

MISSING PATTERN RECOGNIZED DIFFUSION IMPUTATION MODEL FOR MISSING NOT AT RANDOM

Anonymous authors

Paper under double-blind review

ABSTRACT

Missing data frequently arises across diverse domains, including time-series and image domains. In the real world, missing occurrences often depend on the unobservable values themselves, which are referred to as Missing Not at Random (MNAR). To address this, numerous generative models have been proposed, with diffusion models in particular demonstrating strong capabilities in out-of-sample imputation. However, most existing diffusion-based imputation approaches overlook the MNAR setting and instead rely on restrictive assumptions about the missing process, thereby limiting their applicability to practical scenarios. In this work, we introduce the Missing Pattern Recognized Diffusion Imputation Model (PRDIM), a novel framework that explicitly captures the missing pattern and precisely imputes unobserved values. PRDIM iteratively maximizes the likelihood of the joint distribution for observed values and missing mask under an Expectation-Maximization (EM) algorithm. In this sense, we first employ a pattern recognizer, which approximates the underlying missing pattern and provides guidance during every inference toward more plausible imputations with respect to the missing information. In various experimental settings, we demonstrate that PRDIM achieves the state-of-the-art performance compared to previous diffusion imputation approaches under MNAR setting.

1 INTRODUCTION

Missing data imputation aims to recover missing values from partially observed incomplete datasets, and the imputation algorithms serve as a fundamental component in many domains, including healthcare (Goldberger et al., 2000), traffic (Li et al., 2017), and image domain (Xiao et al., 2017). Formally, the imputation goal is to accurately estimate missing values conditioned on observed values. Many recent imputation models assume the missing to be random or to depend on observed values, which are referred to as Missing Completely at Random (MCAR) and Missing at Random (MAR) respectively (Little & Rubin, 1987; Schafer, 1997). However, in real-world scenarios, we have missing values because of some underlying causes such as health deterioration or mortality (Carreras et al., 2021); so the missing tends to have patterns in its occurrences. Therefore, such patterned missing cases are referred to as Missing Not at Random (MNAR), which considered to be more realistic, and which have not been approached by diffusion-based imputation models in the previous works. Consequently, this paper proposes a diffusion-based imputation model to estimate the *missing pattern* to better recover missing under the general MNAR setting.

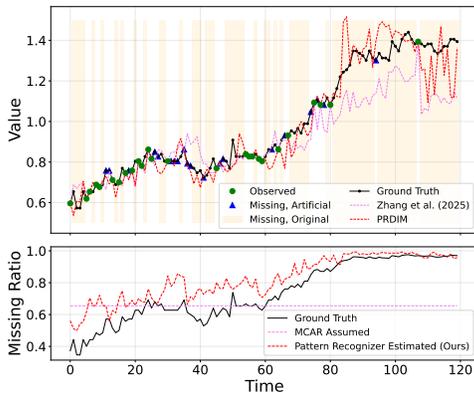


Figure 1: (Top) Comparison of imputation performance on original missing entries (orange) versus artificial missing entries (blue) under observed values (green). (Bottom) Estimated missing ratio by pattern recognizer (red) regard to true missing ratio.

Table 1: Numerical results; imputing original missing presents a more challenging task.

Method	Artificial missing entries			Original missing entries		
	RMSE	MAE	MRE	RMSE	MAE	MRE
Zhang et al. (2025)	0.230	0.146	21.572	1.209	0.782	46.188
PRDIM	0.201	0.124	18.310	1.057	0.663	39.156

To impute high-dimensional data, recent studies have increasingly employed diffusion models (Ho et al., 2020; Song et al., 2021), which provide a powerful framework for capturing complex data distributions under incomplete settings. Existing methods are generally built upon conditional diffusion frameworks (Dhariwal & Nichol, 2021; Ho & Salimans, 2022), in which target entries are artificially masked and treated as missing values for training and evaluation. These methods could be trained on incomplete data; however, direct evaluation on the original missing values is infeasible. Consequently, their performance is typically assessed using artificially masked values. We illustrate this critical distinction between the artificial and original missing entries in the upper part of Figure 1, and argue that the evaluation of the latter is essential for measuring the practical imputation capability. Furthermore, Table 1 demonstrates that imputing original missing entries is a more challenging task. [Additional imputation results and underlying missing probability estimations shown in Figure 1 are provided in Appendix C.2.4 and Appendix C.2.5. A detailed description of *original* versus *artificial* missing entries is further discussed in Section 2.3.](#)

In this work, we target the more realistic and challenging setting of imputing original missing values under MNAR. We introduce **Missing Pattern Recognized Diffusion Imputation Model (PRDIM)**, a diffusion-based imputation model augmented with a *pattern recognizer* that explicitly models the missing pattern with EM framework. By learning the missing pattern, PRDIM provides additional guidance during denoising and provides more precise imputation results than existing methods. Through extensive experiments, we demonstrate that PRDIM achieves state-of-the-art performance on multiple benchmarks, significantly improving imputation accuracy under MNAR conditions.

In summary, our contributions are as follows:

- We propose PRDIM, an EM-based algorithm that maximizes the joint likelihood of observed values and mask information within a diffusion framework. This design enables the model to infer latent missing patterns in incomplete data.
- We theoretically demonstrate that the missing model, termed the pattern recognizer, can provide additional approximate guidance for imputing missing values during the expectation step and inference when the mask information is available.
- We empirically validate PRDIM on incomplete datasets under MNAR setting, showing that it consistently outperforms existing approaches across multiple evaluation metrics.

2 PRELIMINARIES

2.1 MISSING DATA IMPUTATION FOR MNAR

As discussed earlier, accurate estimation of the missing values requires conditioning on both the observed variables X^{obs} and the missing mask M . Let $X = [X_d] \in \mathbb{R}^D$ be a complete instance with D dimensions, and $M \in \{0, 1\}^D$ the missing indicator where $M_d = 1$ if X_d is observed, otherwise 0. We denote the observed and missing subsets by $X^{\text{obs}} = X \odot M$ and $X^{\text{mis}} = X \odot (1 - M)$ under $X = (X^{\text{obs}}, X^{\text{mis}})$, respectively. Given access to the underlying complete data X and its mask M , the ultimate goal is to recover the joint distribution of the observed data, missing mask, and missed data:

$$\max_{\theta, \phi} \mathbb{E}[\log p_{\theta, \phi}(X^{\text{obs}}, X^{\text{mis}}, M)], \quad \text{with } (X^{\text{obs}}, X^{\text{mis}}, M) \sim p_{\text{data}}(X, M).$$

Here, θ and ϕ are the parameters of the conditional distribution, which will be discussed in their individual roles on describing distributions of X and $M|X$, respectively.

In the scenario of missing value imputation on incomplete data, the inference becomes maximizing the joint distribution of only two random variables X^{obs} and M because of an unobservable property of X^{mis} . This likelihood maximization is formulated with the expectation on X^{mis} , which eventually turns the problem into the Expectation-Maximization framework.

$$p_{\theta, \phi}(X^{\text{obs}}, M) = \int_{X^{\text{mis}}} p_{\theta, \phi}(X^{\text{obs}}, X^{\text{mis}}, M) dX^{\text{mis}} = \int_{X^{\text{mis}}} p_{\theta}(X) p_{\phi}(M|X) dX^{\text{mis}}$$

Consequently, principled inference requires joint modeling of $p_{\theta}(X)$ and $p_{\phi}(M|X)$ and optimizing a suitable lower bound of the EM algorithm. Now, the focus becomes how to infer the two distributions: $p_{\theta}(X)$ and $p_{\phi}(M|X)$. Particularly, the inference requirement on $p_{\phi}(M|X)$ becomes different depending on the assumed missing mechanism across different scenarios.

Missing Mechanisms Since the distribution $p_{\theta,\phi}(X, M)$ can be decomposed as $p_{\theta}(X)p_{\phi}(M|X)$, it becomes necessary to explicitly model the generation process of the mask variable M . The following three missing processes (Little & Rubin, 1987) can be modeled as a conditional distribution $p_{\phi}(M|X)$ following the standard taxonomy:

Missing Completely At Random (**MCAR**): $p_{\phi}(M|X) = p_{\phi}(M)$,

Missing At Random (**MAR**): $p_{\phi}(M|X) = p_{\phi}(M|X^{\text{obs}})$,

Missing Not At Random (**MNAR**): $p_{\phi}(M|X) = p_{\phi}(M|X^{\text{obs}}, X^{\text{mis}})$.

Under MCAR/MAR, the likelihood $p_{\theta,\phi}(X^{\text{obs}}, M)$ which is proportional to $p_{\theta}(X^{\text{obs}})$ can be learned while ignoring missing process (Mattei & Frellsen, 2019) with respect to the missing variable X^{mis} . In contrast, the mask also depends on unobserved values under MNAR; makes the missing process non-ignorable (Ipsen et al., 2020). While more realistic scenario comes from MNAR, this new requirement of inferring $p_{\phi}(M|X)$ renders the imputation models under MAR and MCAR to be ineffective and needs to be overhauled significantly.

Missing Model Some previous works have incorporated the missing mask M as a supervised learning target. Originated from Generative Adversarial Network (Goodfellow et al., 2014), GAIN (Yoon et al., 2018a) first proposed that a discriminator $D_{\phi}(X)$ can be trained to approximate $p(M|X)$, with the objective of optimal discriminator D_{ϕ^*} towards the specific missing ratio regardless of data distribution. This design ensures that, when the missing mechanism is independent of the data, the discriminator converges to a uniform prediction over missing and observed variables.

Modified from Variational Autoencoder (Kingma & Welling, 2013), not-MIWAE (Ipsen et al., 2020) extended the discriminator framework to MNAR by directly modeling $p(M|X^{\text{obs}}, X^{\text{mis}})$. Their approach demonstrated that the discriminator loss can be integrated into a variational objective, allowing optimization via minimization of an additional term in the ELBO.

Under the shared assumption adopted by GAIN and not-MIWAE, each missing value indicator $M_d \in \{0, 1\}$ follows Bernoulli distribution conditioned on the entire data (*i.e.* $\log p(M|X) = \sum_{d=1}^D \log p(M_d|X)$); the loss for the missing model D_{ϕ} can be formulated as a binary cross-entropy (BCE) objective:

$$\mathcal{L}(M, X, D_{\phi}) = -M^{\top} \log D_{\phi}(X) - (1 - M)^{\top} \log (1 - D_{\phi}(X)), \quad (1)$$

where optimal $D_{\phi^*}(X)$ predicts the probability whether each entry would be observed or not (*i.e.* $D_{\phi^*}(X) = [p(M_d = 1|X)] \in \mathbb{R}^D$). This formulation provides a flexible way to incorporate the missing mechanism into generative imputation models. In this sense, we hereafter refer to the discriminator D_{ϕ} as the *pattern recognizer*. Correspondingly, the loss \mathcal{L} will be denoted as \mathcal{L}_{PR} .

2.2 DIFFUSION MODELS

Diffusion models have achieved state-of-the-art generation performances in multiple domains, including vision (Esser et al., 2024), audio (Kong et al., 2020), graphs (Jo et al., 2022), and time-series (Coletta et al., 2023). They learn a data distribution by inverting a Markovian noising process. The forward process gradually corrupts a clean sample $X_0 \sim q(X_0)$ with Gaussian noise according to a prescribed variance schedule $\{\beta_t\}_{t=1}^T$:

$$q(X_{1:T}|X_0) := \prod_{t=1}^T q(X_t|X_{t-1}) \text{ where } q(X_t|X_{t-1}) = \mathcal{N}(\sqrt{\alpha_t}, (1 - \alpha_t)\mathbf{I}). \quad (2)$$

Here, $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$. This yields the closed form to sample X_t at time t given X_0 :

$$q(X_t|X_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}X_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (3)$$

The reverse process is modeled as a learned Markov chain that gradually removes noise:

$$p_{\theta}(X_{0:T}) := \prod_{t=1}^T p_{\theta}(X_T)p_{\theta}(X_{t-1}|X_t) \quad (4)$$

where $p_{\theta}(X_T)$ is a standard Gaussian prior, θ is the learnable parameter of a diffusion model.

Diffusion models are trained by maximizing a variational lower bound on the log-likelihood:

$$\mathbb{E}_{X_0 \sim q(X_0)} [\log p_{\theta}(X_0)] \geq \mathbb{E}_{X_0 \sim q(X_0), X_{1:T} \sim q(X_{1:T}|X_0)} \left[\log \frac{p_{\theta}(X_{0:T})}{q(X_{1:T}|X_0)} \right] \quad (5)$$

This objective could be reduced as the data reconstruction (Karras et al., 2022), the noise prediction (Ho et al., 2020), or the score matching objective (Song et al., 2021). At inference time, sampling proceeds by drawing $X_T \sim p_T(X_T)$ and iteratively applying the learned reverse transitions $p_\theta(X_{t-1}|X_t)$ to obtain a clean sample X_0 .

We present the entire methodology of this work from the perspective of data reconstruction (*i.e.* X_0 prediction). Under this view, the diffusion model is trained to reconstruct the clean data X_0 from noisy samples X_t across timesteps. Accordingly, the diffusion loss function for a single sample X_0 is defined as

$$\mathcal{L}_{\text{diff}}(X_0, X_t, t, f_\theta) = \lambda(t) \|f_\theta(X_t, t) - X_0\|_2^2 \quad (6)$$

where $t \sim \mathcal{U}(0, T)$, $\lambda(t) = 1$ in our experiment, f_θ is the X_0 prediction network parameterized by θ , and X_t follows Equation 3.

2.3 REVISITING THE OBJECTIVE OF DIFFUSION IMPUTATION MODELS

In this section, we clarify how the target objective adopted in our work differs in perspective from those used in existing diffusion imputation models, and we justify the validity of our proposed objective. Throughout this paper, the term *original missing* values refers to two situations: (i) naturally occurring missing values in real-world datasets such as PhysioNet (Goldberger et al., 2000) or AirQuality (Zhang et al., 2017) which the underlying values are not accessible for evaluation, and (ii) simulated missing values generated by applying a missing mechanism (under MCAR, MAR, or MNAR) to a complete dataset. To evaluate imputation performance on such original missing values, we construct various incomplete datasets by applying different missing mechanisms to complete benchmarks, and then apply imputation methods to these resulting incomplete data distributions.

For an originally incomplete dataset $(X_0^{\text{obs}}, X_0^{\text{mis}})$, M_O be the corresponding original missing mask. In many prior works (Tashiro et al., 2021; Zhou et al., 2024; Liu et al., 2024), it is common practice to impose an additional mask on X_0^{obs} during training. Let M_A denote this *artificial missing* mask, and $X_0^{\text{obs},A}$ indicates the subset of observed entries remaining after the artificial masking is applied (*i.e.*, $X_0^{\text{obs},A} \equiv X_0^{\text{obs}} \odot M_A$). These artificially masked entries supply the supervision required for both training and quantitative evaluation. We can formulate these works aim to obtain $p_\theta(X_0 | X_0^{\text{obs},A}, M_A)$ and evaluation follows the same protocol by inserting an artificial mask into the test data. From this perspective, our research question emerges from a fundamentally different viewpoint:

If the distribution of the original missing mask M_O differs substantially from that of the artificial mask M_A , can a model trained under $p_\theta(X_0 | X_0^{\text{obs},A}, M_A)$ be expected to perform well under $p_\theta(X_0 | X_0^{\text{obs}}, M_O)$?

This distinction is crucial to the motivation of our work. Our methodology newly introduces a pattern recognizer, D_ϕ , which explicitly learns missing distribution, and it enables to model the distribution $p_{\theta,\phi}(X_0 | X_0^{\text{obs}}, M_O)$ regard to underlying missing mechanism.

3 METHODOLOGY

We decompose PRDIM into two major components: (i) diffusion backbone pre-training and (ii) missing model training with joint distribution fine-tuning. Each component contains practical implementation details that stabilize training and enhance imputation performance. We begin by presenting an overview of PRDIM in Section 3.1, followed by theoretical analysis in subsequent sections.

3.1 PRE-IMPUTATION AND EM ALGORITHM

Figure 2 illustrates the graphical model and overall training procedure of PRDIM. The framework consists of two complementary phases: a diffusion-based pre-imputation stage (Phase 1) and an EM iteration stage (Phase 2). The combination of these two phases enables PRDIM to learn both the data distribution and the missing pattern in a principled manner.

Phase 1: Diffusion model Pre-training and Pre-imputation. Direct optimization of the joint distribution $p_\theta(X_0^{\text{obs}}, X_0^{\text{mis}})$ often suffers from instability and overfitting when missing entries are

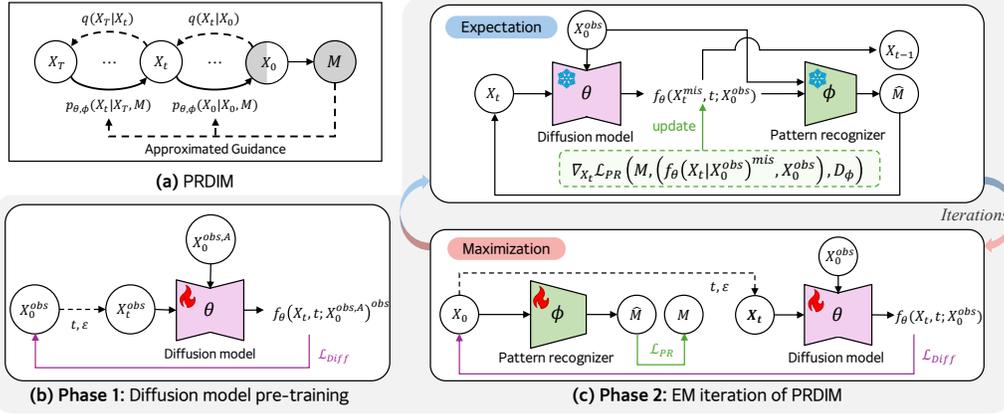


Figure 2: Overall training procedure of PRDIM; (a) Graphical model of PRDIM, (b) Diffusion model pre-training, (c) EM iteration for diffusion model fine-tuning and pattern recognizer training. In expectation step, PRDIM progressively denoising the sample with additional guidance of pattern recognizer (Top). In maximization step, diffusion model and pattern recognizer are trained independently based on X_0^{obs} , X_0^{mis} which generated in the previous expectation step (Bottom).

simply set to zero. Phase 1 of PRDIM pre-trains a diffusion backbone under a Tashiro et al. (2021)’s framework enabling the model to capture plausible data distributions from observed variables. For generic imputation tasks, we extend the diffusion target up to conditional input which is proposed at Du et al. (2023a) as Observed Reconstruction Task, which demonstrates joint distribution learning under a conditional modeling. Therefore, to train the diffusion model for pre-imputation, we need to artificially select missing entries. Unlike CSDI, which adopts artificial masking under the MCAR assumption, we avoid this strategy in our approach. Instead, we introduce an adjacent target masking scheme, where artificial missing entries are placed near original missing values to exploit potential correlations. In multivariate time-series data, the artificial missing values are chosen only across the temporal axis. In image data, they are selected as the top, bottom, left, or right neighboring missing pixels. All subsequent experimental results are conducted under this adjacent target masking setup. We denote this artificial missing mask as A and define the corresponding subset $X_0^{\text{obs},A} := X_0^{\text{obs}} \odot A$. Through this strategy, the diffusion backbone learns to reconstruct plausible imputations while being robust to any missing pattern. Then, the diffusion loss objective can be rewritten as follows:

$$\mathcal{L}_{\text{diff}}(X_t, t, f_\theta) = \|f_\theta(X_t, t; X_0^{\text{obs},A})^{\text{obs}} - X_0^{\text{obs}}\|_2^2 \quad (7)$$

where $f_\theta(X_0, X_t, t; X_0^{\text{obs},A})^{\text{obs}}$ is the subset of X_0 prediction corresponds to observed entries at timestep t . Recent studies have shown that plausible data distributions can be effectively estimated from artificially masked data, as evidenced by He et al. (2022); Peebles & Xie (2023). Furthermore, we demonstrate the performance gain across adjacent target masking and different ratio of MCAR masking schemes in Table 3.

Phase 2: EM Iteration of PRDIM. After diffusion pre-training, PRDIM enters the EM algorithm phase, which is inspired by Zhang et al. (2025). Whereas the previous work trains only joint diffusion model, our approach trains both diffusion model θ and missing pattern recognizer ϕ simultaneously to refine the imputation performance.

In the maximization step, the diffusion model θ is updated to capture the full joint distribution $p_\theta(X_0^{\text{obs}}, X_0^{\text{mis}})$, while the pattern recognizer $D_\phi(X_0) = \hat{M}$ is trained to discriminate the mask variable M as a supervised learning target.

$$\mathcal{L}_{\text{diff}}(X_0, X_t, t, f_\theta) = \|f_\theta(X_t, t; X_0^{\text{obs}}) - X_0\|_2^2 \quad (8)$$

$$\mathcal{L}_{\text{PR}}(M, X_0, D_\phi) = -M^\top \log D_\phi(X_0) - (1 - M)^\top \log (1 - D_\phi(X_0)) \quad (9)$$

This enables the model to explicitly incorporate the missing pattern. Section 3.2 provides the details of maximization objective.

In the expectation step, the diffusion model generates X^{mis} conditioned on M and X^{obs} , while the pattern recognizer provides an additional approximated guidance signal that biases the generation toward imputations consistent with the estimated missing patterns. Importantly, during the early iterations, it is acceptable to use a randomly initialized pattern recognizer as guidance. Since such a recognizer has no discriminative ability, the guidance provides a degenerated signal toward a near-zero vector, yielding a neutral effect on the generation process (Kim et al., 2022). Section 3.3 illustrates the guided generation of X^{mis} for the expectation step.

Furthermore, while Zhang et al. (2025) adopts a soft EM strategy, our framework employs a hard EM variant to enhance the exploration ability to generate X_0^{mis} distribution, which has been theoretically justified in the context of expectation maximization (Samdani et al., 2012). Details of the maximization and expectation steps that are defined in X_0 prediction are summarized in Algorithm 1 and Algorithm 2.

3.2 ELBO OBJECTIVE OF MNAR ON DIFFUSION MODEL FRAMEWORK

To address the Missing Not at Random (MNAR) mechanism, we formulate the evidence lower bound (ELBO) of diffusion imputation model within a diffusion framework. Let $X := X_0 = (X_0^{\text{obs}}, X_0^{\text{mis}})$ denote the true data which can be divided into observed and missing variables, and M the missing mask variable. We derive the ELBO of $\log p_{\theta, \phi}(X_0^{\text{obs}}, M)$ as follows:

Proposition 3.1. *Suppose that the variational distribution and missing process satisfy the following assumptions: (i) Forward process satisfies conditional independence to mask variable given entire data $q(X_{1:T}|X_0, M) = q(X_{1:T}|X_0)$, (ii) Missing process satisfies conditional independence to noisy data $p_{\phi}(M|X_0, X_{1:T}) = p_{\phi}(M|X_0)$.*

Then, the ELBO of joint log-likelihood objective of the observed data and mask can be expressed as

$$\log p_{\theta, \phi}(X_0^{\text{obs}}, M) \geq \mathbb{E}_{X_{1:T}, X_0^{\text{mis}}} \left[\log p_{\phi}(M|X_0) + \log \frac{p_{\theta}(X_T) \prod_{t=1}^T p_{\theta}(X_{t-1}|X_t)}{\prod_{t=1}^T q(X_t|X_{t-1})} \right] \quad (10)$$

$$+ \mathbb{E}_{X_{1:T}} \left[\mathbb{H}(q(X_0^{\text{mis}}|X_0^{\text{obs}}, M)) \right] \quad (11)$$

The proof is given in Appendix B.1. This is the first ELBO derivation of missing patterns under the diffusion setting. Unlike Ipsen et al. (2020), the iterative forward and reverse processes inherent to diffusion models render direct ELBO optimization intractable.

Thus, we utilize the above ELBO in the EM algorithm perspective, so that the generation of X^{mis} becomes an expectation step. In the maximization step, we simultaneously train the missing model ϕ for the missing mask M and the diffusion model θ that estimates the joint probability of $X^{\text{obs}}, X^{\text{mis}}$ by the Equation 10. Subsequently, the expectation step requires conditional generation of X^{mis} given M , which completes the EM loop by requiring another iteration of conditional generation.

In the expectation step, when $q(X_0^{\text{mis}}|X_0^{\text{obs}}, M)$ is replaced with $p_{\theta, \phi}(X_0^{\text{mis}}|X_0^{\text{obs}}, M)$, an additional guidance term with respect to M becomes necessary during the diffusion sampling process. In addition, while the diffusion model is designed to learn the joint distribution, Theorem 1 in Zhang et al. (2025), which is also provided in Appendix B.3, provides a theoretical justification that X^{mis} can be estimated via the joint score function. The guidance mechanism will be discussed in detail in the following subsection.

3.3 FIND BEST X^{MIS} WITH APPROXIMATED GUIDANCE OF PATTERN RECOGNIZER

After training both the missing model and the diffusion model, the expectation step allows us to replace $q(X_0^{\text{mis}}|X_0^{\text{obs}}, M)$ with the parameters θ, ϕ that were optimized in the preceding maximization step. Since the underlying generative process is score-based, the gradient term $\nabla_{X_t} \log p_{\theta, \phi}(X_0|X_0^{\text{obs}}, M)$ can be decomposed into two components: the score function term $\nabla_{X_t} p_{\theta}(X_t|X_0^{\text{obs}})$ corresponding to the joint distribution, and the mask guidance term $\nabla_{X_t} p_{\phi}(M|X_0^{\text{obs}}, X_t^{\text{mis}})$ reflecting the missing pattern. We further show that the mask guidance term can be approximated with the pattern recognizer D_{ϕ} according to Proposition 3.2.

Proposition 3.2. *Suppose that the pattern recognizer D_{ϕ} is optimal which satisfies $D_{\phi^*}(X_0) = [p_{\phi}(M_d|X_0)] \in \mathbb{R}^D$, the score function of the joint log-likelihood with respect to the missing mask*

Algorithm 1 E step - X_0 prediction

```

1: Diffusion model  $\theta$ , pattern recognizer  $\phi$ , observed data  $X^{\text{obs}}$ , and mask data  $M$ .
2: Sample  $X_T \sim \mathcal{N}(0, \mathbf{I})$ 
3: for  $t = T, \dots, 1$  do
4:   Get  $\hat{X}_0 = f_\theta(X_t, t; X_0^{\text{obs}})$ 
5:    $\tilde{X}_0 = \hat{X}_0 \odot M + X_0^{\text{obs}} \odot (1 - M)$ 
6:    $\hat{X}_0 = \tilde{X}_0 - \frac{1 - \bar{\alpha}_t}{\sqrt{\alpha_t}} \nabla_{X_t} \mathcal{L}_{\text{PR}}(M, \tilde{X}_0, D_\phi)$ 
7:   Sample  $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$ 
8:    $X_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{X}_0 + \sqrt{1 - \bar{\alpha}_{t-1}} \varepsilon$ 
9: end for
10: return  $X_0$ 

```

Algorithm 2 M step - X_0 prediction

```

1: Diffusion model  $\theta$ , pattern recognizer  $\phi$ , corresponding learning rate  $\eta_\theta, \eta_\phi$ , imputed data  $X_0$ , mask data  $M$ , and maximization epoch  $N_m$ .
2: for  $i = 1, \dots, N_m$  do
3:   Sample  $t, \varepsilon \sim \mathcal{U}(0, T), \mathcal{N}(0, \mathbf{I})$ 
4:    $X_t = \sqrt{\bar{\alpha}_t} X + \sqrt{1 - \bar{\alpha}_t} \varepsilon$ 
5:    $\theta \leftarrow \theta - \eta_\theta \nabla_\theta \mathcal{L}_{\text{diff}}(X_0, X_t, t, f_\theta)$ 
6:    $\phi \leftarrow \phi - \eta_\phi \nabla_\phi \mathcal{L}_{\text{PR}}(M, X_0, D_\phi)$ 
7: end for
8: return  $\theta, \phi$ 

```

can be approximated as

$$\nabla_{X_t} \log p_{\theta, \phi}(X_t | X_0^{\text{obs}}, M) \simeq \nabla_{X_t} \log p_\theta(X_t | X_0^{\text{obs}}) - \nabla_{X_t} \mathcal{L}_{\text{PR}}(M, \hat{X}_0, D_{\phi^*}) \quad (12)$$

where $\hat{X}_0 := (f_\theta(X_t, t; X_0^{\text{obs}})^{\text{mis}}, X_0^{\text{obs}}) = f_\theta(X_t, t; X_0^{\text{obs}}) \odot (1 - M) + X_0^{\text{obs}} \odot M$.

Here, $f_\theta(X_t, t; X_0^{\text{obs}})$ is X_0 prediction at timestep t , and its superscript means the subset which indicates missing entries (i.e. $f_\theta^{\text{mis}} = f_\theta \odot (1 - M)$). The detailed proof is given in Appendix B.2. This proposition is noteworthy because it steers the more convincing gradient of the intermediate sample X_t with respect to the estimated missing pattern D_ϕ ; it provides a meaningful signal according to the negative missing probability \mathcal{L}_{PR} .

To summarize, within the DDPM framework, given X_t at time step t , the variable X_{t-1} can be obtained through the following three steps. First, using the diffusion model together with Tweedie’s formula (Carlin et al., 2000; Efron, 2011), we can compute the posterior mean and an intermediate estimate $\tilde{X}_{t-1}^{\text{mis}}$ (Ho et al., 2020) only with the diffusion parameter θ :

$$\tilde{X}_{t-1} = \frac{1}{\sqrt{\alpha_t}} (X_t + (1 - \alpha_t) \nabla_{X_t} \log p_\theta(X_t | X_0^{\text{obs}})) + \sigma_t Z \quad \text{where } Z \sim \mathcal{N}(0, \mathbf{I}) \quad (13)$$

$$\hat{X}_0 = \frac{1}{\sqrt{\alpha_t}} (X_t + (1 - \bar{\alpha}_t) \nabla_{X_t} \log p_\theta(X_t | X_0^{\text{obs}})) = f_\theta(X_t, t; X_0^{\text{obs}}) \quad (14)$$

Second, the pattern recognizer evaluates the missing probability for each entry based on the estimated missing values \hat{X}_0^{mis} (Equation 14) and the observed component X_0^{obs} . Finally, according to Proposition 3.2, the denoised sample $X_{t-1} = (X_{t-1}^{\text{obs}}, X_{t-1}^{\text{mis}})$ at time $t-1$ is updated by incorporating the approximated guidance.

$$X_{t-1} = \tilde{X}_{t-1} - \frac{1 - \alpha_t}{\sqrt{\alpha_t}} \nabla_{X_t} \mathcal{L}_{\text{PR}}(M, (f_\theta(X_t, t; X_0^{\text{obs}})^{\text{mis}}, X_0^{\text{obs}}), D_\phi) \quad (15)$$

4 EXPERIMENTS

In this section, we evaluate PRDIM on multiple benchmark datasets under the MNAR setting. The results show that PRDIM achieves consistent improvements over prior methods, confirming the effectiveness of incorporating the missing model into the diffusion framework. The detailed experimental settings are provided in the Appendix C.

4.1 EXPERIMENTAL SETTING

Synthetic MNAR Datasets We evaluate our method on three widely used multivariate time-series datasets and one image dataset: (1) **ETT** (Zhou et al., 2021), which records load and temperature of

Table 2: Overall MAE performance on three benchmark datasets. We report mean \pm std over 5 runs according to each methodology. Best results are in **bold**, and second best results are in underline.

Method	Original / Out-of-Sample			Original / In-Sample		
	ETT	STOCK	PEMS-Bay	ETT	STOCK	PEMS-Bay
Mean	2.034 \pm 0.000	1.949 \pm 0.000	0.813 \pm 0.000	1.486 \pm 0.000	2.039 \pm 0.000	0.789 \pm 0.000
<i>Discriminative models</i>						
TimesNet	1.044 \pm 0.065	1.111 \pm 0.073	0.291 \pm 0.007	1.154 \pm 0.068	1.221 \pm 0.077	0.225 \pm 0.001
TimeMixer++	1.642 \pm 0.025	1.287 \pm 0.239	0.579 \pm 0.018	1.100 \pm 0.032	1.369 \pm 0.260	0.557 \pm 0.020
BRITS	0.992 \pm 0.037	0.627 \pm 0.010	0.278 \pm 0.006	0.491 \pm 0.008	0.701 \pm 0.010	0.182 \pm 0.003
SAITS	0.814 \pm 0.046	0.442 \pm 0.022	0.302 \pm 0.009	0.366 \pm 0.014	0.498 \pm 0.025	0.212 \pm 0.003
<i>Generative models</i>						
GP-VAE	1.511 \pm 0.011	0.902 \pm 0.109	0.345 \pm 0.001	0.896 \pm 0.018	1.010 \pm 0.118	0.292 \pm 0.002
not-MIWAE	1.311 \pm 0.016	0.681 \pm 0.045	0.396 \pm 0.005	0.637 \pm 0.011	0.759 \pm 0.039	0.352 \pm 0.005
<i>Diffusion-based models</i>						
CSDI	1.071 \pm 0.001	0.641 \pm 0.000	0.177 \pm 0.000	0.522 \pm 0.001	0.710 \pm 0.000	0.158 \pm 0.000
MTSCI	0.957 \pm 0.001	0.736 \pm 0.001	0.193 \pm 0.000	0.500 \pm 0.000	0.809 \pm 0.001	0.179 \pm 0.000
cDiffPuter	0.782 \pm 0.000	0.406 \pm 0.000	0.182 \pm 0.000	0.362 \pm 0.000	0.450 \pm 0.000	0.168 \pm 0.000
PRDIM	0.663\pm0.000	0.254\pm0.000	0.170\pm0.000	0.303\pm0.000	0.275\pm0.000	0.154\pm0.000

electricity transformers from 2016 to 2018 and has been a standard benchmark in time-series forecasting and imputation tasks; (2) **STOCK**, which contains historical daily Google stock prices from 2004 to 2019 and reflects complex temporal dynamics with strong non-stationarities; (3) **PEMS-Bay** (Li et al., 2017), which consists of road occupancy rates collected from highway sensors in the Bay Area, exhibiting highly correlated spatial-temporal patterns, and (4) **Fashion-MNIST** (FMNIST) (Xiao et al., 2017), which consists of gray-scale images of clothing items across 10 categories, widely used as a benchmark for evaluating imputation models.

These datasets are originally complete, and we arbitrarily mask missing values using a MNAR mechanism. As discussed in the introduction, all baseline models are trained on the resulting incomplete data, while evaluation is performed on the imputation of unobserved ground-truth values. The missing mechanism in our main experiment is proposed in not-MIWAE (Ipsen et al., 2020). A detailed description of the employed missing mechanism is provided in the Appendix C.

Baselines We compare our approach PRDIM against 12 representative imputation methods. **Mean** serves as a traditional statistical baseline. **TimesNet** (Wu et al., 2022), **TimeMixer++** (Wang et al., 2024), **BRITS** (Cao et al., 2018), and **SAITS** (Du et al., 2023a) are discriminative models, known to achieve strong performance for time-series imputations. **GP-VAE** (Fortuin et al., 2020) and **not-MIWAE** (Ipsen et al., 2020) are VAE-based imputation model that could be implemented on incomplete data. Among diffusion-based methods, we reproduce **CSDI** (Tashiro et al., 2021), **MTSCI** (Zhou et al., 2024), and **DiffPuter** (Zhang et al., 2025) which shows robust performance on various datasets. Specifically, We modified DiffPuter into conditional diffusion framework, which denote as **cDiffPuter**. In image domain, we reproduce **misGAN** (Li et al., 2019) and **MCFlow** (Richardson et al., 2020) which are GAN, Flow based generative imputation model respectively. The details of baselines can be found in Appendix A.

Evaluation Metrics We consider two types of evaluation: (i) *original in-sample* imputation and *original out-of-sample* imputation, both targeting the recovery of true missing values. As highlighted in the introduction, our main goal is to impute original missing entries, as opposed to artificially masked ones. We report the three error-based metrics (i) RMSE, (ii) MAE, and (iii) MRE which are defined in Appendix C, throughout this section.

4.2 OVERALL PERFORMANCE

We first evaluate PRDIM against representative discriminative, generative, and diffusion-based baselines on three multivariate time-series datasets (ETT, STOCK, PEMS-Bay) and FMNIST dataset. As reported in Table 2 and detailed in the Appendix C.1, PRDIM consistently achieves the best performance across all metrics and datasets. In particular, the gains are most pronounced on out-of-sample imputation tasks, while in-sample results remain competitive, suggesting that PRDIM generalizes well to unseen missing values.

432 Figure 3 shows qualitative comparisons on FMNIST, further
 433 illustrating that PRDIM generates semantically more
 434 consistent reconstructions than other approaches. In addition,
 435 the corresponding FID scores demonstrate that our model
 436 achieves outstanding generative quality, corroborating the
 437 visual improvements with quantitative evidence. These results
 438 highlight the importance of explicitly modeling the missing
 439 process.
 440

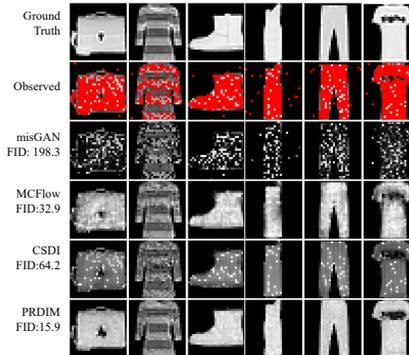


Figure 3: Imputation on FMNIST. Second row shows observed inputs, where red pixels indicate missing.

442 4.3 CONTROLLED PARAMETER ANALYSIS

443 To better understand the dynamics of PRDIM, we analyze
 444 the training behavior of the EM procedure. Figure 4
 445 shows the convergence of the pattern recognizer’s loss,
 446 where red curves indicate the ability to distinguish miss-
 447 ing values and blue curves correspond to observed values.
 448 The results on both the ETT and STOCK datasets demonstrate
 449 that the pattern recognizer effectively captures the missing
 450 pattern, thereby providing informative guidance during
 451 generation. Furthermore, Figure 5 illustrates the evolution
 452 of MRE across EM iterations on ETT and STOCK datasets,
 453 revealing a consistent improvement in imputation accuracy
 454 as the number of EM epochs increases.

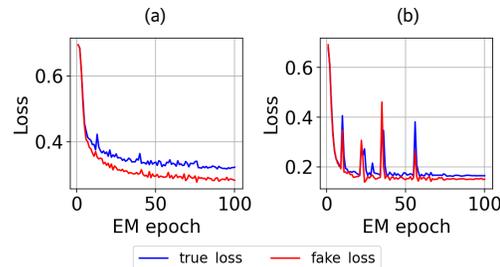
The results on both the ETT and STOCK datasets demonstrate that the pattern recognizer effectively captures the missing pattern, thereby providing informative guidance during generation. Furthermore, Figure 5 illustrates the evolution of MRE across EM iterations on ETT and STOCK datasets, revealing a consistent improvement in imputation accuracy as the number of EM epochs increases.

455 Moreover, according to the findings on a previous work (Ho & Salimans, 2022), increasing the
 456 weight of guidance generally leads to the better generation. As shown in Figure 6, the imputation
 457 performance consistently improves as the guidance scale increases, demonstrating the effective-
 458 ness of guidance weighting. These findings confirm that EM refinement is a critical component of
 459 PRDIM, substantially enhancing its capacity to model the joint distribution of the data and the miss-
 460 ing pattern. Despite the strong performance of our proposed PRDIM framework, it should be noted
 461 that introducing an additional guidance term in diffusion models inherently incurs extra computa-
 462 tional cost (Kim et al., 2022; Chung et al., 2023). To clarify this trade-off, we provide supplementary
 463 results in Table 8 of Appendix C, where training and inference times are compared across different
 464 diffusion-based imputation models.

462 4.4 ABLATION STUDIES

463 We further conduct ablation experiments to quantify the contribution of each component in PRDIM.
 464 Table 3 reports results when either the pattern recognizer is removed or hard EM is replaced with
 465 soft EM. Both modifications lead to a significant performance drop, indicating that explicit missing
 466 modeling and iterative EM updates are indispensable for exploring the missing data distribution.
 467 Furthermore, we investigate the impact of the artificial missing mask A . While previous works
 468 typically adopt A under the MCAR, we report performance for missing rates of 10%, 50%, and
 469 90% of A . Overall, the missing rate of A has limited influence on imputation accuracy. Specifically,
 470 extremely high missing rate degrades the quality of pre-imputation in Phase 1, thereby requiring
 471 more EM iterations for convergence.

472 Tables 4 and 5 investigate robustness under different missing mechanisms by applying MNAR and
 473 MCAR masks to the ETT dataset. Implementation details for the corresponding mechanisms are
 474



484 Figure 4: Convergence of the pattern recognizer’s loss during EM training. (a) ETT (b)
 485 STOCK.

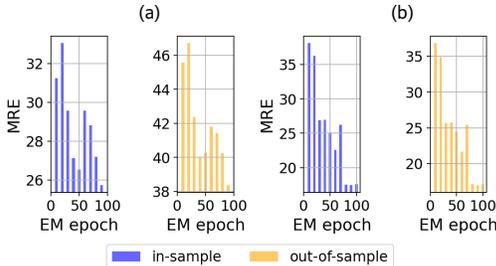


Figure 5: Evolution of MRE across EM epochs. (a) ETT (b) STOCK.

Table 3: Ablation study results on the **STOCK** dataset. We report the average over 5 runs. $n\%$ missing means utilizing artificial missing masking A under MCAR mechanism with $n\%$ instead of adjacent target masking defined in Section 3.1.

Method	original out-of-sample			original in-sample		
	RMSE	MAE	MRE	RMSE	MAE	MRE
PRDIM	0.599	0.254	16.794	0.633	0.275	17.150
10% missing of A	0.590	0.258	17.057	0.630	0.282	17.632
50% missing of A	0.624	0.259	17.118	0.667	0.281	17.563
90% missing of A	0.631	0.293	19.396	0.677	0.323	20.146
w/o PR	0.650	0.306	20.233	0.691	0.339	21.171
w/o PR and hard EM	0.734	0.406	26.878	0.778	0.450	28.064

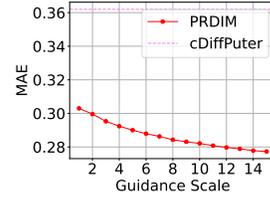


Figure 6: Effect of guidance scale on STOCK.

described in Appendix C. PRDIM consistently outperforms baselines under MNAR, whereas under MCAR the advantage diminishes, as the pattern recognizer effectively learns randomly in this scenario. Together, these ablation studies confirm the necessity of PRDIM’s design choices and its robustness across varying missing conditions.

Table 4: Imputation performance on the **ETT** dataset in different **MNAR** pattern. We report the average RMSE / MAE / MRE over 5 runs.

Method	out-of-sample	in-sample
CSDI	0.345 / 0.218 / 13.825	0.233 / 0.159 / 13.259
cDiffPuter	0.264 / 0.177 / 11.168	0.281 / 0.178 / 14.836
PRDIM	0.282 / 0.171 / 10.837	0.197 / 0.130 / 10.816

Table 5: Imputation performance on the **ETT** dataset in **MCAR** mechanism. We report the average RMSE / MAE / MRE over 5 runs.

Method	out-of-sample	in-sample
CSDI	0.237 / 0.160 / 15.550	0.182 / 0.128 / 16.752
cDiffPuter	0.251 / 0.162 / 15.737	0.202 / 0.136 / 17.795
PRDIM	0.225 / 0.147 / 14.267	0.172 / 0.118 / 15.485

4.5 APPLICATION OF PRDIM

We conducted two additional experiments to investigate the importance of Phase 1 initialization and the generalization capability of PRDIM. First, because the main objective of PRDIM can operate even without an explicit Phase 1, we evaluated several initialization strategies (MEAN, BRITS, SAITS) under identical conditions to examine how different pre-imputation methods influence PRDIM’s performance. As shown in Table 6, initialization with CSDI yields the best results. However, PRDIM consistently outperforms the baseline regardless of the chosen initialization method, demonstrating that it does not heavily rely on a particular Phase 1 strategy.

Second, we assessed post-imputation classification accuracy on the FMNIST dataset to evaluate how well the imputed data support downstream tasks. This experiment examines whether higher-quality imputations translate into improved task performance. As reported in Table 7, PRDIM surpasses Flow-based, GAN-based, and prior diffusion-based approaches, indicating that PRDIM is capable of restoring semantically meaningful information even in the image domain.

Table 6: Out-of-sample MAE under different Phase 1 methods.

Initialization Method	ETT	STOCK	PEMS-Bay
MEAN (close to w/o phase 1)	0.766	0.326	0.207
BRITS	0.824	0.300	0.188
SAITS	0.774	0.307	0.180
CSDI (PRDIM)	0.663	0.254	0.170

Table 7: Post-imputation classification results on FMNIST.

Method	Accuracy (%)
Clean Data	92.59
PRDIM	91.14
MCFLOW	90.49
CSDI	87.41
misGAN	84.34

5 CONCLUSION

We presented PRDIM, a diffusion-based imputation framework that incorporates an additional discriminator denoted pattern recognizer under an EM algorithm to explicitly estimate missing patterns. We shows that the guidance which understands missing pattern could be helpful for generating missing values precisely in diffusion imputation model. Our theoretical derivation and extensive experiments demonstrate that PRDIM consistently improves imputation compared to existing methods.

REFERENCES

- 540
541
542 Juan Miguel Lopez Alcaraz and Nils Strodthoff. Diffusion-based time series imputation and fore-
543 casting with structured state space models. *arXiv preprint arXiv:2208.09399*, 2022.
- 544
545 Wei Cao, Dong Wang, Jian Li, Hao Zhou, Lei Li, and Yitan Li. Brits: Bidirectional recurrent
546 imputation for time series. *Advances in neural information processing systems*, 31, 2018.
- 547
548 Bradley P Carlin, Thomas A Louis, et al. *Bayes and empirical Bayes methods for data analysis*.
Chapman & Hall/CRC, 2000.
- 549
550 Giulia Carreras, Guido Miccinesi, Andrew Wilcock, Nancy Preston, Daan Nieboer, Luc Deliens,
551 Mogensm Groenvold, Urska Lunder, Agnes van der Heide, Michela Baccini, et al. Missing not at
552 random in end of life care studies: multiple imputation and sensitivity analysis on data from the
553 action study. *BMC medical research methodology*, 21(1):13, 2021.
- 554
555 Hyungjin Chung, Jeongsol Kim, Michael T McCann, Marc L Klasky, and Jong Chul Ye. Diffu-
556 sion posterior sampling for general noisy inverse problems. In *11th International Conference on*
Learning Representations, ICLR 2023, 2023.
- 557
558 Andrea Coletta, Sriram Gopalakrishnan, Daniel Borrajo, and Svitlana Vyetrenko. On the con-
559 strained time-series generation problem. *Advances in Neural Information Processing Systems*,
560 36:61048–61059, 2023.
- 561
562 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
in neural information processing systems, 34:8780–8794, 2021.
- 563
564 Wenjie Du, David Côté, and Yan Liu. Saits: Self-attention-based imputation for time series. *Expert*
Systems with Applications, 219:119619, 2023a.
- 565
566 Wenjie Du, Yiyuan Yang, Linglong Qian, Jun Wang, and Qingsong Wen. PyPOTS: A Python Toolkit
567 for Machine Learning on Partially-Observed Time Series. *arXiv preprint arXiv:2305.18811*,
568 2023b.
- 569
570 Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Associa-*
571 *tion*, 106(496):1602–1614, 2011.
- 572
573 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
574 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers
575 for high-resolution image synthesis. In *Forty-first international conference on machine learning*,
576 2024.
- 577
578 Vincent Fortuin, Dmitry Baranchuk, Gunnar Rätsch, and Stephan Mandt. Gp-vae: Deep probabilis-
579 tic time series imputation. In *International conference on artificial intelligence and statistics*, pp.
1651–1661. PMLR, 2020.
- 580
581 Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G
582 Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank,
583 physiotookit, and physionet: components of a new research resource for complex physiologic
584 signals. *circulation*, 101(23):e215–e220, 2000.
- 585
586 Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
587 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information*
processing systems, 27, 2014.
- 588
589 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
590 toencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*
vision and pattern recognition, pp. 16000–16009, 2022.
- 591
592 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
593 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*
neural information processing systems, 30, 2017.

- 594 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint*
595 *arXiv:2207.12598*, 2022.
- 596
- 597 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
598 *neural information processing systems*, 33:6840–6851, 2020.
- 599
- 600 Niels Bruun Ipsen, Pierre-Alexandre Mattei, and Jes Frellsen. not-miwae: Deep generative mod-
601 elling with missing not at random data. *arXiv preprint arXiv:2006.12871*, 2020.
- 602
- 603 Daniel Jarrett, Bogdan C Cebere, Tennison Liu, Alicia Curth, and Mihaela van der Schaar. Hy-
604 perimpute: Generalized iterative imputation with automatic model selection. In *International*
Conference on Machine Learning, pp. 9916–9937. PMLR, 2022.
- 605
- 606 Jaehyeong Jo, Seul Lee, and Sung Ju Hwang. Score-based generative modeling of graphs via the
607 system of stochastic differential equations. In *International conference on machine learning*, pp.
10362–10383. PMLR, 2022.
- 608
- 609 Mehmed Kantardzic. *Data mining: concepts, models, methods, and algorithms*. John Wiley &
610 Sons, 2011.
- 611
- 612 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
613 based generative models. *Advances in neural information processing systems*, 35:26565–26577,
2022.
- 614
- 615 Dongjun Kim, Yeongmin Kim, Se Jung Kwon, Wanmo Kang, and Il-Chul Moon. Refining gen-
616 erative process with discriminator guidance in score-based diffusion models. *arXiv preprint*
arXiv:2211.17091, 2022.
- 617
- 618 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
619 *arXiv:1312.6114*, 2013.
- 620
- 621 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile
622 diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- 623
- 624 Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive
625 facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition*
(CVPR), 2020.
- 626
- 627 Steven Cheng-Xian Li, Bo Jiang, and Benjamin Marlin. Misgan: Learning from incomplete data
628 with generative adversarial networks. *arXiv preprint arXiv:1902.09599*, 2019.
- 629
- 630 Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural net-
631 work: Data-driven traffic forecasting. *arXiv preprint arXiv:1707.01926*, 2017.
- 632
- 633 Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*. New York: Wiley,
1987.
- 634
- 635 Yixin Liu, Thalaiyasingam Ajanthan, Hisham Husain, and Vu Nguyen. Self-supervision improves
636 diffusion models for tabular data imputation. In *Proceedings of the 33rd ACM International*
Conference on Information and Knowledge Management, pp. 1513–1522, 2024.
- 637
- 638 Yukai Liu, Rose Yu, Stephan Zheng, Eric Zhan, and Yisong Yue. Naomi: Non-autoregressive
639 multiresolution sequence imputation. *Advances in neural information processing systems*, 32,
2019.
- 640
- 641 Chao Ma and Cheng Zhang. Identifiable generative models for missing not at random data imputa-
642 tion. *Advances in Neural Information Processing Systems*, 34:27645–27658, 2021.
- 643
- 644 Pierre-Alexandre Mattei and Jes Frellsen. Miwae: Deep generative modelling and imputation of
645 incomplete data sets. In *International conference on machine learning*, pp. 4413–4423. PMLR,
2019.
- 646
- 647 Boris Muzellec, Julie Josse, Claire Boyer, and Marco Cuturi. Missing data imputation using optimal
transport. In *International Conference on Machine Learning*, pp. 7130–7140. PMLR, 2020.

- 648 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
649 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 650
- 651 Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising
652 diffusion models for multivariate probabilistic time series forecasting. In *International conference*
653 *on machine learning*, pp. 8857–8868. PMLR, 2021.
- 654 Trevor W Richardson, Wencheng Wu, Lei Lin, Beilei Xu, and Edgar A Bernal. Mcflow: Monte
655 carlo flow models for data imputation. In *Proceedings of the IEEE/CVF conference on computer*
656 *vision and pattern recognition*, pp. 14205–14214, 2020.
- 657
- 658 Rajhans Samdani, Ming-Wei Chang, and Dan Roth. Unified expectation maximization. In *Proceed-*
659 *ings of the 2012 Conference of the North American Chapter of the Association for Computational*
660 *Linguistics: Human Language Technologies*, pp. 688–698, 2012.
- 661 Joseph L Schafer. *Analysis of incomplete multivariate data*. CRC press, 1997.
- 662
- 663 Maximilian Seitzer. pytorch-fid: FID Score for PyTorch. [https://github.com/mseitzer/](https://github.com/mseitzer/pytorch-fid)
664 [pytorch-fid](https://github.com/mseitzer/pytorch-fid), August 2020. Version 0.3.0.
- 665 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
666 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-*
667 *ional Conference on Learning Representations*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=PXTIG12RRHS)
668 [forum?id=PXTIG12RRHS](https://openreview.net/forum?id=PXTIG12RRHS).
- 669
- 670 Daniel J Stekhoven and Peter Bühlmann. Missforest—non-parametric missing value imputation for
671 mixed-type data. *Bioinformatics*, 28(1):112–118, 2012.
- 672
- 673 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. CsdI: Conditional score-based dif-
674 fusion models for probabilistic time series imputation. *Advances in neural information processing*
675 *systems*, 34:24804–24816, 2021.
- 676
- 677 Stef Van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equa-
678 tions in r. *Journal of statistical software*, 45:1–67, 2011.
- 679
- 680 Shiyu Wang, Jiawei Li, Xiaoming Shi, Zhou Ye, Baichuan Mo, Wenzhe Lin, Shengtong Ju, Zhixuan
681 Chu, and Ming Jin. Timemixer++: A general time series pattern machine for universal predictive
682 analysis. *arXiv preprint arXiv:2410.16032*, 2024.
- 683
- 684 Haixu Wu, Tengge Hu, Yong Liu, Hang Zhou, Jianmin Wang, and Mingsheng Long. Timesnet: Tem-
685 poral 2d-variation modeling for general time series analysis. *arXiv preprint arXiv:2210.02186*,
686 2022.
- 687
- 688 Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmark-
689 ing machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 690
- 691 Jinsung Yoon, James Jordon, and Mihaela Schar. Gain: Missing data imputation using generative
692 adversarial nets. In *International conference on machine learning*, pp. 5689–5698. PMLR, 2018a.
- 693
- 694 Jinsung Yoon, William R Zame, and Mihaela Van Der Schar. Estimating missing data in temporal
695 data streams using multi-directional recurrent neural networks. *IEEE Transactions on Biomedical*
696 *Engineering*, 66(5):1477–1490, 2018b.
- 697
- 698 Xinyu Yuan and Yan Qiao. Diffusion-ts: Interpretable diffusion for general time series generation.
699 *CoRR*, 2024.
- 700
- 701 Hengrui Zhang, Liancheng Fang, Qitian Wu, and Philip S Yu. Diffputer: Empowering diffusion
models for missing data imputation. In *The Thirteenth International Conference on Learning*
Representations, 2025.
- 702
- 703 Shuyi Zhang, Bin Guo, Anlan Dong, Jing He, Ziping Xu, and Song Xi Chen. Cautionary tales on
air-quality improvement in beijing. *Proceedings of the Royal Society A: Mathematical, Physical*
and Engineering Sciences, 473(2205):20170457, 2017.

702 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai Zhang.
703 Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings*
704 *of the AAAI conference on artificial intelligence*, volume 35, pp. 11106–11115, 2021.
705
706 Jianping Zhou, Junhao Li, Guanjie Zheng, Xinbing Wang, and Chenghu Zhou. Mtsci: A conditional
707 diffusion model for multivariate time series consistent imputation. In *Proceedings of the 33rd*
708 *ACM International Conference on Information and Knowledge Management*, pp. 3474–3483,
709 2024.
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A RELATED WORKS

Traditional imputation approaches typically rely on simple statistical heuristics such as mean, median, or last observation carried forward to fill in missing entries of multivariate time-series data (Kantardzic, 2011). These strategies often fail to capture the complex temporal dynamics and cross-feature dependencies in them. To improve upon these methods, more sophisticated methods such as MICE (Van Buuren & Groothuis-Oudshoorn, 2011), which iteratively applies the EM algorithm, and MissForest (Stekhoven & Bühlmann, 2012), which leverages random forests for iterative refinement, have been proposed. Although these techniques provide more better imputation performances than naive statistical rules, their capacity remains limited when handling high-dimensional time-series data with intricate dependencies.

Imputation with Deep Learning Early attempts to exploit deep learning for time-series imputation include mRNN (Yoon et al., 2018b), which leverages recurrent neural networks to capture temporal dependencies and model complex patterns in partially observed sequences. BRITS (Cao et al., 2018) further improves upon RNN-based imputers by introducing a bidirectional structure, allowing information to flow forward and backward across time to enhance estimation accuracy. Then NAOMI (Liu et al., 2019) combines multi-resolution RNNs with adversarial training strategies to refine imputations at different time scales. SAITS (Du et al., 2023a) introduces a self-attention mechanism to better capture long-range temporal dependencies. Most of them could be reproduced with the released Python toolkit (Du et al., 2023b).

Beyond time-series-specific architectures, several general-purpose imputation frameworks have also shaped the development of recent methods. GAIN (Yoon et al., 2018a), although not tailored to time-series data, was the first to introduce an adversarial discriminator to imputation, providing a novel mechanism to distinguish observed values and missing values. Successively, misGAN (Li et al., 2019) provided multiple generators and discriminators system for stable training across varied missing patterns. Flow-based approaches such as MCFLOW (Richardson et al., 2020) further demonstrated that the EM algorithm can be combined with invertible generative models to jointly optimize flow parameters and missing entries. not-MIWAE (Ipsen et al., 2020) addressed the MNAR scenario by explicitly optimizing a missing model within the ELBO objective, which can inference missing values by missing model weighted importance sampling. Although originally proposed for general imputation tasks, these frameworks have significantly influenced subsequent advances in time-series imputation by highlighting the value of probabilistic, adversarial, and likelihood-based modeling.

Diffusion-based Approaches for Imputation In the time-series domain, TimeGrad (Rasul et al., 2021) applies diffusion to probabilistic forecasting, though its design mainly focuses on forecasting task. For the imputation task, CSDI (Tashiro et al., 2021) introduces a conditional diffusion framework with masking to handle arbitrary missing. Building on this line of research, methods such as SSSD (Alcaraz & Strodthoff, 2022) and Diffusion-TS (Yuan & Qiao, 2024) incorporate additional regularization losses tailored to time-series characteristics, thereby enhancing the interpretability of the imputed sequences. On the other hand, MTSCI (Zhou et al., 2024) integrates a contrastive loss to maximize mutual information between observed variable and missed variable which improves generated sample consistency. More recently, DiffPuter (Zhang et al., 2025) further improves probabilistic imputation with the EM algorithm, which progressively refines the missing values.

B PROOFS

B.1 PROOF OF PROPOSITION 3.1

Proposition 3.1. *Suppose that the variational distribution and missing process satisfy the following assumptions: (i) Forward process satisfies conditional independence to mask variable given entire data $q(X_{1:T}|X_0, M) = q(X_{1:T}|X_0)$, (ii) Missing process satisfies conditional independence to noisy data $p_\phi(M|X_0, X_{1:T}) = p_\phi(M|X_0)$.*

Then, the ELBO of joint log-likelihood objective of the observed data and mask can be expressed as

$$\log p_{\theta, \phi}(X_0^{obs}, M) \geq \mathbb{E}_{X_{1:T}, X_0^{mis}} \left[\log p_{\phi}(M|X_0) + \log \frac{p_{\theta}(X_T) \prod_{t=1}^T p_{\theta}(X_{t-1}|X_t)}{\prod_{t=1}^T q(X_t|X_{t-1})} \right] \quad (10)$$

$$+ \mathbb{E}_{X_{1:T}} \left[\mathbb{H}(q(X_0^{mis}|X_0^{obs}, M)) \right] \quad (11)$$

Proof.

$$\log p_{\theta, \phi}(X^{obs}, M) = \log \int_Z \int_{X^{mis}} p_{\theta, \phi}(X^{obs}, X^{mis}, M, Z) dX^{mis} dZ \quad (16)$$

$$= \log \int_Z \int_{X^{mis}} \frac{p_{\theta, \phi}(X^{obs}, X^{mis}, M, Z)}{q(X^{mis}, Z|X^{obs}, M)} q(X^{mis}, Z|X^{obs}, M) dX^{mis} dZ \quad (17)$$

$$= \log \int_Z \int_{X^{mis}} \frac{p_{\phi}(M|X) p_{\theta}(X|Z) p_{\theta}(Z)}{q(Z|X) q(X^{mis}|X^{obs}, M)} q(Z|X) q(X^{mis}|X^{obs}, M) dX^{mis} dZ \quad (18)$$

$$\geq \mathbb{E}_{Z \sim q(Z|X), X^{mis} \sim q(X^{mis}|X^{obs}, M)} \left[\log p_{\phi}(M|X) + \log \frac{p_{\theta}(X|Z) p_{\theta}(Z)}{q(Z|X)} \right] \quad (19)$$

$$+ \mathbb{E}_{Z \sim q(Z|X)} \left[\mathbb{H}(q(X^{mis}|X^{obs}, M)) \right] \quad (20)$$

Equation 19 defines the loss objective between the true parameters and the corresponding variational distribution. Replacing the latent variable Z with the diffusion latents $X_{1:T}$, can be formulated as equation 21 under the Markov property of the diffusion process.

$$\mathbb{E}_{X_{1:T} \sim q(X_{1:T}|X), X^{mis} \sim q(X^{mis}|X^{obs}, M)} \left[\log p_{\phi}(M|X) + \log \frac{p_{\theta}(X|X_1) p_{\theta}(X_{1:T})}{q(X_{1:T}|X)} \right] \quad (21)$$

□

B.2 PROOF OF PROPOSITION 3.2

Proposition 3.2. Suppose that the pattern recognizer D_{ϕ} is optimal which satisfies $D_{\phi^*}(X_0) = [p_{\phi}(M_d|X_0)] \in \mathbb{R}^D$, the score function of the joint log-likelihood with respect to the missing mask can be approximated as

$$\nabla_{X_t} \log p_{\theta, \phi}(X_t|X_0^{obs}, M) \simeq \nabla_{X_t} \log p_{\theta}(X_t|X_0^{obs}) - \nabla_{X_t} \mathcal{L}_{PR}(M, \hat{X}_0, D_{\phi^*}) \quad (12)$$

where $\hat{X}_0 := (f_{\theta}(X_t, t; X_0^{obs})^{mis}, X_0^{obs}) = f_{\theta}(X_t, t; X_0^{obs}) \odot (1 - M) + X_0^{obs} \odot M$.

Proof. From the graphical model of Figure 2 (a), conditional independence for missing process satisfies $p_{\phi}(M|X_0, \hat{X}_t) = p_{\phi}(M|X_0)$ at Equation 23. The approximation 24 came from the Thm 1 of (Chung et al., 2023).

$$\nabla_{X_t} \log p_{\theta, \phi}(X_t|X_0^{obs}, M) = \nabla_{X_t} \log p_{\theta}(X_t|X_0^{obs}) + \nabla_{X_t} \log p_{\theta, \phi}(M|X_t, X_0^{obs}) \quad (22)$$

$$= \nabla_{X_t} \log p_{\theta}(X_t|X_0^{obs}) + \nabla_{X_t} \log \int p_{\phi}(M|X_0) p_{\theta}(X_0^{mis}|X_t, X_0^{obs}) dX_0^{mis} \quad (23)$$

$$\simeq \nabla_{X_t} \log p_{\theta}(X_t|X_0^{obs}) + \nabla_{X_t} \log p_{\phi}(M|f_{\theta}(X_t, t; X_0^{obs})^{mis}, X_0^{obs}) \quad (24)$$

$$= \nabla_{X_t} \log p_{\theta}(X_t|X_0^{obs}) + \nabla_{X_t} M \log D_{\phi^*}(f_{\theta}(X_t, t; X_0^{obs})^{mis}, X_0^{obs}) \quad (25)$$

$$+ \nabla_{X_t} (1 - M) \log \left\{ 1 - D_{\phi^*}(f_{\theta}(X_t, t; X_0^{obs})^{mis}, X_0^{obs}) \right\} \quad (26)$$

$$= \nabla_{X_t} \log p_{\theta}(X_t|X_0^{obs}) - \nabla_{X_t} \mathcal{L}_{PR}(M, (f_{\theta}(X_t, t; X_0^{obs})^{mis}, X_0^{obs}), D_{\phi^*}) \quad (27)$$

Since d -th element of $D_{\phi^*}(X_0)$ converges to $p_{\phi}(M_d = 1|X_0)$ (Ipsen et al., 2020; Ma & Zhang, 2021), entire probability of mask variable M follows:

$$\begin{aligned} \log p_{\phi}(M|f_{\theta}(X_t, t; X_0^{\text{obs}})^{\text{mis}}, X_0^{\text{obs}}) &= \sum_{i=1}^D \log p_{\phi}(M_d|f_{\theta}(X_t, t; X_0^{\text{obs}})^{\text{mis}}, X_0^{\text{obs}}) \\ &= M \log D_{\phi^*}\left(f_{\theta}(X_t, t; X_0^{\text{obs}})^{\text{mis}}, X_0^{\text{obs}}\right) + (1 - M) \log \left\{1 - D_{\phi^*}\left(f_{\theta}(X_t, t; X_0^{\text{obs}})^{\text{mis}}, X_0^{\text{obs}}\right)\right\} \end{aligned} \tag{28}$$

$$\tag{29}$$

□

B.3 REWRITTEN THEOREM 1 OF ZHANG ET AL. (2025)

Theorem 1. Let X_T be a sample from the prior distribution $p_{\theta}(X_T) = \mathcal{N}(0, \mathbf{I})$, X be the data to impute, and the known entries of X are denoted by $X^{\text{obs}} = X_0^{\text{obs}}$. The score function $\nabla_{X_t} \log p(X_t)$ could be parameterized by neural network $f_{\theta}(X_t, t; X_0^{\text{obs}})$. Applying forward and reverse process of the diffusion model iteratively from $t = T \gg 0$ until $t = 0$ with $\Delta t \rightarrow 0$, then \hat{X}_0 is a sample from $p_{\theta}(X)$, under the condition that its observed entries $\hat{X}_0^{\text{obs}} = X_0^{\text{obs}}$. Formally,

$$\hat{X}_0 \sim p_{\theta}(X|X^{\text{obs}} = X_0^{\text{obs}}) \tag{30}$$

This section presents a rewritten version of Theorem 1 from Zhang et al. (2025), and we refer the reader to the original paper for the detailed proof. The theorem establishes that when learning the joint probability of X^{obs} and X^{mis} , the missing values can be inferred by conditioning on the observed values of a given sample. We adopt the same line of reasoning in Section 3.3 to support our theoretical development.

C EXPERIMENT DETAILS

Model Configuration In our implementation, the pattern recognizer is designed as a lightweight multi-layer perceptron (MLP) to minimize additional model complexity. As summarized in Table 8, the overall model size remains comparable to other diffusion-based approaches. The higher inference time of PRDIM, relative to competing methods, arises from the use of autograd of the input with respect to the outputs of the pattern recognizer during the inference process. This design choice, while incurring additional computational cost, enables the model to provide more informative guidance for imputing missing values.

Table 8: Detailed model configuration. Model Size indicates the number of parameter of each diffusion-based model, and $+n$ in PRDIM shows the number of pattern recognizer. The inference time (s) is measured based on the a single inference required to impute the entire out-of-sample data.

Method	ETT			STOCK		
	Model Size	Training Time (s)	Inference Time (s)	Model Size	Training Time (s)	Inference Time (s)
CSDI	164025	366	39	164017	138	11
MTSCI	162321	585	21	146969	354	12
cDiffPuter	163769	981	28	163761	516	11
PRDIM	163769+17376	1812	47	163761+17359	1052	28

EM Configuration Since PRDIM is built upon the EM framework, both training time and performance are affected by the design of EM configuration. In this section, we analyze the impact of different EM settings on the STOCK dataset and show that the PRDIM can overcome its inherent training time consumption with respect to imputation performance. We denote $1\mathbf{E} N_m \mathbf{M}$ as the number of training epochs for the maximization step per expectation step, and N as the total number of EM iterations. In the main

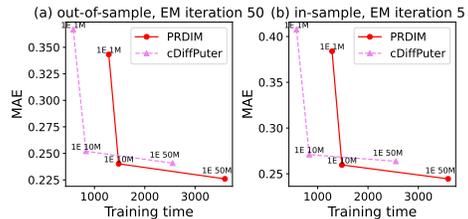


Figure 7: Training time and imputation performance (MAE) under different EM configurations on the STOCK dataset.

experiment, both cDiffPuter and PRDIM were trained with $N = 100$ EM iterations and **1E 1M** configuration, i.e., one epoch for the maximization step per expectation step.

When varying the number of training epochs in the maximization step, we observe that the training time increases proportionally to the epoch size. As the number of training epochs in each maximization step increases, the relative computational cost of the expectation step becomes negligible. Thus, the training time of PRDIM converges to that of cDiffPuter as the maximization step grows. Nevertheless, as shown in Figures 7 (a) and 7 (b), PRDIM with the configuration of **1E 10M** and a total of 50 EM iterations achieves superior performance and requires less training time compared to cDiffPuter whose configuration is **1E 50M** and 50 EM iterations. This observation highlights the importance of appropriate hyperparameter tuning of the EM configuration, and further suggests that PRDIM provides strong scalability and practical utility for imputation tasks on real-world datasets.

Dataset Configuration Table 9 summarizes the dataset configurations employed in our experiments across (train / test / valid) set. For the three time-series datasets (ETT, STOCK, and PEMS-Bay), the data size are represented as *time length* \times *feature dimension*, while for FMNIST, the dimensions are denoted as *width* \times *height*. The reported missing ratios correspond to the proportion of original missing values observed after applying the MNAR mechanism to generate incomplete data for training, thereby reflecting the intrinsic difficulty of the imputation task.

Table 9: Dataset configuration used in the main experiments. For time-series datasets (ETT, STOCK, PEMS-Bay), the data size is denoted as *time length* \times *feature dimension*, whereas for FMNIST, it is expressed as *width* \times *height*. The missing ratios represent the proportion of original missing values after applying the MNAR mechanism to construct incomplete data for training.

Dataset	ETTM1	STOCK	PEMS-Bay	FMNIST
Data size	24×7	24×6	12×325	28×28
# of Samples	3861 / 983 / 959	2418 / 622 / 622	5788 / 1448 / 1448	60000 / 5000 / 5000
Missing ratio (%)	21.4 / 43.9 / 14.0	21.2 / 20.0 / 20.9	13.5 / 13.0 / 14.1	25.8 / 25.8 / 25.8

Missing Mechanisms Table 2 reports results obtained under the following MNAR mechanism. Inspired by the MNAR mechanism of not-MIWAE (Ipsen et al., 2020), we design the MNAR mechanism such that the probability of a missing entry increases exponentially with its value:

$$p(M_d = 1|X_d) = \frac{1}{1 + e^{-\text{logits}}}, \quad \text{logits} = W(X_d - b), \quad (31)$$

where $M_d \in \{0, 1\}$ denotes the mask variable of entry X_d , W controls the slope, and b is a bias term. This mechanism ensures that entries with values larger than the mean are more likely to be missing, thus faithfully mimicking MNAR conditions. In main experiments on section 4.2, we set $W = 5$ and $b = 0.8$ for all time-series dataset while $W = 7$ and $b = 0.6$ for FMNIST.

To verify coherent results under MCAR and MNAR mechanisms, we follow the missingness simulation procedures of Hyperimpute (Jarrett et al., 2022) and MissingOT (Muzellec et al., 2020), as reported in Table 5 and Table 4. In MCAR setting, each value is excluded following the Bernoulli random variable with a fixed parameter. In our implementation, we randomly assign 10% of the entries as missing to maintain MCAR property. Otherwise, to construct the MNAR setting, we employ a quantile-based mechanism distinct from the previous logistic approach. A subset of variables is randomly selected, and missing values are generated within the q -quantile. In our implementation, a missing rate of around 30% was created using a bilateral 25% quantile.

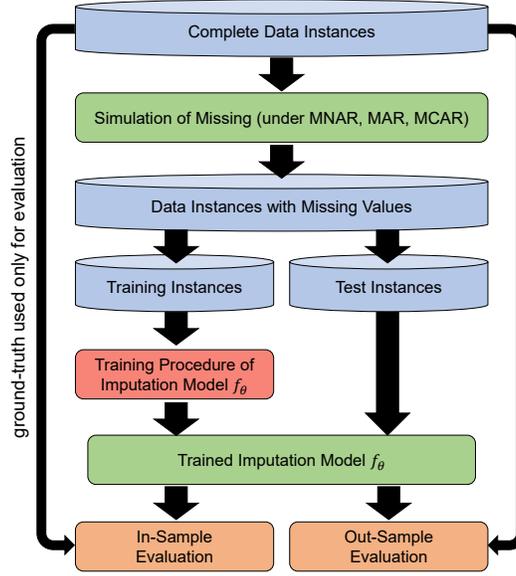


Figure 8: Overall information flow between the in-sample imputation and the out-of-sample imputation.

Evaluation Metrics We report three error-based metrics:

$$\text{RMSE} = \sqrt{\frac{\sum_{n=1, d=1}^{N, D} (X_d^n - \hat{X}_d^n)^2 \times M_d^n}{\sum_{n=1, d=1}^{N, D} M_d^n}}, \quad (32)$$

$$\text{MAE} = \frac{\sum_{n=1, d=1}^{N, D} |X_d^n - \hat{X}_d^n| \times M_d^n}{\sum_{n=1, d=1}^{N, D} M_d^n}, \quad (33)$$

$$\text{MRE} = \frac{\sum_{n=1, d=1}^{N, D} |X_d^n - \hat{X}_d^n| \times M_d^n}{\sum_{n=1, d=1}^{N, D} |X_d^n| \times M_d^n} \times 100(\%), \quad (34)$$

where X_d^n denotes the ground-truth value of dimension d of n th sample and \hat{X}_d^n its imputed counterpart. These complementary measures assess squared error, absolute error, and relative error, providing a comprehensive evaluation of imputation performance.

In-sample and Out-of-sample imputation In our experiments, In-sample imputation refers to evaluating imputation performance on the same dataset used for training, whereas out-of-sample imputation evaluates the model on held-out test splits containing unseen data. The simulated MNAR pattern distributions are generated consistently for each split to ensure complete evaluation. Figure 8 summarizes the distinction between the in-sample imputation process and the out-of-sample imputation process.

C.1 OVERALL PERFORMANCE

In this section, we provide additional results on the overall imputation performance across all time-series datasets in Table 10, Table 11, and Table 12. These results complement the main findings reported in the paper and further validate the effectiveness of our approach.

Table 10: Overall performance on the **ETT** dataset. We report $mean \pm std$ over 5 runs according to each methodology. Best results are in **bold**. Second best results are in underline.

Method	Original / Out-of-Sample			Original / In-Sample		
	RMSE	MAE	MRE	RMSE	MAE	MRE
Mean	2.307 \pm 0.000	2.034 \pm 0.000	120.233 \pm 0.000	1.618 \pm 0.000	1.486 \pm 0.000	127.379 \pm 0.000
<i>Discriminative models</i>						
TimesNet	1.393 \pm 0.038	1.044 \pm 0.065	69.040 \pm 4.305	1.485 \pm 0.044	1.154 \pm 0.068	72.065 \pm 4.222
TimeMixer++	1.965 \pm 0.012	1.642 \pm 0.025	97.093 \pm 0.015	1.283 \pm 0.015	1.100 \pm 0.032	94.319 \pm 0.028
BRITS	1.461 \pm 0.048	0.992 \pm 0.037	58.600 \pm 0.022	0.850 \pm 0.020	0.491 \pm 0.008	42.067 \pm 0.007
SAITS	1.247 \pm 0.069	0.814 \pm 0.046	48.119 \pm 0.027	0.626 \pm 0.018	0.366 \pm 0.014	31.417 \pm 0.012
<i>Generative models</i>						
GP-VAE	1.915 \pm 0.006	1.511 \pm 0.011	89.315 \pm 0.638	1.147 \pm 0.008	0.896 \pm 0.018	76.809 \pm 1.507
not-MIWAE	1.781 \pm 0.012	1.311 \pm 0.016	77.512 \pm 0.972	0.945 \pm 0.010	0.637 \pm 0.011	54.643 \pm 0.918
<i>Diffusion-based models</i>						
CSDI	1.658 \pm 0.001	1.071 \pm 0.001	63.254 \pm 0.038	0.822 \pm 0.001	0.522 \pm 0.001	44.733 \pm 0.049
MTSCI	1.335 \pm 0.001	0.957 \pm 0.001	56.574 \pm 0.036	0.730 \pm 0.000	0.500 \pm 0.000	42.827 \pm 0.018
cDiffPuter	<u>1.209</u> \pm 0.001	<u>0.782</u> \pm 0.000	<u>46.188</u> \pm 0.020	<u>0.612</u> \pm 0.001	<u>0.362</u> \pm 0.000	<u>31.069</u> \pm 0.020
PRDIM	1.057 \pm 0.000	0.663 \pm 0.000	39.156 \pm 0.009	0.538 \pm 0.000	0.303 \pm 0.000	25.986 \pm 0.015

Table 11: Overall performance on the **STOCK** dataset. We report $mean \pm std$ over 5 runs according to each methodology. Best results are in **bold**. Second best results are in underline.

Method	Original / Out-of-Sample			Original / In-Sample		
	RMSE	MAE	MRE	RMSE	MAE	MRE
Mean	2.079 \pm 0.000	1.949 \pm 0.000	128.903 \pm 0.000	2.168 \pm 0.000	2.039 \pm 0.000	127.313 \pm 0.000
<i>Discriminative models</i>						
TimesNet	1.415 \pm 0.054	1.111 \pm 0.073	73.528 \pm 0.049	1.509 \pm 0.057	1.221 \pm 0.077	76.237 \pm 0.048
TimeMixer++	1.490 \pm 0.223	1.287 \pm 0.239	85.153 \pm 0.158	1.569 \pm 0.239	1.369 \pm 0.260	85.456 \pm 0.162
BRITS	0.953 \pm 0.016	0.627 \pm 0.010	41.478 \pm 0.006	1.020 \pm 0.016	0.701 \pm 0.010	43.757 \pm 0.006
SAITS	0.743 \pm 0.021	0.442 \pm 0.022	29.115 \pm 0.015	0.801 \pm 0.023	0.498 \pm 0.025	31.071 \pm 0.016
<i>Generative models</i>						
GP-VAE	1.239 \pm 0.118	0.902 \pm 0.109	59.684 \pm 7.220	1.333 \pm 0.123	1.010 \pm 0.118	63.046 \pm 7.371
not-MIWAE	1.028 \pm 0.043	0.681 \pm 0.045	45.039 \pm 0.296	1.114 \pm 0.043	0.759 \pm 0.039	47.368 \pm 0.243
<i>Diffusion-based models</i>						
CSDI	0.932 \pm 0.000	0.641 \pm 0.000	42.393 \pm 0.004	0.995 \pm 0.000	0.710 \pm 0.000	44.330 \pm 0.001
MTSCI	0.988 \pm 0.002	0.736 \pm 0.001	48.629 \pm 0.009	1.056 \pm 0.001	0.809 \pm 0.001	50.485 \pm 0.041
cDiffPuter	<u>0.734</u> \pm 0.000	<u>0.406</u> \pm 0.000	<u>26.878</u> \pm 0.100	<u>0.778</u> \pm 0.000	<u>0.450</u> \pm 0.000	<u>28.064</u> \pm 0.008
PRDIM	0.599 \pm 0.001	0.254 \pm 0.000	16.794 \pm 0.027	0.633 \pm 0.000	0.275 \pm 0.000	17.150 \pm 0.008

C.2 ADDITIONAL EXPERIMENTS

C.2.1 FASHION-MNIST

In the image domain, missing values are not limited to the MNAR mechanism demonstrated in the main experiment. As a representative example in the image dataset, we additionally exhibit imputation results under the *block missing* mechanism. Experiments are conducted on the FMNIST dataset, and we further include comparisons with representative GAN-based and Flow-based imputation approaches, namely misGAN (Li et al., 2019) and MCFlow (Richardson et al., 2020).

As illustrated by the following Figure 9, among the three methods, our proposed PRDIM most effectively captures the underlying object structure and achieves the most faithful reconstructions. These findings demonstrate that PRDIM can generalize beyond MNAR to handle other types of missingness, such as block-MAR, while retaining its ability to generate semantically plausible imputations.

Table 12: Overall performance on the **PEMS-Bay** dataset. We report $mean \pm std$ over 5 runs according to each methodology. Best results are in **bold**. Second best results are in underline.

Method	Original / Out-of-Sample			Original / In-Sample		
	RMSE	MAE	MRE	RMSE	MAE	MRE
Mean	0.901 \pm 0.000	0.813 \pm 0.000	119.064 \pm 0.000	0.868 \pm 0.000	0.789 \pm 0.000	119.066 \pm 0.000
<i>Discriminative models</i>						
TimesNet	0.481 \pm 0.009	0.291 \pm 0.007	42.579 \pm 0.010	0.392 \pm 0.004	0.225 \pm 0.001	33.970 \pm 0.002
TimeMixer++	0.684 \pm 0.013	0.579 \pm 0.018	84.816 \pm 0.026	0.652 \pm 0.015	0.557 \pm 0.020	84.113 \pm 0.030
BRITS	0.503 \pm 0.007	0.278 \pm 0.006	40.758 \pm 0.008	0.342 \pm 0.004	0.182 \pm 0.003	27.490 \pm 0.004
SAITS	0.481 \pm 0.011	0.302 \pm 0.009	44.266 \pm 0.013	0.356 \pm 0.004	0.212 \pm 0.003	31.970 \pm 0.005
<i>Generative models</i>						
GP-VAE	0.537 \pm 0.001	0.345 \pm 0.001	50.561 \pm 0.208	0.470 \pm 0.003	0.292 \pm 0.002	44.039 \pm 0.367
not-MIWAE	0.623 \pm 0.006	0.396 \pm 0.005	57.510 \pm 0.774	0.608 \pm 0.004	0.352 \pm 0.005	53.181 \pm 0.746
<i>Diffusion-based models</i>						
CSDI	<u>0.338</u> \pm 0.002	<u>0.177</u> \pm 0.000	<u>25.912</u> \pm 0.017	0.302 \pm 0.000	<u>0.158</u> \pm 0.000	<u>23.910</u> \pm 0.005
MTSCI	0.349 \pm 0.000	0.193 \pm 0.000	28.289 \pm 0.011	0.322 \pm 0.000	0.179 \pm 0.000	27.017 \pm 0.003
cDiffPuter	0.349 \pm 0.007	0.182 \pm 0.000	26.714 \pm 0.011	0.330 \pm 0.000	0.168 \pm 0.000	25.377 \pm 0.005
PRDIM	0.334 \pm 0.002	0.170 \pm 0.000	24.966 \pm 0.015	<u>0.306</u> \pm 0.000	0.154 \pm 0.000	23.304 \pm 0.006

Furthermore, to highlight the general imputation ability of our model under the MNAR mechanism from the main experiment, we also present additional qualitative results in Figure 10. For both experiments, we evaluate the quality of generated samples using the Fréchet Inception Distance (FID) (Heusel et al., 2017), which is computed with the released Python library (Seitzer, 2020).

C.2.2 CELEBA-HQ

To verify the scalability of PRDIM on high-dimensional data, we conducted an additional imputation experiment on the RGB image benchmark dataset named CelebA-HQ (Lee et al., 2020). We compared our method with a vanilla diffusion model trained under the CSDI objective. Each image in CelebA-HQ is accompanied by corresponding annotation mask vectors that label facial attributes such as eyes, nose, mouth, and hair. To design an incomplete dataset under an MNAR pattern, we utilized the annotation masks of eyes, nose, and mouth to construct a missing-value mask. Specifically, for each facial attribute, we introduced missing pixels with an 80% probability within the annotated regions, forming an MNAR missing mechanism for the experiment.

To efficiently manage the training time of the EM-based PRDIM model trained from a scratch, both images and their corresponding mask vectors were resized to a resolution of 64 \times 64. The original CelebA-HQ dataset consists of 1024 \times 1024 images and 512 \times 512 annotation masks.

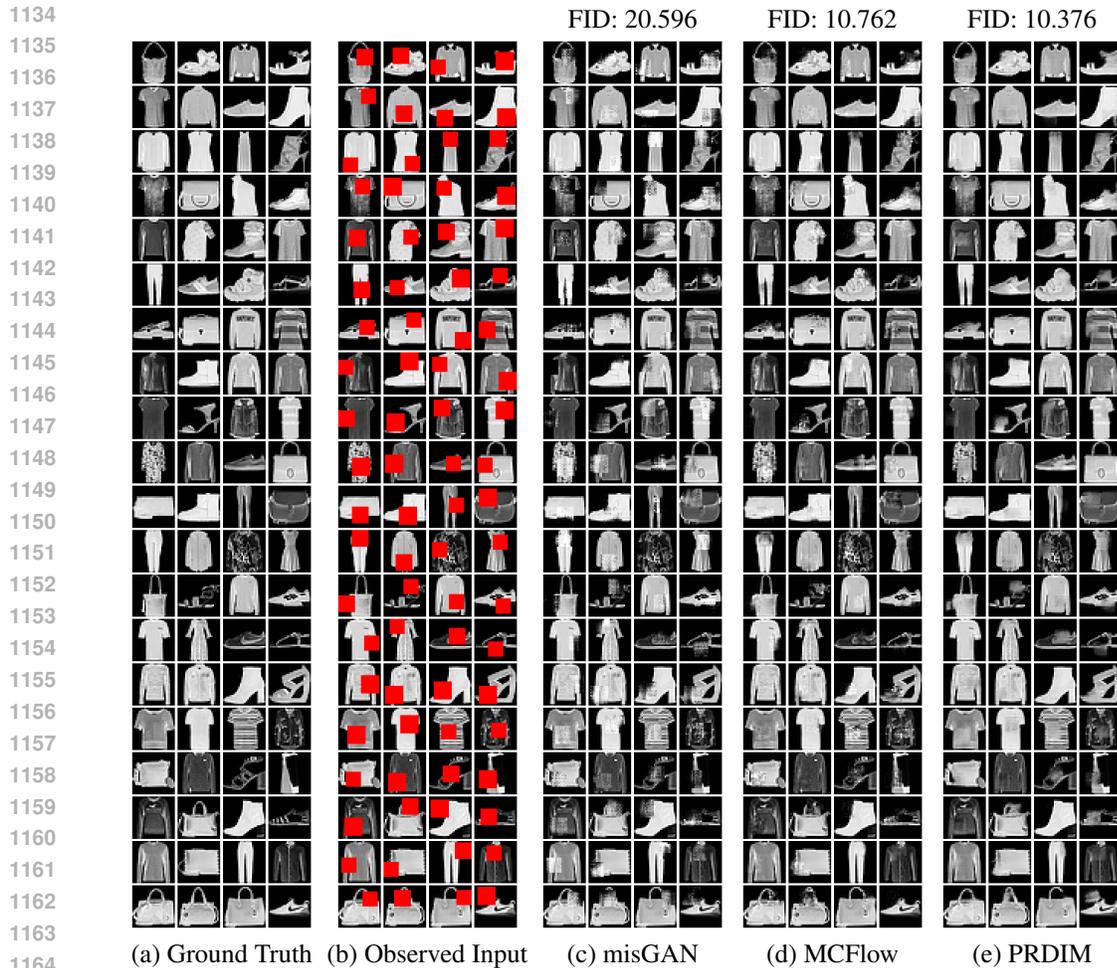
Qualitative results are shown in Figure 11. The vanilla diffusion model trained with the CSDI objective tends to fill in missing areas with averaged color tones around the missing regions, resulting in naive reconstructions and relatively high FID scores despite a moderate missing ratio. In contrast, PRDIM generates more detailed and realistic facial structures, accurately reconstructing attribute boundaries and color variations, which leads to significantly improved perceptual quality and lower FID values.

C.2.3 TABULAR DATA

To further demonstrate the generalization capability of PRDIM, we reproduced the official implementation of DiffPuter¹ and compared its performance with PRDIM under the MNAR setting.

In this experiment, we followed the practical implementation of DiffPuter with default configuration regardless to dataset, which differs from the main experiments in that incomplete samples were not used as conditional information during imputation.

¹<https://github.com/hengruizhang98/DiffPuter>



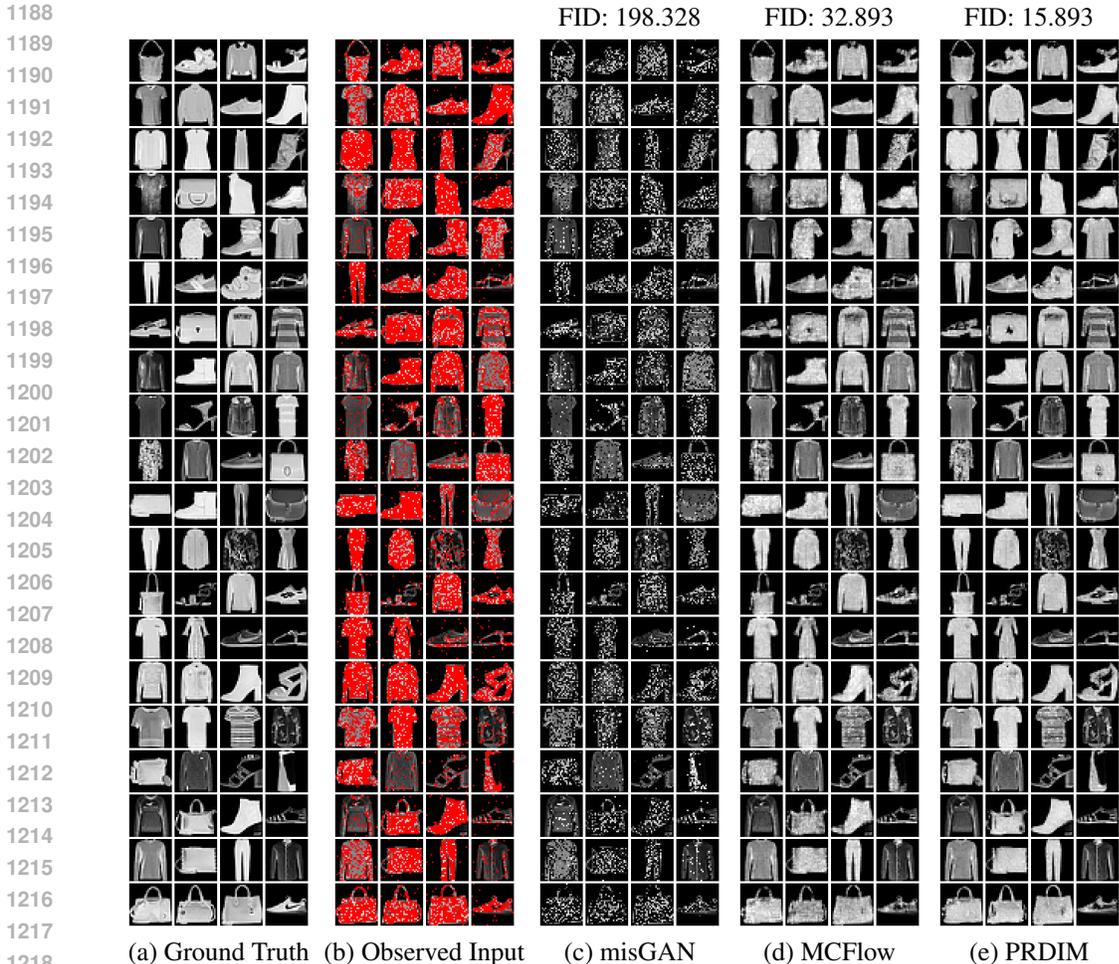
1165 Figure 9: Comparison of out-of-sample imputation results under block missing pattern.
1166
1167
1168

1169 Table 13: In-sample imputation results across tabular datasets. We report the mean \pm std over five
1170 different training/test split combinations.
1171

	adult	bean	default	gesture	magic
DiffPuter	0.497 \pm 0.013	0.240 \pm 0.069	0.374 \pm 0.092	0.391\pm0.023	0.539 \pm 0.088
PRDIM	0.474\pm0.012	0.199\pm0.058	0.336\pm0.074	0.394 \pm 0.029	0.490\pm0.083

1172
1173
1174
1175
1176
1177
1178
1179 We selected 5 different tabular datasets available from the UCI Machine Learning Repository.
1180 Among them, bean, gesture, and magic consist solely of continuous features, while adult and default
1181 contain both continuous and discrete features. The discrete attributes were label-encoded to preserve
1182 the original data dimensionality, and the corresponding mask vectors were designed to match this
1183 structure.

1184 Tables 13 and 14 present the in-sample and out-of-sample imputation performance of DiffPuter and
1185 PRDIM, respectively. Across most datasets, PRDIM achieves more accurate imputations, validating
1186 its robustness and adaptability across different data modalities.
1187



1219 Figure 10: Comparison of out-of-sample imputation results under MNAR missing pattern.
1220

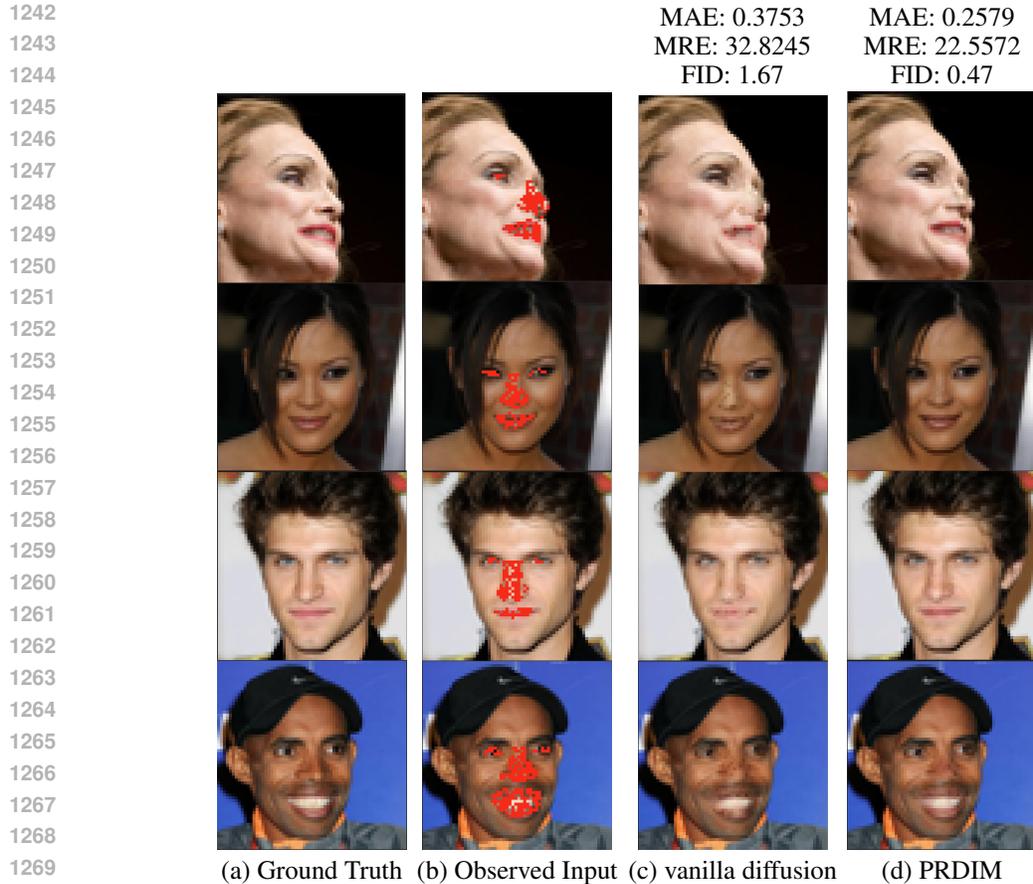
1221
1222 Table 14: Out-of-sample imputation results across tabular datasets. We report the mean \pm std over
1223 five different training/test split combinations.

	adult	bean	default	gesture	magic
DiffPuter	0.504 \pm 0.012	0.219 \pm 0.053	0.315 \pm 0.040	0.353\pm0.007	0.539 \pm 0.049
PRDIM	0.482\pm0.022	0.199\pm0.053	0.279\pm0.039	0.371 \pm 0.052	0.488\pm0.047

1224
1225
1226
1227
1228
1229
1230
1231 C.2.4 INTERPRETABILITY OF THE PATTERN RECOGNIZER WITH CASE STUDY

1232
1233 To evaluate whether the pattern recognizer D_ϕ trained under the EM iterations has effectively
1234 learned the underlying missing mechanism, we conducted a case study on three time-series datasets:
1235 ETT, STOCK, and PEMS-Bay. Specifically, we randomly sampled instances from the ETT dataset
1236 and plotted both the true missing ratio for each entry and the corresponding output of the trained
1237 pattern recognizer, $D_\phi(\hat{X}_0)$, where \hat{X}_0 represents the imputed sample generated through approxi-
1238 mate guided generation by PRDIM. This visualization allows us to examine whether the learned D_ϕ
1239 accurately captures and mimics the MNAR missing patterns inherent in high dimensional data.

1240 For each dataset, we sampled time intervals of length 72. The ETT and STOCK datasets contain
1241 7 and 6 features, respectively, and we visualized all features for completeness. In the case of the
PEMS-Bay dataset, which has a total of 325 feature dimensions, only the first 10 features were used



1271 Figure 11: Qualitative results of PRDIM on celebA-HQ (64×64 downsized) test dataset compared
1272 to vanilla diffusion model under custom MNAR missing pattern. Red regions in observed inputs
1273 denote missing entries under the MNAR setting.

1274
1275 for visualization due to its high dimensionality. The experimental results for ETT, STOCK, and
1276 PEMS-Bay are shown in Figure 13, Figure 14, and Figure 15, 16, respectively.

1278 Overall, the results demonstrate that while the pattern recognizer tends to slightly overestimate miss-
1279 ing entries, it nonetheless captures the overall tendency and structure of the true missing pattern
1280 remarkably well, indicating its strong capability to model MNAR mechanisms. This interpretability
1281 analysis provides empirical evidence that the pattern recognizer contributes meaningful guidance
1282 during the generation phase.

1284 C.2.5 QUALITATIVE RESULTS OF PRDIM IMPUTATION

1286 To further illustrate the behavior of PRDIM under the out-of-sample imputation setting, we pro-
1287 vide qualitative visualizations across ETT, STOCK, and PEMS-Bay on Figure 17, 18, and 19
1288 respectively. For each dataset, we randomly sample 4 test instances and display (i) the locations of
1289 missing values as yellow points and (ii) the corresponding imputation results produced by CSDI,
1290 MTSCI, cDiffPuter, MTSI, and PRDIM respectively. These visualization results allow for a direct
1291 visual comparison of reconstruction quality, highlighting the degree to which each model captures
1292 temporal structure and recovers unseen missing values.

1293 To examine whether PRDIM can operate on data with naturally occurring missing values where
1294 ground-truth values for the missing entries are not available, we additionally conducted experiments
1295 on the PhysioNet (Goldberger et al., 2000) dataset. Figure 20 visualizes the out-of-sample imputa-
tion results of CSDI, cDiffPuter, and PRDIM on PhysioNet.

One notable observation is that both cDiffPutter and PRDIM involve a joint optimization procedure over the latent missing variables X_0^{mis} and the observed variables X_0^{obs} during the EM iterations. As a consequence, when the natural missing rate is extremely high (approaching nearly 80% in the PhysioNet dataset), the imputed values may become biased toward zero (i.e. initial imputed value). This highlights an inherent limitation of EM-based diffusion imputation methods under severe natural missingness.

D DETAILED DESCRIPTION OF DATA PROCESSING AND OBJECTIVE

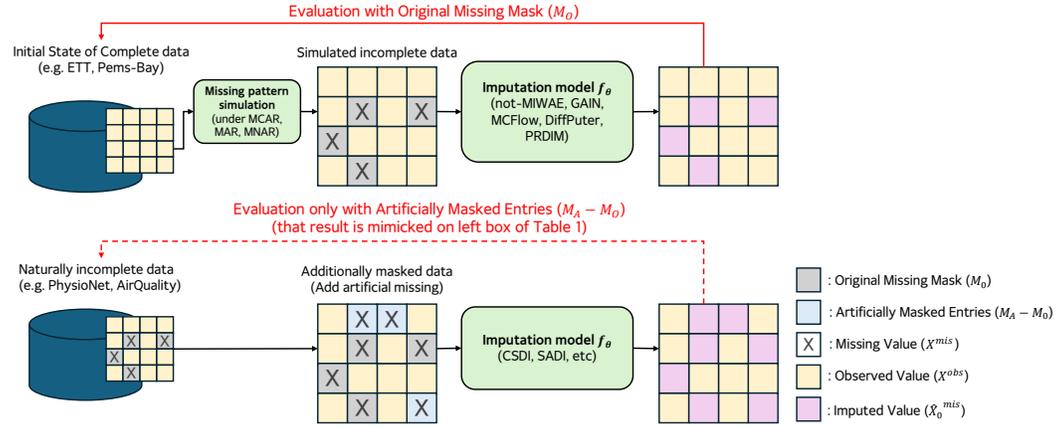


Figure 12: Overview of the data processing pipeline and the distinction between two classes of imputation objectives.

In this section, we aim to clarify the rationale behind our choice of datasets, draw theoretical connections to the EM-based training procedure adopted in PRDIM, and contrast our evaluation protocol with that of prior studies that share a similar experimental framework. Figure 12 provides an overview of two distinct classes of objectives used in existing imputation research, highlighting why directly evaluating models on naturally incomplete datasets such as PhysioNet (Goldberger et al., 2000) or AirQuality (Zhang et al., 2017) can be problematic.

Imputation applicable diffusion models including CSDI (Tashiro et al., 2021), SSSD (Alcaraz & Strodthoff, 2022), and Diffusion-TS (Yuan & Qiao, 2024), generally rely on one of two strategies. (i) introducing artificial missingness into a complete dataset so that ground-truth values are available during training, or (ii) injecting additional artificial missingness into already incomplete datasets, thereby increasing the overall missing ratio and using the resulting data as model input. A key commonality between the two imputation paradigms is that the ground-truth values employed for evaluation are implicitly utilized during model training.

Let M_O denote the original missing mask of the incomplete dataset X_0^{obs} , and let M_A denote the mask obtained after applying additional artificial missingness. The missingness distributions induced by these two masks differ intrinsically, which can be formalized as $p(M_O|X_0) \neq p(M_A|X_0, M_O)$. Consequently, the imputed results generated under these differing mask conditions also diverge $p_\theta(X_0|M_O, X_0^{\text{obs}}) \neq p_\theta(X_0|M_A, X_0^{\text{obs},A})$. Such discrepancies indicate that the imputation task inevitably involves a latent missing variable X_0^{mis} , whose distribution cannot be directly inferred from artificially masked data alone. This observation motivates the necessity of adopting an Expectation–Maximization (EM) training framework, wherein the missing entries are treated as latent variables and iteratively refined during model optimization.

E THE USE OF LARGE LANGUAGE MODELS

We’ve got some help by Large Language Models (LLMs) only in the areas of translation and grammar examination. The core research ideation and all theoretical statements are our own work.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

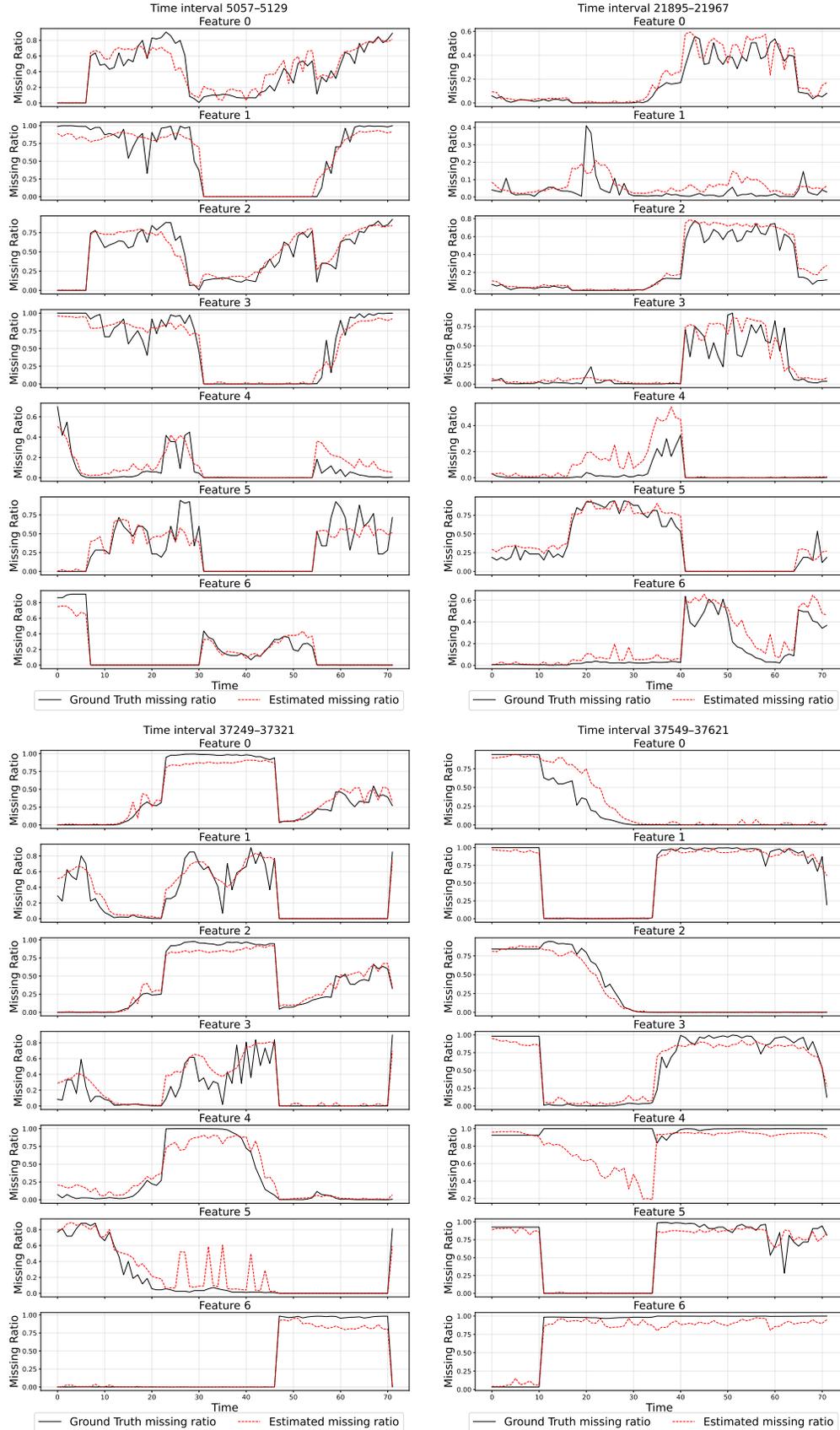


Figure 13: Randomly sampled 4 ETT segments with time length 72: Ground-truth missing ratio (black) versus Pattern Recognizer-estimated missing ratio $D_\phi(\hat{X}_0)$ across 7 features.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

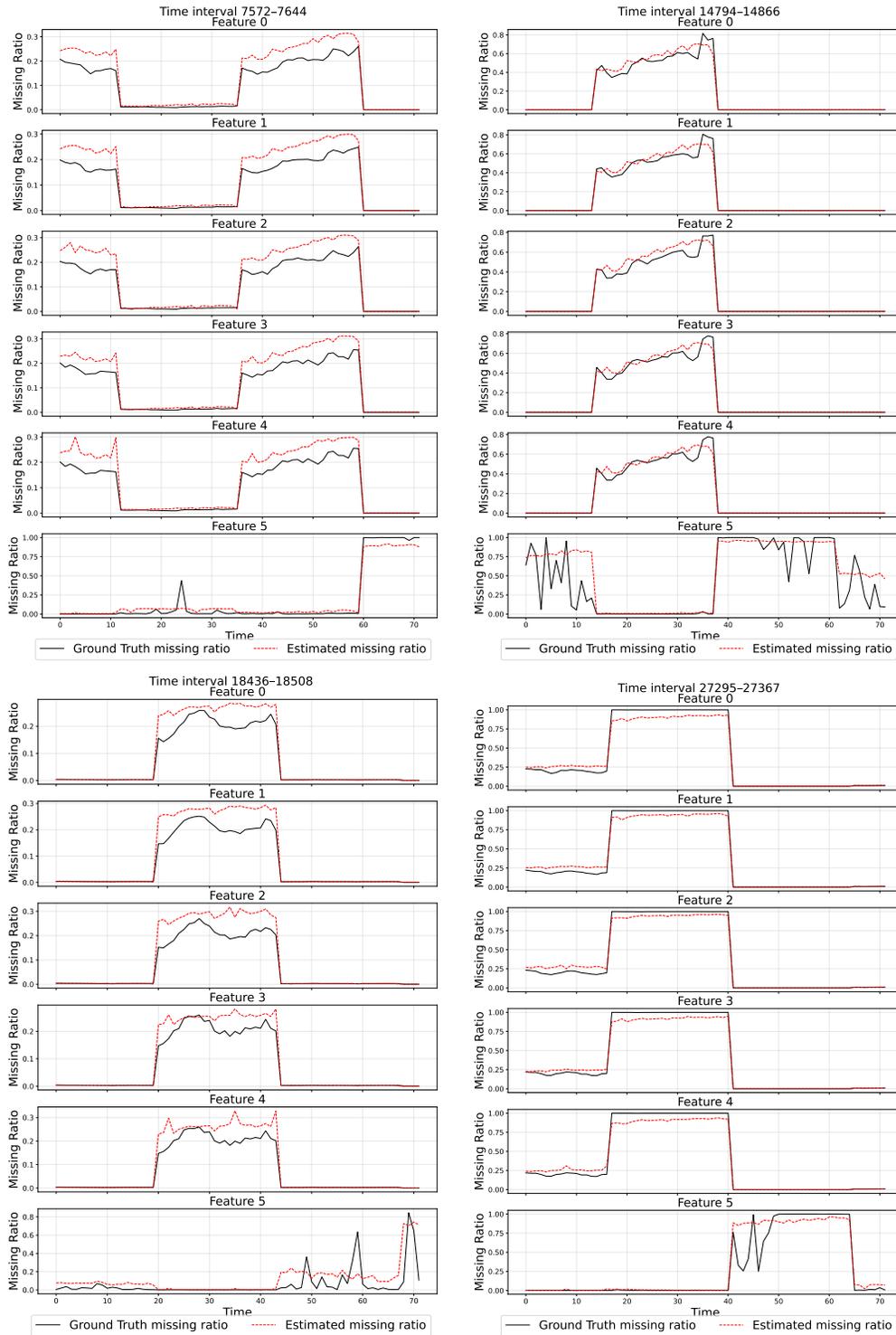


Figure 14: Randomly sampled 4 STOCK segments with time length 72: Ground-truth missing ratio (black) versus Pattern Recognizer-estimated missing ratio $D_\phi(\hat{X}_0)$ across 6 features.

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

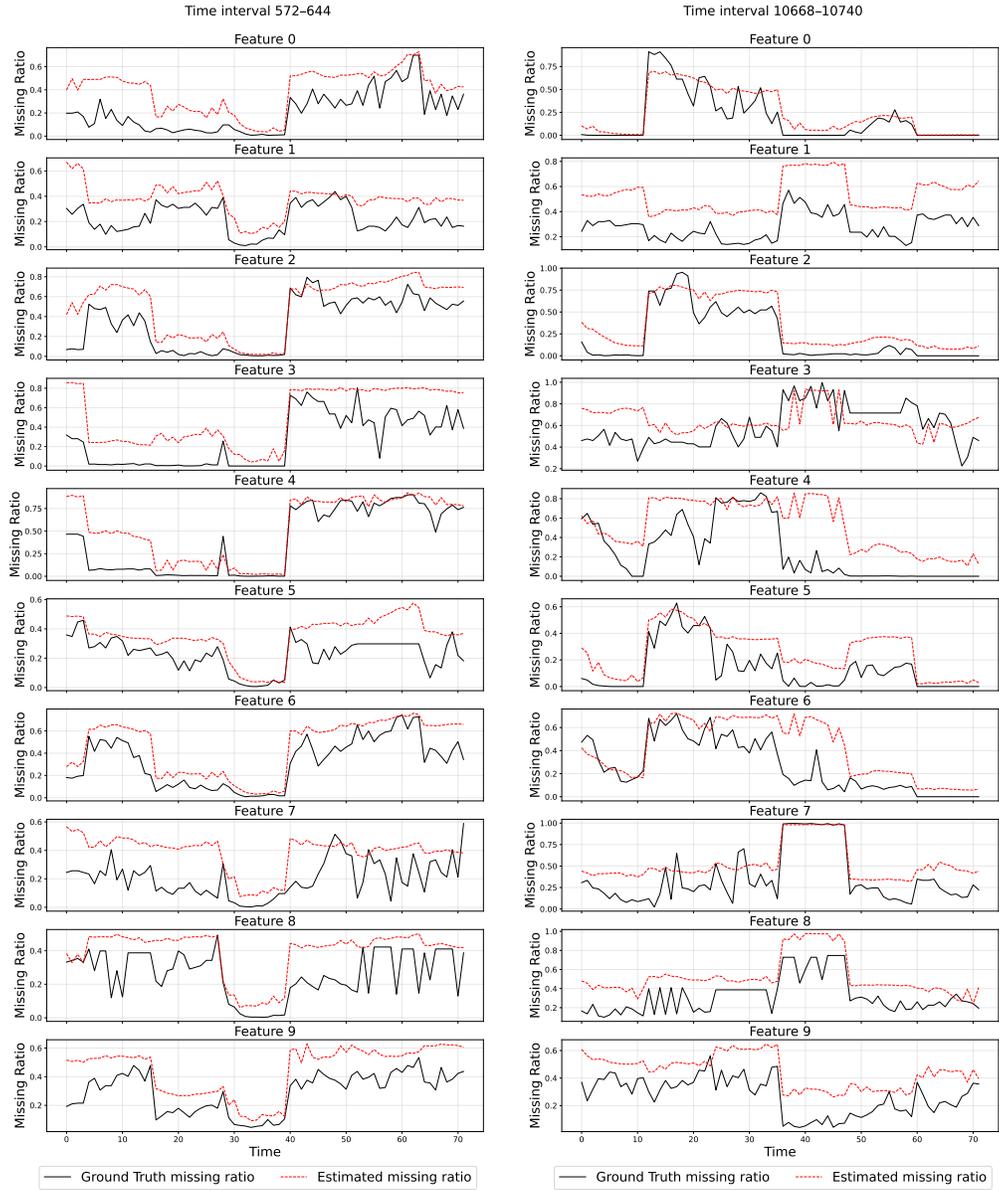


Figure 15: Randomly sampled 2 PEMS-Bay segments with time length 72: Ground-truth missing ratio (black) versus Pattern Recognizer-estimated missing ratio $D_\phi(\hat{X}_0)$ across 10 of 325 features.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

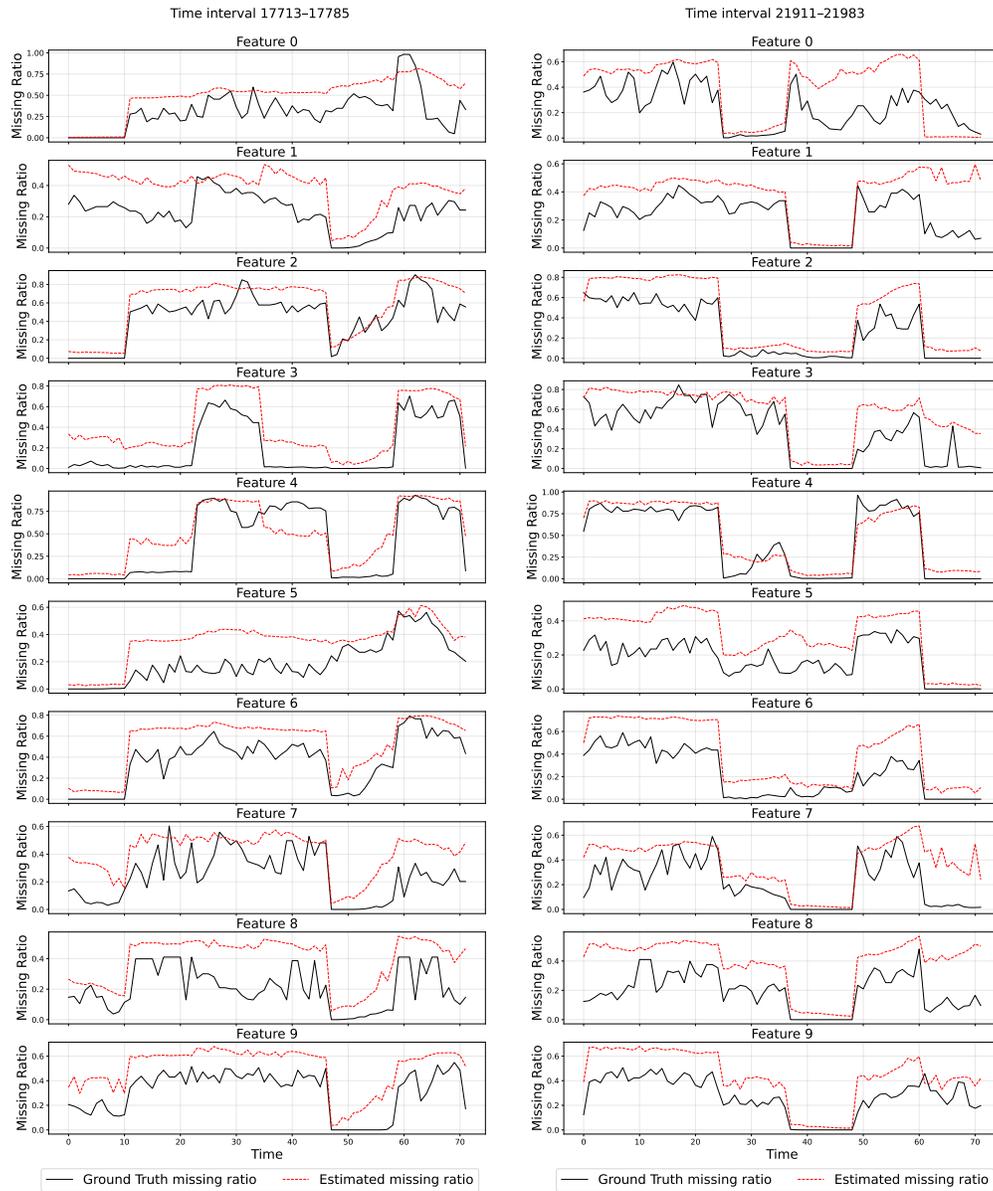


Figure 16: Additional 2 random samples of 2 PEMS-Bay segments with time length 72: Ground-truth missing ratio (black) versus Pattern Recognizer-estimated missing ratio $D_\phi(\hat{X}_0)$ across 10 of 325 features.

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

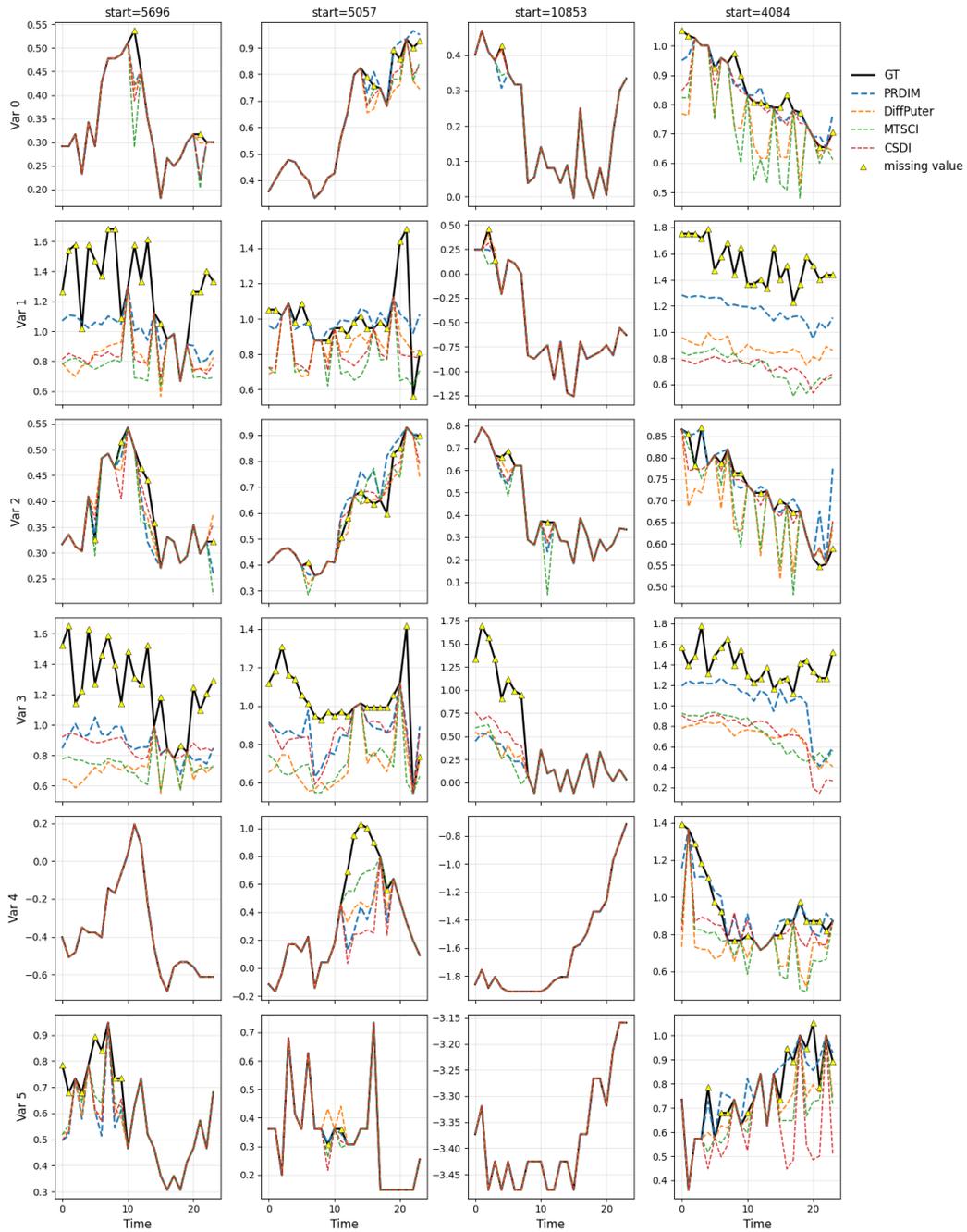


Figure 17: Qualitative results of PRDIM compared to other diffusion imputation models. 4 randomly selected imputed out-of-samples from the ETT dataset are visualized. Each panel labeled start = n corresponds to the time interval $(n, n+24)$.

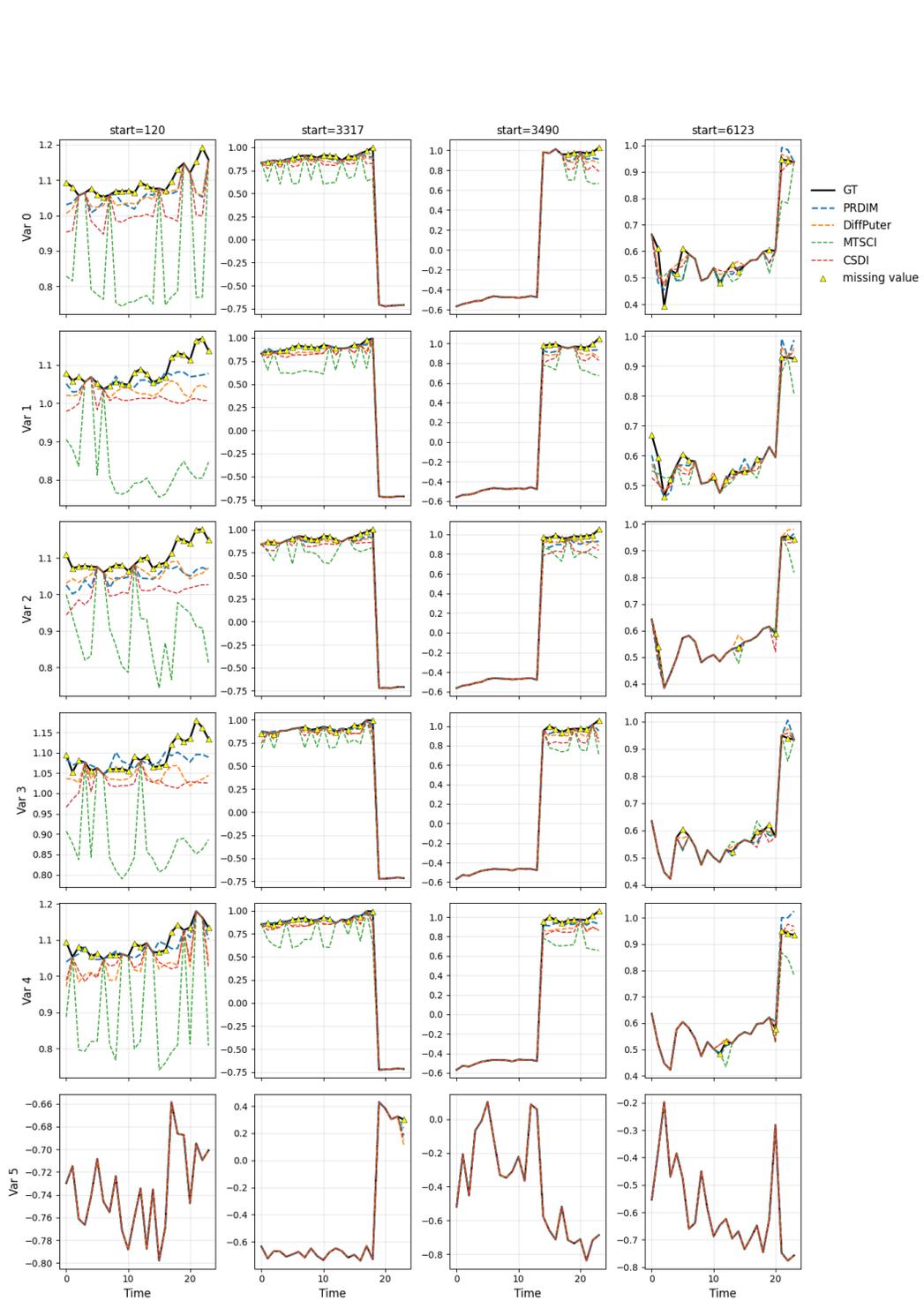


Figure 18: Qualitative results of PRDIM compared to other diffusion imputation models. 4 randomly selected imputed out-of-samples from the STOCK dataset are visualized. Each panel labeled start = n corresponds to the time interval $(n, n+24)$.

1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704
 1705
 1706
 1707
 1708
 1709
 1710
 1711
 1712
 1713
 1714
 1715
 1716
 1717
 1718
 1719
 1720
 1721
 1722
 1723
 1724
 1725
 1726
 1727

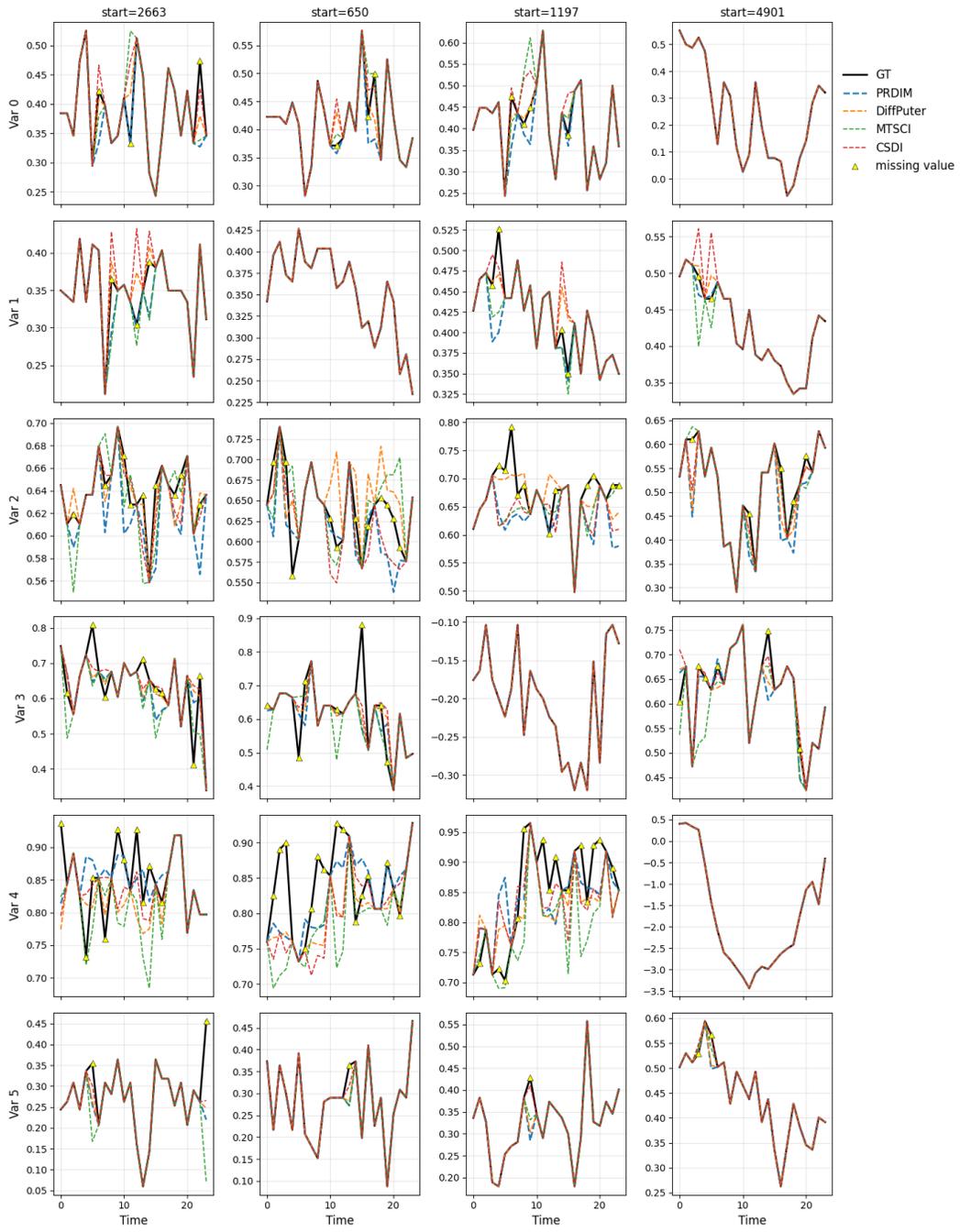


Figure 19: Qualitative results of PRDIM compared to other diffusion imputation models. 4 randomly selected imputed out-of-samples from the PEMS-Bay dataset are visualized. Each panel labeled start = n corresponds to the time interval $(n, n+24)$.

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

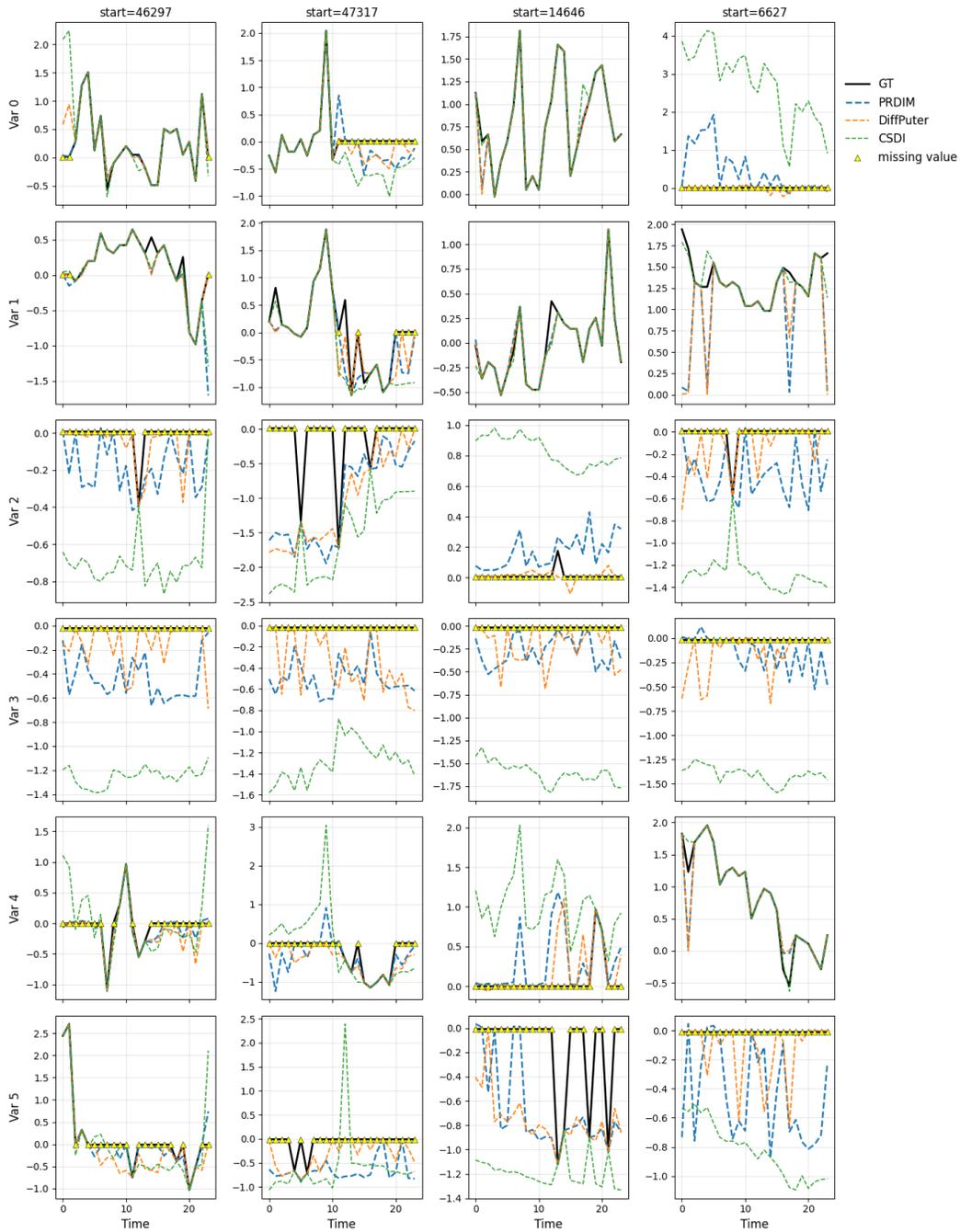


Figure 20: Qualitative results of PRDIM compared to other diffusion imputation models. 4 randomly selected imputed out-of-samples from the PhysioNet dataset are visualized. Each panel labeled start = n corresponds to the time interval $(n, n+24)$.