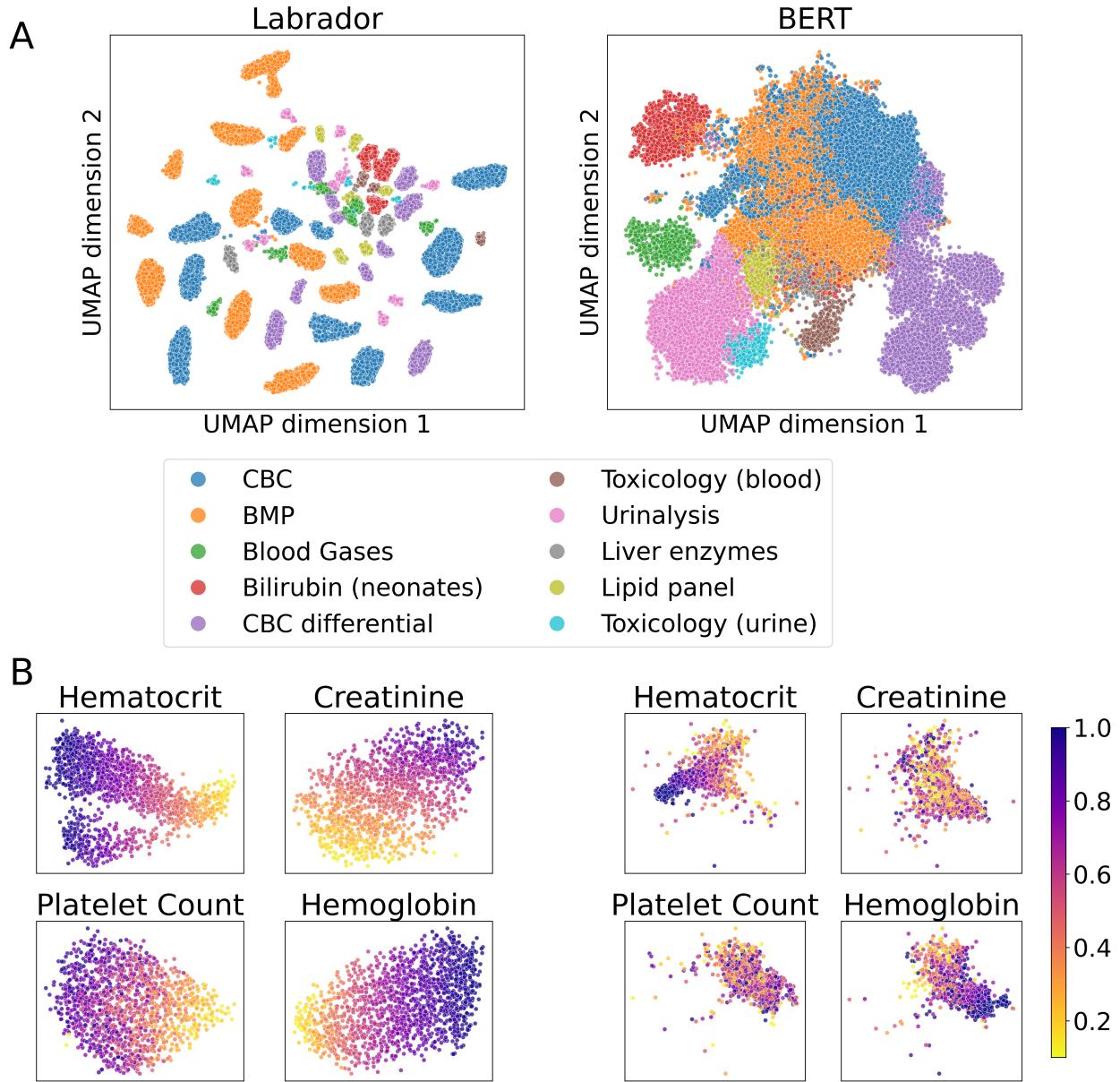


Figure 2: UMAP Visualization of LABRADOR and BERT Embeddings. **A.** A global view of the embedding space structure for LABRADOR (left) and BERT (right). The 70 most frequently ordered lab tests are shown and colored according to the panel of tests they are typically ordered with (see Appendix I for all panel definitions). All labs on this figure are from the test split and were not seen during pre-training. **B.** Visualization of embeddings for four routinely collected laboratory measurements colored by lab value and scaled to the interval [0, 1]. LABRADOR appears to encode the measured lab value in a much more natural way with a smooth gradient for lab value compared to BERT.

## Appendix D. Imputation Evaluation of Labrador and BERT

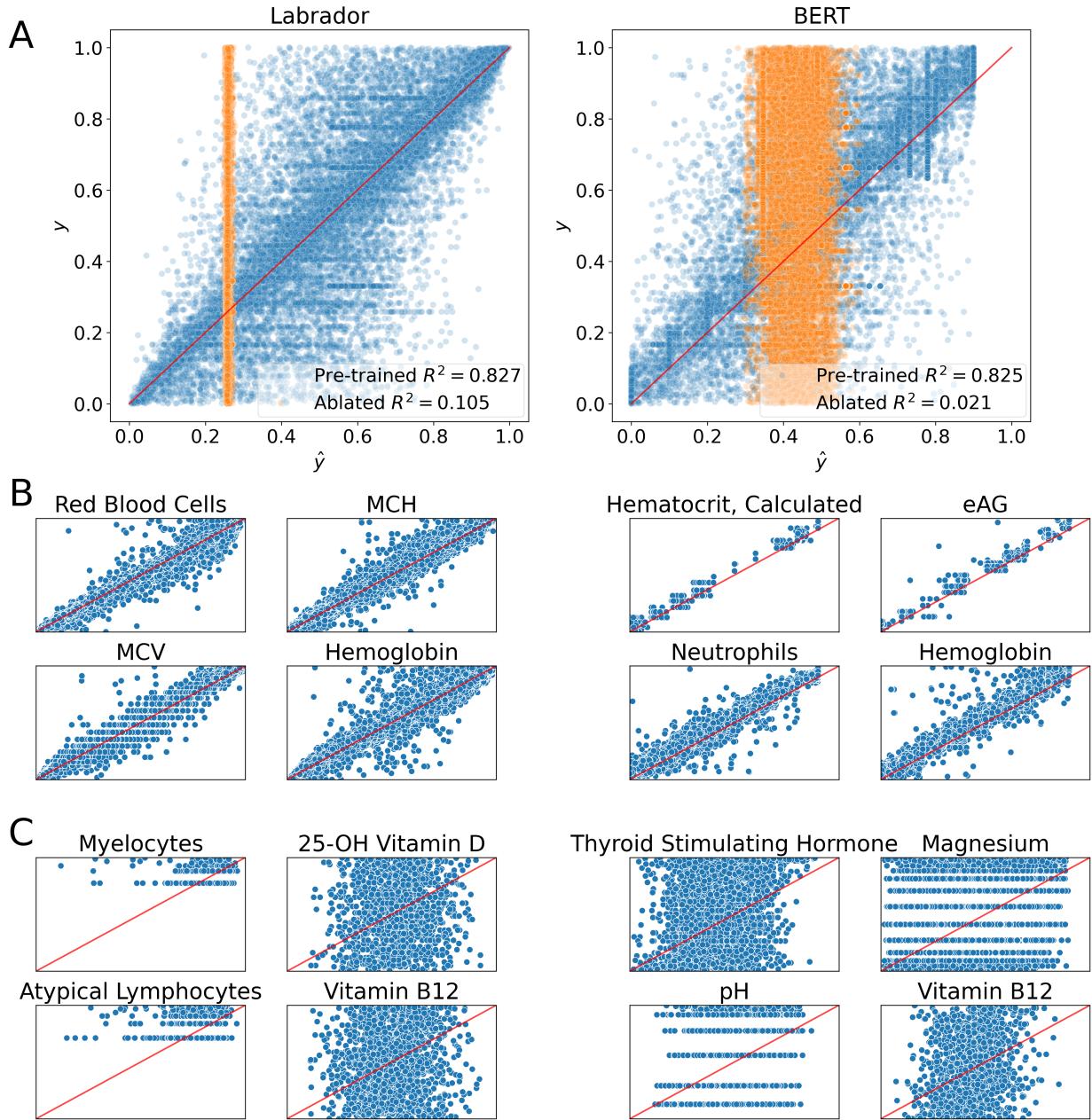


Figure 3: Intrinsic Evaluation of LABRADOR and BERT by lab value imputation. **A.** We masked a random lab from each bag in the test set of the pre-training data and imputed these values using pre-trained LABRADOR (left) and BERT (right). Both LABRADOR and BERT achieve a Pearson correlation  $r^2 > 0.8$ , in contrast to their ablations (orange). **B.** Imputations for the four best lab tests as measured by Pearson correlation. **C.** Imputations for the four worst lab tests as measured by Pearson correlation.

## Appendix E. Parameter-matched BERT experiments

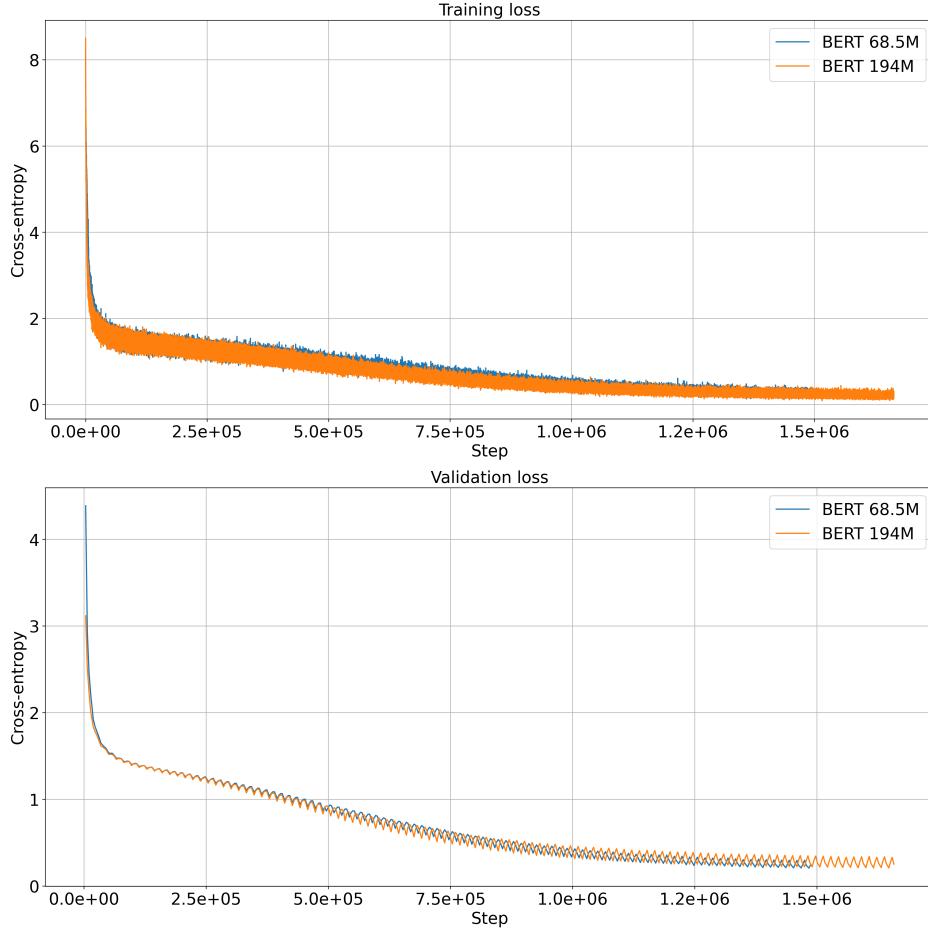


Figure 4: Pre-training loss for BERT with 68.5M ( $d_k = \lfloor \frac{d_{\text{model}}}{h} \rfloor \rfloor$ ) versus 194M parameters ( $d_k = d_{\text{model}}$ ). **Top:** The training loss for both BERT models. **Bottom:** The validation loss.

To ensure that our choice for LABRADOR’s key dimension,  $d_k = d_{\text{model}}$  (instead of the more common  $d_k = \frac{d_{\text{model}}}{h}$ , where  $h$  is the number of attention heads) does not alter our results or conclusions, we pre-trained a second BERT model with the same  $d_k$  as LABRADOR and repeated our evaluations.

BERT with  $d_k = \frac{d_{\text{model}}}{h}$  has 68.5M parameters, whereas the BERT model with the larger  $d_k$  has 194M parameters. The latter matches LABRADOR’s parameter count of 196M as closely as possible, since their backbone structure is identical besides the additional 2M weights in LABRADOR’s continuous embedding layer and continuous prediction head. Figure 4 compares the pre-training loss curves for BERT with 68.5M versus 194M parameters. Figure 5 compares their intrinsic lab value imputation performance. The imputation performance of BERT 194M is slightly better than BERT 68.5M, despite their pre-training loss curves appearing nearly identical.

Most importantly, we also include the larger BERT’s fine-tuning performance on the downstream outcome prediction tasks, both with pre-training and without. Table 3 shows that the smaller BERT model outperforms the larger BERT on 7 out of the 8 evaluations, but in general the difference is not statistically significant.

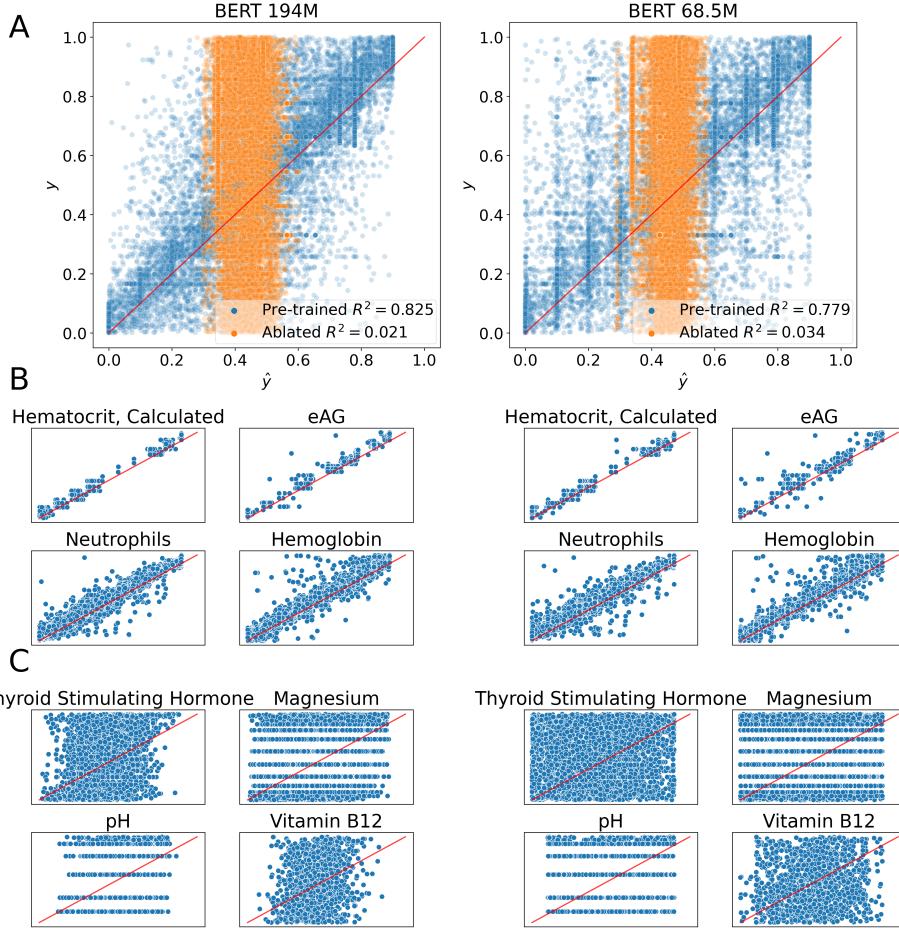


Figure 5: Lab value imputations from BERT with 68.5M ( $d_k = \lfloor \frac{d_{\text{model}}}{h} \rfloor$ ) versus 194M parameters ( $d_k = d_{\text{model}}$ ). **A.** Imputations from both pre-trained BERT models (blue) as well as their ablations (orange) on the test split of the pre-training data. **B.** Imputations for the four best lab tests as measured by Pearson correlation. **C.** Imputations for the four worst lab tests as measured by Pearson correlation.

Table 3: Comparison of fine-tuning performance between BERT with 68.5M versus 194M parameters. **Left:** Performance of each model size after pre-training. **Right:** Performance of each model size starting from randomly initialized weights.

	Pre-trained		Random weights	
	BERT 68M ( $\downarrow$ )	BERT 194M ( $\downarrow$ )	BERT 68M ( $\downarrow$ )	BERT 194M ( $\downarrow$ )
Sepsis mortality	0.406 (0.368, 0.431)	<b>0.400</b> (0.368, 0.424)	<b>0.404</b> (0.378, 0.418)	0.406 (0.379, 0.419)
Cancer diagnosis	<b>1.01</b> (0.859, 1.127)	1.08 (0.978, 1.20)	<b>0.921</b> (0.774, 1.054)	0.972 (0.884, 1.10)
COVID-19 diagnosis	<b>0.441</b> (0.391, 0.513)	0.484 (0.400, 0.544)	<b>0.455</b> (0.412, 0.511)	0.470 (0.435, 0.489)
Alcohol consumption	<b>8.34</b> (5.80, 12.71)	8.68 (6.57, 12.41)	<b>7.38</b> (4.93, 10.43)	7.49 (4.39, 11.28)

## Appendix F. Further experimental details

### F.1. UMAP

After performing several scaling experiments with UMAP reduction on the full test split of the pre-training data, we observed that the results were not significantly altered if we used a random subsample of the data. We randomly selected bags of labs totalling 50,000 (lab code, lab value) pairs from the test split. These bags separately underwent LABRADOR and BERT’s forward pass to obtain the final layer embeddings from each Transformer. The UMAP algorithm was run with `n_neighbors=600`, `min_dist=0.99` and a random seed of 3141592. In preliminary experiments, we found that increasing the number of neighbors and increasing the minimum distance between points in the reduction provided a better global structure of the embedding space. In consultation with the MIMIC database authors, we manually labeled lab codes according to their panel membership, such as the Complete Blood Count (see Appendix I for all panels). This allowed us to color Figure 2 accordingly.

### F.2. Intrinsic imputations

For intrinsic imputations, we randomly sampled approximately 2 million bags of labs from the pre-training test split and randomly masked one (lab code, lab value) pair. These masked bags were run through the forward pass of LABRADOR and BERT in order to assess their predictions for the masked lab value. LABRADOR’s continuous prediction head outputs a predicted lab value, which we compared directly to the true masked lab value.

Since BERT only predicts tokens, we experimented with two methods for converting these categorical predictions into continuous values. First, we tried an argmax method, which simply extracts the token with the greatest logit in BERT’s output distribution and set the predicted lab value equal to the lower bound of that token’s defined decile. Second, we tried a weighted quantile method, which takes a weighted average of the lower bounds of each token’s decile where the weights are the re-normalized probabilities for that lab code’s subset of the output distribution. We found that this weighted quantile method works far better than the argmax method, which is likely due to the fact that the model is able to hedge its probabilities across several deciles.

We present the Pearson correlation between the true lab values and each model’s predictions in Figure 3. We also calculated Pearson correlations for each unique lab code and used this to rank lab codes from best to worst in terms of their imputation accuracy. Finally, we performed an ablation experiment to highlight what was learned during pre-training. The ablation was conducted by randomly initializing the weights for each model prior to performing the steps described above.

### F.3. Pre-processing the evaluation data

#### F.3.1. COVID-19 COHORT

For COVID-19 diagnosis, we used a public dataset ([Cabitza et al., 2020, 2021](#)) that contains blood work for 1,624 patients (52% COVID-19 positive) admitted at San Raphael Hospital from February to May 2020. The data contain 34 features and a binary target variable for COVID-19 diagnosis determined by RT-PCR at a later point in time. Features ‘CK’ (creatinine kinase) and ‘UREA’ (blood urea nitrogen) were dropped from the data because their rates of missingness were too high (60% and 38.5%, respectively). Any patients that were missing greater than 25% of the remaining 32 features were dropped as well since we believe that the imputation quality for these patients would be too poor. This resulted in 1,312 patients.

The features contain 26 lab tests that exist in the MIMIC-IV vocabulary of LABRADOR and BERT. We manually verified the scale of each lab test to ensure that the units of measurement were identical. One test, ‘CA’ (calcium), was initially in mmol/L so we converted it to mg/dL, which is the unit of measurement used in MIMIC-IV for this test (itemid 50893). Our pre-processed version of the COVID-19 dataset can be found in our codebase on GitHub.

During evaluation of the baseline methods, any missing values were imputed using the column mean in the current split during k-fold cross-validation. In contrast, both Transformers have natural ways of handling missing values without requiring imputation. LABRADOR uses a special null token for missing values and BERT has a special token in its vocabulary for missing values for each lab code. Finally, we used the stored eCDFs from pre-training to convert all lab values from the 26 in-vocabulary features to the interval [0, 1]. LABRADOR was fine-tuned directly on these values. For BERT, an additional step

is required which maps the eCDF values to their respective tokens in the model’s vocabulary.

### F.3.2. SEPSIS COHORT

We downloaded the MIMIC-IV sepsis cohort from Google BigQuery under `physionet-data/mimic-derived/sepsis3`. Accessing the MIMIC-IV data (Johnson et al., 2020, 2023) requires signing a data use agreement and completing a short online training, so we cannot share this dataset directly. However, after obtaining MIMIC-IV access on PhysioNet, you can find instructions on accessing the data using Google BigQuery.

The sepsis dataset contains information on 34,678 septic patients, such as their patient ID, hospital stay ID, and several other variables. The variable `sofa_time`, which is the timestamp for the patient’s most recent SOFA score, was used to link patients to their other MIMIC records. These patients were linked to their admission record to determine whether they survived their ICU stay or not. We created a binary outcome variable equal to 1 if the patient died by the end of their stay and 0 otherwise. We also linked patients to their lab test results. We filtered the labs to only include those that were measured in the first 24 hours of the patient’s stay in the ICU. Patients with fewer than 2 lab tests during the first 24 hours were dropped from the analysis, which left 34,136 remaining patients.

In order to conduct a fair comparison with baseline methods, we opted to create a tabular dataset out of this ragged dataset of lab tests. To do so, we selected the 23 most frequently occurring lab tests in the first 24 hours. We chose the top 23 because this produced a tabular dataset for all 34,136 patients with under 5% missingness. Finally, we converted all lab values to the interval [0, 1] using the eCDF method described in Section 3.3 and mapped these values into BERT’s token vocabulary. The code to recreate this dataset is included in our codebase.

### F.3.3. CANCER DIAGNOSIS COHORT

For cancer diagnosis, we used a public dataset shared by Tsai et al. (2022) that contains blood work for 1,336 patients. Briefly, each patient received a diagnosis of either cystitis, bladder cancer, kidney cancer, uterine cancer, or prostate cancer. The original dataset contained 39 features in addition to the 5-class target variable. Five of these 39 features were

removed for having too much missingness (urine epithelium, A/G ratio, urine ketone, urine glucose, and strip WBC). Among the remaining 34 features, 24 were lab tests within the MIMIC-IV vocabulary, 2 lab tests were out of vocabulary, and there were 8 social variables (age, gender, smoking status, diabetes status, etc.). We manually verified that each in-vocabulary lab test was using the same units of measurement as in MIMIC-IV. For categorical variables that had missing values, we used the missing indicator method. We transformed all lab values using the eCDF method and missing lab values were handled in the same manner as for the COVID-19 dataset (Section F.3.1).

### F.3.4. ALCOHOL CONSUMPTION COHORT

We obtained the alcohol consumption dataset from the UCI Machine Learning Repository<sup>2</sup>. This small dataset contains 345 patients, 5 lab test features and 1 continuous outcome variable. Briefly, the outcome is the number of half-pint equivalents of alcoholic beverages drunk per day. The 5 lab test features are all thought to be sensitive to liver disorders that might arise from excessive alcohol consumption. There is no missingness in this dataset therefore no pre-processing was required for the baseline methods. For LABRADOR and BERT, the lab values were transformed with the eCDF method.

## F.4. Transformer fine-tuning

In order to fine-tune LABRADOR and BERT on the four downstream outcome prediction tasks, we created a simple wrapper class for each model that could handle a combination of in-vocabulary lab features as well as previously unseen features. The pseudocode for LABRADOR’s wrapper class’ call method is as follows:

```
def LabradorFinetuneWrapper(
    inputs: Dict[Tensor]) -> Tensor:
    x = base_model(inputs['lab_features'])
    x = pool(x)

    if 'non_mimic_features' in inputs.keys():
        non_mimic_features = dense(
            inputs['non_mimic_features'])
        x = concat([x, non_mimic_features],
                   axis=1)
```

---

2. <https://doi.org/10.24432/C54G67>

```

x = dense(x)
x = dropout(x)
x = dense(x)
return x

```

The lab features are passed to the base model as usual. The base model’s final embeddings are averaged into a single vector of size `hidden_dim` (i.e. 1024 in our experiments). If there are other features, these are passed through a 1-to-1 dense layer with a ReLU activation followed by concatenation with the pooled embedding vector. This combined representation is passed through an additional 1-to-1 ReLU-activated dense layer prior to the final dense layer that projects this vector down to the size required by the outcome prediction task. This final dense layer also has an activation function that is dependent on the outcome prediction task (e.g. linear for alcohol consumption prediction, sigmoid for sepsis mortality prediction). The pseudocode is the same for BERT.

Across all fine-tuning experiments, we froze the parameters of the base model and performed hyperparameter tuning across a grid consisting of the learning rate, dropout rate, number of fine-tuning epochs and batch size. More specifically, we use the following grid of hyperparameters:

- Number of epochs: [30, 60, 90]
- Batch size: [16, 32, 64]
- Learning rate: [1e-4, 3e-4, 5e-4, 1e-3]
- Dropout rate: [0.1, 0.3, 0.5, 0.7]

In our preliminary experiments, we found that this grid spans all plausible regions of interest in hyperparameter space. We also conducted preliminary experiments with unfrozen base models and the same hyperparameter grid as the original authors of BERT, but we found that this resulted in worse performance.

## F.5. Baseline methods for outcome prediction

We include a set of simple baseline methods for comparison in our fine-tuning experiments. This included linear regression, logistic regression, Random Forest (classification and regression) and XGBoost (classification and regression). For all logistic regression experiments, we tuned  $C$ , the coefficient for the strength of  $L_2$ -regularization, using the grid: [0.0001, 0.001, 0.01, 0.1]. Otherwise, we allowed 200 iterations

for the solver to converge and used the lbfgs solver. For linear regression, we used ordinary least squares. For the random forest baseline and XGBoost, the hyperparameter grids can be seen below.

Random forest:

- Number of estimators: [30, 100, 200]
- Max depth: [3, 5, 10]
- Max features: [All, sqrt]

XGBoost grid:

- Number of estimators: [30, 100, 200]
- Subsample: [0.6, 0.8, 1.0]
- Column sample by tree: [0.6, 0.8, 1.0]
- Max depth: [3, 5]

During fine-tuning, continuous features were normalized using the mean and standard deviation of the current fold in K-fold cross-validation.

## Appendix G. Pre-training loss curves

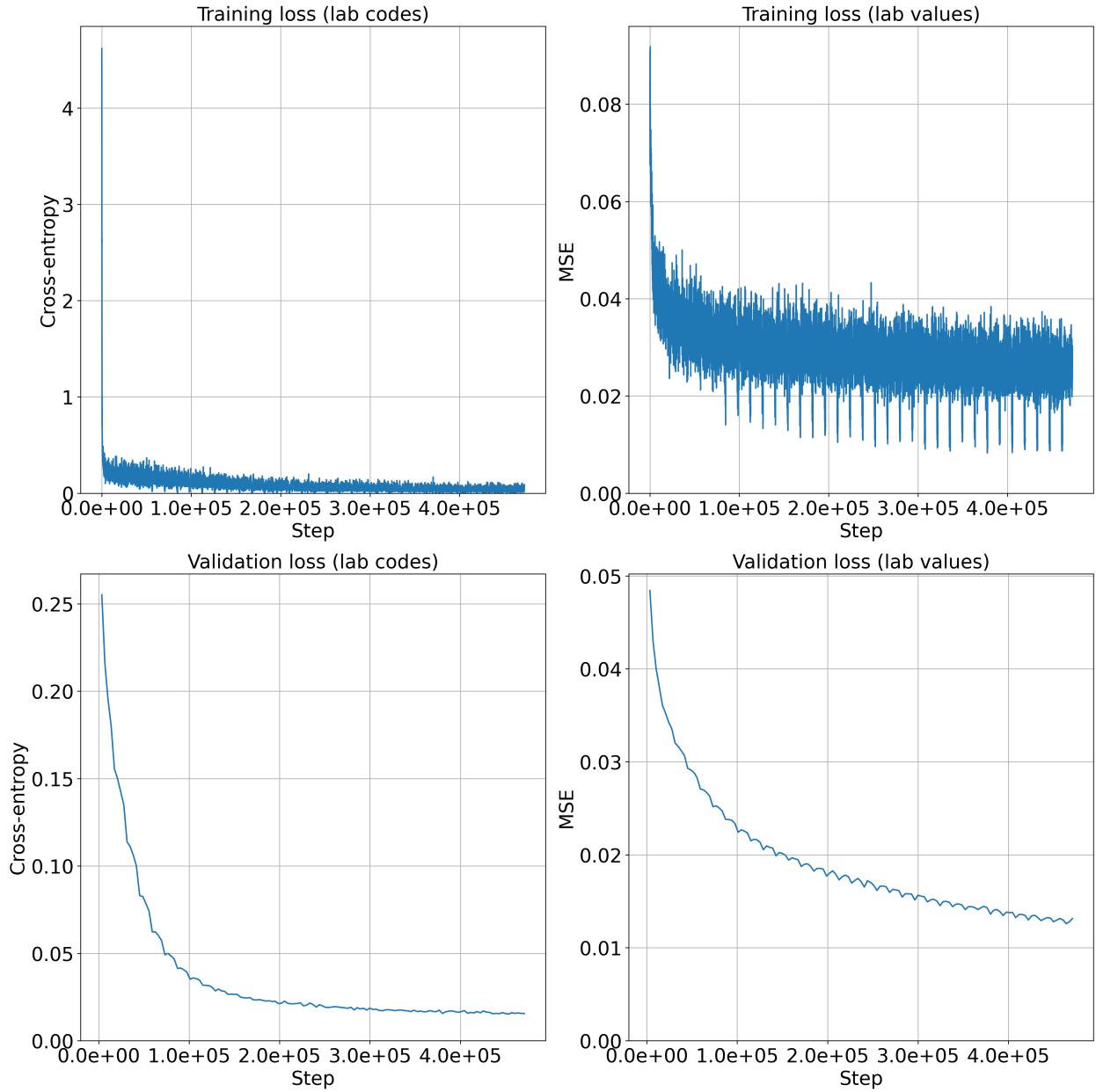


Figure 6: Pre-training loss curves for LABRADOR, the continuous Transformer. **Top:** The training loss for both the lab code (cross-entropy) and lab value (MSE) prediction heads. **Bottom:** The validation cross-entropy and MSE.

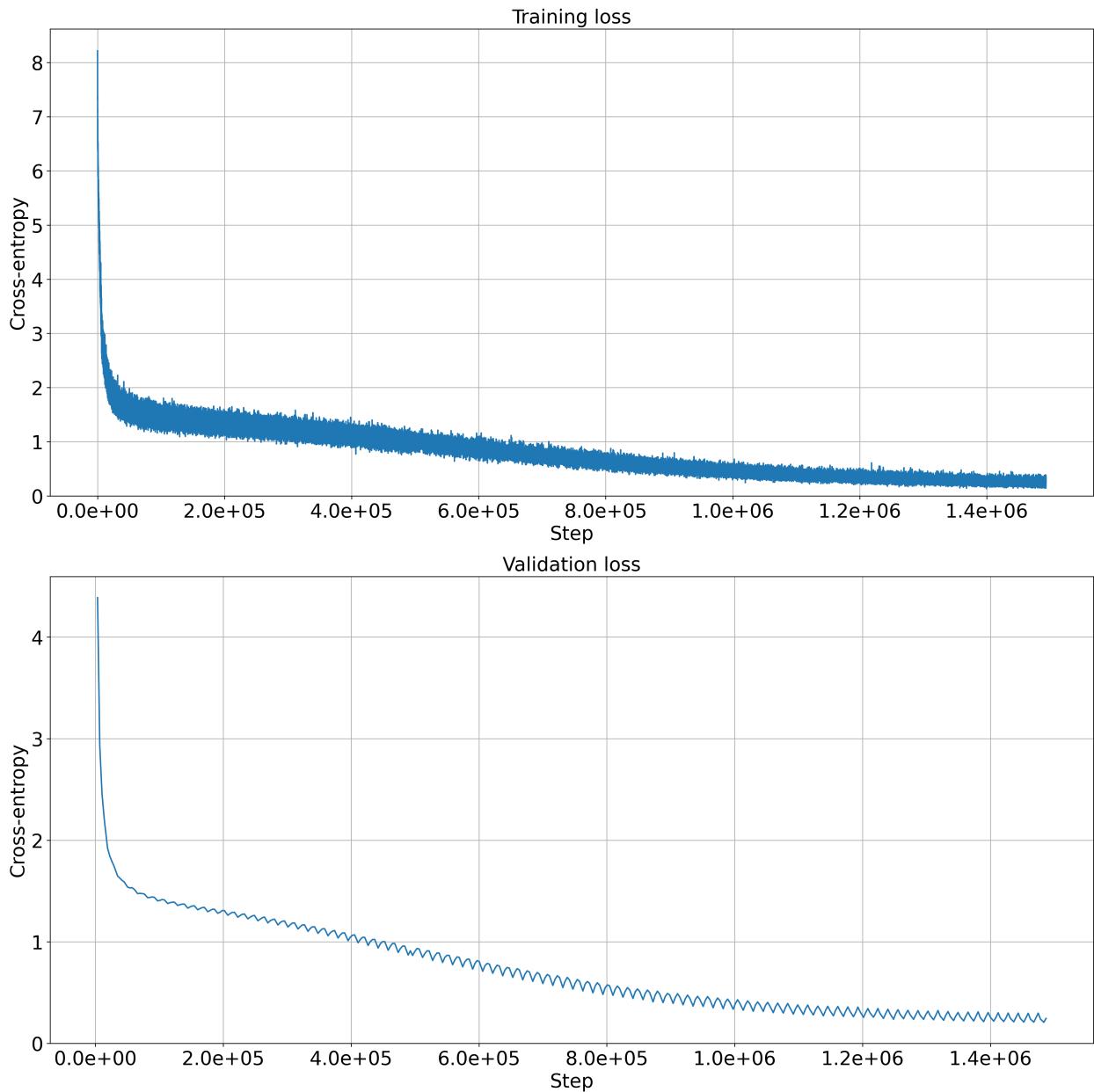


Figure 7: Pre-training loss curves for BERT. **Top:** The training cross-entropy. **Bottom:** The validation cross-entropy.

## Appendix H. Lab value correlation plot

Figure 8 is a heatmap of the Pearson correlation coefficients between each pair of the 70 most frequently ordered lab tests in MIMIC-IV. Labs are grouped by their panel membership and black lines separate these panels. For example, the top left box contains the correlation coefficients within the Basic Metabolic Panel (BMP), whereas the box at row 6, column 3 represents the inter-panel correlation between lipid tests and blood gas tests. This allows us to compare intra-panel (main diagonal) to inter-panel correlation structure (off-diagonal).

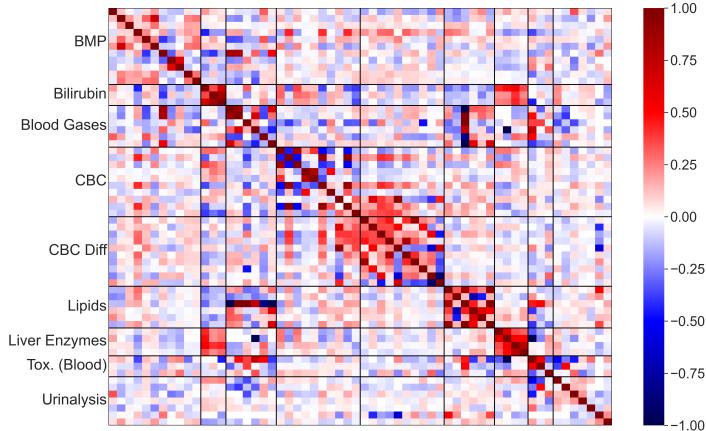


Figure 8: Lab value correlation matrix. A  $70 \times 70$  matrix of Pearson correlation coefficients between each pair of the 70 most frequently ordered lab tests in MIMIC-IV. The 70 tests are grouped by the panel that they are most frequently ordered with (see Appendix I for panel definitions). Black lines divide the matrix into panel-by-panel correlation blocks.

## Appendix I. Lab panel definitions

Table 4: Definition of the MIMIC-IV bilirubin panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
50885	Bilirubin, Total	Blood	Chemistry	1975-2
50883	Bilirubin, Direct	Blood	Chemistry	1968-7
50884	Bilirubin, Indirect	Blood	Chemistry	1971-1

Table 5: Definition of the MIMIC-IV blood gas panel. These tests do not have LOINC codes.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category
50820	pH	Blood	Blood Gas
50821	pO2	Blood	Blood Gas
50802	Base Excess	Blood	Blood Gas
50804	Calculated Total CO2	Blood	Blood Gas
50818	pCO2	Blood	Blood Gas
50813	Lactate	Blood	Blood Gas

Table 6: Definition of the MIMIC-IV basic metabolic panel (BMP).

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
50912	Creatinine	Blood	Chemistry	2160-0
51006	Urea Nitrogen	Blood	Chemistry	3094-0
50971	Potassium	Blood	Chemistry	2823-3
50983	Sodium	Blood	Chemistry	2951-2
50902	Chloride	Blood	Chemistry	2075-0
50882	Bicarbonate	Blood	Chemistry	1963-8
50868	Anion Gap	Blood	Chemistry	1863-0
50931	Glucose	Blood	Chemistry	6777-7
50893	Calcium, Total	Blood	Chemistry	2000-8
50960	Magnesium	Blood	Chemistry	2601-3
50970	Phosphate	Blood	Chemistry	2777-1

Table 7: Definition of the MIMIC-IV complete blood count (CBC) panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
51221	Hematocrit	Blood	Hematology	4544-3
51265	Platelet Count	Blood	Hematology	777-3
51222	Hemoglobin	Blood	Hematology	718-7
51301	White Blood Cells	Blood	Hematology	804-5
51249	MCHC	Blood	Hematology	786-4
51279	Red Blood Cells	Blood	Hematology	789-8
51250	MCV	Blood	Hematology	787-2
51248	MCH	Blood	Hematology	785-6
51277	RDW	Blood	Hematology	788-0
52172	RDW-SD	Blood	Hematology	N/A

Table 8: Definition of the MIMIC-IV complete blood count with differential panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
51256	Neutrophils	Blood	Hematology	761-7
51244	Lymphocytes	Blood	Hematology	731-0
51254	Monocytes	Blood	Hematology	742-7
51146	Basophils	Blood	Hematology	704-7
51200	Eosinophils	Blood	Hematology	711-2
52075	Absolute Neutrophil Count	Blood	Hematology	751-8
52073	Absolute Eosinophil Count	Blood	Hematology	711-2
52074	Absolute Monocyte Count	Blood	Hematology	742-7
51133	Absolute Lymphocyte Count	Blood	Hematology	731-0
52069	Absolute Basophil Count	Blood	Hematology	704-7

Table 9: Definition of the MIMIC-IV lipid panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
51000	Triglycerides	Blood	Chemistry	1644-4
50907	Cholesterol, Total	Blood	Chemistry	2093-3
50904	Cholesterol, HDL	Blood	Chemistry	2086-7
50903	Cholesterol Ratio (Total/HDL)	Blood	Chemistry	9322-9
50905	Cholesterol, LDL, Calculated	Blood	Chemistry	2090-9
50906	Cholesterol, LDL, Measured	Blood	Chemistry	N/A

Table 10: Definition of the MIMIC-IV liver function test (LFT) panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
50861	Alanine Aminotransferase (ALT)	Blood	Chemistry	1742-6
50878	Aspartate Aminotransferase (AST)	Blood	Chemistry	1920-8
50863	Alkaline Phosphatase	Blood	Chemistry	6768-6
50927	Gamma Glutamyltransferase	Blood	Chemistry	2324-2

Table 11: Definition of the MIMIC-IV toxicology (blood) panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
50922	Ethanol	Blood	Chemistry	5642-4
50856	Acetaminophen	Blood	Chemistry	3297-9
50981	Salicylate	Blood	Chemistry	4023-8

Table 12: Definition of the MIMIC-IV toxicology (urine) panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
51092	Opiate Screen, Urine	Urine	Chemistry	N/A
51079	Cocaine, Urine	Urine	Chemistry	N/A
51071	Amphetamine Screen, Urine	Urine	Chemistry	N/A
51090	Methadone, Urine	Urine	Chemistry	N/A
51074	Barbiturate Screen, Urine	Urine	Chemistry	3377-9
51989	Oxycodone	Urine	Chemistry	N/A
51089	Marijuana	Urine	Chemistry	N/A

Table 13: Definition of the MIMIC-IV urinalysis panel.

MIMIC-IV ItemID	Test Name	MIMIC-IV Test Fluid	MIMIC-IV Test Category	LOINC Code
51506	Urine Appearance	Urine	Hematology	5767-9
51498	Specific Gravity	Urine	Hematology	5811-5
51491	pH	Urine	Hematology	5803-2
51478	Glucose	Urine	Hematology	5792-7
51492	Protein	Urine	Hematology	5804-0
51484	Ketone	Urine	Hematology	5797-6
51464	Bilirubin	Urine	Hematology	5770-3
51487	Nitrite	Urine	Hematology	5802-4
51514	Urobilinogen	Urine	Hematology	5818-0
51087	Length of Urine Collection	Urine	Chemistry	N/A