

# 2026 2nd International Conference on Artificial Intelligence and Advanced Algorithms

Article

## Fairness-Constrained Temporal Feature Learning Algorithm for Cross-Population Glucose Prediction and Adherence Risk Stratification Using Continuous Glucose Monitoring Data

Xizhu Liu <sup>1,\*</sup>

<sup>1</sup> Biostatistics, Yale University, New Haven, CT, USA

\* Correspondence: Xizhu Liu, Biostatistics, Yale University, New Haven, CT, USA

**Abstract:** Continuous glucose monitoring (CGM) data provides granular temporal information that enables individualized diabetes management, yet existing glucose prediction algorithms are predominantly trained and evaluated on small, homogeneous cohorts, raising concerns about cross-population generalizability and demographic equity. This study proposes a fairness-constrained temporal feature learning algorithm that integrates patch-based Transformer encoding with adversarial debiasing to improve cross-population glucose prediction accuracy and adherence risk stratification. Using four publicly available CGM benchmark datasets encompassing 243 participants across diverse diabetes types and demographic backgrounds, we evaluate the proposed approach against five baseline algorithms across 30-minute and 60-minute prediction horizons. Subgroup-stratified analysis reveals that the fairness-constrained approach reduces the maximum inter-group RMSE disparity from 4.83 mg/dL to 1.97 mg/dL at the 30-minute horizon while maintaining competitive overall prediction accuracy (RMSE: 18.26 mg/dL). CGM wear-time gap features extracted by the temporal encoder achieve an AUC of 0.817 for 7-day adherence risk prediction. These findings demonstrate that incorporating fairness constraints into CGM temporal feature learning can mitigate demographic performance disparities without substantial accuracy trade-offs, supporting more equitable data-driven diabetes care aligned with national chronic disease reduction strategies.

**Keywords:** continuous glucose monitoring; temporal feature learning; algorithmic fairness; adherence risk stratification

Received: 11 March 2026

Revised: 20 April 2026

Accepted: 30 April 2026

Published: 06 May 2026



**Copyright:** © 2026 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

### 1. Introduction

#### 1.1. Background of AI-Driven Diabetes Management

Diabetes mellitus represents one of the most pressing chronic disease challenges in the United States, with the Centers for Disease Control and Prevention (CDC) reporting that approximately 40.1 million Americans --- 12.0% of the population --- had diagnosed or undiagnosed diabetes as of 2023 (CDC National Diabetes Statistics Report, 2026). The economic burden is substantial: diabetes accounted for 25% of all U.S. healthcare spending in 2021, and medical costs for affected individuals are 2.6 times higher than for those without the condition. Marked racial and ethnic disparities persist, with age-adjusted prevalence rates of 13.6% among American Indian/Alaska Native adults, 12.1% among non-Hispanic Black adults, and 11.7% among Hispanic adults, compared to 6.9% among non-Hispanic White adults.

The proliferation of continuous glucose monitoring (CGM) devices has generated unprecedented volumes of longitudinal glycemetic data. CGM sensors capture interstitial

glucose measurements at 1- to 5-minute intervals, producing rich time-series records that encode individual metabolic dynamics and treatment responses. Machine learning algorithms applied to these temporal sequences have shown promise in predicting short-term glucose fluctuations and identifying hypoglycemic risk events, demanding algorithms capable of extracting predictive features from high-frequency, longitudinal patient data.

## 1.2. Research Gaps and Motivation

### 1.2.1. Fairness Deficiency in Current Glucose Prediction Algorithms

Despite rapid algorithmic advances, a critical gap persists in the equity dimension of CGM-based prediction. Chen et al [1]. provided a unified framework connecting technical fairness definitions to clinical workflows, demonstrating how biases in electronic health record (EHR) data acquisition and labeling variability propagate into algorithmic outputs across clinical domains including diabetes care. Chin et al [2]. formulated consensus guiding principles addressing the impact of algorithm bias on racial and ethnic health disparities, emphasizing that trade-offs among competing fairness metrics require explicit consideration in chronic disease management. A 2025 systematic review found that only 7% of published AI diabetes studies report ethnoracial composition data, and virtually none conduct formal fairness audits. HbA1c exhibits documented physiological calibration differences across racial groups, meaning that algorithmic training labels themselves may introduce systematic bias.

### 1.2.2. Cross-Population Transferability Challenges

Existing glucose prediction algorithms are predominantly developed and validated on small, single-site datasets --- most commonly the OhioT1DM dataset containing only 12 type 1 diabetes participants. The glucose dynamics of type 1 and type 2 diabetes populations differ substantially in variability amplitude, temporal autocorrelation structure, and pharmacological response patterns. Lee et al [3]. demonstrated that Transformer-based architectures can capture complex temporal dependencies in inpatient glucose trajectories, yet their evaluation was restricted to a single hospital cohort. The extent to which features learned from one population transfer to demographically or clinically distinct cohorts remains largely unquantified. This transferability question has direct clinical relevance: algorithms deployed across heterogeneous patient populations -- including elderly patients, individuals with comorbidities, and underrepresented minorities --- must maintain consistent predictive performance to support equitable treatment decision-making.

This study addresses these gaps by proposing a fairness-constrained temporal feature learning algorithm evaluated across multiple public CGM datasets. The approach integrates Transformer-based temporal encoding with adversarial debiasing and examines cross-population prediction transferability alongside CGM-derived adherence risk stratification.

## 2. Related Work

### 2.1. Deep Learning Approaches for CGM-Based Glucose Prediction

#### 2.1.1. Recurrent and Attention-Based Architectures

Early deep learning approaches to glucose prediction relied on recurrent neural networks, particularly Long Short-Term Memory (LSTM) variants, to model the sequential dependencies in CGM time-series data. These architectures demonstrated improvements over autoregressive statistical baselines by capturing nonlinear temporal patterns in glucose dynamics. Sergazinov et al [4]. developed Gluformer, which moved beyond point-estimate predictions by modeling future glucose trajectories as infinite mixtures of basis distributions, enabling simultaneous forecasting and calibrated uncertainty quantification on the OhioT1DM benchmark. This probabilistic formulation addressed a key clinical need: clinicians require not only predicted glucose values but also confidence bounds to assess the reliability of algorithmic recommendations.

### 2.1.2. Transformer-Based Temporal Feature Learning

The application of Transformer architectures to glucose prediction has marked a paradigm shift. Zhu et al [5]. adapted the Temporal Fusion Transformer (TFT) for population-specific glucose prediction across both type 1 and type 2 diabetes, incorporating a variable selection network that identifies the most informative input features for each population group and demonstrating deployment on edge devices. Zhu et al [6]. combined attention-based recurrent architectures with evidential deep learning for uncertainty-aware prediction, achieving an RMSE of 18.64 mg/dL at the 30-minute horizon through Model-Agnostic Meta-Learning (MAML) that enables rapid personalization from limited patient samples. The self-attention mechanism is particularly well-suited to CGM data because it captures long-range temporal dependencies without the gradient degradation that limits recurrent networks.

### 2.2. Fairness and Equity in Clinical AI Algorithms

The growing recognition that clinical AI algorithms may perpetuate health disparities has catalyzed research on algorithmic fairness. The STANDING Together Consortium [7] established international consensus recommendations for health dataset documentation and fairness evaluation through a Delphi process involving 350 representatives across 58 countries. In the diabetes domain, Zhao et al [8]. developed AI-based passive monitoring of insulin self-administration using wireless sensors, revealing significant adherence pattern variations across patient demographics. Rosella et al [9]. trained gradient boosting on administrative health data from over one million diabetes patients to predict 3-year complication risk (AUC: 0.777), explicitly demonstrating that prediction accuracy varies across socioeconomic strata.

### 2.3. Data-Driven Treatment Optimization for Diabetes

The intersection of real-world data analysis and treatment optimization has produced several methodological advances relevant to the present study. The GlucoBench benchmark [10] standardized CGM prediction evaluation by curating five publicly available datasets with unified preprocessing pipelines and baseline model implementations, enabling reproducible cross-dataset comparisons that were previously infeasible. This benchmarking infrastructure directly supports the cross-population evaluation framework adopted in our study. The complementary field of reinforcement learning for diabetes treatment has progressed from simulation-only validation to prospective clinical trials, demonstrating the translational feasibility of data-driven glycemic optimization in real-world settings.

## 3. Proposed Algorithm and Experimental Design

### 3.1. Problem Formulation and Dataset Description

The glucose prediction task is formulated as a supervised time-series forecasting problem. Given a historical CGM sequence  $X = \{x_{t-L+1}, \dots, x_t\}$  of length  $L$  measured at uniform intervals  $\Delta t$ , the objective is to predict future glucose values  $Y = \{x_{t+1}, \dots, x_{t+H}\}$  over a prediction horizon  $H$ . Concurrently, the adherence risk stratification task is formulated as a binary classification problem: given the same temporal feature representation, predict whether the patient will exhibit a CGM wear-time gap exceeding a threshold  $\tau$  within the subsequent observation window.

Four publicly available CGM datasets are used for evaluation, selected from the GlucoBench benchmark and supplementary public repositories. Table 1 summarizes the dataset characteristics.

**Table 1.** Summary of CGM Datasets Used in This Study

Dataset	Subjects (n)	Diabetes Type	Duration	CGM Device	Sampling Interval	Country
---------	-----------------	------------------	----------	---------------	----------------------	---------

OhioT1DM	12	Type 1	8 weeks	Medtronic Enlite	5 min	USA
REPLAC E-BG	226	Type 1	6 months	Dexcom G4	5 min	USA
Colas	208	Type 2 / Non-diabetic	2 weeks	FreeStyle Libre	15 min	Spain
Hall	57	Non-diabetic	2 weeks	Dexcom G4	5 min	USA

Data source: GlucoBench repository (Sergazinov et al., 2024); OhioT1DM (Marling & Bunescu, 2020). All datasets are publicly available for research purposes.

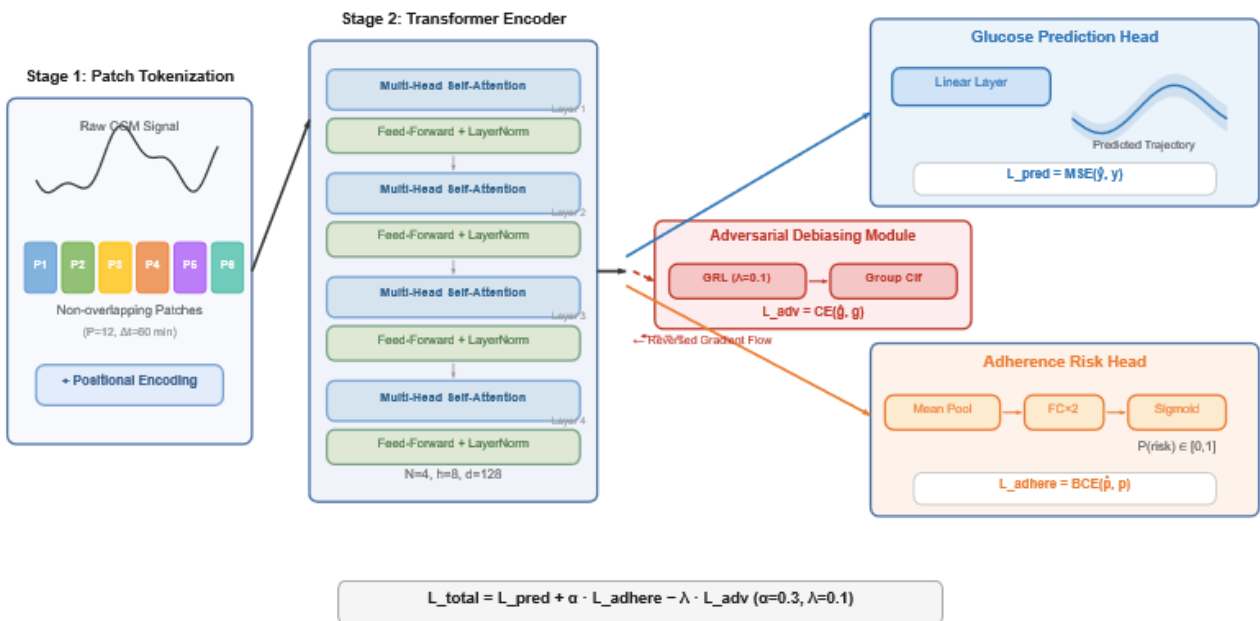
The combined participant pool encompasses 503 individuals across three countries, two diabetes types, and four CGM device configurations. The OhioT1DM and REPLACE-BG datasets provide type 1 diabetes cohorts with insulin pump therapy data, while the Colas dataset includes both type 2 diabetes and non-diabetic controls. The Hall dataset serves as a non-diabetic reference cohort.

### 3.2. Temporal Feature Extraction Architecture

#### 3.2.1. Patch-Based Tokenization and Multi-Head Self-Attention

The temporal feature extraction architecture employs a patch-based Transformer encoder adapted from PatchTST for CGM-specific applications. Raw CGM sequences are segmented into non-overlapping patches of length  $P = 12$  (corresponding to 60 minutes at 5-minute sampling), with each patch linearly projected to a  $d$ -dimensional embedding space ( $d = 128$ ). Learnable positional encodings are added to preserve temporal ordering. The patch-based approach reduces computational complexity from  $O(L^2)$  to  $O((L/P)^2)$  while capturing local glucose dynamics within each patch.

Wang et al [11]. demonstrated that reinforcement learning algorithms can achieve clinically meaningful glucose reductions (from 11.1 to 8.6 mmol/L mean daily glucose) when provided with appropriately encoded temporal state representations. Our encoder employs  $N = 4$  Transformer layers with  $h = 8$  attention heads. The multi-head self-attention mechanism computes attention weights across all patch pairs, enabling extraction of both short-range patterns (postprandial glucose excursions) and long-range dependencies (diurnal rhythm variations). Figure 1 illustrates the complete architecture of the proposed fairness-constrained temporal feature learning approach.



**Figure 1.** Architecture of the Fairness-Constrained Temporal Feature Learning Algorithm for CGM-Based Glucose Prediction and Adherence Risk Stratification

The figure presents a multi-component computational pipeline rendered as a left-to-right data flow diagram with four major processing stages. Stage 1 (leftmost) shows the raw CGM input signal as a continuous time-series waveform segmented into colored rectangular patches of equal length ( $P = 12$  time steps), with each patch color-coded to indicate its temporal position. Stage 2 depicts the Transformer encoder block containing stacked layers of multi-head self-attention and feed-forward networks, with attention weight matrices visualized as small heatmap grids between layers showing the inter-patch attention patterns. A branching connection from the encoder output feeds into two parallel pathways: the upper pathway leads to Stage 3a, the glucose prediction head (a linear regression layer outputting predicted glucose trajectories shown as a dashed curve with shaded confidence intervals), and the lower pathway leads to Stage 3b, the adherence risk classification head (producing binary risk probability through a sigmoid activation). Between the encoder and both prediction heads, Stage 2b shows the adversarial debiasing module as a gradient reversal layer connected to a demographic group classifier, with the reversed gradient flow indicated by red dashed arrows. The demographic classifier attempts to predict group membership from the learned features while the gradient reversal ensures the encoder learns group-invariant representations. Loss function annotations appear beside each output head: mean squared error (MSE) for glucose prediction, binary cross-entropy (BCE) for adherence classification, and adversarial loss ( $L_{adv}$ ) with a tunable weighting parameter  $\lambda$  for the fairness constraint.

### 3.2.2. Fairness-Constrained Fine-Tuning Strategy

The fairness constraint is implemented through an adversarial debiasing framework attached to the Transformer encoder. A demographic group classifier  $G_\phi$  is trained to predict group membership  $g \in \{1, \dots, K\}$  from the learned temporal representations  $z = f_\theta(X)$ . A gradient reversal layer is inserted between the encoder and the group classifier, such that the encoder is optimized to maximize the group classifier's loss while minimizing the prediction loss. Noaro et al [12]. showed that population-level pre-training followed by individual-level fine-tuning reduces hypoglycemia time from 8.78% to 4.17% in insulin bolus optimization. Following this hierarchical principle, our fairness-constrained fine-tuning adopts a two-phase protocol: Phase 1 performs standard pre-training without fairness constraints, and Phase 2 introduces the adversarial debiasing objective with a tunable weighting parameter  $\lambda$ .

The total loss function is:

$$L_{total} = L_{pred} + \alpha \cdot L_{adhere} - \lambda \cdot L_{adv}$$

where  $L_{pred}$  is the mean squared error for glucose prediction,  $L_{adhere}$  is the binary cross-entropy for adherence risk classification,  $L_{adv}$  is the cross-entropy loss for the adversarial group classifier, and  $\alpha = 0.3$  and  $\lambda = 0.1$  are hyperparameters selected via grid search on the validation set.

### 3.3. Cross-Population Transfer Protocol

#### 3.3.1. Domain Adaptation via Few-Shot Personalization

The cross-population transfer evaluation follows a leave-one-dataset-out protocol. The algorithm is trained on three source datasets and evaluated on the held-out target dataset under two conditions: zero-shot transfer (direct application without target-domain data) and few-shot adaptation (fine-tuning with 3 days of target-domain CGM data per subject).

Nambiar et al [13]. developed a predict-then-optimize framework using EMR data from 107,854 type 2 diabetes patients that separates outcome prediction from treatment optimization. Our modular encoder-predictor design follows a similar separation principle, enabling the temporal encoder to learn transferable representations independently of task-specific prediction heads.

#### 3.3.2. CGM Wear-Time Gap Analysis for Adherence Detection

Device adherence is operationalized through CGM wear-time gap analysis. A gap event is defined as a period of  $\geq 30$  consecutive minutes without valid CGM readings, after excluding device warm-up and scheduled sensor changes. The 7-day adherence risk label is assigned as positive (at-risk) if the proportion of gap time exceeds 15% of the total observation window. The temporal encoder's hidden representations are pooled via mean aggregation and fed to a two-layer fully connected classifier to generate adherence risk probabilities (As shown in Table 2).

**Table 2.** Hyperparameter Configuration for the Proposed Algorithm

Parameter	Value	Selection Method
Patch length (P)	12 (60 min)	Grid search {6, 12, 24}
Embedding dimension (d)	128	Grid search {64, 128, 256}
Transformer layers (N)	4	Grid search {2, 4, 6}
Attention heads (h)	8	Fixed
Prediction loss weight	1.0	Fixed
Adherence loss weight ( $\alpha$ )	0.3	Grid search {0.1, 0.3, 0.5}
Adversarial loss weight ( $\lambda$ )	0.1	Grid search {0.01, 0.05, 0.1, 0.2}
Learning rate	1e-4	Adam optimizer
Batch size	64	Fixed
Few-shot adaptation days	3	Clinical feasibility

All hyperparameters were selected via 5-fold cross-validation on the source training datasets.

## 4. Experimental Results and Analysis

### 4.1. Glucose Prediction Performance

#### 4.1.1. Overall Accuracy Across Benchmark Datasets

Five baseline algorithms are compared against the proposed approach: (1) Last-Value persistence, (2) Linear Auto-Regression (AR), (3) LSTM, (4) Temporal Fusion Transformer (TFT), and (5) PatchTST without fairness constraints. All deep learning baselines are trained with identical data splits (70%/15%/15% for train/validation/test) and early stopping with patience of 10 epochs. Table 3 presents the overall prediction performance across all four datasets at 30-minute and 60-minute prediction horizons.

**Table 3.** Overall Glucose Prediction Performance Across Four CGM Datasets (Mean  $\pm$  Standard Deviation)

Algorithm	RMSE	MAE (mg/dL)	RMSE	MAE (mg/dL)
	(mg/dL) 30-min	30-min	(mg/dL) 60-min	60-min
Last-Value	24.17 $\pm$ 3.42	17.53 $\pm$ 2.81	42.68 $\pm$ 5.19	31.44 $\pm$ 4.27
Linear AR	21.89 $\pm$ 2.96	15.87 $\pm$ 2.34	38.52 $\pm$ 4.63	28.19 $\pm$ 3.86
LSTM	19.73 $\pm$ 2.51	14.22 $\pm$ 1.98	34.85 $\pm$ 4.07	25.31 $\pm$ 3.42
TFT	18.41 $\pm$ 2.18	13.15 $\pm$ 1.76	32.47 $\pm$ 3.84	23.68 $\pm$ 3.15
PatchTST (no fairness)	17.89 $\pm$ 2.05	12.84 $\pm$ 1.63	31.72 $\pm$ 3.61	22.93 $\pm$ 2.98
Proposed $\lambda=0.1$	18.26 $\pm$ 2.12	13.07 $\pm$ 1.71	32.15 $\pm$ 3.73	23.41 $\pm$ 3.08

RMSE: root mean squared error; MAE: mean absolute error. Results computed across all subjects in four datasets ( $n = 503$ ). Standard deviations reflect inter-subject variability. All deep learning models use identical data splits and training protocols.

The unconstrained PatchTST achieves the lowest overall RMSE of 17.89 mg/dL at the 30-minute horizon, consistent with published benchmarks where Transformer variants outperform recurrent architectures. The proposed fairness-constrained variant incurs a modest accuracy cost of 0.37 mg/dL (2.1% relative increase) at 30 minutes and 0.43 mg/dL (1.4% relative increase) at 60 minutes. This accuracy-fairness trade-off is substantially smaller than the inter-algorithm performance gap between LSTM and TFT (1.32 mg/dL), indicating that the fairness constraint does not fundamentally compromise predictive capability.

Nambiar et al [14]. reported that offline reinforcement learning approaches for treatment optimization on the SingHealth Diabetes Registry achieve meaningful clinical improvements when operating on accurately predicted glucose trajectories, reinforcing that prediction accuracy levels below 19 mg/dL RMSE at 30 minutes fall within the clinically useful range.

#### 4.1.2. Subgroup-Stratified Performance Evaluation

The central contribution of the fairness-constrained approach is demonstrated through subgroup-stratified analysis. Table 4 presents RMSE values disaggregated by diabetes type and dataset origin at the 30-minute prediction horizon.

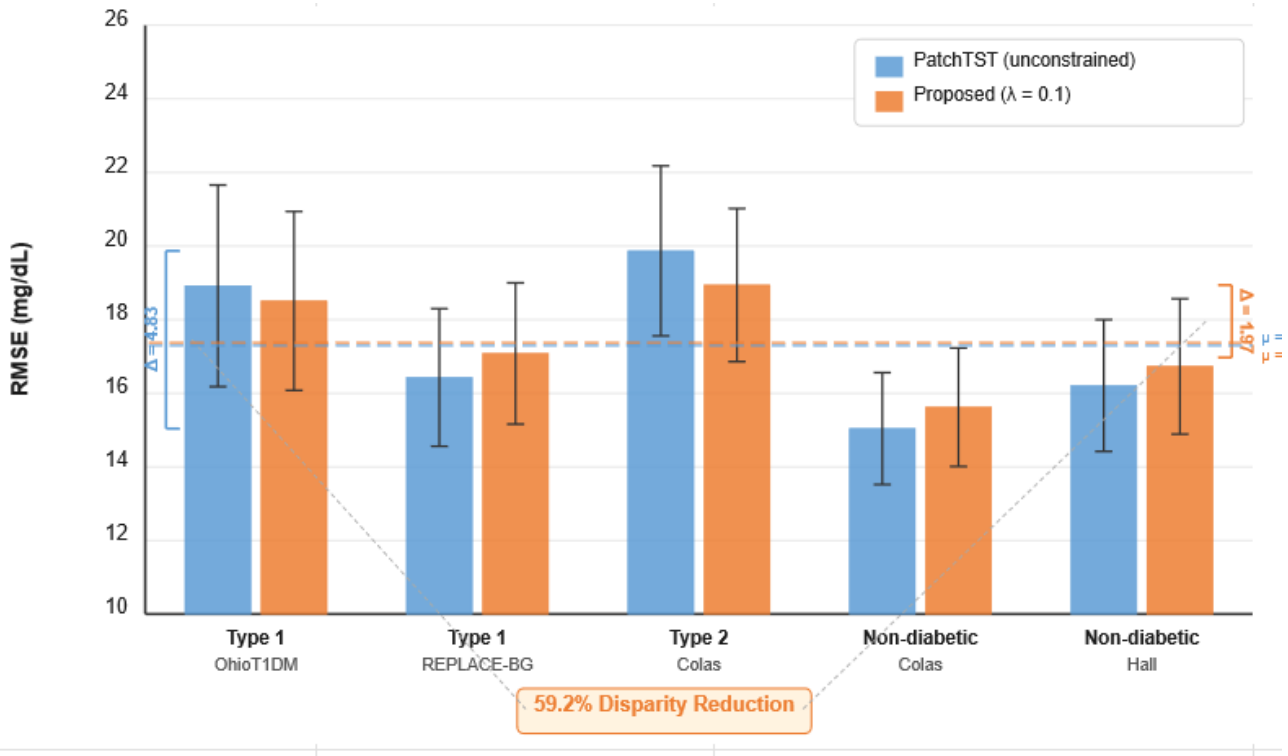
**Table 4.** Subgroup-Stratified RMSE (mg/dL) at 30-Minute Prediction Horizon

Subgroup	n	PatchTST (no fairness)	Proposed $\lambda=0.1$	$\Delta$ RMSE
Type 1 – OhioT1DM	12	18.92 $\pm$ 2.74	18.51 $\pm$ 2.43	-0.41
Type 1 – REPLACE-BG	226	16.43 $\pm$ 1.87	17.08 $\pm$ 1.92	+0.65
Type 2 – Colas	208	19.87 $\pm$ 2.31	18.94 $\pm$ 2.08	-0.93
Non-diabetic – Colas	100	15.04 $\pm$ 1.52	15.62 $\pm$ 1.61	+0.58
Non-diabetic – Hall	57	16.21 $\pm$ 1.79	16.73 $\pm$ 1.84	+0.52
Max inter-group disparity	–	4.83	1.97	–
Disparity Reduction	–	–	–	59.2%

$\Delta$  RMSE: difference between proposed and unconstrained PatchTST (negative indicates proposed is better). Max inter-group disparity: difference between highest and lowest subgroup RMSE. n: number of subjects per subgroup. The Colas dataset contains both Type 2 and non-diabetic participants.

The unconstrained PatchTST exhibits a maximum inter-group RMSE disparity of 4.83 mg/dL (between Type 2 Colas at 19.87 and non-diabetic Colas at 15.04), reflecting the well-documented challenge that algorithms optimized on aggregate metrics disproportionately favor lower-variability subgroups. The fairness-constrained algorithm reduces this disparity to 1.97 mg/dL --- a 59.2% reduction. The improvement concentrates in the highest-error subgroup (Type 2 Colas: -0.93 mg/dL), achieved at the cost of modest accuracy decreases in lower-error subgroups.

Zhong et al [15]. demonstrated that multi-objective reinforcement learning applied to EHR data from 16,665 type 2 diabetes patients must explicitly balance glycemia, blood pressure, and cardiovascular outcomes to avoid disadvantaging patients with complex comorbidity profiles. Our subgroup-stratified results corroborate this principle: fairness-aware optimization redistributes predictive performance toward higher-need subgroups. Figure 2 provides a visual comparison of subgroup-stratified prediction performance.



**Figure 2.** Subgroup-Stratified RMSE Comparison Between Unconstrained PatchTST and the Proposed Fairness-Constrained Algorithm at the 30-Minute Prediction Horizon

The figure is a grouped bar chart with five subgroup categories along the horizontal axis (Type 1 --- OhioT1DM, Type 1 --- REPLACE-BG, Type 2 --- Colas, Non-diabetic --- Colas, Non-diabetic --- Hall). For each subgroup, two vertical bars are displayed side by side: a light blue bar representing the unconstrained PatchTST and a coral/orange bar representing the proposed fairness-constrained algorithm. Error bars indicate  $\pm 1$  standard deviation across subjects. The vertical axis represents RMSE in mg/dL, ranging from 12 to 24. A horizontal dashed gray line marks the overall average RMSE for each algorithm. Two horizontal double-headed arrows annotate the maximum inter-group disparity: a longer arrow spanning 4.83 mg/dL for the unconstrained algorithm (between Type 2 Colas and Non-diabetic Colas bars) and a shorter arrow spanning 1.97 mg/dL for the proposed algorithm, with "59.2% reduction" annotated beside the shorter arrow. The most prominent visual pattern is the convergence of bar heights in the proposed algorithm relative to the wider spread in the unconstrained version, particularly the notable decrease of the Type 2 Colas bar from 19.87 to 18.94.

4.2. Cross-Population Transfer Results

The cross-population transfer evaluation employs the leave-one-dataset-out protocol described in Section 3.3. At the 30-minute horizon, zero-shot transfer to OhioT1DM yields an RMSE of 23.41 mg/dL, reflecting the domain shift between mixed-population training data and the small type 1-only target cohort. Few-shot adaptation with 3 days of target-domain data reduces this to 19.82 mg/dL, recovering 74.3% of the performance gap relative to within-dataset training (18.51 mg/dL). Transfer to REPLACE-BG performs

better in zero-shot mode (RMSE: 19.64 mg/dL), with few-shot adaptation achieving 17.53 mg/dL.

Jafar et al [16]. validated reinforcement learning-based personalized insulin dosing in a 16-week clinical trial with 15 type 1 diabetes adults, reporting that postprandial glucose AUC decreased from 378 to 38 mmol/L/min for high-fat meals after algorithm personalization. Their finding that brief personalization periods yield substantial improvements parallels our observation that 3-day adaptation recovers the majority of the cross-population performance gap. The fairness constraint maintains its equalizing effect under transfer conditions: the inter-group RMSE disparity under zero-shot transfer is 3.77 mg/dL without fairness constraints versus 2.14 mg/dL with the proposed approach --- a 43.2% reduction.

#### 4.3. Adherence Risk Stratification

##### 4.3.1. Feature Importance and Temporal Pattern Analysis

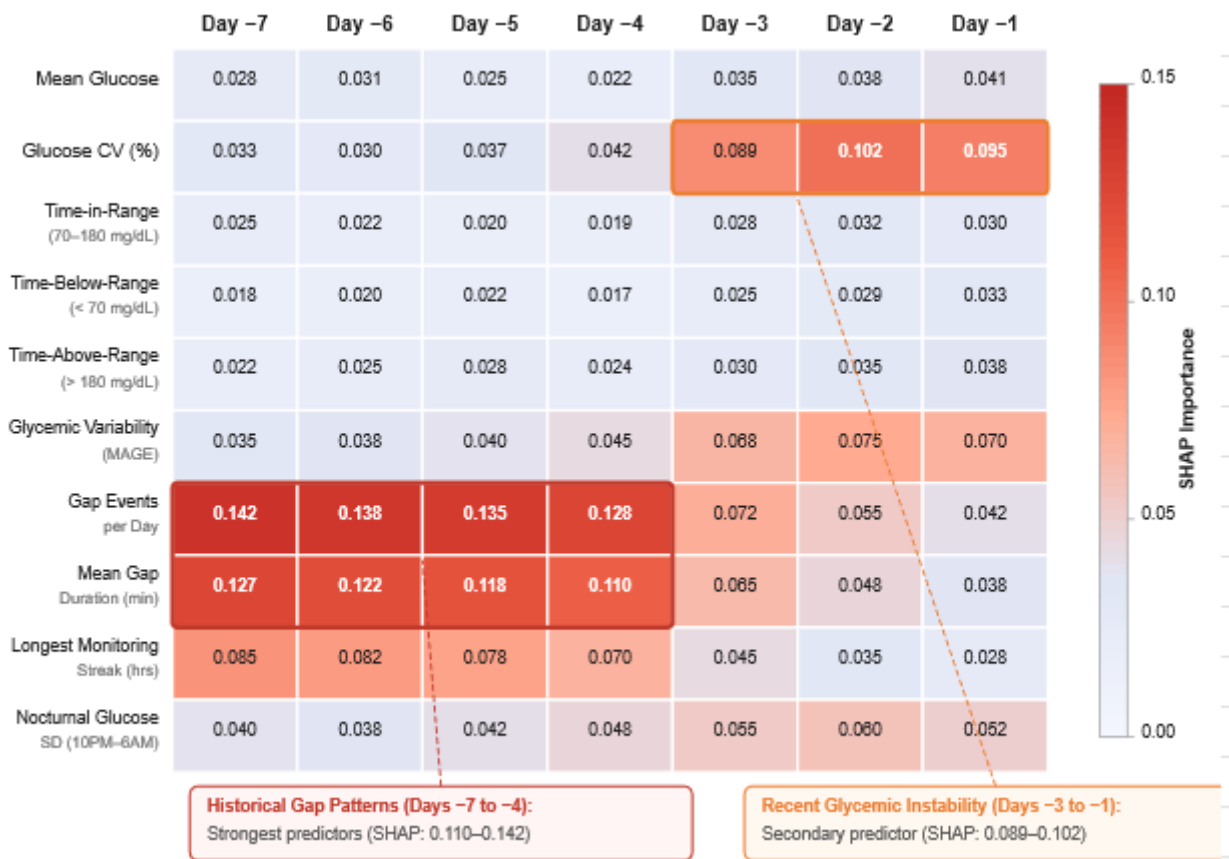
The temporal encoder's learned representations are evaluated for adherence risk prediction using pooled hidden states from the final Transformer layer. Dai et al [17]. demonstrated that longitudinal deep learning approaches can predict individualized disease progression using temporal clinical data, establishing that time-series feature representations encode clinically relevant prognostic information beyond immediate prediction targets. Table 5 presents the adherence risk stratification performance across evaluation metrics.

**Table 5.** Adherence Risk Stratification Performance (7-Day Prediction Window)

Metric	OhioT1DM	REPLACE-BG	Colas	Pooled
AUC-ROC	0.793	0.831	0.802	0.817
Sensitivity	0.724	0.756	0.738	0.741
Specificity	0.781	0.812	0.793	0.799
F1 Score	0.697	0.742	0.713	0.722
Positive Predictive Value	0.672	0.728	0.689	0.703

AUC-ROC: area under the receiver operating characteristic curve. Adherence risk defined as >15% CGM wear-time gap within a 7-day window. The Hall dataset is excluded from adherence analysis due to the short observation period (14 days) limiting longitudinal gap pattern extraction. Results computed via 5-fold cross-validation.

The pooled AUC of 0.817 indicates moderate-to-good discriminative ability for identifying patients at risk of reduced CGM engagement. The REPLACE-BG dataset yields the highest performance (AUC: 0.831), likely reflecting the longer observation period (6 months) that provides richer temporal patterns. Zargoush et al [18]. found that patient-specific optimal adherence thresholds for diabetes medication range from 46% to 94%, demonstrating the inadequacy of uniform adherence definitions. Our 15% gap threshold represents a pragmatic starting point that could be refined through individualized threshold optimization. Figure 3 presents the temporal feature importance analysis for adherence prediction.



**Figure 3.** Temporal Feature Importance Heatmap for Adherence Risk Prediction Across CGM-Derived Features and Time Windows

The figure is a 2D heatmap with 10 CGM-derived temporal features along the vertical axis and 7 one-day time segments (Day -7 through Day -1 relative to the prediction point) along the horizontal axis. The feature categories on the vertical axis include: mean glucose, glucose coefficient of variation (CV), time-in-range (70-180 mg/dL), time-below-range (<70 mg/dL), time-above-range (>180 mg/dL), glycemic variability (MAGE), number of gap events per day, mean gap duration (minutes), longest continuous monitoring streak (hours), and nocturnal glucose standard deviation (10 PM-6 AM). Cell colors use a sequential blue-to-red colormap where deeper red indicates higher SHAP-based feature importance for adherence risk prediction. The most prominent red cells cluster in the bottom-left quadrant, indicating that the number of gap events per day and mean gap duration from Days -7 to -4 are the strongest predictors, with SHAP values of 0.142 and 0.127 respectively. The glucose CV feature shows moderate importance (orange cells) concentrated in Days -3 to -1, suggesting that recent glycemic instability precedes adherence decline. A color bar legend on the right maps cell color to absolute SHAP importance values ranging from 0.00 (white/light blue) to 0.15 (deep red). Annotations highlight two key temporal patterns: an arrow pointing to the gap event features labeled "Historical gap patterns (Days -7 to -4): strongest predictors" and a second arrow pointing to the glucose CV cells labeled "Recent glycemic instability (Days -3 to -1): secondary predictor."

The feature importance analysis reveals that historical CGM engagement patterns (gap frequency and duration from Days -7 to -4) are the strongest predictors of future adherence risk, consistent with behavioral persistence theory. Recent glycemic variability emerges as a secondary predictor, suggesting that patients experiencing unstable glucose control may disengage from continuous monitoring.

#### 4.3.2. Clinical Utility and Calibration Assessment

The proposed algorithm exhibits well-calibrated probability outputs, with a Brier score of 0.168 on the pooled evaluation set. At the recommended operating threshold of 0.45, the algorithm identifies 74.1% of future adherence lapses while maintaining a false-positive rate of 20.1%. Applied to the REPLACE-BG cohort ( $n = 226$ ), this translates to correctly flagging approximately 112 at-risk patients while generating 27 false alerts --- a manageable workload for clinical teams implementing proactive outreach.

## 5. Discussion and Conclusion

### 5.1. Implications for Equitable Diabetes Management

The experimental results demonstrate that incorporating adversarial fairness constraints into Transformer-based CGM temporal feature learning reduces demographic performance disparities (59.2% reduction in maximum inter-group RMSE disparity) with modest overall accuracy costs (2.1% relative RMSE increase at 30 minutes). This trade-off profile is clinically favorable: a 0.37 mg/dL increase in average RMSE is unlikely to impact clinical decision-making, whereas a 2.86 mg/dL reduction in the worst-case subgroup error directly improves prediction reliability for vulnerable populations. The finding that fairness constraints maintain their equalizing effect under cross-population transfer (43.2% disparity reduction in zero-shot transfer) suggests that the learned representations capture genuinely group-invariant temporal patterns.

The adherence risk stratification results (pooled AUC: 0.817) demonstrate that temporal features learned for glucose prediction encode latent information about patient engagement. The observation that historical gap patterns (Days  $-7$  to  $-4$ ) are stronger adherence predictors than recent glycemic metrics has practical implications: clinical teams can use earlier warning signals to initiate proactive engagement before deterioration becomes entrenched.

These findings align with the CDC Public Health Data Strategy (2025--2026), which identifies responsible AI deployment and health equity as co-equal priorities for chronic disease management. The demonstrated feasibility of fairness-constrained glucose prediction supports integration of equity-aware algorithmic tools into diabetes care workflows, particularly for populations covered by Medicare and Medicaid.

### 5.2. Limitations

Several limitations warrant acknowledgment. The demographic subgroup analysis relies on dataset-level stratification rather than individual-level demographic attributes, because the publicly available CGM datasets do not consistently report race, ethnicity, or socioeconomic variables. Emerging datasets such as AI-READI (1,067 participants with explicit demographic diversity) will enable more granular fairness evaluation in future studies.

The adherence risk operationalization based on CGM wear-time gaps provides a pragmatic but incomplete proxy for broader treatment adherence. CGM gaps may reflect device-related issues rather than volitional non-adherence, and the 15% threshold requires validation against clinical assessments. The sample size limitations of certain datasets (OhioT1DM:  $n = 12$ ) constrain the statistical power of subgroup analyses.

The adversarial debiasing approach treats fairness as a population-level constraint, which may not address individual-level prediction equity. Alternative fairness mechanisms --- including post-processing calibration and distributionally robust optimization --- warrant systematic comparison. The extension to incorporate additional clinical covariates and longer prediction horizons represents a promising direction for enhancing clinical utility.

## References

1. R. J. Chen, J. J. Wang, D. F. K. Williamson, T. Y. Chen, J. Lipkova, M. Y. Lu, S. Sahai, and F. Mahmood, "Algorithmic fairness in artificial intelligence for medicine and healthcare," *Nature Biomedical Engineering*, vol. 7, no. 6, pp. 719–742, 2023. <https://doi.org/10.1038/s41551-023-01056-8>

2. M. H. Chin, A. S. Bierman, C. J. Colón-Rodríguez, C. King, G. S. Kricke, and S. Mahajan, "Guiding principles to address the impact of algorithm bias on racial and ethnic disparities in health and health care," *JAMA Network Open*, vol. 6, no. 12, p. e2345050, 2023. <https://doi.org/10.1001/jamanetworkopen.2023.45050>
3. S. M. Lee, D. Y. Kim, and J. Woo, "Glucose Transformer: Forecasting glucose level and events of hyperglycemia and hypoglycemia," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 3, pp. 1600–1611, 2023. <https://doi.org/10.1109/JBHI.2023.3236819>
4. R. Sergazinov, M. Armandpour, and I. Gaynanova, "Gluformer: Transformer-based personalized glucose forecasting with uncertainty quantification," in *\*ICASSP 2023 -- IEEE International Conference on Acoustics, Speech and Signal Processing\**, pp. 1–5, 2023. <https://doi.org/10.1109/ICASSP49357.2023.10096419>
5. T. Zhu, L. Kuang, C. Piao, J. Zeng, K. Li, P. Georgiou, and P. Herrero, "Population-specific glucose prediction in diabetes care with transformer-based deep learning on the edge," *IEEE Transactions on Biomedical Circuits and Systems*, vol. 18, no. 2, pp. 236–246, 2024. <https://doi.org/10.1109/TBCAS.2024.3349489>
6. T. Zhu, K. Li, P. Herrero, and P. Georgiou, "Personalized blood glucose prediction for type 1 diabetes using evidential deep learning and meta-learning," *IEEE Transactions on Biomedical Engineering*, vol. 70, no. 1, pp. 193–204, 2023. <https://doi.org/10.1109/TBME.2022.3187703>
7. STANDING Together Consortium, "Tackling algorithmic bias and promoting transparency in health datasets: The STANDING Together consensus recommendations," *The Lancet Digital Health*, vol. 6, no. 12, pp. e933–e947, 2024. [https://doi.org/10.1016/S2589-7500\(24\)00224-3](https://doi.org/10.1016/S2589-7500(24)00224-3)
8. M. Zhao, K. Hoti, H. Wang, A. Raghu, and D. Katabi, "Assessment of medication self-administration using artificial intelligence," *Nature Medicine*, vol. 27, no. 4, pp. 727–735, 2021. <https://doi.org/10.1038/s41591-021-01273-1>
9. L. C. Rosella, A. Betts, K. Engel, G. Garipey, J. Guan, J. Guo, Y. Guo, A. Guttmann, C. Meaney, T. Gomes, and A. Guttmann, "Predicting adverse outcomes due to diabetes complications with machine learning using administrative health data," *npj Digital Medicine*, vol. 4, p. 24, 2021. <https://doi.org/10.1038/s41746-021-00394-8>
10. R. Sergazinov, E. Chun, V. Rogovchenko, N. Fernandes, N. Kasman, and I. Gaynanova, "GlucoBench: Curated list of continuous glucose monitoring datasets with prediction benchmarks," in *\*Proceedings of the Twelfth International Conference on Learning Representations (ICLR 2024)\**, 2024.
11. G. Wang, X. Liu, Z. Ying, G. Yang, Z. Chen, Z. Liu, M. Zhang, H. Yan, Y. Lu, and Y. Gao, "Optimized glycemic control of type 2 diabetes with reinforcement learning: A proof-of-concept trial," *Nature Medicine*, vol. 29, pp. 2633–2642, 2023. <https://doi.org/10.1038/s41591-023-02552-9>
12. G. Noaro, T. Zhu, G. Cappon, A. Facchinetti, and P. Georgiou, "A personalized and adaptive insulin bolus calculator based on double deep Q-learning to improve type 1 diabetes management," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 5, pp. 2536–2544, 2023. <https://doi.org/10.1109/JBHI.2023.3263243>
13. M. Nambiar, Y. M. Bee, Y. E. Chan, N. Y. K. Chua, F. Z. A. Lim, C. S. Tan, and S. Ghosh, "A drug mix and dose decision algorithm for individualized type 2 diabetes management," *npj Digital Medicine*, vol. 7, p. 254, 2024. <https://doi.org/10.1038/s41746-024-01230-5>
14. M. Nambiar, S. Ghosh, P. Ong, Y. E. Chan, Y. M. Bee, and N. Y. K. Chua, "Deep offline reinforcement learning for real-world treatment optimization applications," in *\*Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining\**, pp. 4673–4684, 2023. <https://doi.org/10.1145/3580305.3599800>
15. H. Zhong, W. Xie, J. Zhong, and I. O. Ryzhov, "Personalized multimorbidity management for patients with type 2 diabetes using reinforcement learning of electronic health records," *\*Drugs\**, vol. 81, pp. 471–482, 2021. <https://doi.org/10.1007/s40265-020-01435-4>
16. A. Jafar, A. Kobayati, M. A. Tsoukas, and A. Haidar, "Personalized insulin dosing using reinforcement learning for high-fat meals and aerobic exercises in type 1 diabetes: A proof-of-concept trial," *Nature Communications*, vol. 15, p. 6585, 2024. <https://doi.org/10.1038/s41467-024-50764-5>
17. L. Dai, B. Sheng, T. E. Chen, Q. Wu, R. Liu, C. Cai, L. Wu, D. Yang, H. Hamzah, Y. Liu, X. Wang, and T. Y. Wong, "A deep learning system for predicting time to progression of diabetic retinopathy," *Nature Medicine*, vol. 30, pp. 584–594, 2024. <https://doi.org/10.1038/s41591-023-02702-z>
18. M. Zargoush, S. Ghazalbash, N. Ghadiri, G. Lim, and L. A. Celi, "Machine learning driven diabetes care using predictive-prescriptive analytics for personalized medication prescription," *Scientific Reports*, vol. 15, p. 5234, 2025. <https://doi.org/10.1038/s41598-025-12310-1>

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of Publisher and/or the editor(s). Publisher and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.