

# A Multi-target Bayesian Transformer Framework for Predicting Cardiovascular Disease Biomarkers during Pandemics

Trusting Inekwe      Emmanuel Agu  
*Computer Science Dept*  
*Worcester Polytechnic Institute*  
 Worcester, MA 01609  
 {toinekwe, emmanuel}@wpi.edu

Winnie Mkandawire      Andres Colubri  
*Genomics and Computational Biology*  
*University of Massachusetts Chan Medical School*  
 Worcester, MA 01605  
 {Winnie.Mkandawire, Andres.Colubri}@umassmed.edu

**Abstract**—The COVID-19 pandemic disrupted healthcare systems worldwide, disproportionately impacting individuals with chronic conditions such as cardiovascular disease (CVD). These disruptions—through delayed care and behavioral changes—affected key CVD biomarkers, including LDL cholesterol (LDL-C), HbA1c, BMI, and systolic blood pressure (SysBP). Accurate modeling of these changes is crucial for predicting disease progression and guiding preventive care. However, prior work has not addressed multi-target prediction of CVD biomarker from Electronic Health Records (EHRs) using machine learning (ML), while jointly capturing biomarker interdependencies, temporal patterns, and predictive uncertainty.

In this paper, we propose MBT-CB, a Multi-target Bayesian Transformer (MBT) with pre-trained BERT-based transformer framework to jointly predict LDL-C, HbA1c, BMI and SysBP CVD biomarkers from EHR data. The model leverages Bayesian Variational Inference to estimate uncertainties, embeddings to capture temporal relationships and a DeepMTR model to capture biomarker inter-relationships. We evaluate MBT-CT on retrospective EHR data from 3,390 CVD patient records (304 unique patients) in Central Massachusetts during the Covid-19 pandemic. MBT-CB outperformed a comprehensive set of baselines including other BERT-based ML models, achieving an MAE of 0.00887, RMSE of 0.0135 and MSE of 0.00027, while effectively capturing data and model uncertainty, patient biomarker inter-relationships, and temporal dynamics via its attention and embedding mechanisms. MBT-CB’s superior performance highlights its potential to improve CVD biomarker prediction and support clinical decision-making during pandemics.

**Index Terms**—bayesian neural networks, transformers, uncertainty, variational inference, multi-target regression

## I. INTRODUCTION

**Motivation:** The COVID-19 pandemic caused unprecedented disruptions to healthcare systems worldwide, disproportionately affecting individuals with chronic illnesses such as cardiovascular diseases (CVDs) [1]. These disruptions included reduced access to medical services [2], delays in routine preventive care [3], and

pandemic-induced lifestyle changes [4]. Such interruptions had a profound impact on patient outcomes, particularly affecting critical biomarkers such as LDL cholesterol (LDL-C), Glycated hemoglobin (HbA1c), BMI and Systolic Blood Pressure (SysBP) [5]. The ability to model and predict such pandemic-related biomarker changes can enhance preventive care strategies and early detection of diseases.

**Challenges:** *EHR data has temporal structure and can be recorded at irregular intervals*, introducing uncertainty [6], which complicates predictive modeling. Irregular visits and missing or erroneous entries lead to *aleatoric* (data noise) and *epistemic* (model-related) uncertainty. Estimating these is challenging in deep learning and requires approaches such as *Bayesian inference* [7], which is rarely integrated into standard attention models [8]. Additionally, *CVD biomarkers—such as HbA1c and LDL-C—often exhibit heteroscedasticity*, with non-constant variance (Appendix Figures 11, 12). Addressing this requires models that learn across noisy outputs. *CVD biomarkers are also interdependent* [9], reflecting shared physiological and clinical factors. Thus, models that jointly model temporal dynamics, output dependencies, and uncertainty are critical for reliable prediction during crisis periods like pandemics.

**Limitations of prior research:** While prior studies reported the pandemic’s adverse effects on CVD patients [4], EHR-based models [10] overlooked three key challenges. First, they predicted biomarkers independently rather than jointly, ignoring interdependencies and clinical coherence [11]. Second, they lacked temporal modeling, missing trends and trajectories [12]. Third, they failed to quantify uncertainty in data and predictions [13], limiting clinical reliability [14].

**Our approach:** We propose MBT-CB NN, a novel transformer-based architecture for modeling temporal relationships in patient EHR data, multi-target biomarker

prediction, and uncertainty estimation. Building on the ClinicalBERT model [15], MBT-CB integrates Bayesian Variational Inference (BVI) to capture aleatoric and epistemic uncertainty [16], a Deep Multi-Target Regression (DeepMTR) layer [17] to learn shared and target-specific representations, and positional embeddings [18] for temporal encoding. Trained on 3,390 EHRs (304 unique patients) from Central Massachusetts between Jan 2019–June 2021, MBT-CB is tailored for CVD biomarker prediction during the pandemic.

**Novelty of Work:** While transformers have been used for CVD detection [19], multi-target prediction [11], and uncertainty-aware COVID-19 models [20], no prior work combines BVI, multi-target prediction, and temporal attention in a unified transformer framework for longitudinal CVD biomarker prediction. Our model delivers robust, confidence/uncertainty-aware predictions suitable for high-stakes clinical use.

#### Our Contributions:

- 1) *Bayesian Multi-Target Transformer:* We propose MBT-CB, a novel transformer architecture with BVI that jointly predicts LDL-C, HbA1c, BMI, and SysBP from EHR data while modeling uncertainty/confidence, temporal dynamics, and biomarker interdependencies.
- 2) *Systematic Rigorous evaluation:* We benchmark MBT-CB against other models. MBT-CB outperforms traditional ML, DL baselines, and pretrained transformers, achieving MAE 0.00887, RMSE 0.0135, and MSE 0.00027, showing high accuracy and generalizability under pandemic-related variability.
- 3) *Explainability via attention and uncertainty visualization:* MBT-CB reveals dependencies such as HbA1c–BMI, highlights aleatoric and epistemic uncertainty, supporting transparent, risk-aware predictions.

## II. RELATED WORK

### A. Transformers for COVID-19 and CVD prediction

**CVD detection using Transformers on EHR & Image datasets:** Antikainen et al. [19] compared BERT and XLNet for mortality prediction in 23,000 cardiac patients using EHR time series data. XLNet slightly outperformed BERT, with an AUC of 76.0% vs. 75.5%, and showed a 9.8% higher recall, indicating better identification of at-risk patients. Similarly, Kilimci et al [21] explored deep vision transformers (Google-ViT, Microsoft-Beit, and Swin-Tiny) for heart disease detection from ECG images, achieving high accuracies Swin-Tiny (95.9%), Microsoft-Beit (95.5%) and Google-ViT (94.3%).

**Transformers and Multi-target Prediction:** Poulain et al. [11], enhanced BERT with a Deep Multi-Target Regression (DeepMTR) module to predict 11 modifiable CVD risk factors. The model significantly improved RMSE (by 0.053) for targets with over 80% missing

data and reduced MAE by an average of 12.6% across all targets compared to baselines.

**COVID detection using transformers and Bayesian NN):** Chen et al., [20] applied a hybrid transformer model (CBAM, ViT, Swin Transformer) to chest X-rays for COVID-19 detection, adding a Bayesian NN layer to capture weight uncertainty. Unlike their approach that adds Bayesian methods externally, our study embeds Bayesian inference directly into the attention mechanism to quantify uncertainty in CVD prediction from EHR data.

**Transformers and Bayesian approach:** Sankararaman [22] investigated a Bayesian transformer approach by treating learnable parameters as random distributions. In contrast, our approach selectively applies Bayesian methods to attention weights, treating them as random variables while keeping other parameters constant.

**Research gap addressed** Although Bayesian and multi-target transformers have been applied to clinical prediction, no prior work has unified multi-target learning, temporal modeling, and uncertainty estimation within a single framework to predict the impact of COVID-19 on key CVD biomarkers (HbA1c, LDL, BMI, and SysBP). Addressing these gaps is essential for reliable deployment in high-stakes healthcare settings affected by pandemic-related care disruptions.

## III. METHODOLOGY

### A. Overview of Methodology

Figure 1 outlines MBT-CB, our framework for predicting CVD biomarker trajectories during the COVID-19 pandemic. We utilized EHR data from 304 CVD patients treated at UMass Chan and Memorial Hospital, spanning pre-pandemic (Jan 2018–Dec 2019) and pandemic (June 2020–June 2021) periods. The dataset included demographics (age, gender, income, race) and biomarkers (HbA1c, LDL, BMI, SysBP). To predict biomarker values at the first pandemic-era visit, biomarker data were encoded using ClinicalBERT representations. Variational inference was applied to attention weights, with embeddings and outputs processed by the DeepMTR FFNN model.

**Data Pre-processing** First, patients with CVD were identified using ICD-10 codes (I63, I67, I48, I73, I25, I65, I50). Records with nulls, outliers, implausible values, or no hospital visits during the COVID-19 pandemic were removed, reducing the dataset from 80,917 to 3,390 records across 304 unique patients. Visit frequency distribution is shown in Appendix Figure 5.

**Feature Extraction and Engineering** HbA1c, LDL-C, BMI, and SysBP were selected as target variables. Socioeconomic status—estimated via median household income by zip code [23]—was included due to its impact on health outcomes [24]. Categorical features

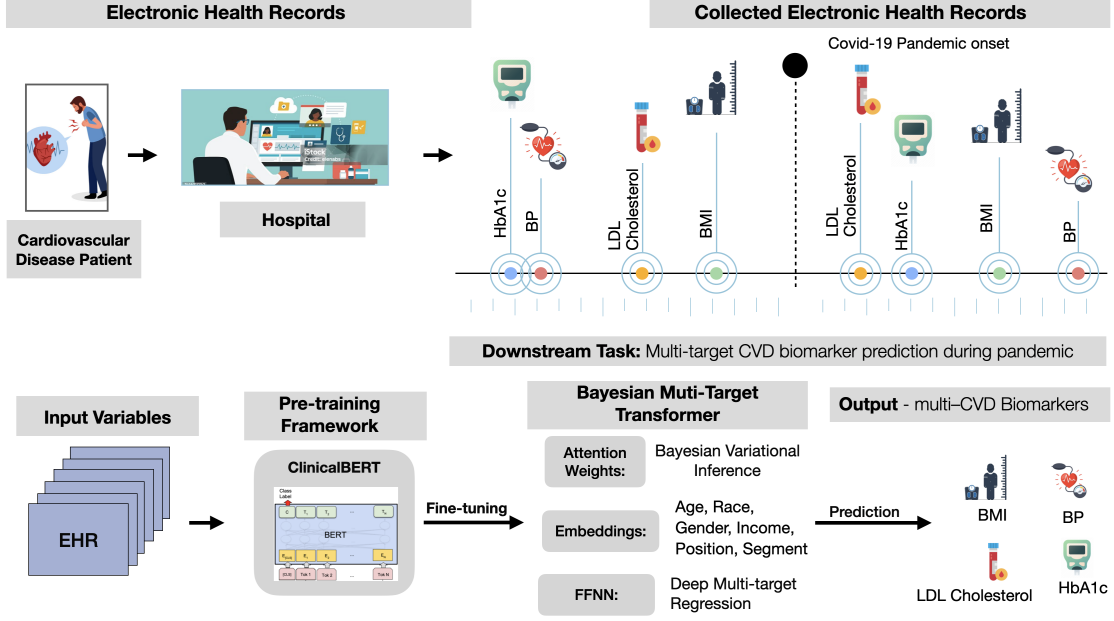


Fig. 1: Overview of Methodology

(gender, race, income) were one-hot encoded, numerical features were normalized, and skewed distributions log-transformed.

The dataset comprised four biomarkers: HbA1c (4% to 14%), LDL-C (20–370 mg/dL), systolic blood pressure(SysBP) (84–196 mmHg), and BMI (15–100). Demographics: income level (38.2% upper-middle, 36.8% middle, 24.3% lower-middle, 0.7% upper), age (45–96 years), race (88.2% White, 7.2% Black, 4.6% Asian), and gender (61.5% male, 38.5% female). These variables capture both physiological and social determinants relevant to cardiovascular health. Data was split 60/20/20 train/validation/test at patient level, preventing leakage and preserving longitudinal visit structure.

**MBT-CB Implementation for CVD Prediction** details of the transformer model are in Appendix 1.

#### 1) Input sequence and patient representation:

As shown in Appendix Figure 3, each patient’s EHR is structured as a chronological sequence of clinical visits, forming a time series of four key CVD biomarkers: SysBP, BMI, HbA1c, and LDL-C. For patient  $p$ , the visit sequence is denoted as  $\mathbf{V}_p = \{v_p^1, v_p^2, \dots, v_p^{n_p}\}$ , where  $n_p$  is the number of visits. Each visit  $v_p^i$  is transformed into a structured sentence:  $T_p^i = \text{“sys: } x_{\text{sys}}^i; \text{bmi: } x_{\text{bmi}}^i; \text{hba1c: } x_{\text{HbA1c}}^i; \text{chol: } x_{\text{chol}}^i\text{”}$  where each  $x^i$  is a normalized scalar biomarker value. These text sequences are input into a pre-trained transformer for contextual representation learning.

#### 2) Transfer Learning:

**Pre-training:** To leverage existing domain-specific knowledge, a transfer learning strategy using Clinical-

BERT [15], [25] is adopted. ClinicalBERT is a language model pre-trained on over 3million structured EHR records using the masked language modeling objective on over 1.2billion words of clinical text.

**Fine-Tuning:** Each input sentence  $T_p^i$  is tokenized using the ClinicalBERT’s tokenizer and passed through the encoder to produce a contextualized embedding from the [CLS] token:  $\mathbf{h}_p^i = \text{Model}_{\text{pre}}(T_p^i)_{[\text{CLS}]} \in \mathbb{R}^d$ . This contextual embedding is then projected into a task-specific representation using a linear transformation:  $\mathbf{z}_p^i = \mathbf{W}_{\text{proj}} \cdot \mathbf{h}_p^i + \mathbf{b}_{\text{proj}}$ . To enhance these representations, we append several auxiliary embeddings:

- **Positional embeddings**  $\text{pos}_p = \{0, 1, \dots, n_p - 1\}$  to capture visit order,
- **Segment embeddings**  $\text{seg}_p = \{s_p^1, \dots, s_p^{n_p}\}$  where  $s_p^i \in \{0, 1\}$  indicates whether the visit occurred before or after the onset of COVID-19,
- **Demographic identifiers:** gender  $g_p$ , race  $r_p$ , and income class  $i_p$ , encoded as categorical indices.

To enable the transformer model to jointly learn from temporal visit patterns, biomarker trajectories, and relevant demographic context, each patient  $p$ ’s complete input is represented as  $\text{Input}_p = (\{T_p^i\}_{i=1}^{n_p}, \text{pos}_p, \text{seg}_p, g_p, r_p, i_p)$ .

3) **Variational Self-Attention Mechanism:** To capture epistemic uncertainty, the fixed projection weights are replaced with a Gaussian distribution parameterized by a learnable mean  $\mu$  and log standard deviation  $\sigma$ . Weights are sampled during each forward pass:  $W = \mu + \exp(\log \sigma) \cdot \varepsilon$ ,  $\varepsilon \sim \mathcal{N}(0, I)$ , Keys are transformed as

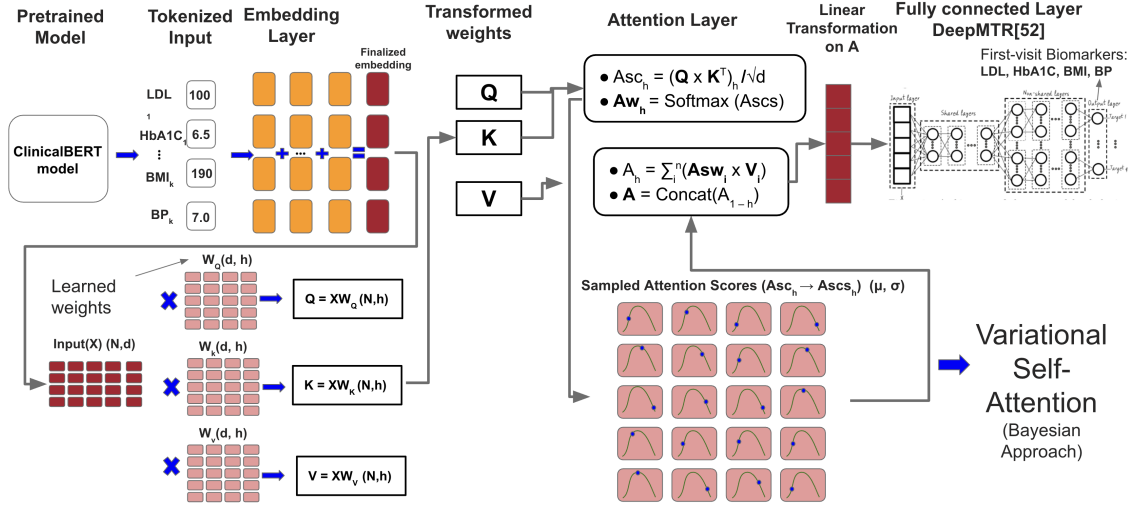


Fig. 2: Our proposed MBT-CB framework based on a Transformer with Variational Self Attention. Patient’s EHR biomarker values for  $k$  visits (1 to  $k$ ) are passed as input variables to the MBT-CB model. Prediction is on the 1<sup>st</sup> visit, where  $k < n$ . DeepMTR image from [17]

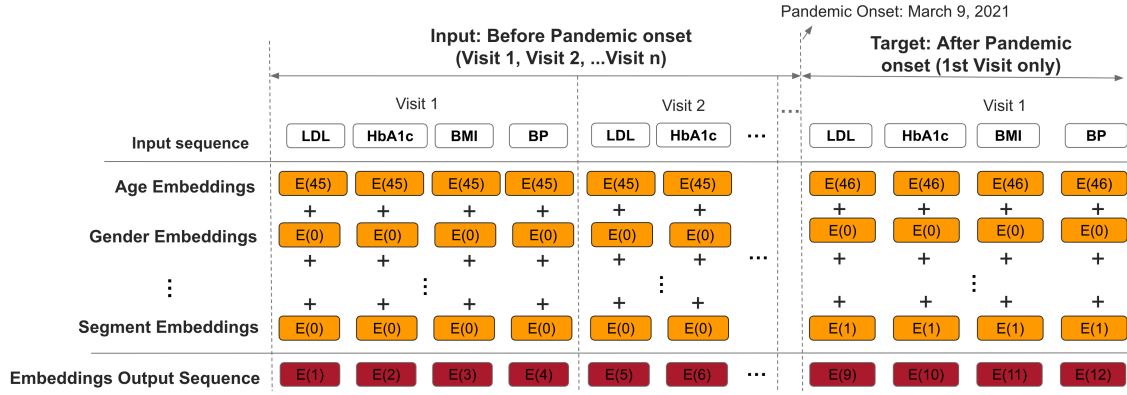


Fig. 3: Modification of a patient’s EHR record. Each row represents a set of chronological biomarker values from an individual clinical visit, suitable for transformer-based modeling.

$k' = kW$ , and attention scores computed via:  $S = \frac{q \cdot k'^T}{\sqrt{d}}$ . More details in Appendix Sections 2 and 3 respectively

4) **Bayesian Prior and Objective Function:** A standard normal prior is placed on  $W$ , regularized via KL divergence. The training loss is  $\mathcal{L}_{\text{total}} = \text{MSE}(y, \hat{y}) + \lambda \cdot \text{KL}(q(W) \parallel p(W))$  where  $\text{KL} = \sum [-\log \sigma + \frac{1}{2} (\sigma^2 + \mu^2 - 1)]$ .

5) **DeepMTR:** To capture both shared and biomarker-specific patterns, the attention output feeds into DeepMTR, which includes shared layers and target-specific heads for each biomarker.

**The MBT-CB Model Pipeline** restructures patient visit histories into biomarker sentences, tokenized to generate contextual embeddings using ClinicalBERT. Positional, segment, and demographic embeddings are added to enrich the representation, and Variational Self-Attention applied to model epistemic uncertainty. The

resultant representations are input to a DeepMTR head for multi-biomarker prediction, with uncertainty quantified via multiple stochastic forward passes.

#### IV. EVALUATION

**Evaluation metrics:** used to assess MBT-CB’s performance on a test set included MAE, MSE and RMSE across multiple CVD biomarkers. The mathematical expressions of these metrics are shown below:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i| \quad (1)$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (2)$$



$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2} \quad (3)$$

**Baseline Models** were selected from the following categories: Bio\_ClinicalBERT, MedBERT, BERT BERT-based models (default, and also integrated with BVI and DeepMTR), a FFNN and a Linear regression. All models were configured for multi-target prediction and are described briefly.

- 1) *Bio\_ClinicalBERT*: derived from BioBERT [26], was trained on clinical notes from MIMIC III EHR dataset [27] and has approximately 880 million words.
- 2) *MedBERT* [28]: derived from Bio\_ClinicalBERT, was pre-trained on over 57.46 million tokens from multiple medical datasets including the N2C2, BioNLP project4 and Wikipedia corpora.
- 3) *BERT* [29]: is a language model pre-trained on a large corpus of 11,038 unpublished books and Wikipedia in English. BERT-base-uncased, trained on 110 million words, was utilized.
- 4) *Traditional ML models*: including a multi-target FFNN and Linear Regression models.

In the results section, performance is reported across all four pre-trained model variants. The best-performing variant is selected for final evaluation and uncertainty quantification.

An *ablation study* was conducted study to evaluate the contributions of the DeepMTR and BVI modules to model performance.

## V. RESULTS

### A. MBT-CB prediction performance and comparisons

As shown in Table I, the complete (proposed) MBT-CB model outperformed all baselines in multi-target prediction of CVD biomarkers during the pandemic with a lowest mean value for MAE (0.00887), RMSE (0.0135) and MSE (0.00027).

**1. Comparisons to other integrated pre-trained transformer models (bayesian and DeepMTR):** MBT-CB outperformed MedBERT, Bio\_ClinicalBERT and BERT integrated pretrained transformer, achieving an MAE of 0.0051, 0.0071, 0.0173, and 0.0059 for SysBp, BMI, HbA1c, and LDL-C respectively. Corresponding RMSE values were 0.0065, 0.0094, 0.0298, and 0.0082, while MSE values remained consistently low across all targets (SysBp: 0.000042, BMI: 0.000088, HbA1c: 0.000890, LDL-C: 0.000067), indicating minimal error variance and strong prediction stability. Bio\_ClinicalBERT was the next best performing model with a mean MAE of 0.0138, MSE(0.00044) and RMSE (0.0185). MedBERT followed with a mean MAE of 0.1450, MSE(0.00048) and RMSE (0.0199). Lastly, with

a mean MAE of 0.0169, MSE (0.00055) and RMSE (0.0217), BERT exhibited the lowest performance in this group. These results reveal that incorporating Bayesian inference and the DeepMTR architecture into pre-trained clinical language models enhances accuracy, particularly in predicting multi-target biomarker trajectories.

**2. Comparisons to non-integrated transformer-based baselines:** While the non-integrated pretrained models performed reasonably well across all biomarkers, they were still outperformed by our proposed MBT-CB model. This suggests that there is an added benefit in using a fully end-to-end Bayesian Transformer architecture with integrated multi-target regression capabilities to enhance predictive accuracy and generalization across multiple CVD biomarkers.

**3. Comparisons with non-transformer models:** such as multi-target linear regression and multi-target FFNN revealed that they were unable to generalize, with MAE values exceeding 0.5 for most targets and MSE values as high as 1.02. This highlights the limitations of non-transformer approaches in handling complex, longitudinal EHR data with interrelated targets.

**4. Ablative analyses:** The full MBT-CB model outperformed all ablation variants across all biomarkers. Removing BVI led to higher MAEs, except for LDL-C, which saw slight improvements (MAE: 0.0039, RMSE: 0.00576, MSE: 0.000033). Still, the complete model achieved the best overall performance, highlighting the value of Bayesian attention. Excluding the DeepMTR head resulted in a substantial performance drop (mean MAE: 0.0245, RMSE: 0.033, MSE: 0.0011), confirming its critical role in modeling biomarker interdependencies.

**5. Training dynamics and attention visualization:** As shown in Appendix, Figure 9, the MBT-CB model’s training and validation losses steadily decreased over 50 epochs without overfitting. Appendix Figure 10 visualizes the attention pattern of MBT-CB’s Bayesian attention mechanism. The attention scores display token-level dependencies across visits and biomarkers, suggesting that MBT-CB learns nuanced patterns, such as the relationship between HbA1c and BMI trajectories, structured, temporal input.

**Summary of findings:** The MBT-CB architecture—with BVI self-attention and a multi-target regression head—accurately predicted multiple CVD biomarker trajectories and substantially outperformed all baselines, while supporting interpretation via attention patterns. Details on hardware and computational resources used are in Appendix A1

### B. Interpreting Uncertainty in Biomarker Predictions

Figure 4 presents MBT-CB’s predictions with ground truth and uncertainty bands for four CVD biomarkers. SysBP and LDL-C show narrow intervals overall, with

TABLE I: MAE, RMSE, and MSE Performance of Various Models for Multi-target Biomarker Prediction (all values in decimal form, normalized to exponential level  $1e-2$ )

Model	MAE					RMSE					MSE				
	SysBp	BMI	HbA1c	LDL	Mean	SysBp	BMI	HbA1c	LDL	Mean	SysBp	BMI	HbA1c	LDL	Mean
<b>Multi-target</b>															
FFNN	73.8	30.3	51.8	66.0	55.4	101.0	41.8	70.3	83.0	74.1	103.0	17.5	49.4	68.9	59.6
Linear Regression	59.9	30.2	55.7	73.8	55.0	81.5	42.2	81.6	89.9	73.8	66.4	17.8	66.6	80.8	58.0
<b>Standard &amp; Multi-target FFNN</b>															
MedBERT	0.537	0.977	3.620	0.653	1.450	0.725	1.280	4.940	0.876	1.960	0.0053	0.0163	0.2450	0.0077	0.0685
BERT	0.583	1.030	3.090	0.623	1.380	0.894	1.430	4.620	1.090	2.060	0.0080	0.0204	0.2130	0.0118	0.0634
Bio_ClinicalBERT	0.630	1.290	3.140	0.477	1.380	1.120	1.730	4.630	0.778	2.060	0.0125	0.0301	0.2140	0.0060	0.0656
<b>Bayesian &amp; DeepMTR</b>															
MedBERT	0.906	1.330	2.660	0.918	1.450	1.230	1.740	3.630	1.170	1.990	0.0151	0.0301	0.1320	0.0137	0.0477
BERT	1.140	1.510	2.840	1.060	1.690	1.410	1.920	3.800	1.350	2.170	0.0200	0.0367	0.1450	0.0182	0.0549
Bio_ClinicalBERT	1.020	1.160	2.540	0.792	1.380	1.350	1.510	3.520	1.030	1.850	0.0181	0.0229	0.1240	0.0105	0.0439
<b>Bayesian + DeepMTR (Proposed)</b>															
ClinicalBERT	<b>0.511</b>	<b>0.714</b>	<b>1.730</b>	0.592	<b>0.887</b>	<b>0.652</b>	<b>0.936</b>	<b>2.980</b>	0.818	<b>1.350</b>	<b>0.0042</b>	<b>0.0088</b>	<b>0.0890</b>	0.0067	<b>0.0272</b>
<b>Ablation Study</b>															
w/o Bayesian	0.580	1.000	3.170	<b>0.390</b>	1.280	0.769	1.330	4.590	<b>0.576</b>	1.810	0.0059	0.0176	0.2100	<b>0.0033</b>	0.0593
w/o DeepMTR	2.000	2.290	3.470	2.050	2.450	2.610	2.970	4.490	2.910	3.300	0.0682	0.0885	0.2020	0.0844	0.1110

Note: These normalized errors can be converted to raw clinical units using the biomarker ranges in Section III-A and Appendix Table II. They all correspond to values within clinically meaningful thresholds.

localized spikes in epistemic uncertainty where predictions deviate—indicating model uncertainty in less familiar regions. BMI and HbA1c display broader, aleatoric-dominated bands, suggesting higher intrinsic variability or measurement noise.

## VI. DISCUSSION

**MBT-CB outperforms other pre-trained transformer models** enhanced with Bayesian and DeepMTR components. ClinicalBERT’s superior performance may be attributed to its richer pre-training on over 1.2 billion clinical tokens and 3M EHRs, compared to smaller corpora used by Bio\_ClinicalBERT (880M), BERT (110M), and MedBERT (57.5M). This broader clinical language exposure improved generalization to our heteroscedastic dataset, aligning with findings in [30].

**Ablation results highlight the value of MBT-CB’s components**, particularly Bayesian self-attention and DeepMTR. Bayesian attention enhanced generalization by modeling epistemic uncertainty, improving robustness to sparse or noisy data. DeepMTR captured shared and target-specific patterns across correlated biomarkers, enabling the model to learn biomarker interdependencies and improve predictive accuracy on complex EHR data.

**Strengths of the MBT-CB Architecture** MBT-CB integrates pretrained ClinicalBERT, Bayesian self-attention, and a DeepMTR architecture to address challenges in longitudinal EHR-based biomarker prediction. Benefiting from ClinicalBERT’s large(1.2 bil-

lion)biomedical and EHR pretraining, the model encodes temporally ordered biomarker prompts using visit-level tokenization, with positional and segment embeddings to track visit order and pandemic phases.

Bayesian self-attention (4 heads, latent dim 128) models attention weights as Gaussian distributions, with variational inference learning  $\mu$  and  $\log \sigma$ , and KL regularization ( $1 \times 10^{-4}$ ) to capture epistemic uncertainty and reduce overfitting on sparse, irregular data.

For multi-target prediction, DeepMTR employs a shared fully connected block ( $512 \rightarrow 128 \rightarrow 64$ ) and four output heads for SysBP, BMI, HbA1c, and LDL-C. Demographic, positional, and segment embeddings support personalized, temporal modeling. MBT-CB was trained for 50 epochs (batch size 4, weight decay 0.01) using Hugging Face’s Trainer.

It consistently outperformed baseline transformers, demonstrating the value of combining pretrained models with Bayesian reasoning and multi-target learning for robust, personalized clinical prediction. These findings underscore the potential of hybrid transformer-based models for advancing personalized prediction in real-world healthcare settings.

**Uncertainty-aware biomarker prediction during the pandemic** MBT-CB uses Bayesian self-attention with variational inference to generate interpretable, patient-level uncertainty estimates. Epistemic uncertainty spikes—most notably for LDL-C and SysBP—occur when predictions diverge from ground truth, reflecting

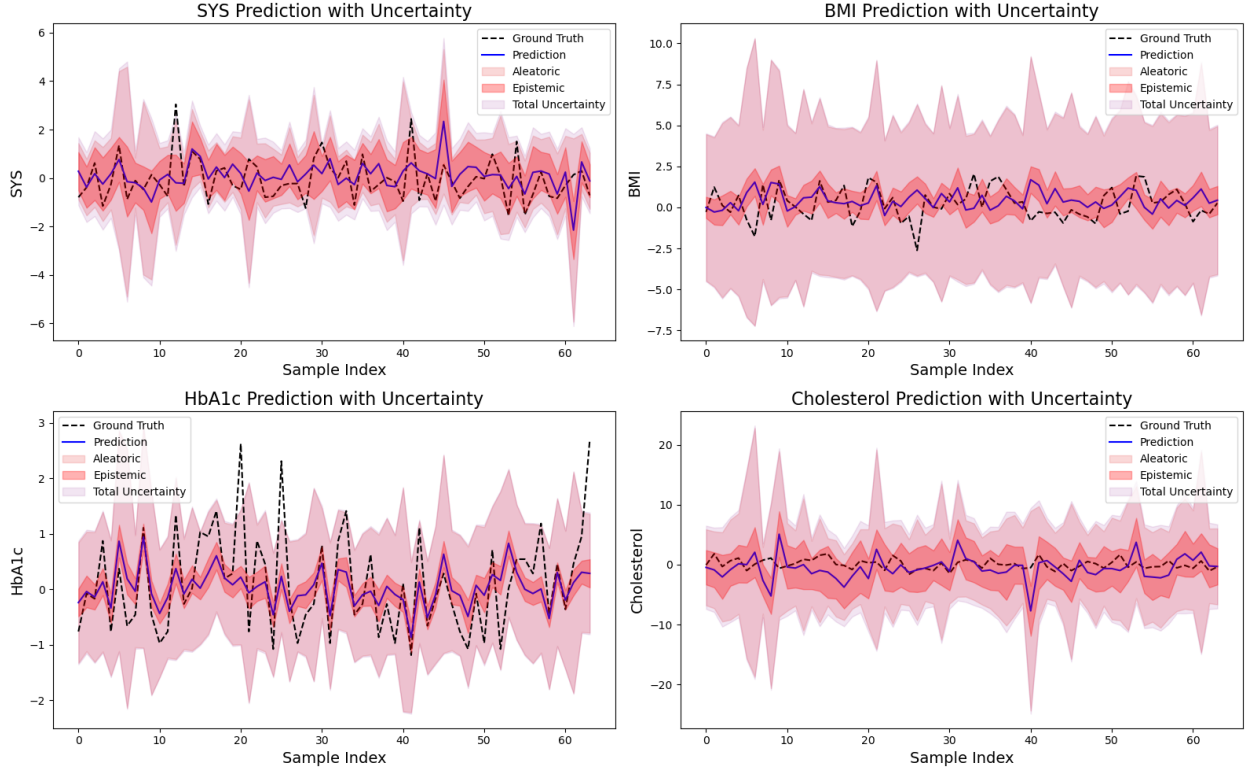


Fig. 4: Uncertainty Estimation of MBT-CB

sparse training or unmodeled factors such as care gaps.

During the pandemic, distinct biomarker-specific patterns emerged: SysBP and LDL-C remained mostly stable with occasional confidence drops, while BMI and HbA1c showed consistently broader, aleatoric-driven uncertainty—indicating variability in weight [31] and glycemic control. By disentangling and localizing uncertainty, MBT-CB enables risk-aware decisions amid care disruptions and data shifts in pandemic contexts.

**MBT-CB achieved strong, balanced performance across multiple biomarkers and generalized well.** It captured complex interactions between patient demographics, temporal visit data, and biomarker interdependencies. Its resilience to noise and data irregularities—especially under COVID-19-related disruptions—demonstrates its suitability for highly variable, real-world, clinical scenarios.

**Methodological Advancements** By integrating Bayesian self-attention and multi-target regression into a transformer, our approach captures non-linear, sequential biomarker relationships while quantifying both aleatoric and epistemic uncertainty. This enables per-sample confidence estimation, supporting more informed and risk-aware decisions—vital when data conditions change unexpectedly, as in pandemics.

**Clinical Implications and Healthcare Impact** Our

findings highlight the broader value of uncertainty-aware, multi-target prediction for CVD risk management beyond COVID-19. MBT-CB can enable proactive identification of at-risk patients. In low-resource or crisis settings, combining risk prediction with confidence estimates supports more targeted and timely interventions, advancing AI-driven clinical decision support.

## VII. LIMITATIONS AND FUTURE DIRECTIONS

The training data, drawn from two hospitals in Central Massachusetts with 89% White patients, limits generalizability of our findings. Future work will include multi-site, diverse patients. While Bayesian attention enhances uncertainty estimation, its stochastic nature complicates interpretation. Integrating SHAP values or counterfactual reasoning can improve explainability.

Future directions include: (1) scaling to larger, more heterogeneous EHR datasets; (2) replacing ClinicalBERT with a custom transformer trained on broader clinical corpora; (3) incorporating SHAP-based attribution to clarify feature contributions to COVID-related biomarker shifts; and (4) benchmarking MBT-CB against deep learning models auto-generated using Network Architecture Search (NAS), and CVD-specific models.

## VIII. CONCLUSION

We presented MBT-CB, a transformer-based model for predicting LDL-C, HbA1c, BMI, and SysBP biomarkers during the COVID-19 pandemic. Leveraging ClinicalBERT, Bayesian self-attention for uncertainty, multi-target regression and temporal patterns, MBT-CB outperformed all baselines. It produced clear uncertainty bands—narrow with epistemic spikes for SysBP and LDL-C, broader and aleatoric-dominated for BMI and HbA1c—highlighting its ability to distinguish model uncertainty from data noise. This enhances interpretability at the patient level. Future work includes real-time deployment, broader generalization, and SHAP-based explanations.

## REFERENCES

- [1] R. Verity *et al.*, “Estimates of severity of coronavirus disease 2019: a model-based analysis,” *Lancet Infectious Diseases*, vol. 20, no. 6, pp. 669–677, 2020.
- [2] T. C. Tsai *et al.*, “Association of community-level social vulnerability with us acute care hospital intensive care unit capacity during covid-19,” *Healthcare*, vol. 10, p. 100611, Mar. 2022.
- [3] N. Bilgin Dogan and E. Ozel, “The missing stems and lifestyle changes during covid-19 pandemic,” *APJPH*, vol. 33, no. 2-3, pp. 296–298, 2021.
- [4] Mattioli *et al.*, “Covid-19 pandemic: effects of quarantine on cardiovascular risk,” *EJCN*, vol. 74, pp. 852–855, Jun. 2020.
- [5] Aparisi *et al.*, “Low-density lipoprotein cholesterol levels are associated with poor clinical outcomes in covid-19,” *Nutrition, Metabolism, and Cardiovascular Diseases*, vol. 31, pp. 2619–2627, Aug. 2021.
- [6] C. Xiao, E. Choi, and J. Sun, “Opportunities & challenges in developing deep learning models using ehr data: a systematic review,” *JAMIA*, vol. 25, no. 10, pp. 1419–1428, 2018.
- [7] J. Liu *et al.*, “Simple and principled uncertainty estimation with deterministic deep learning via distance awareness,” *Adv. NeurIPS*, vol. 33, pp. 7498–7512, 2020.
- [8] E. Kostenok, D. Cherniavskii, and A. Zaytsev, “Uncertainty estimation of transformers’ predictions via topological analysis of the attention matrices,” *arXiv preprint arXiv:2308.11295*, 2023.
- [9] C.-H. Wu, W.-J. Yao, F.-H. Lu, J.-S. Wu, and C.-J. Chang, “Relationship between glycosylated hemoglobin, blood pressure, serum lipid profiles and body fat distribution in healthy chinese,” *Atherosclerosis*, vol. 137, no. 1, pp. 157–165, 1998.
- [10] T. Inekwe, W. Mkandawire, B. Wee, E. Agu, and A. Colubri, “Biomarker trajectory prediction and causal analysis of the impact of the covid-19 pandemic on cvd patients using machine learning,” in *Proc. IEEE/ACM CHASE*, pp. 1–12, 2024.
- [11] R. Poulain, M. Gupta, R. Foraker, and R. Beheshti, “Transformer-based multi-target regression on electronic health records for primordial prevention of cardiovascular disease,” in *Proc IEEE BIBM*, pp. 726–731, IEEE, 2021.
- [12] Y.-Q. Cai *et al.*, “Pitfalls in developing machine learning models for predicting cardiovascular diseases: Challenge and solutions,” *J. Medical Internet Research*, vol. 26, p. e47645, 2024.
- [13] J. Kim, “Bert: Pre-training of deep bidirectional transformers for language understanding,” 2022.
- [14] D. Sam *et al.*, “Bayesian neural networks with domain knowledge priors,” *arXiv preprint arXiv:2402.13410*, 2024.
- [15] X. Liu *et al.*, “A generalist medical language model for disease diagnosis assistance,” *Nature Medicine*, pp. 1–11, 2025.
- [16] P. Lauret, E. Fock, *et al.*, “Bayesian neural network approach to short time load forecasting,” *Energy conversion and mgmt.*, vol. 49, no. 5, pp. 1156–1166, 2008.
- [17] O. Reyes and S. Ventura, “Performing multi-target regression via a parameter sharing-based deep network,” *International journal of neural systems*, vol. 29, no. 09, p. 1950014, 2019.
- [18] A. Vaswani, “Attention is all you need,” *Advances in Neural Information Processing Systems*, 2017.
- [19] E. Antikainen *et al.*, “Transformers for cardiac patient mortality risk prediction from heterogeneous electronic health records,” *Scientific Reports*, vol. 13, no. 1, p. 3517, 2023.
- [20] H. Chen, J.-Y. Hsieh, H.-Y. Hsu, and Y.-F. Chang, “Covid-19 detection based on chest x-ray images using attention mechanism modules and weight uncertainty in bayesian neural networks,” in *Int’l Conf. IoT and Health*, pp. 104–115, Springer, 2023.
- [21] Z. H. Kilimci *et al.*, “Heart disease detection using vision-based transformer models from ecg images,” *arXiv preprint arXiv:2310.12630*, 2023.
- [22] K. A. Sankararaman *et al.*, “Bayesformer: Transformer with uncertainty estimation,” *arXiv preprint arXiv:2206.00826*, 2022.
- [23] UnitedStatesZipCodes, “Us zip codes,” 2023. [Online]. Available: <https://www.unitedstateszipcodes.org>. [Accessed: Dec. 13, 2023].
- [24] R. B. . o. Hawkins, “Socio-economic status & covid-19-related cases and fatalities,” *Public health*, vol. 189, pp. 129–134, 2020.
- [25] G. Wang *et al.*, “Optimized glycemic control of type 2 diabetes with reinforcement learning: a proof-of-concept trial,” *Nature Medicine*, vol. 29, no. 10, pp. 2633–2642, 2023.
- [26] J. Lee *et al.*, “Biobert: a pre-trained biomedical language representation model for biomedical text mining,” *Bioinformatics*, vol. 36, no. 4, pp. 1234–1240, 2020.
- [27] A. Johnson *et al.*, “Mimic-iii, a freely accessible critical care database, sci,” *Data*, vol. 3, no. 1, pp. 1–9, 2016.
- [28] C. Vasantharajan *et al.*, “Medbert: A pre-trained language model for biomedical named entity recognition,” in *Proc. APSIPA ASC*, pp. 1482–1488, 2022.
- [29] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Proc NAACL-HLT*, pp. 4171–4186, 2019.
- [30] K. Huang *et al.*, “Clinicalbert: Modeling clinical notes & predicting hosp. readmission,” *arXiv preprint arXiv:1904.05342*, 2019.
- [31] E. A. Knapp, Y. Dong, A. L. Dunlop, J. L. Aschner, J. B. Stanford, T. Hartert, S. L. Teitelbaum, M. L. Hudak, K. Carroll, T. G. O’Connor, *et al.*, “Changes in bmi during the covid-19 pandemic,” *Pediatrics*, vol. 150, no. 3, p. e2022056552, 2022.
- [32] Y. Xiao *et al.*, “Bayesian variational transformer: A generalizable model for rotating machinery fault diagnosis,” *Mech. Systems and Signal Proc.*, vol. 207, p. 110936, 2024.



## APPENDIX

TABLE II: Dataset Variables with Descriptions and Distributions

Feature	Description	Range / Distribution (After Outlier Removal)
Glycated Hemoglobin (HbA1c)	Average blood sugar level for the past 2-3 months	$4\% < \text{HbA1c} \leq 14\%$
LDL Cholesterol	A type of lipoprotein that carries cholesterol in the blood	$20 \text{ mg/dL} < \text{LDL} \leq 370 \text{ mg/dL}$
Blood Pressure (BP)	Force of blood against artery walls (systolic and diastolic)	Systolic: $84 \leq \text{BP} \leq 196$
Body Mass Index (BMI)	A measure of weight in relation to height	$15 \leq \text{BMI} \leq 100$
Socioeconomic Status (Income)	Position in social and economic hierarchy	38.2% upper middle class, 36.8% middle class, 24.3% lower middle class, 0.7% upper class
Age	Age distribution in the dataset	$45 \leq \text{Age} \leq 96$
Race	Race distribution in the dataset	88.2% White, 7.2% Black/African American, 4.6% Asian
Gender	Biological sex recorded at time of visit	61.5% Male, 38.5% Female

1) *The Vanilla Encoder-only Transformer*: Transformers can be either encoder-based, for example the Bi-directional Encoder Representations Transformer (BERT) model or decoder-based, such as Generative Pre-trained Transformer (GPT). Figure 6 illustrates the high-level structure of an encoder-based transformer model. From the image, biomarker sequences for each patient are tokenized and mapped into a high-dimensional embedding space. Each token is then transformed into three vectors (Query, Key, and Value)—via learned linear projections. These vectors are used within the attention layer to determine which parts of the sequence are most relevant to each other. Identified relevancies (attention weights) are eventually multiplied by a value matrix and passed to the fully connected layer for further learning and prediction.

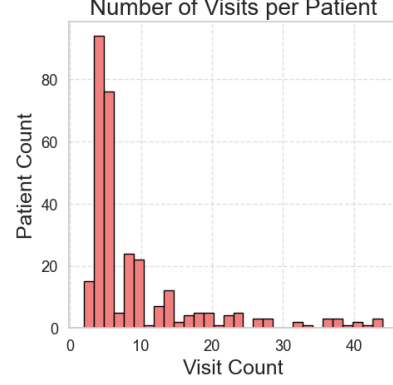


Fig. 5: Visit frequency of patients

A significant aspect of the Transformer architecture is the multi-head self attention [32] which comprises multiple parallel self attention heads, each consisting of linear transformations and dot-product operations (see Figure 8). Given a set of learned query weights  $W^q \in \mathbb{R}^{d_m \times d_k}$ , learned key weights  $W^k \in \mathbb{R}^{d_m \times d_k}$ , learned values weights  $W^v \in \mathbb{R}^{d_m \times d_v}$  and input embedding  $I \in \mathbb{R}^{m \times d_m}$ , the transformer input values (queries, keys, and values) fed into the attention layer are obtained through linear transformations: queries  $q = IW^q \in \mathbb{R}^{m \times d_k}$ , keys  $k = IW^k \in \mathbb{R}^{m \times d_k}$  and values  $v = IW^v \in \mathbb{R}^{m \times d_v}$ . As a first step, q and k are multiplied via the dot product to get matrix  $\varphi$ :  $\varphi = f_{dot}(q, k) = qk^T / (d_k)^{1/2} \in \mathbb{R}^{m \times m}$ .

$(d_k)^{1/2}$  is used to scale matrix values. Next, matrix scores are normalized using the Softmax function on  $\varphi$ ,  $A = \text{softmax}(\varphi)$  to derive attention weights. These normalized weights are then multiplied by the value matrix to obtain attention results  $O = Av \in \mathbb{R}^{m \times d_v}$ . Finally, the attention results from all the self attention heads are concatenated to obtain the final attention value  $O_{final} = \text{Concat}(O_1, O_2, O_3, \dots, O_H)W^O \in \mathbb{R}^{m \times d_m}$ . The concatenated attention value matrix from the multi-head self attention are then passed into a linear layer to transform the large matrix into form suitable for input into a NN layer for more learning and prediction. More details is illustrated in Appendix Figure 2. In our BMT model, we extend the basic transformer architecture by introducing Bayesian self-attention where deterministic attention weights are replaced with stochastic distributions. However, to better understand this concept, in the next section we explain bayesian methods in DL.

2) *Bayesian Methods in DL*: Bayesian methods in DL incorporate uncertainty by treating the weights of the network as probability distributions rather than fixed/deterministic values. This probabilistic treatment allows the model to account for uncertainty in both the data and the learned parameters, resulting in more robust and informative predictions. An example of this

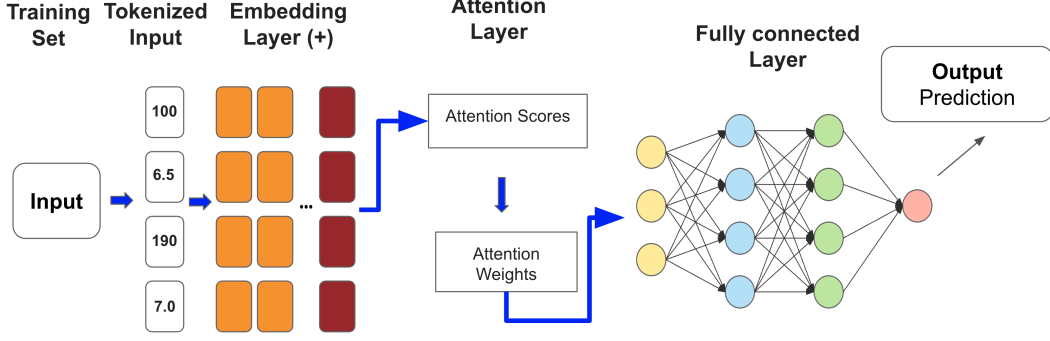


Fig. 6: A simplified example illustration of an Encoder-only Transformer model

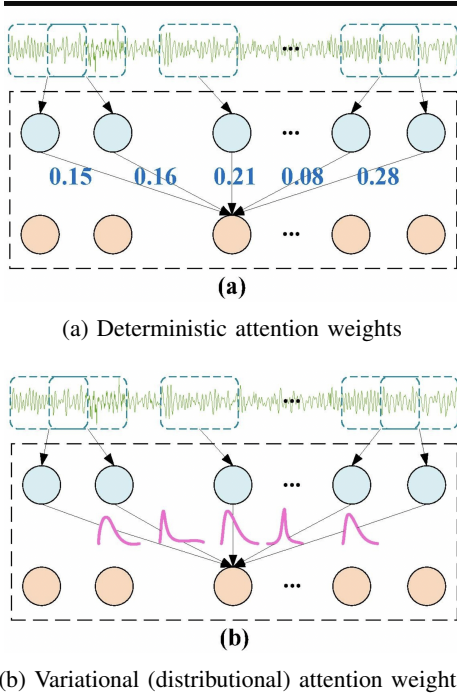


Fig. 7: Comparison between deterministic and variational attentions. (a) Deterministic attention; (b) Variational attention. Image adapted from [32].

transformation is as shown in Appendix, Figure 7b. In the bayesian method for NN, weights  $w$  are initially assumed to follow a distribution called prior distribution  $p(w)$ , which is then multiplied by the likelihood  $p(D|w)$  and divided by the marginal likelihood  $p(D)$ . The result gives you a posterior distribution  $p(w|D)$ , which is the posterior distribution gotten after training on the data (via likelihood), minimizing loss and leveraging the prior knowledge/distribution. This process is governed by Bayes' theorem shown in equation 4:

$$p(w|D) = \frac{p(D|w)p(w)}{p(D)} \quad (4)$$

Here,  $p(w)$  is the prior distribution over weights,  $p(D|w)$  is the likelihood of the data given the weights,  $p(D)$  is the marginal likelihood (also known as evidence), and  $p(w|D)$  is the resulting posterior after observing data. While the posterior distribution  $p(w|D)$  captures the complete uncertainty in the model's parameters, it is often intractable to estimate considering the large number of weights that can be present in a NN model, hence various approximation approaches are used including variational inference or Markov Chain Monte Carlo (MCMC) methods.

3) **Variational Inference for Attention Weights:** In variational inference, we aim to approximate the true posterior distribution  $p(w | D)$  by introducing a simpler, tractable distribution  $q(w)$  that is close enough to it. The objective is to minimize the Kullback-Leibler (KL) divergence between the two distributions:

$$\text{KL}(q(w) || p(w | D)) = \int q(w) \log \left( \frac{q(w)}{p(w | D)} \right) dw \quad (5)$$

Minimizing this divergence allows us to efficiently estimate the posterior while still capturing meaningful uncertainty. This forms the foundation of *Bayesian NN*, which not only generate predictions but also provide well-calibrated uncertainty estimates—an essential feature in high-stakes applications such as healthcare. In our model, we apply variational inference specifically to the *attention weights*, treating them as latent random variables. The goal is to approximate the true posterior distribution over attention weights  $p(A | x, y)$  using a variational distribution  $q_\theta(A)$ . According to Bayes' theorem:

$$p(A | x, y) = \frac{p(x, y | A) p(A)}{p(x, y)} \quad (6)$$

However, since the marginal likelihood  $p(x, y)$  is typically intractable, we instead minimize the KL di-

vergence between the variational approximation and the true posterior:

$$\begin{aligned}
KL(q_\phi(A)||p(A|x,y)) &= \int q_\phi(A) \log \frac{q_\phi(A)}{p(A|x,y)} dA \\
&= \int q_\phi(A) \log \frac{q_\phi(A)}{p(x,y)p(x,y|A)p(A)} dA \\
&= \log p(x,y) - \underbrace{\int q_\phi(A) \log \frac{p(x,y|A)p(A)}{q_\phi(A)} dA}_{L(x,y)}
\end{aligned}$$

This formulation reveals that minimizing the KL divergence is equivalent to maximizing the *Evidence Lower Bound (ELBO)*, defined as:

$$\mathcal{L}(x,y) = \mathbb{E}_{q_\theta(A)} [\log p(x,y | A)] - KL(q_\theta(A)||p(A))$$

By maximizing the ELBO during training, we encourage the learned attention weight distribution  $q_\theta(A)$  to approximate the true posterior  $p(A | x, y)$  as closely as possible, enabling principled uncertainty quantification in our Bayesian Multi-Target Transformer model.

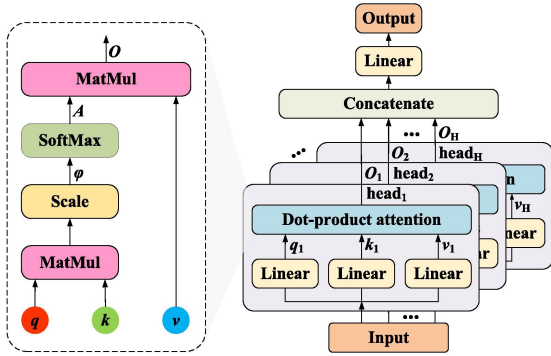


Fig. 8: Multi-head self-attention mechanism. Image from [32]

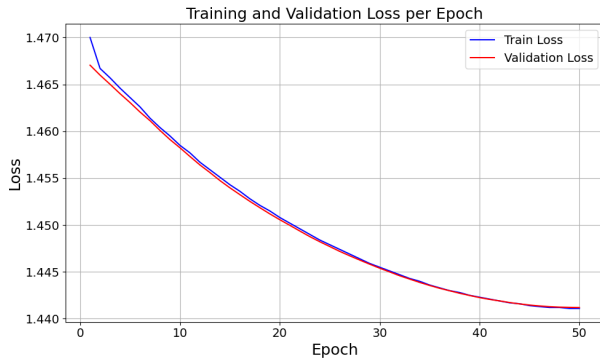


Fig. 9: Loss curve of MBT-CB model

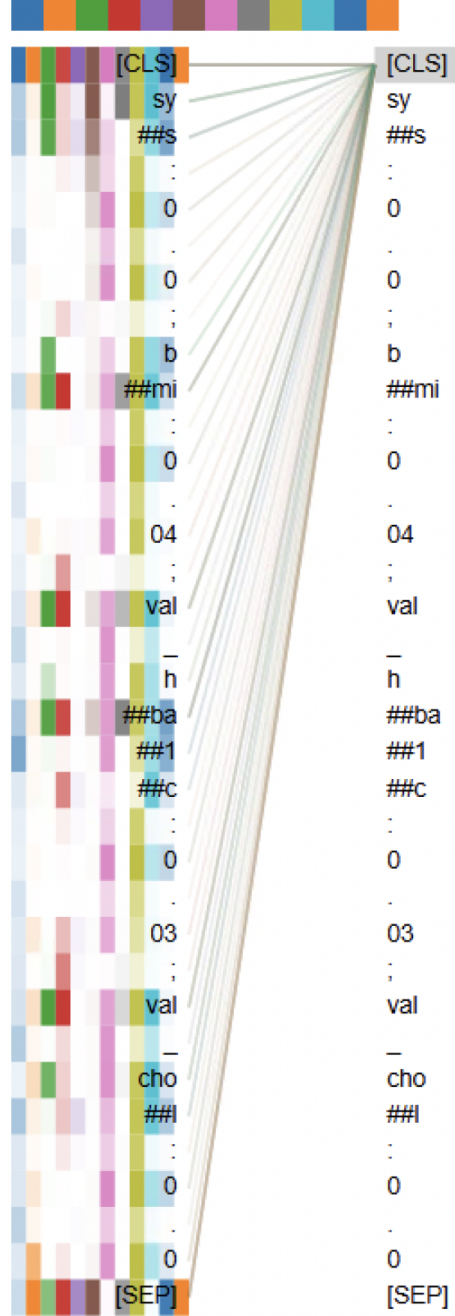


Fig. 10: Attention relationship visualization

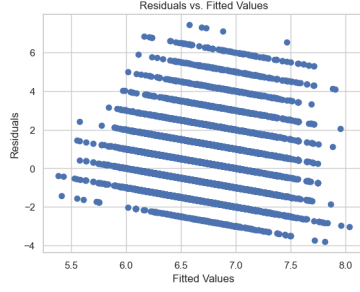


Fig. 11: Heteroscedasticity of HbA1c showing non-constant residual variance

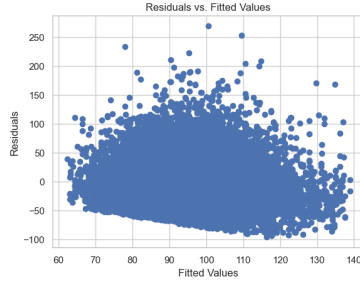


Fig. 12: Heteroscedasticity of LDL-C showing non-constant residual variance

#### A. Heteroscedasticity

1) *Computational Complexity and Feasibility:* All experiments were run on Amazon WorkSpaces (EC2 m7i-flex.2xlarge) with Microsoft Windows Server 2019, 8 vCPUs (Intel Xeon Platinum 8375C @ 2.4 GHz), and 32 GB RAM, without GPU acceleration. MBT-CB trained for 50 epochs in approximately 6 hours, and inference per patient sequence required 150–200 ms, indicating that the model is computationally efficient even on CPU-only infrastructure. For deployment in low-resource settings, model compression techniques (e.g., distillation and quantization) can further reduce runtime and resource usage, making MBT-CB practical for real-world clinical environments.