
Volatility-Aware Masking Improves Performance and Efficiency of Pretrained EHR Foundation Models

Rajna Fani^{1,2,*}, Rafi Al Attrach^{1,2,*}, David Restrepo³, Yugang Jia¹,
Leo Anthony Celi^{1,6,7,†}, Peter Schüffler^{4,5,†}

¹Massachusetts Institute of Technology (MIT), USA

²Technical University of Munich (TUM), Germany

³MICS, CentraleSupélec – Université Paris-Saclay, France

⁴Institute of Pathology, Technical University of Munich, Germany

⁵Munich Center for Machine Learning (MCML), Germany

⁶Harvard Medical School, USA

⁷Beth Israel Deaconess Medical Center, USA

{rajnaf, rafiaa, yugang, lceli}@mit.edu, david.restrepo@centralesupelec.fr,
peter.schueffler@tum.de

Abstract

Masked autoencoder (MAE) models are increasingly applied to electronic health records (EHR) as a pre-training method to learn general-purpose representations that support diverse downstream clinical tasks. However, existing approaches typically rely on uniform random masking, implicitly assuming that all clinical features are equally predictable. In practice, laboratory tests exhibit substantial heterogeneity in temporal volatility: certain biomarkers (e.g., sodium) remain relatively stable, whereas others (e.g., lactate) fluctuate considerably and are more challenging to model. To address this limitation, we propose Volatility-Aware Masking strategy (CV-Masking), a pretraining strategy that adaptively adjusts masking probabilities according to the intrinsic variability of each feature. Our experiments on a large panel of laboratory tests demonstrate that CV-Masking consistently outperforms both random and variance-based masking strategies, yielding improved downstream predictive performance and faster convergence.

1 Introduction

Foundation models are transforming healthcare research by enabling the learning of general-purpose representations from large-scale electronic health records (EHR) [11, 15]. Depending on the architecture and application, these models are pretrained using diverse strategies and architectures, including state-space models [4], generative modeling [18, 14], contrastive learning [7], or zero-shot transfer [16]. Among these approaches, masked autoencoders (MAEs) have emerged as a powerful framework for representation learning in EHR, reconstructing masked inputs from partially observed sequences [17]. For laboratory test modeling, MAEs are particularly relevant since they can accurately reconstruct missing values, not only reflecting representation quality but also enabling clinically meaningful applications such as decision support and risk prediction [6, 2, 17].

However, existing MAE pretraining strategies almost exclusively adopt uniform random masking, implicitly assuming that all features are equally predictable. This is especially critical in clinical data

*Shared first authors: Rajna Fani and Rafi Al Attrach

†Shared corresponding authors: Leo Anthony Celi and Peter Schüffler

26 such as lab values, where biomarkers differ substantially in their temporal volatility. For example,
27 sodium levels remain tightly regulated, whereas lactate varies considerably during acute illness.
28 Ignoring this variability can waste model capacity, slow convergence, and limit the clinical utility of
29 learned representations.

30 To address this gap, we propose Volatility-Aware Masking (CV-Masking), a pretraining strategy
31 that adapts masking probabilities to the intrinsic variability of each laboratory test. Inspired by
32 curriculum learning [3, 5] and informed masking policies [10, 12], CV-Masking focuses learning
33 on less predictable signals, improving both efficiency and representation quality. Similarly, for
34 tabular imputation tasks, recent work adjusts masking based on observed missingness proportions [9],
35 though our temporal approach prioritizes volatile features in sequences rather than static missingness
36 patterns.

37 Our contributions are threefold:

- 38 1. **CV-Masking Strategy.** A principled masking policy guided by the coefficient of variation
39 (CV). By prioritizing inherently volatile laboratory values, CV-Masking creates a natural
40 curriculum that directs learning capacity toward clinically uncertain signals.
- 41 2. **Value-Only Masked Autoencoder Objective (VO-MAE).** We adapt the masked autoen-
42 coder framework to a VO-MAE, where lab identifiers and timestamps remain visible while
43 only results are masked. This design mirrors real-world clinical workflows—orders are
44 observed, but outcomes are unknown—and encourages the model to learn meaningful value
45 representations in temporal context.
- 46 3. **Comprehensive Empirical Validation.** Using 100 high-frequency laboratory tests from
47 MIMIC-IV [8], we show that CV-Masking (i) improves reconstruction accuracy on 71% of
48 labs, (ii) enhances downstream prediction of in-ICU mortality, in-hospital mortality, and
49 30-day readmission, and (iii) achieves up to 50% faster convergence compared to random
50 masking. Perturbation analysis further demonstrates that CV-Masking promotes deeper
51 reliance on patient-specific temporal context.

52 These results demonstrate that integrating clinical volatility into masking strategies substantially
53 improves the efficiency, robustness, and clinical relevance of MAE-based EHR foundation models.³

54 2 Methods

55 2.1 Data and Preprocessing

56 The experiments use data from the MIMIC-IV v3.1 [8] critical care database, structured in the MEDS
57 format [1]. This representation organizes a patient’s history into a sequence of (time, code, value)
58 triplets, where each triplet represents a laboratory measurement with its timestamp, test identifier,
59 and numeric result. To ensure robust statistical comparisons, our evaluation focuses on a fixed set of
60 100 target laboratory tests representing the most clinically relevant and frequently ordered tests in
61 critical care settings. This selection follows established practices in clinical MAE literature [17, 2]
62 and avoids statistical noise from rare or infrequently ordered laboratories.

63 2.2 Architecture

64 Our model follows an asymmetric encoder-decoder MAE framework, with a transformer encoder
65 processing unmasked events and a lightweight transformer decoder reconstructing masked values
66 (see Appendix Figure 3 for complete architecture details). All three masking strategies (random,
67 variance-based, and CV-Masking) employ the identical VO-MAE architecture, differing only in the
68 masking probability distribution assigned to each laboratory test, ensuring performance differences
69 are attributable solely to the masking strategy.

70 The key innovation is our **value-only masking objective**: given input sequences of (time, code,
71 value) triplets, we mask only the value components while preserving temporal and categorical context,
72 enabling the model to focus learning on the challenging value prediction task. During training, each
73 laboratory value is masked independently with probability proportional to its assigned weight (0.8 for

³Code available at <https://github.com/rajna-fani/meds-triplet-mae>

74 high-CV labs, 0.2 for low-CV labs), maintaining an overall target masking ratio of approximately
75 25% across the full sequence.

76 We train with mean squared error using a joint objective:

$$\mathcal{L} = \mathcal{L}_{\text{masked}} + \lambda \mathcal{L}_{\text{unmasked}} \quad (1)$$

77 where $\mathcal{L}_{\text{masked}}$ is the primary reconstruction loss on masked positions, $\mathcal{L}_{\text{unmasked}}$ regularizes visible
78 positions to prevent representation collapse, and $\lambda = 0.1$ balances the objectives (selected via grid
79 search over $\{0.05, 0.1, 0.2\}$ on validation data).

80 2.3 Masking Policies

81 We implement and compare three distinct masking strategies:

- 82 • **Random Masking (Baseline):** Each lab value has a uniform masking probability, treating
83 all tests as equally learnable.
- 84 • **Variance-Based Masking:** Masking probability proportional to raw variance (σ^2). This
85 scale-sensitive approach overweights labs with large numerical ranges regardless of relative
86 clinical volatility.
- 87 • **CV-Based Masking (Proposed):** Our proposed principled approach assigns masking
88 weights based on the 75th percentile threshold of CV values: laboratories with CV above the
89 threshold receive weight 0.8, while those below receive weight 0.2. The CV is calculated
90 as $CV = \sigma/\mu$, a dimensionless measure of relative variability well-suited for weighting
91 heterogeneous lab tests. This binary assignment enables the masking policy to prioritize truly
92 volatile labs requiring sophisticated temporal modeling. Masking weights are precomputed
93 from training statistics and remain fixed.

94 3 Experiments and Results

95 3.1 Intrinsic Evaluation: CV-Based Masking Systematically Improves Reconstruction

96 To establish that principled masking strategies are needed, we analyzed which laboratory character-
97 istics predict imputation difficulty. Evaluating reconstruction performance (R^2 score) against CV
98 across 100 labs revealed that CV serves as a meaningful predictor of imputation difficulty (Pearson’s
99 $r = -0.486$, $p < 0.000001$), with CV explaining 23.6% of variance in reconstruction performance.
100 While the relationship shows expected variability across diverse laboratory types, this provides
101 statistical support for using CV to guide masking strategies.

102 CV-based masking consistently outperforms both baselines, achieving superior reconstruction on
103 **71.0%** of laboratories compared to random masking and **68.0%** compared to variance-based masking,
104 with large effect sizes (Cohen’s $d = 0.73$) and statistically significant improvements ($p < 0.000009$
105 after Bonferroni correction, Wilcoxon signed-rank test) exceeding chance expectation. Figure 1
106 demonstrates both the statistical foundation and systematic improvements across laboratory types.

107 Win rates: CV vs. Random (71.0%), Variance vs. Random (65.0%), CV vs. Variance (68.0%), all
108 statistically significant after Bonferroni correction.

109 3.2 Extrinsic Evaluation: Downstream Tasks Evaluation

110 We evaluated pretrained encoders on three high-stakes downstream prediction tasks using linear
111 probes on frozen representations to isolate representation quality from downstream model capacity
112 effects. Following the evaluation protocol established in MEDS-Torch [13], we assess performance
113 on in-ICU mortality, in-hospital mortality, and 30-day readmission prediction tasks. These tasks
114 represent clinically meaningful outcomes with varying prediction horizons and class imbalance
115 characteristics.

116 The CV-based model achieves the highest performance across all tasks and metrics (Table 1).
117 The CV-based approach achieves an in-ICU mortality AUROC of 0.713, representing meaningful
118 improvements over random masking (0.682) and variance-based masking (0.694). The 0.031 AUROC
119 improvement represents meaningful clinical impact, potentially enabling earlier interventions and

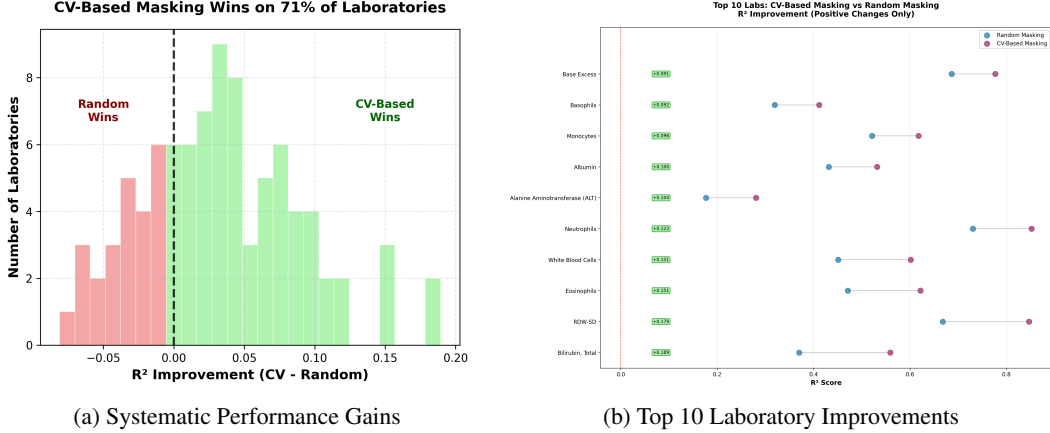


Figure 1: **CV-Based Masking Systematically Improves Reconstruction Performance.** (a) Systematic performance gains show CV-based masking wins on 71% of all 100 laboratories, demonstrating systematic rather than random improvements. (b) Top 10 laboratory improvements with most significant R^2 gains, sorted by magnitude.

improved outcomes. The substantial gains in AUPRC, a metric sensitive to minority class performance, are particularly significant, suggesting learned representations better capture subtle patterns indicative of adverse clinical outcomes.

Table 1: Downstream Task Performance on Clinical Prediction (linear probes on frozen encoders). For each masking strategy, we trained one encoder and evaluated it with 3 independent downstream probe runs using different random seeds. Standard deviations capture probe-training variability and are reported rounded to 3 significant figures.

Task	Metric	Random	Variance	CV-Based
In-ICU Mortality	AUROC	0.682 ± 0.017	0.694 ± 0.017	0.713 ± 0.017
	AUPRC	0.083 ± 0.009	0.091 ± 0.009	0.107 ± 0.009
In-Hospital Mortality	AUROC	0.657 ± 0.014	0.668 ± 0.014	0.691 ± 0.014
	AUPRC	0.124 ± 0.007	0.131 ± 0.007	0.149 ± 0.007
30-Day Readmission	AUROC	0.618 ± 0.016	0.627 ± 0.016	0.648 ± 0.016
	AUPRC	0.156 ± 0.008	0.162 ± 0.008	0.173 ± 0.008

3.3 Mechanistic Analysis: CV-Masking Promotes Deeper Contextual Learning

Analysis of reconstruction error against patient history reveals that CV-based masking more effectively leverages available context (see Appendix Figure 5). While all methods improve with more historical data, the performance gap between CV-based models and baselines widens as more patient history becomes available, indicating superior contextual learning.

To validate that CV-based masking learns meaningful patient-specific temporal patterns rather than simple memorization, we implemented a controlled perturbation experiment. We corrupted historical laboratory values with Gaussian noise while preserving target predictions and temporal context. CV-based models exhibited $2.1 \times$ greater performance degradation compared to random masking (9.8% vs 4.7% MAE increase, Figure 2), providing causal evidence that CV-based masking learns meaningful temporal patterns across diverse laboratory categories including hematology markers (Basophils, Monocytes), electrolytes (Sodium, Potassium), and metabolic indicators (Triglycerides). See Appendix 4.3 for detailed methodology.

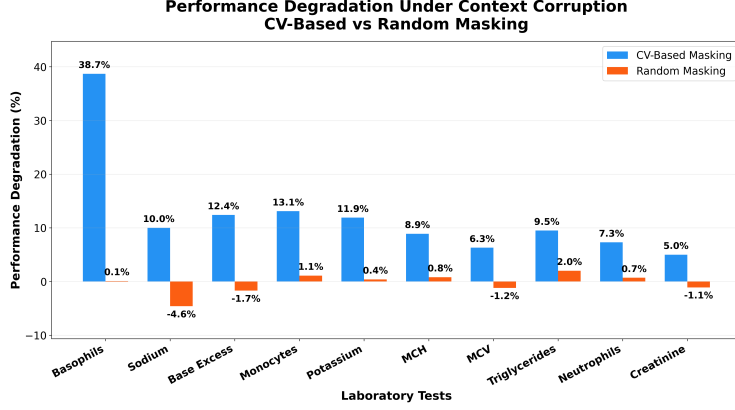


Figure 2: **Perturbation analysis demonstrates deeper contextual learning in CV-based models.** When historical laboratory values are artificially corrupted, CV-based models show 2.1× greater performance degradation across representative laboratory types (n=10). Higher degradation indicates stronger reliance on historical context (desirable for learning patient-specific patterns). Notable effects observed in clinically critical markers: hematology (Basophils, Monocytes), electrolytes (Sodium, Potassium), and metabolic indicators (Triglycerides).

4 Discussion and Conclusion

Our results demonstrate that clinically-informed masking curricula are essential for EHR foundation models operating in complex healthcare environments. CV-Masking achieves consistent improvements across reconstruction, downstream prediction, and mechanistic analyses, establishing clinical volatility as a meaningful signal for guiding pretraining strategies. This work advances domain-aware self-supervised learning by moving beyond naive architectural adaptations toward principled integration of medical expertise into foundational model training.

While our focus was on CV-Masking’s clinical impact, future studies should consider comprehensive ablation analyses to disentangle value-only masking contributions from CV-Masking policy effects. Such investigations could further illuminate the complementary roles of architectural choices and masking strategies across diverse clinical domains.

CV-Masking offers a principled pathway toward more efficient and robust EHR foundation models. The CV-Masking principle is architecture-agnostic and could be adapted to other EHR foundation models through loss reweighting or prioritized sampling strategies. As the field moves toward deployment-ready clinical AI systems, incorporating clinical knowledge into self-supervised objectives appears essential for developing models that are both technically sound and clinically meaningful.

References

- [1] Bert Arnrich, Edward Choi, Jason A. Fries, Matthew B. A. McDermott, Jungwoo Oh, Tom J. Pollard, Nigam Shah, Ethan Steinberg, Michael Wornow, and Robin van de Water. Medical event data standard (meds): Facilitating machine learning for health. In *ICLR 2024 Workshop on Learning from Time Series for Health*, 2024.
- [2] David R Bellamy, Bhawesh Kumar, Cindy Wang, and Andrew Beam. Labrador: Exploring the limits of masked language modeling for laboratory data. *arXiv preprint arXiv:2312.11502*, 2023.
- [3] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.
- [4] Adibvafa Fallahpour, Mahshid Alinoori, Wenqian Ye, Xu Cao, Arash Afkanpour, and Amrit Krishnan. Ehrmamba: Towards generalizable and scalable foundation models for electronic health records. *arXiv preprint arXiv:2405.14567*, 2024.

- [5] Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *international conference on machine learning*, pages 1311–1320. Pmlr, 2017.
- [6] Sujeong Im, Jungwoo Oh, and Edward Choi. Labtop: A unified model for lab test outcome prediction on electronic health records. *arXiv preprint arXiv:2502.14259*, 2025.
- [7] Hyewon Jeong, Nassim Oufattole, Aparna Balagopalan, Matthew B. McDermott, Payal Chandak, Marzyeh Ghassemi, and Collin Stultz. Event-based contrastive learning for medical time series, 2023.
- [8] Alistair EW Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gayles, Ayad Shammout, Steven Horng, Tom J Pollard, Sicheng Hao, Benjamin Moody, Brian Gow, et al. Mimic-iv, a freely accessible electronic health record dataset. *Scientific data*, 10(1):1, 2023.
- [9] Jungkyu Kim, Kibok Lee, and Taeyoung Park. To predict or not to predict? proportionally masked autoencoders for tabular data imputation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 17886–17894, 2025.
- [10] Yoav Levine, Barak Lenz, Opher Lieber, Omri Abend, Kevin Leyton-Brown, Moshe Tennenholtz, and Yoav Shoham. Pmi-masking: Principled masking of correlated spans. *arXiv preprint arXiv:2010.01825*, 2020.
- [11] Yikuan Li, Shishir Rao, José Roberto Ayala Solares, Abdelaali Hassaine, Rema Ramakrishnan, Dexter Canoy, Yajie Zhu, Kazem Rahimi, and Gholamreza Salimi-Khorshidi. Behrt: transformer for electronic health records. *Scientific reports*, 10(1):7155, 2020.
- [12] Neelu Madan, Nicolae-Cătălin Ristea, Kamal Nasrollahi, Thomas B Moeslund, and Radu Tudor Ionescu. Cl-mae: Curriculum-learned masked autoencoders. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2492–2502, 2024.
- [13] Nassim Oufattole, Teya Bergamaschi, Pawel Renc, Aleksia Kolo, Matthew BA McDermott, and Collin Stultz. Meds-torch: An ml pipeline for inductive experiments for ehr medical foundation models. In *NeurIPS Workshop on Time Series in the Age of Large Models*, 2024.
- [14] Chao Pang, Xinzhuo Jiang, Nishanth Parameshwar Pavinkurve, Krishna S Kalluri, Elise L Minto, Jason Patterson, Linying Zhang, George Hripcsak, Gamze Gürsoy, Noémie Elhadad, et al. Cehr-gpt: Generating electronic health records with chronological patient timelines. *arXiv preprint arXiv:2402.04400*, 2024.
- [15] Laila Rasmy, Yang Xiang, Ziqian Xie, Cui Tao, and Degui Zhi. Med-bert: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction. *npj Digital Medicine*, 4:86, 2021. doi: 10.1038/s41746-021-00455-y.
- [16] Pawel Renc, Yugang Jia, Anthony E Samir, Jaroslaw Was, Quanzheng Li, David W Bates, and Arkadiusz Sitek. Zero shot health trajectory prediction using transformer. *NPJ digital medicine*, 7(1):256, 2024.
- [17] David Restrepo, Chenwei Wu, Yueran Jia, Jaden K Sun, Jack Gallifant, Catherine G Bielick, Yugang Jia, and Leo A Celi. Representation learning of lab values via masked autoencoder. *arXiv preprint arXiv:2501.02648*, 2025.
- [18] Zhichao Yang, Avijit Mitra, Weisong Liu, Dan Berlowitz, and Hong Yu. Transformehr: transformer-based encoder-decoder generative model to enhance prediction of disease outcomes using electronic health records. *Nature communications*, 14(1):7857, 2023.

Appendix

4.1 Architecture Overview

Training Details: Models are trained with a masking ratio of 25%, using AdamW optimizer with learning rate $1e-4$ and weight decay 0.01, for up to 100 epochs with early stopping (patience = 10

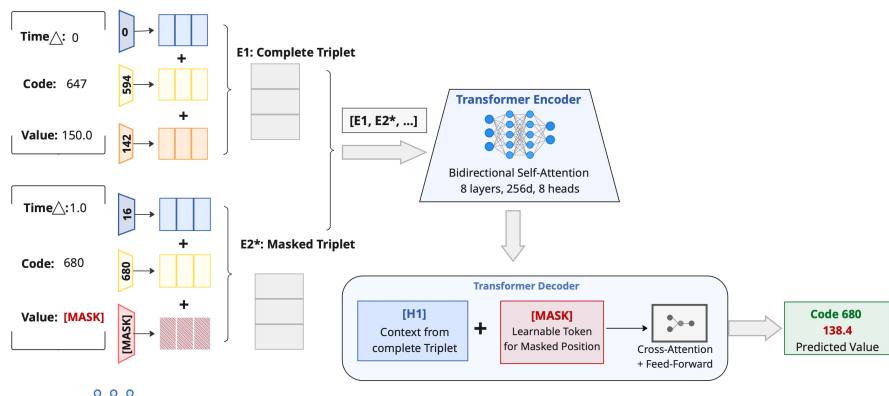


Figure 3: **Value-Only MAE Architecture with CV-Based Masking.** The complete workflow from MEDS triplet input through pretraining to downstream fine-tuning. During pretraining, only laboratory values are masked while preserving temporal and categorical context. The lightweight decoder reconstructs masked values using cross-attention between learnable mask tokens and encoded visible representations. For downstream tasks, the frozen encoder provides contextualized representations to task-specific classifiers.

epochs). We use a batch size of 32 and normalize laboratory values using z-score standardization computed per-lab across the training set. The model architecture consists of 8 encoder layers with 256 embedding dimensions and 8 attention heads.

4.2 Computational Efficiency

While variance-based masking required 100 epochs and random masking required approximately 67 epochs to converge, CV-based masking achieved optimal performance in just 33 epochs - representing a 67% and 50% reduction in training time respectively. This efficiency gain likely stems from the principled curriculum learning effect [3, 5], where the model focuses computational resources on truly challenging prediction tasks rather than easy, low-variance tests.

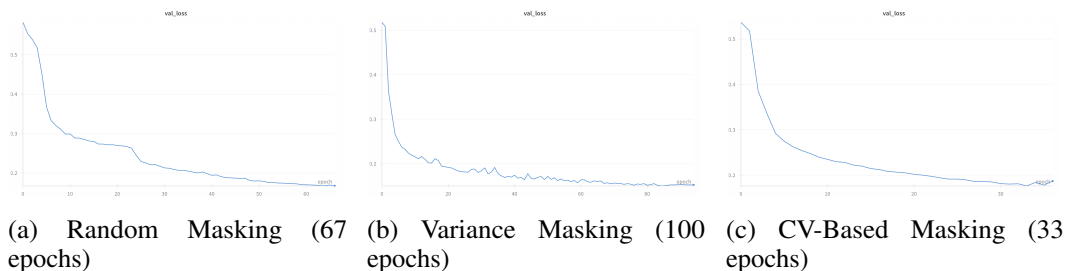


Figure 4: **Training Efficiency Comparison.** Training efficiency comparison across masking strategies. CV-Masking converges in 33 epochs, representing a 50% reduction in training time compared to random masking.

This efficiency gain likely stems from the principled curriculum learning effect of CV-based masking, where the model focuses computational resources on truly challenging prediction tasks rather than easy, low-variance laboratories.

4.3 Detailed Mechanistic Analysis

Historical Context Utilization: We analyzed reconstruction error as a function of available patient history, measuring performance across varying amounts of lab-specific history (prior occurrences of the same test) and general patient context (timeline length). CV-based masking consistently

shows superior contextual learning, with performance gaps widening as more historical data becomes available.

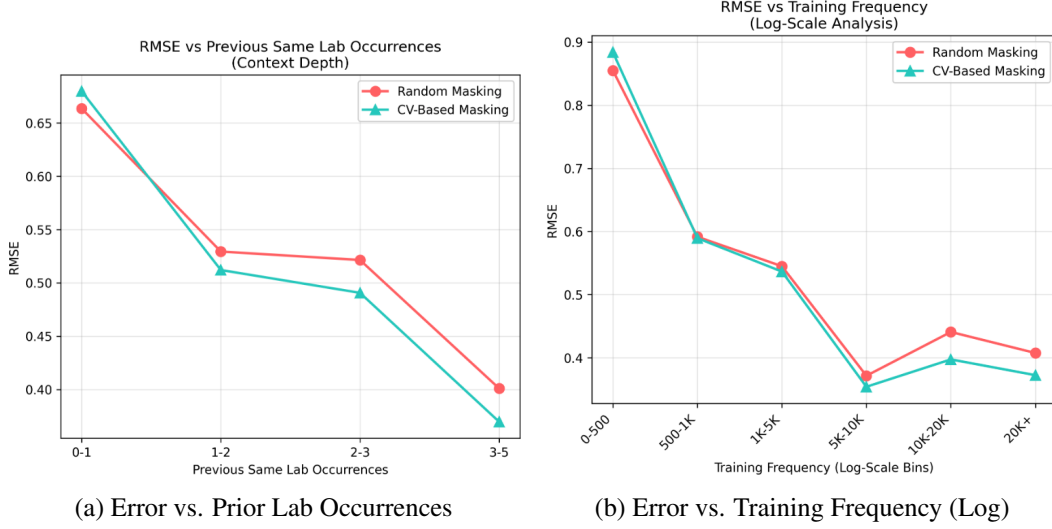


Figure 5: **CV-Based Masking Shows Superior Learning from Patient Context.** (a) Performance improves with more lab-specific history for all models, with CV-based model benefiting most. (b) CV-based model maintains superior performance across all data availability levels.

Perturbation Stress Test: To validate meaningful learning vs. memorization, we implemented a controlled experiment corrupting historical laboratory values while preserving prediction targets.

Experimental Protocol:

- **Corruption:** Historical values of the same lab type corrupted with adaptive Gaussian noise: $v'_i = v_i + \mathcal{N}(0, \sigma_{\text{adaptive}})$
- **Noise scaling:** $\sigma_{\text{adaptive}} = \text{std}(v_{\text{prior}}) \times 0.6 \times 1.5$
- **Controls:** Identical corruption (50% intensity, fixed seeds) for fair comparison
- **Preservation:** Target values and temporal/categorical context unchanged

Performance Measurement:

$$\text{Degradation} = \frac{\text{MAE}_{\text{corrupted}} - \text{MAE}_{\text{original}}}{\text{MAE}_{\text{original}}} \times 100\%$$

Results: CV-based models showed $2.1 \times$ greater degradation (9.8% vs 4.7%), indicating stronger reliance on patient-specific temporal patterns across diverse laboratory categories including hematology, electrolytes, and metabolic markers.

4.4 Per-Laboratory Reconstruction Example

As a concrete example of the reconstruction improvements, Figure 6 provides a visual example of reconstruction performance on Red Cell Distribution Width (RDW-SD), a moderate-volatility biomarker. The scatter plots demonstrate progressive tightening of predictions around the ideal $y=x$ line as the masking strategy improves from Random to Variance-based to CV-Based, with R^2 increasing dramatically from 0.668 to 0.775 to 0.847.

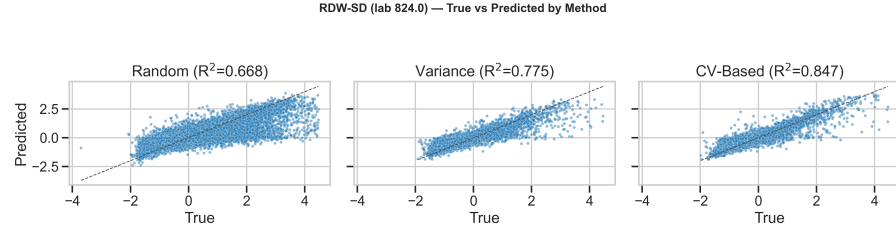


Figure 6: **True vs. Predicted scatter plots for Red Cell Distribution Width (RDW-SD).** Each point represents a single masked value prediction. The scatter becomes progressively tighter and more aligned with the ideal $y=x$ line (dashed) as the masking strategy improves from Random (left, $R^2=0.668$) to Variance-based (center, $R^2=0.775$) to CV-Based (right, $R^2=0.847$). This visual improvement demonstrates the substantial benefit of volatility-aware curriculum on moderately predictable laboratories.