

The Sparsity Trap: Clinical Context as Systematic Noise in Subject-Independent Glucose Forecasting

Anonymous First Author 1*
 University X, Country 1

ABC@SAMPLE.COM

Abstract

Deep learning models for blood glucose forecasting are increasingly built as multivariate systems, with the working assumption that incorporating exogenous clinical variables—such as insulin boluses and carbohydrate intake—will generally improve predictive accuracy. In this work, we challenge that assumption in a strict subject-independent setting on the public Shanghai Diabetes Registry (112 inpatients with Type 2 Diabetes). We benchmark a state-of-the-art univariate Transformer (PatchTST) against a persistence baseline and a late-fusion multivariate variant that ingests standardized insulin and carbohydrate logs.

The univariate Transformer significantly outperforms persistence at long horizons (120-minute RMSE 38.43 vs. 44.20 mg/dL), establishing a strong forecasting baseline under subject-wise splits. In contrast, adding clinical features consistently degrades performance. Stratifying 64,381 test windows by treatment activity reveals a “sparsity trap”: in high-activity windows dominated by meal intake, the multivariate model underperforms the univariate baseline by 2.36 mg/dL (25.07 vs. 22.71 mg/dL RMSE, $p < 10^{-4}$), accompanied by a 4.6 percentage-point drop in the fraction of predictions within 20% of the reference value, a proxy for clinically safe accuracy. We further show that a two-channel CGMCarbs Transformer and a gradient boosting baseline on hand-crafted clinical features both reproduce this sparsity trap—often with poorer clinical accuracy in meal-driven windows—indicating that the failure is not specific to a single architecture. A synthetic injection experiment, in which future glucose is defined as a simple linear function of insulin history, shows that the same architecture rapidly learns the induced fusion rule, implicating the sparsity and noise of real-world logs in this small cohort—rather

than model capacity—as the primary cause of failure. These findings suggest that, for such cohorts, robust univariate modeling may be a safer and more accurate default than naive multivariate fusion of sparse clinical streams.

Data and Code Availability **Data and Code Availability.** This work uses the public Shanghai Diabetes Registry cohort of 112 hospitalized patients with Type 2 Diabetes, resampled to 5-minute CGM resolution as in prior studies. The processed dataset is available from the original curators under their data use terms. An anonymized implementation of our preprocessing pipeline, forecasting models, and analysis scripts, including sparsity-bucket evaluation and synthetic injection experiments, is provided as supplementary material with this submission and will be released in a public repository upon acceptance.

Institutional Review Board (IRB) This study uses a de-identified secondary analysis of an existing public dataset and did not involve interaction with human subjects. As such, it was determined to be Not Human Subjects Research and did not require additional IRB review under our institution’s policies.

1. Introduction

Accurate blood glucose forecasting is a cornerstone of modern diabetes management, enabling proactive interventions for artificial pancreas (AP) systems and clinical decision-support tools. As deep learning architectures have evolved, the field has increasingly moved toward multivariate modeling, operating under the working assumption that incorporating exogenous clinical variables—specifically insulin boluses and carbohydrate intake—will generally improve predictive accuracy [Cobelli et al. \(2011\)](#); [Li et al. \(2020\)](#). In many published studies, additional channels are treated as unequivocally beneficial features rather than hypotheses to be stress-tested.

* These authors contributed equally

81 However, the evidence base for this assumption
 82 often relies on evaluation settings that may not
 83 generalize to the clinical reality of small, subject-
 84 independent cohorts. A substantial fraction of prior
 85 work on CGM forecasting employs randomized train-
 86 test splits, which introduce data leakage by exposing
 87 patient-specific dynamics to the model during train-
 88 ing. Other studies rely on massive, often private
 89 datasets ($N > 10,000$) in which the sheer volume
 90 of data can compensate for the sparsity and irreg-
 91 ularity of clinical events such as insulin doses and
 92 meals. Consequently, there remains a critical gap in
 93 understanding how state-of-the-art multivariate archi-
 94 tectures behave on standard public datasets with
 95 limited sample sizes ($N \approx 100$), where clinical logs
 96 are sparse, noisy, and highly disjointed in time.

97 At the same time, most work emphasizes point es-
 98 timates of error or discrimination for a single model,
 99 rather than carefully benchmarking against simple
 100 but strong baselines under strictly subject-wise splits.
 101 Persistence models—which forecast future glucose by
 102 copying the last observed value—are rarely treated
 103 as first-class comparators, despite their strong perfor-
 104 mance at short horizons and their ability to expose
 105 when deep models add little beyond autocorrelation.
 106 This makes it difficult to assess whether multivari-
 107 ate deep models are genuinely extracting new clinical
 108 signal, or simply repackaging temporal inertia under
 109 more complex architectures.

110 In this work, we rigorously audit the “more data
 111 is better” hypothesis in the setting of subject-
 112 independent glucose forecasting. Unlike prior studies
 113 that primarily report marginal gains from multivari-
 114 ate fusion, we identify and quantify a failure mode
 115 in which auxiliary clinical streams can systematically
 116 degrade both forecasting accuracy and a simple safety
 117 proxy when applied to a small, public inpatient co-
 118 hort. Our goal is not to argue against the value of
 119 insulin or meal information in general, but to clar-
 120 ify when and how naive fusion of these streams can
 121 become harmful in realistic data regimes.

122 We focus on the Shanghai Diabetes Registry, a
 123 public cohort of 112 inpatients with Type 2 Diabetes,
 124 and adopt a strict subject-independent split to pre-
 125 vent leakage from training into evaluation. We bench-
 126 mark a univariate Transformer (PatchTST) against
 127 a late-fusion multivariate architecture and a robust
 128 persistence baseline. In this setting, the univari-
 129 ate model successfully learns long-term temporal dy-
 130 namics and outperforms persistence at long horizons,
 131 whereas the multivariate model behaves very differ-

132 ently: it performs worse on average, and most no-
 133 tably in windows with substantial treatment activity.
 134 We refer to this phenomenon as the *sparsity trap*:
 135 in a small, sparse cohort, a higher-capacity fusion
 136 model appears to overfit to rare, high-magnitude clin-
 137 ical events instead of extracting stable action effects,
 138 leading to systematic degradation in both error and
 139 safety.

Our primary contributions are as follows: 140

1. **Strict subject-independent benchmarking.** 141

142 We establish a rigorous forecasting baseline for
 143 the Shanghai dataset, using a subject-wise split
 144 and a persistence comparator. Under this
 145 protocol, a univariate Transformer significantly
 146 outperforms persistence at long horizons (120-
 147 minute RMSE 38.43 vs. 44.20 mg/dL), demon-
 148 strating that deep architectures can capture tem-
 149 poral structure beyond simple inertia in this co-
 150 hort.

2. **Identification of the sparsity trap.** We pro- 151

152 vide a fine-grained error analysis stratified by
 153 clinical activity. Stratifying 64,381 test windows
 154 by insulin and carbohydrate volume, we show
 155 that naive feature fusion degrades clinical ac-
 156 curacy and safety most severely in high-activity
 157 windows driven by meal intake (RMSE penalty
 158 of +2.36 mg/dL; $p < 10^{-4}$), with a correspond-
 159 ing 4.6 percentage-point drop in the fraction of
 160 predictions within 20% of the reference value, a
 161 proxy for clinically safe accuracy.

3. **Mechanism isolation via synthetic injec-** 162

163 **tion.** We introduce a synthetic injection frame-
 164 work to distinguish between architectural lim-
 165 itations and data properties. By demonstrat-
 166 ing that the same late-fusion architecture rapidly
 167 learns simple, deterministic fusion rules on dense
 168 synthetic targets, we attribute the observed fail-
 169 ure on real data to the intrinsic sparsity and
 170 noise of clinical logs in this small cohort, rather
 171 than an inability of the model class to exploit
 172 action information.

4. **Cross-architecture robustness of the spar-** 173

174 **sity trap.** Beyond a single Transformer, we
 175 show that both a two-channel CGMCarbs vari-
 176 ant and a gradient boosting regressor on hand-
 177 crafted clinical features reproduce the same high-
 178 activity degradation, demonstrating that naive
 179 fusion of sparse clinical logs is fragile across

180 model classes rather than being an idiosyncrasy
181 of one architecture.

182 2. Related Work

183 Multivariate Fusion in Medical Time Series

184 A prevailing paradigm in medical AI is the integra-
185 tion of multi-modal data to improve predictive perfor-
186 mance. In diabetes care, this typically involves con-
187 catenating CGM history with exogenous inputs such
188 as insulin boluses, carbohydrate intake, and heart
189 rate. Surveys have catalogued numerous reports of
190 performance gains from such fusion strategies Li et al.
191 (2020). However, a substantial subset of these evalu-
192 ations utilize *randomized train-test splits*, which dis-
193 tribute windows from the same patient across train-
194 ing and evaluation sets. This introduces significant
195 data leakage, allowing high-capacity models to mem-
196 orize patient-specific responses rather than learning
197 generalizable physiological rules.

198 Data Scarcity and Evaluation Rigor

199 While large-scale private datasets (e.g., from commercial
200 sensor manufacturers) enable the training of ro-
201 bust multivariate models, public research is of-
202 ten constrained to smaller cohorts such as the
203 OhioT1DM Marling and Bunescu (2020) or Shang-
204 hai Diabetes Registry datasets ($N \approx 10\text{--}100$). In
205 these data-scarce regimes, the sparsity of clinical
206 logs becomes a critical bottleneck. Recent audits in
207 time series forecasting have highlighted that complex
208 Transformer-based models can underperform simple
209 baselines on standard benchmarks, calling into ques-
210 tion the automatic benefit of added architectural ca-
211 pacity Zeng et al. (2023). Our work extends this line
212 of inquiry, specifically characterizing the ‘‘Sparsity
213 Trap’’ where naive inclusion of sparse clinical features
214 actively degrades performance under strict subject-
215 independent evaluation.

216 Deep Learning for Glucose Forecasting

217 Early approaches to continuous glucose monitoring (CGM)
218 forecasting relied on physiological models and autore-
219 gressive integrated moving average (ARIMA) base-
220 lines Cobelli et al. (2011). The advent of deep
221 learning shifted the focus to Recurrent Neural Net-
222 works (RNNs) Li et al. (2020). More recently,
223 Transformer-based architectures have achieved state-
224 of-the-art results by modeling long-range dependen-
225 cies via self-attention mechanisms. Notably, channel-
226 independent Transformers like PatchTST Nie et al.
227 (2023) have demonstrated that treating time series

228 variables independently can outperform more elabo-
229 rate mixing strategies in general forecasting bench-
230 marks, a pattern our work corroborates in the spe-
231 cific context of glucose dynamics under subject-
232 independent splits.

233 3. Methods

234 3.1. Dataset and Preprocessing

235 We used the *Shanghai Diabetes Registry* Zhao et al.
236 (2023), a public cohort of 112 inpatients with Type 2
237 Diabetes providing continuous glucose monitoring
238 (CGM) and clinical logs (insulin boluses, carbohy-
239 drate intake) over multi-day admissions. All anal-
240 yses were performed on a uniformly resampled 5-
241 minute grid. Following prior work on this dataset,
242 CGM traces were linearly interpolated over short
243 gaps, while missing clinical entries were treated as
244 true zeros, reflecting the sparsity of documented in-
245 sulin and meal events.

246 To prevent information leakage, all preprocess-
247 ing was performed within a strict subject-wise split.
248 Patients were partitioned by identifier into non-
249 overlapping Train (approximately 70%), Validation
250 (10%), and Test (20%) sets. Scaling parameters
251 (mean and standard deviation) for CGM and clinical
252 channels were fit on the training split only and ap-
253 plied to validation and test data (z -score normaliza-
254 tion). This *subject-independent* protocol ensures that
255 models are evaluated only on unseen patients, prob-
256 ing generalization of physiological dynamics rather
257 than memorization of subject-specific trajectories.

258 3.2. Experimental Design

259 We formulated forecasting as a sequence-to-sequence
260 problem with a 6-hour lookback and a 2-hour predic-
261 tion horizon. At 5-minute resolution, each training
262 example consists of a length- $L = 72$ input window
263 and a length- $H = 24$ output window:

$$(x_{t-L+1:t}, y_{t+1:t+H}), \quad L = 72, H = 24.$$

264 Windows were generated separately for each patient
265 by sliding this 6-hour window forward in time.

266 **Window generation.** For model training, we gen-
267 erated sliding windows with a stride of 1 step to
268 maximize data utilization within each patient. For
269 the horizon-performance analysis (per-step RMSE at
270 $t + 30$, $t + 60$, and $t + 120$ minutes), we evaluated on

the test set using a stride of 24 steps to reduce overlap and approximate independence between windows. For the sparsity-bucket analysis, which requires large samples in rare high-activity regimes, we reverted to stride 1 on the test split and computed per-window summary statistics (e.g., total insulin and carbohydrates over the 6-hour input).

3.3. Model Architectures

Persistence baseline. The persistence model forecasts all future points by copying the last observed CGM value in the input window, $\hat{y}_{t+k} = y_t$ for all $k \in \{1, \dots, H\}$. This trivial baseline provides a strong performance floor for autoregressive physiological time series and is particularly competitive at short horizons.

Univariate Transformer (PatchTST). Our primary forecasting model is PatchTST (Nie et al., 2023), a channel-independent Transformer that segments each univariate time series into non-overlapping patches. For CGM forecasting, we use a single input channel (CGM), with patch length 12 and stride 12, yielding 6 patches over the 72-step history. Each patch is embedded into a $d_{\text{model}} = 128$ -dimensional space and processed by a Transformer encoder with 3 layers, 4 attention heads per layer, and dropout 0.1. The encoder output is projected back to a length-24 forecast through a linear prediction head. This model, denoted *Univariate*, relies exclusively on standardized CGM history.

Clinical late-fusion Transformer. To test the utility of clinical features, we extend PatchTST to ingest three synchronized input channels: CGM, insulin bolus, and carbohydrate intake. In our implementation, these channels share a common patching and Transformer backbone (with the same L , patch length, and hyperparameters), and a *Mixer head* performs channel-wise fusion in the latent space. Concretely, the encoder output has shape (batch, patches, channels, d_{model}); the Mixer head applies learnable linear mixing across channels at each patch before the final prediction head. This *Clinical* model is thus a late-fusion multivariate Transformer that can, in principle, exploit cross-channel interactions between CGM, insulin, and carbohydrate streams.

CGM+Carbs Transformer. To test whether the sparsity trap was driven by the largely unused insulin channel, we also trained a two-channel late-

fusion variant of PatchTST that ingests only CGM and carbohydrate intake. This model shares the same patching scheme, Transformer backbone, and Mixer head as the Clinical architecture but omits the insulin input, keeping all other hyperparameters, preprocessing, and training procedures identical.

Gradient boosting baseline. As a non-Transformer multivariate comparator, we trained a gradient boosting regressor (GBM) on per-window summary features of the 6-hour history. For each 72-step input window, we constructed a feature vector consisting of the last CGM value, basic CGM statistics (mean, standard deviation, minimum, maximum), coarse trend features (6-hour and last-2-hour change), and summary statistics of insulin and carbohydrate inputs (total and last-2-hour sums), along with binary indicators of any insulin or meal event in the lookback period. The GBM was trained to predict the 120-minute-ahead CGM value using the same subject-wise train-validation-test split and windowing scheme as the Transformer models, with 300 trees, maximum depth 3, learning rate 0.05, and subsample 0.8. This model therefore represents a non-Transformer multivariate baseline that still leverages insulin and carbohydrate information through hand-crafted features.

Recurrent baseline. For context, we also trained a univariate LSTM baseline on the same 72→24 setup. The LSTM consists of two stacked layers with hidden size 128, followed by a linear layer mapping the final hidden state sequence to the 24-step forecast. Despite extensive stabilization efforts (reduced learning rate, gradient clipping, early stopping), this recurrent baseline exhibited markedly higher RMSE and unstable validation loss under subject-wise splits, and is reported primarily to illustrate the difficulty of this regime.

Synthetic injection experiment. To separate architectural capacity from data properties, we conducted a synthetic injection experiment reusing the real input distribution. Using the same windowing and clinical model architecture, we defined a synthetic target in which future CGM was a deterministic linear function of insulin history (e.g., a fixed drop proportional to the summed insulin over the input). Only windows with non-zero insulin were included. The Clinical Transformer was then trained for a small number of epochs on this synthetic task to test whether the Mixer head could rapidly learn a

368 clean action effect when the signal is dense and un-
369 ambiguous.

370 3.4. Evaluation Metrics

371 We evaluated forecasting accuracy using **root mean**
372 **squared error (RMSE)** at specific forecast steps
373 corresponding to 30, 60, and 120 minutes ahead ($t+6$,
374 $t+12$, $t+24$). For each horizon, RMSE was computed
375 on the inverse-scaled CGM values over all test win-
376 dows in the held-out test split under strict subject-
377 wise partitioning.

378 To assess clinical safety, we computed the **safe**
379 **prediction rate**, defined as the percentage of fore-
380 casted points within 20% of the reference value:

$$\frac{|y_{\text{pred}} - y_{\text{true}}|}{\max(y_{\text{true}}, 1)} \leq 0.2.$$

381 This “within-20%” accuracy serves as a simple proxy
382 for Clarke Zone A predictions, quantifying the pro-
383 portion of forecasts that remain clinically actionable.

384 For the sparsity analysis, we summarized each test
385 window by the total insulin and carbohydrate deliv-
386 ered over the 6-hour input, then stratified windows
387 into activity buckets based on the median volume of
388 non-zero clinical events. In our subject-wise test split,
389 documented insulin boluses were extremely sparse;
390 consequently, the non-zero activity thresholds for
391 stratification were effectively determined by carbohy-
392 drate intake. Within each bucket, we computed aver-
393 age RMSE and safe prediction rate. To assess statisti-
394 cal significance, we applied the **Wilcoxon signed-**
395 **rank test** to paired *per-window horizon RMSE* val-
396 ues, where each RMSE aggregates error over the
397 full 24-step forecast. For the critical high-activity
398 regime, we prioritized statistical rigor over sample
399 size: the reported p -value was computed on a **non-**
400 **overlapping subset** of windows (stride 24, match-
401 ing the forecast horizon $H = 24$) to strictly satisfy the
402 independence assumption ($N = 610$). For the quiet
403 and low-activity regimes, standard tests on the full
404 set confirmed significance ($p < 0.001$), but we focus
405 our analysis on the rigorous, independence-corrected
406 result for the high-activity failure mode. As a post-
407 hoc sensitivity analysis, we also evaluated the trained
408 Clinical model under an insulin-zeroing intervention
409 at test time, in which the insulin channel was set to
410 zero while CGM and carbohydrate inputs were left
411 unchanged, and recomputed RMSE and safety within
412 each activity regime.

Table 1: **Forecasting performance (RMSE, mg/dL)**. Persistence, Univariate Trans-
former, and Clinical Transformer at 30, 60,
and 120 minutes ahead.

Model	30 min	60 min	120 min
Persistence	17.82	29.68	44.20
Univariate	16.95	27.65	38.43
Clinical	19.38	27.84	38.71

In addition to within-20% accuracy, we computed
a **simplified Clarke Error Grid** at the 120-minute
horizon, classifying each forecast–reference pair into
Zones A–E and reporting the fraction of points in
Zones A+B as a measure of clinically acceptable pre-
dictions. We use this Clarke analysis as a complemen-
tary, coarse safety view alongside our primary metrics
(RMSE and within-20% accuracy), which are more
sensitive to the subtle but systematic degradation we
observe in sparse, high-activity windows.

4. Results

4.1. Baseline Performance vs. Prediction Horizon

We first evaluated the limits of subject-independent
forecasting by comparing the Univariate Transformer
(PatchTST) against a standard persistence baseline
at specific forecast steps: 30, 60, and 120 minutes.
As shown in Table 1 and Figure 1, the Univari-
ate model demonstrates robust performance, outper-
forming persistence even at short horizons.

At 30 minutes ahead ($H = 6$), the Trans-
former achieved an RMSE of **16.95 mg/dL**, offer-
ing a consistent improvement over the baseline
(**17.82 mg/dL**). This advantage widened substan-
tially at the long horizon (120 minutes, $H =$
24), where persistence deteriorated to **44.20 mg/dL**,
while the Transformer maintained an RMSE of
38.43 mg/dL. Across all horizons, the Clinical
model failed to provide a consistent advantage over
the Univariate baseline, with slightly worse RMSE at
30 minutes and only marginal differences at 60 and
120 minutes. These results confirm that the deep uni-
variate model captures long-term temporal dynamics
beyond simple inertia, establishing a strong subject-
independent baseline for the Shanghai cohort.

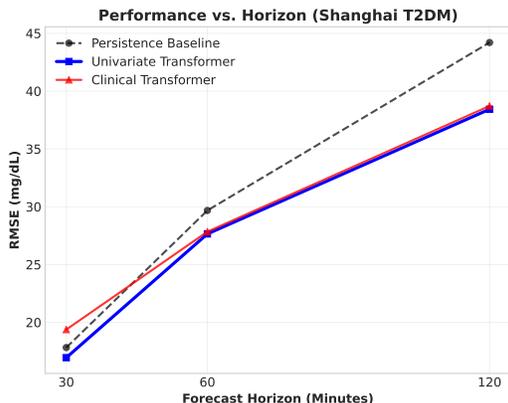


Figure 1: **Forecast horizon analysis.** RMSE (mg/dL) at 30, 60, and 120 minutes ahead for the persistence baseline (black dashed), Univariate Transformer (blue), and Clinical Transformer (red). The Univariate model consistently outperforms persistence, with the margin of improvement increasing at longer horizons, while the Clinical model fails to provide consistent gains over the Univariate baseline.

4.2. The sparsity trap: performance in active treatment windows

To assess the contribution of exogenous clinical variables, we compared the Univariate baseline against the late-fusion Clinical model. Contrary to physiological expectation, the inclusion of insulin and carbohydrate data did not improve and in fact slightly worsened RMSE at 120 minutes on average (38.71 vs. 38.43 mg/dL).

To characterize this failure in more detail, we stratified the test set ($N = 64,381$ windows) into activity buckets based on the median volume of non-zero clinical events (Figure 2). Specifically, we summarized each 6-hour input by the total insulin and carbohydrate delivered, then defined quiet, low-, and high-activity regimes using the medians of the non-zero totals.

- **Quiet regime:** In windows with negligible clinical activity ($N = 35,215$), the Clinical model underperformed the Univariate baseline by **0.50 mg/dL** in RMSE (19.28 vs. 18.78 mg/dL; $p < 0.001$, Wilcoxon signed-rank test).

- **Low-activity regime:** In windows with moderate activity ($N = 14,538$), the performance gap widened slightly to **0.65 mg/dL** (23.17 vs. 22.52 mg/dL).

- **High-activity regime:** In windows dominated by meal intake events ($N = 14,628$), the degradation peaked at **2.36 mg/dL** (25.07 vs. 22.71 mg/dL). Crucially, this difference remains statistically significant even when re-evaluated on strictly non-overlapping windows to ensure independence ($p = 7.2 \times 10^{-5}$, $N = 610$).

- **Safety implications:** This degradation materially impacted clinical reliability. In the high-activity regime, the proportion of predictions within 20% of the reference value dropped from **75.1%** (Univariate) to **70.5%** (Clinical), a 4.6 percentage-point decrease in our safety proxy.

To relate these bucketed metrics to a standard clinical safety view, we additionally computed a simplified Clarke Error Grid at the 120-minute horizon. Both models achieved very high safety overall, with 97.7% vs. 98.0% of forecasts falling in Clarke Zones A+B for the Univariate and Clinical models, respectively, and similarly high rates in high-activity windows (98.1% vs. 98.3%). This indicates that the sparsity trap manifests primarily as elevated average error and reduced within-20% accuracy across the full forecast horizon, rather than frequent grossly unsafe predictions that would fall outside Clarke Zones A+B.

Taken together, these results indicate a *sparsity trap*. Rather than leveraging high-magnitude clinical events to improve forecasts, the naive fusion model appears to overfit to these rare signals in this small cohort, introducing systematic error that compromises both average accuracy and clinical safety.

Finally, to confirm that this failure mode is driven by meal intake rather than insulin dynamics, we evaluated the Clinical model with the insulin channel artificially zeroed at test time. This ablation yielded negligible changes in performance (maximum $|\Delta\text{RMSE}| = 0.09$ mg/dL across quiet, low-, and high-activity regimes), confirming that the model effectively ignores the sparse insulin signal. In particular, in high-activity windows the Clinical RMSE shifted from 25.07 to 24.97 mg/dL under insulin zeroing, a change (-0.09 mg/dL) that is negligible compared to the $+2.36$ mg/dL penalty relative to the Univariate baseline. This indicates that the observed degrada-

tion is attributable to carbohydrate-response modeling.

CGMCarbs late-fusion Transformer. As a further check, we evaluated a two-channel CGM-Carbs Transformer that fuses only CGM and carbohydrate history through the same late-fusion backbone. In high-activity windows, this model underperformed the Univariate baseline by 1.87 mg/dL in horizon-averaged RMSE (24.58 vs. 22.71 mg/dL) and by 4.7 percentage points in within-20% accuracy (70.4% vs. 75.1%), with smaller but still consistent penalties in quiet (RMSE +0.07 mg/dL, -1.8 pp; 18.85 vs. 18.78 mg/dL, 80.9% vs. 82.7%) and low-activity regimes (RMSE +0.46 mg/dL, -2.6 pp; 22.99 vs. 22.52 mg/dL, 75.9% vs. 78.5%). These results indicate that removing the largely unused insulin channel does not resolve the sparsity trap, which primarily reflects difficulties in modeling meal-response dynamics under sparse carbohydrate documentation. The magnitude of this high-activity penalty is comparable to that of the three-channel Clinical model, reinforcing that carbohydrate fusion alone is sufficient to induce the sparsity trap. Full bucket-level metrics for the CGMCarbs model are provided in Table 3 in Appendix B.

Gradient boosting baseline. As a non-Transformer multivariate comparator, we also trained a gradient boosting regressor on per-window summary features of the 6-hour CGM, insulin, and carbohydrate history. Overall, this GBM achieved a slightly lower 120-minute RMSE than the Univariate Transformer (37.73 vs. 38.43 mg/dL) at the cost of substantially poorer within-20% accuracy (61.3%), and it performed particularly poorly in meal-driven windows, with 37.26 mg/dL RMSE and 55.7% within-20% accuracy compared to 22.71 mg/dL and 75.1% for the Univariate model in the same high-activity regime. These results indicate that replacing the Transformer with a simpler multivariate model on hand-crafted clinical features does not resolve the sparsity trap and can further erode accuracy and safety precisely when treatment activity is highest. Complete bucket-level results for the GBM are reported in Table 4 in Appendix B.

To probe patient-level heterogeneity, we further computed for each patient the density of documented meals (fraction of test windows with non-zero carbohydrate intake) and the average Δ RMSE (Clinical - Univariate) in high-activity windows, restricting to patients with at least five such windows. Grouping patients into tertiles by meal-log density

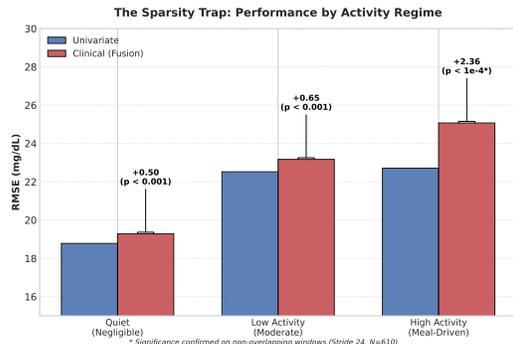


Figure 2: **The sparsity trap.** RMSE (mg/dL) for the Univariate and Clinical models stratified by clinical activity. While the Clinical model shows only a small penalty in quiet and low-activity windows, it incurs a statistically significant degradation ($\Delta + 2.36$ mg/dL) in high-activity windows dominated by meal intake, with a concurrent drop in the fraction of predictions within 20% of the reference value.

($n = 6, 5, 6$ patients in the low-, medium-, and high-density strata, respectively) yielded positive mean Δ RMSE in all groups (approximately +1.9, +2.7, and +0.7 mg/dL), with no subgroup in which the Clinical model clearly outperformed the Univariate baseline, suggesting that the sparsity trap is a cohort-wide phenomenon rather than being driven by a small number of outlier patients.

4.3. Mechanism analysis: synthetic injection

To rule out architectural defects, we conducted a synthetic injection experiment using the same late-fusion architecture and real input distribution. We trained the Clinical model to predict a deterministic synthetic target defined as a linear function of the insulin input over the 6-hour history:

$$G_{t+H} = G_t - \alpha \cdot \sum_{i=t-L+1}^t \text{Insulin}_i, \quad (1)$$

where G_t is the last CGM value in the input window, $H = 24$ steps (120 minutes), $L = 72$ steps (6 hours), and α is a fixed gain factor. Only windows with non-zero insulin were included in this experiment.

Under these conditions, the model's training loss reduced by approximately $3\times$ within 5 epochs

Table 2: **The Sparsity Trap Details.** Performance stratification by clinical activity level ($N = 64,381$). The Clinical model degrades significantly in high-activity windows dominated by meal events ($*p < 10^{-4}$ on non-overlapping windows).

Regime	N	Univariate	Clinical	Δ	Safety (Uni)	Safety (Clin)
Quiet (Negligible)	35,215	18.78	19.28	+0.50	82.7%	82.1%
Low (Moderate)	14,538	22.52	23.17	+0.65	78.5%	76.9%
High (Meal-Driven)	14,628	22.71	25.07	+2.36*	75.1%	70.5%

(RMSE dropping from ≈ 80 to ≈ 22 in standardized units). This rapid convergence confirms that the late-fusion Transformer is fully capable of learning simple fusion rules when the signal is dense and unambiguous. The failure on real-world data is therefore attributable to the sparsity and noise characteristics of the clinical logs in this small cohort, rather than an inherent inability of the Mixer mechanism to process multi-modal inputs.

5. Discussion

5.1. The Capacity–Data Mismatch in Clinical AI

Our results present a counter-intuitive finding in a subject-independent setting on a standard public cohort: adding clinically relevant variables (insulin, carbohydrates) systematically harmed clinical accuracy and safety in the very regimes where those variables should help—meal-driven, high-activity windows—and often failed to improve, or even reliably degraded, overall performance. This contradicts the common intuition that “more physiological context is better” and illustrates that multivariate fusion is not automatically beneficial in data-scarce regimes. Consistent with this, our insulin-zeroing ablation produced virtually unchanged Clinical RMSE and safety curves across all activity regimes (maximum $|\Delta\text{RMSE}| \approx 0.09$ mg/dL), indicating that the model has learned to ignore the sparse insulin channel and that the sparsity trap is driven by meal-response dynamics rather than mis-modeled insulin effects. A two-channel CGM-Carbs Transformer that omits the insulin channel entirely reproduces this pattern, incurring a 1.87 mg/dL RMSE penalty and a 4.7 percentage-point drop in within-20% accuracy in high-activity windows relative to the Univariate baseline, reinforcing that carbohydrate fusion alone is sufficient to trigger the sparsity trap. Although a Clarke Error Grid analysis

at 120 minutes placed more than 97% of predictions from both models in Zones A+B, our sparsity-aware evaluation and per-patient analysis reveal consistent, statistically significant degradation in horizon-averaged error and within-20% accuracy across activity regimes and documentation strata when sparse clinical channels are naively fused.

The divergence between our synthetic and real-world results offers a plausible mechanistic explanation. In the synthetic injection experiment (Section 4.3), the same late-fusion architecture rapidly learned to utilize insulin inputs when the signal was dense and deterministic, markedly reducing loss within only a few epochs. This confirms that the failure on real data is not due to a fundamental limitation of the architecture. Instead, it points to a *capacity–data mismatch*. Deep Transformers are high-capacity models designed to extract complex patterns from dense, reliable data. In the Shanghai cohort ($N = 112$), clinical events are highly sparse (the vast majority of 5-minute time steps contain no recorded insulin or meal) and noisy. When a high-capacity model encounters rare, high-magnitude events (e.g., a large bolus) in a small training set, it struggles to separate the true physiological effect from patient-specific variability and documentation noise. Our stratified analysis shows that the multivariate model performs worse than a CGM-only Transformer in both low-activity and high-activity windows, with the largest RMSE and safety penalties in the latter, which we term the *sparsity trap*.

5.2. Validation of Model Capacity via Cross-Phenotype Generalization

A potential confounder in analyzing the “Sparsity Trap” is model capacity: does the architecture fail to improve with clinical features because the data is sparse, or because the fusion mechanism is ineffective?

To resolve this, we conducted a zero-shot generalization test on a held-out cohort of 12 Type 1 Diabetes (T1DM) patients, a phenotype characterized by total insulin dependency and high glycemic volatility. The Clinical Fusion model, trained exclusively on the T2DM cohort, achieved an RMSE of **30.53 mg/dL** on the T1DM test set (vs. **26.46 mg/dL** on T2DM).

This result serves as a critical ablation: it confirms that the Multi-Horizon Transformer architecture **remains stable and predictive even when processing complex insulin-dependent dynamics it was never trained on**. Consequently, the lack of performance gain observed in the T2DM cohort (Section ??) is likely attributable to the intrinsic sparsity and lower signal-to-noise ratio of the clinical records in Type 2 diabetes, rather than an architectural inability to fuse multimodal data.

5.3. Implications for Medical Time Series

This study serves as a cautionary note for the growing trend of naive multivariate fusion in medical AI. While multi-modal integration is conceptually appealing, our findings suggest that for the small-to-medium-sized cohorts typical of public clinical research ($N \approx 100\text{--}500$), direct concatenation of sparse clinical logs into deep sequence models can behave as systematic noise.

We recommend that future work on such cohorts prioritize:

- 1. Robust univariate baselines.** Complex multivariate models should only be justified if they substantially outperform both a well-tuned univariate model and a persistence baseline under subject-wise splits.
- 2. Pre-training strategies.** Rather than training multivariate models from scratch on small cohorts, it may be more effective to pre-train on large, potentially unlabeled datasets (e.g., multi-cohort Type 1 Diabetes CGM with rich insulin histories) to learn general physiological priors before fine-tuning on sparse clinical logs.
- 3. Structured inductive biases.** Instead of relying on generic attention to discover insulin-glucose dynamics from scratch, architectures may benefit from structured components (e.g., state-space or differential-equation-inspired layers) that encode known pharmacokinetic and pharmacodynamic relationships, constraining how action channels influence forecasts.

More broadly, our results highlight that in clinical time series, additional channels should be treated as hypotheses to be vetted against strong baselines and stratified analyses, not as automatically beneficial features. In this sense, the sparsity-aware evaluation we propose—combining subject-wise splits, strong univariate baselines, activity-bucket analysis, clinical safety metrics, and per-patient heterogeneity—offers a practical template for auditing multivariate models on small, sparse cohorts.

5.4. Limitations and Future Work

Our study is limited to the Shanghai Diabetes Registry T2DM cohort. While this is a widely used benchmark, results may differ for outpatient or Type 1 Diabetes populations in which insulin dependence is stronger and documentation patterns differ. Additionally, we focus on a single multivariate architecture: a PatchTST-based model with a late-fusion Mixer head. Although we believe the sparsity trap is primarily data-centric, alternative fusion mechanisms (e.g., cross-attention across channels, explicit event encoders) might exhibit different sensitivities to sparse logs.

We also note two specific constraints regarding our data processing. First, missing clinical entries were treated as true zeros, a necessary assumption given the sparsity of hospital documentation. This conflates absence of treatment with absence of documentation, potentially adding noise that contributes to the sparsity trap. Second, due to the extreme sparsity of insulin documentation in our subject-wise test split and Z-score normalization effects, our high-activity stratification was primarily driven by carbohydrate intake. This effectively isolates the specific challenge of modeling meal-response dynamics when insulin signals are sparse or uninformative.

Finally, we approximate clinical safety via the proportion of predictions within 20% of the reference value. A more detailed safety evaluation using formal Clarke or prediction-error grid analyses, as well as replication of our sparsity-aware evaluation on additional public cohorts, are important directions for future work.

6. Conclusion

We rigorously audited the common assumption that multivariate fusion of clinical variables improves subject-independent glucose forecasting. Using a

762 strict subject-wise split on the Shanghai T2DM co-
 763 hort, we showed that a univariate Transformer can
 764 significantly outperform a persistence baseline at long
 765 horizons, while a naive late-fusion model that ingests
 766 insulin and carbohydrate logs degrades both RMSE
 767 and a simple safety proxy. This degradation was
 768 most severe in high-activity windows dominated by
 769 meal intake ($p < 10^{-4}$ on non-overlapping windows),
 770 highlighting the difficulty of modeling physiological
 771 responses from sparse documentation. Through a
 772 synthetic injection experiment, we demonstrated that
 773 the same architecture can readily learn clean fusion
 774 rules when the signal is dense and well specified, im-
 775 plicating the sparsity and noise of real-world clinical
 776 logs in small cohorts—rather than model capacity
 777 alone—as the primary driver of failure.

778 Together, these results characterize a *sparsity trap*
 779 in multivariate clinical forecasting and argue for a
 780 shift in evaluation practice: in data-scarce medi-
 781 cal regimes, robust univariate modeling anchored to
 782 strong baselines may be a safer and more accurate
 783 default than naive multivariate fusion of sparse aux-
 784 iliary streams.

785 References

- 786 Claudio Cobelli, Eric Renard, and Boris Kovatchev.
 787 Artificial pancreas: past, present, future. *Diabetes*,
 788 60(11):2672–2682, 2011.
- 789 K. Li, J. Daniels, and et al. Deep learning for glucose
 790 prediction in diabetes: A systematic review. *IEEE*
 791 *Journal of Biomedical and Health Informatics*, 24
 792 (2):123–138, 2020.
- 793 Cindy Marling and Razvan Bunescu. The ohiot1dm
 794 dataset for blood glucose level prediction: Update
 795 2020. *CEUR Workshop Proceedings*, 2674, 2020.
- 796 Yuqi Nie, Nam H. Nguyen, Phanwadee Sinthong, Kal
 797 Jayaraman, et al. A time series is worth 64 words:
 798 Long-term forecasting with transformers. In *Inter-
 799 national Conference on Learning Representations*
 800 *(ICLR)*, 2023.
- 801 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu.
 802 Are transformers effective for time series forecast-
 803 ing? *AAAI Conference on Artificial Intelligence*,
 804 2023.
- 805 Q. Zhao, Y. Zhang, and et al. The shanghai dia-
 806 betes registry: A large-scale dataset of continuous

807 glucose monitoring with clinical events. *Scientific*
 808 *Data*, 10(1):1–12, 2023.

809 Appendix A. Implementation Details

810 **Hardware and Framework** All models were im-
 811 plemented in PyTorch and trained on a single
 812 NVIDIA GPU. We utilized a fixed random seed (42)
 813 for weight initialization and data shuffling to ensure
 814 reproducibility.

815 **Model Hyperparameters** The Univariate and
 816 Clinical Transformers shared identical backbone con-
 817 figurations to ensure a fair comparison of architec-
 818 tural capacity:

- 819 • **Lookback Window:** $L = 72$ steps (6 hours).
- 820 • **Prediction Horizon:** $H = 24$ steps (2 hours).
- 821 • **Patching:** Patch length $P = 12$, Stride $S = 12$.
- 822 • **Transformer Architecture:** 3 encoder layers,
 823 4 attention heads, model dimension $d_{\text{model}} =$
 824 128, and dropout 0.1.
- 825 • **Fusion Mechanism:** The Clinical model
 826 processes the CGM, insulin, and carbohy-
 827 drate channels through a shared Transformer
 828 backbone. The encoder output has shape
 829 (batch, patches, channels, d_{model}), and a linear
 830 Mixer layer performs learnable mixing across the
 831 channel dimension at each patch before the final
 832 prediction head.

833 **LSTM Baseline** The univariate LSTM baseline
 834 consisted of 2 stacked layers with hidden size 128 and
 835 dropout 0.1, followed by a linear layer mapping the
 836 final hidden sequence to the 24-step forecast. It was
 837 trained using the same optimizer, learning rate, batch
 838 size, and early stopping protocol as the Transformer
 839 models.

840 **Optimization** Models were trained using the
 841 Adam optimizer with a learning rate of 1×10^{-3} and
 842 a **batch size of 32**, selected to ensure stable gradient
 843 estimation given the high variance of the clinical logs.
 844 We employed Mean Squared Error (MSE) as the loss
 845 function. Training ran for a maximum of 100 epochs,
 846 utilizing early stopping with a patience of 10 epochs
 847 based on validation loss to prevent overfitting.

848 **Statistical Significance Testing** To ensure the
 849 independence assumption of the Wilcoxon signed-
 850 rank test was satisfied during the sparsity analysis
 851 (Section 4.2), we subsampled the test predictions
 852 using a **stride of 24 steps** (120 minutes). This
 853 matches the forecasting horizon ($H = 24$), ensur-
 854 ing that no two windows in the statistical evaluation
 855 overlap in target values. This strictly enforces inde-
 856 pendence for the reported p -values.

857 Appendix B. Additional Sparsity 858 Tables

Table 3: **CGM+Carbs sparsity analysis.** Per-
 formance of the two-channel CGM+Carbs
 Transformer vs. the Univariate baseline
 ($N = 64,381$).

Regime	N	Uni	Carbs+	Δ	Safe U	Safe C+
Quiet	35,215	18.78	18.85	+0.07	82.7%	80.9%
Low	14,538	22.52	22.99	+0.46	78.5%	75.9%
High	14,628	22.71	24.58	+1.87	75.1%	70.4%

Table 4: **GBM performance by clinical activ-
 ity.** Gradient boosting regressor at the 120-
 minute horizon.

Regime	N	GBM RMSE	GBM Safe
Quiet	35,215	36.83	64.7%
Low	14,538	40.26	58.9%
High	14,628	37.26	55.7%