

# DIFFPM: DIFFUSION-BASED GENERATIVE FRAMEWORK FOR TIME SERIES SYNTHESIS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Generative models for time series often fail to reconcile local accuracy with global structure. Autoregressive models accumulate errors over long horizons, while standard diffusion approaches can degrade long-range dependencies, resulting in samples with phase drift and weakened correlations. We introduce **DiffPM**, a non-autoregressive diffusion framework that resolves this tension by factorizing the generation process. DiffPM learns to model time series by explicitly separating them into low-frequency trends and high-frequency residuals, training two specialized, window-conditioned diffusion models. At inference, the models generate short, overlapping windows for each component in parallel. The individual segments are subsequently reassembled into a full sequence via a position-aware stitching mechanism that enforces inter-window consistency. This modular, decompose-and-recombine architecture allows specialized models to excel at local generation, while guaranteeing global coherence in the final synthesis. Our extensive evaluations demonstrate that DiffPM holds a performance advantage over existing methods as well as being considerably faster at inference. This advancement is quantified by marked improvements in metrics for distributional fidelity, such as Contextual FID, and enhanced temporal coherence over long-horizon benchmarks.

## 1 INTRODUCTION

Synthetic time-series data provides a powerful solution for augmenting limited datasets, validating analytical systems, and enabling privacy-preserving data sharing in sectors from finance to health-care. A formidable barrier, however, prevents the full realization of this potential: the difficulty of generating samples that are both locally accurate and globally consistent. An effective generative model must replicate transient, short-term phenomena like signal spikes while simultaneously sustaining long-term structural properties such as seasonality, gradual trends, and multi-variate dependencies. This requirement to satisfy both local and global constraints exposes the inherent brittleness of many established generative techniques. In augmentation settings, practitioners require full-length, unconditional samples—not short imputations or forecast snippets. Most recent diffusion baselines are trained and evaluated on windowed tasks, making whole-series synthesis slow and operationally awkward.

The landscape of early deep generative techniques for sequential data was principally shaped by adversarial and variational paradigms. On one hand, GAN-based approaches, including TimeGAN (Yoon et al., 2019) and COT-GAN (Xu et al., 2020), sought to enforce temporal consistency through specialized loss functions, yet could not fully overcome the notorious challenges of mode collapse and unstable training dynamics. On the other hand, VAE-based models such as TimeVAE (Desai et al., 2021) offered a more reliable optimization process but frequently yielded reconstructions that lacked high-frequency detail due to oversmoothing. This fundamental tension between sample quality and training stability has led the field to explore diffusion models as a promising alternative. Despite their ability to generate high-fidelity samples without adversarial training, ensuring the integrity of long-range temporal dependencies in these models introduces a distinct set of research problems.

While powerful, applying diffusion models to time series reveals a core conflict: a single denoising network struggles to simultaneously learn low-frequency global structure and high-frequency local

054 details. This cross-scale interference can cause a model to sacrifice one for the other. Recent work  
 055 has attempted to mitigate this in several ways. One line of research guides a single model with aux-  
 056 iliary losses or built-in decompositions, as in DiffusionTS (Yuan & Qiao, 2024) or multi-resolution  
 057 schemes (Shen et al., 2024; Fan et al., 2024), but still risks interference within a monolithic architec-  
 058 ture. Another approach offloads structural guidance to external data via retrieval (Liu et al., 2024),  
 059 sacrificing the goal of a purely generative, self-contained model. A third path uses auto-regressive  
 060 or block-wise dependencies (Hoogetboom et al., 2022; Arriola et al., 2025), trading full parallelism  
 061 for sequential consistency. These are all partial solutions that highlight the need for a framework  
 062 that is both fully generative and explicitly structured to avoid scale interference.

063 Beyond fidelity, **wall-clock feasibility** matters. When sampling entire sequences—particularly in  
 064 **high-dimensional** datasets—window-wise or auto-regressive baselines require sequential passes  
 065 and cross-window dependencies, leading to **substantially slower** generation for the same number  
 066 of indices.

067 **DiffPM** tackles time-series generation by explicitly factorizing the process: we decompose each  
 068 sequence into a low-frequency trend and a high-frequency residual, train two independent diffusion  
 069 models specialized for these components, and at inference generate short overlapping windows for  
 070 both in parallel before recombining them with position-aware stitching and smooth cross-fades.  
 071 This decompose→generate→recombine pipeline avoids cross-scale interference that plagues single-  
 072 network designs, delivering high local fidelity from the residual model while the trend model  
 073 preserves long-range structure; within windows the process is fully non-autoregressive, and across  
 074 windows it is embarrassingly parallel, enabling fast synthesis of long, globally coherent sequences  
 075 without retrieval or extra supervision. Across diverse benchmarks, **DiffPM** improves distributional  
 076 fit (e.g., contextual FID) and temporal coherence (cross-variable correlation) as well as sampling  
 077 speed, establishing a strong new baseline for unconditional time-series generation.

078 Our contributions are as follows:

- 079 • We introduce **DiffPM**, a model that decomposes time series into trend and residual compo-  
 080 nents and learns them with two specialized diffusion models. This explicitly mitigates the  
 081 cross-scale interference common in monolithic architectures.
- 082 • We propose a window-based synthesis strategy where each local window is generated con-  
 083 ditioned on its absolute position. This position-aware mechanism enables the model to  
 084 achieve both high-fidelity local detail and long-range global coherence.
- 085 • Our framework is built for speed: all windows are generated independently and in paral-  
 086 lel and stitched once, removing sequential bottlenecks and keeping whole-series synthesis  
 087 feasible in high-dimensional settings.

## 090 2 RELATED WORK

091 The landscape of time-series generation has rapidly evolved from classical statistical models to  
 092 sophisticated deep generative frameworks. Our work builds on this progression, but its core con-  
 093 tribution—explicitly factorizing the generation process to resolve cross-scale interference—is best  
 094 understood by tracing the challenges encountered by prior paradigms. We structure our review  
 095 around this central theme.

### 098 2.1 EARLY DEEP GENERATIVE MODELS: THE EMERGENCE OF THE MULTI-SCALE 099 PROBLEM

100 Early deep generative architectures for time series synthesis were largely dominated by two  
 101 paradigms: Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs), both  
 102 offering compelling data-driven alternatives to traditional methods. GAN-based frameworks, ex-  
 103 emplified by TimeGAN (Yoon et al., 2019), leveraged a combination of adversarial and supervised  
 104 learning signals within the latent space to enforce temporal consistency. Other works in this do-  
 105 main sought to better respect the underlying sequential structure through mechanisms like causal  
 106 optimal transport (Xu et al., 2020). Concurrently, VAEs provided a more stable training framework.  
 107 For example, TimeVAE (Desai et al., 2021) advanced the VAE approach by explicitly designing

108 the model with inductive biases for trend and seasonality, achieving robust performance even in  
109 low-data scenarios.

110 The successes of these foundational frameworks were consistently shadowed by a core methodolog-  
111 ical tension. On one hand, the adversarial training of GANs, while capable of generating high-  
112 fidelity details, is notoriously brittle and prone to issues like mode collapse. On the other hand,  
113 the stable optimization of the ELBO in VAEs frequently leads to an undesirable oversmoothing of  
114 high-frequency components, prioritizing global structural integrity at the expense of local precision.  
115 This persistent conflict between local detail generation and global structural modeling indicates the  
116 inherent difficulty for a singular, uniform model to capture the multi-scale temporal dynamics of  
117 time series data. Consequently, this observation strongly motivates the development of architectures  
118 designed to explicitly address disparate temporal scales.

## 119 2.2 DIFFUSION MODELS: A NEW PARADIGM AND A FAMILIAR CHALLENGE

120 Denoising diffusion models have recently emerged as the state of the art in high-fidelity synthesis for  
121 many domains, including audio waveforms (Kong et al., 2021) and images Rombach et al. (2022);  
122 Kasaei et al. (2025). Their application to time series was a natural next step. In contrast to GANs,  
123 they offer stable training and high-quality sample generation without adversarial objectives. At their  
124 core, these models operate non-autoregressively, learning to denoise an entire sequence in parallel  
125 at each step of a reverse process.  
126

127 However, this parallel denoising process re-introduces the multi-scale challenge in a new form. A  
128 significant challenge arises when a single diffusion network is tasked with concurrently predicting  
129 noise across all frequency bands. This simultaneous prediction can lead to a phenomenon known  
130 as cross-scale interference, where the model’s learning process is skewed. For instance, the opti-  
131 mization may favor the reconstruction of pronounced low-frequency trends over the preservation of  
132 subtle high-frequency information, or the inverse. Such an imbalance often manifests as artifacts  
133 in the generated signal, including phase discrepancies or a reduction in expected volatility. Recent  
134 work in time-series diffusion can be viewed as a collection of strategies to mitigate this very issue:  
135

136 Many approaches retain a single diffusion network but attempt to guide its learning process.  
137 Diffusion-TS (Yuan & Qiao, 2024) established a strong baseline for unconditional generation  
138 by incorporating internal decomposition losses and frequency-domain penalties. Multi-resolution  
139 schemes denoise from coarse to fine scales (Shen et al., 2024; Fan et al., 2024), forcing the model  
140 to attend to different levels of granularity. While often effective, these methods still rely on a single  
141 set of network weights to manage all scales, risking residual interference.

142 Alternative strategies seek to impose external structure on the diffusion process. For instance, au-  
143 toregressive diffusion models (ARDMs) (Hoogeboom et al., 2022) and their block-wise counterparts  
144 (Arriola et al., 2025) generate sequences one segment at a time. This design enforces sequential con-  
145 sistency at the cost of full parallelism. Another approach, retrieval-augmented diffusion (Liu et al.,  
146 2024), leverages examples from a reference database to steer the generation process, anchoring the  
147 output to relevant precedents.

## 148 2.3 POSITIONING.

149 Unlike methods that guide a *single* denoiser or rely on external retrieval, DiffPM *factorizes*  
150 generation into two independent diffusion models (trend/residual) and reassembles via position-  
151 conditioned, coverage-normalized overlap-add. This avoids cross-scale interference while remain-  
152 ing fully generative and non-autoregressive at the window level, with embarrassingly parallel sam-  
153 pling. See 4 for details.  
154

## 155 3 PRELIMINARIES

156 We build on denoising diffusion probabilistic models (DDPMs) for unconditional generation (Ho  
157 et al., 2020). A fixed forward (noising) process gradually corrupts a clean sample  $x_0 \in \mathbb{R}^{T \times D}$  over  
158  $S$  steps:

$$159 q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t) \mathbf{I}), \quad t = 1, \dots, S,$$

where  $\{\beta_t\}_{t=1}^S$  is a variance schedule with  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ . Marginalizing the chain yields a convenient closed form

$$q(x_t | x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) \mathbf{I}),$$

so one can obtain a noisy  $x_t$  directly via

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}).$$

The reverse (denoising) process is parameterized by a neural network  $\epsilon_\theta(x_t, t, c)$  that predicts the noise  $\epsilon$  added at step  $t$ . The reverse transition is Gaussian,

$$p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t, c), \sigma_t^2 \mathbf{I}),$$

with mean computed from the noise predictor

$$\mu_\theta(x_t, t, c) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t, c) \right),$$

and variance  $\sigma_t^2$  chosen per schedule (e.g., fixed, learned, or  $\beta_t$ ). Training minimizes the (weighted) MSE between true and predicted noise

$$\mathcal{L}(\theta) = \mathbb{E}_{x_0, \epsilon, t} \left[ w_t \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t, c)\|_2^2 \right],$$

where, in theory,  $t$  is drawn uniformly from  $\{1, \dots, S\}$  and  $w_t$  is a per-step weight (e.g., constant or cosine/SNR-aware) (Nichol & Dhariwal, 2021b). Sampling starts from  $x_S \sim \mathcal{N}(0, \mathbf{I})$  and iteratively applies the reverse transitions  $S \rightarrow S-1 \rightarrow \dots \rightarrow 0$ . Importantly, diffusion denoises *all* dimensions in parallel at each step, making it non-autoregressive in the spatial/temporal sense.

**Scope.** We adopt the standard DDPM setup above and specialize it in our method: a windowed, position-conditioned formulation with a two-branch factorization (trend/residual) and a stitching procedure for whole-series synthesis. Implementation details appear in Sec. 4; supporting analyses of schedules/SNR and stitching are deferred to Appendix. A and C.

## 4 METHOD

### 4.1 OVERVIEW

**DiffPM** is a non-autoregressive, *unconditional* diffusion framework that follows a *decompose*  $\rightarrow$  *generate*  $\rightarrow$  *recombine* pipeline. Given a multivariate series  $\mathbf{x}_{1:T} \in \mathbb{R}^{T \times D}$ , we split it into a low-frequency *trend*  $\tau_{1:T}$  and a high-frequency *residual*  $\mathbf{r}_{1:T}$ . We then train two *window-conditioned* diffusion models—one for  $\tau$ , one for  $\mathbf{r}$ —on short, overlapping windows taken from each component. At inference, both models synthesize overlapping windows *in parallel* from Gaussian noise; windows are stitched by overlap-aware averaging to produce  $\hat{\tau}_{1:T^*}$  and  $\hat{\mathbf{r}}_{1:T^*}$ , and we form the final sample by

$$\hat{\mathbf{x}}_{1:T^*} = \hat{\tau}_{1:T^*} + \hat{\mathbf{r}}_{1:T^*}.$$

Training and sampling are configured per branch (trend/residual may use different cutoffs) to reduce gradient variance and align the learned score field with inference; see Appendix A for details on this branch-specific tuning. An architectural overview is shown in Fig. 1.

### 4.2 TREND-RESIDUAL DECOMPOSITION

Let  $\mathbf{x}_{1:T} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$  with  $\mathbf{x}_t \in \mathbb{R}^D$ . We compute a smooth trend by a centered moving average of width  $K$  (odd, for clarity). For  $t \in \{1, \dots, T\}$ ,

$$\tau_t = \frac{1}{K} \sum_{j=t-(K-1)/2}^{t+(K-1)/2} \mathbf{x}_j, \quad (1)$$

implemented as a same-length discrete convolution with a uniform kernel (out-of-range indices treated as zeros). The residual is the high-frequency remainder

$$\mathbf{r}_t = \mathbf{x}_t - \tau_t. \quad (2)$$

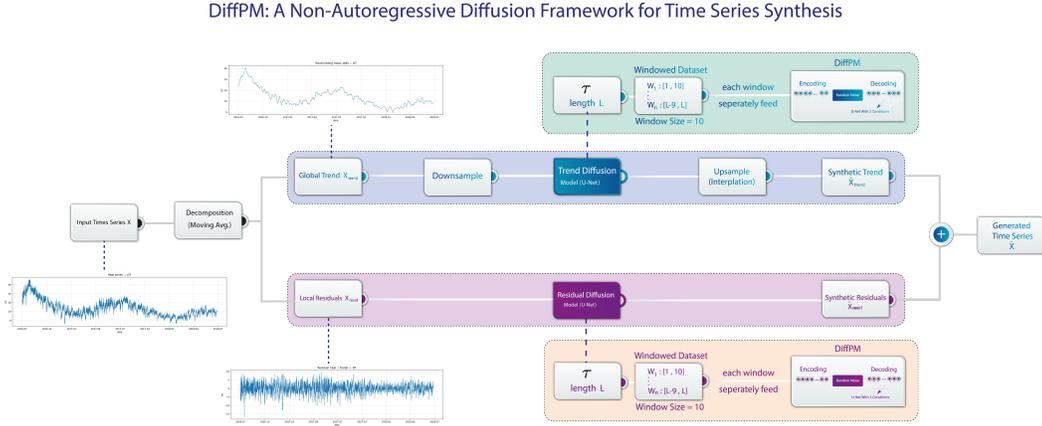


Figure 1: **DiffPM: decompose**  $\rightarrow$  **generate**  $\rightarrow$  **recombine**. *Left*: a centered moving average (width  $K$ ) produces a low-frequency trend  $\tau$  and a zeromean residual  $r = x - \tau$ . *Middle*: two window-conditioned diffusion models (1D UNets with timestep and position embeddings) are trained independently on sliding windows of  $\tau$  and  $r$ . *Right*: at inference, overlapping windows are sampled *in parallel* and stitched with overlapaware averaging, yielding  $\hat{\tau}$ ,  $\hat{r}$ , and the final sample  $\hat{x} = \hat{\tau} + \hat{r}$ .

Both sequences retain the original length and dimensionality. From this point onward, the *trend branch is trained on a downsampled trend* obtained by deterministic decimation (factor  $K$ ); for notational simplicity we continue to denote this downsampled trend by  $\tau$ .

This separation reduces cross-scale interference during learning: the trend model focuses on low-frequency, long-range structure while the residual model specializes in local variability. This deterministic analysis operator also underpins our likelihood: under this decomposition the training objective decomposes into a sum of two diffusion ELBOs (and admits an optional hierarchical residual); see Appendix B, Secs. B.3–B.4.

### 4.3 WINDOW-CONDITIONED DIFFUSION MODELS

We segment each series into fixed-length, overlapping windows. For length  $T$ , window size  $L$ , and stride  $\Delta$  ( $1 \leq \Delta < L$ ), the start set is  $\mathcal{S} = \{s \in \{1, 1 + \Delta, 1 + 2\Delta, \dots\} : s \leq T - L + 1\}$ . From each  $s \in \mathcal{S}$  we extract windows for the trend and residual:  $\tau_{s:s+L-1}, r_{s:s+L-1} \in \mathbb{R}^{L \times D}$  (recall that the trend branch uses the *downsampled* trend; for notational simplicity we still write  $\tau$ ). Each window is conditioned on its absolute position via  $c = (s, T)$ , embedded and injected into the denoiser to preserve global context. We train on all windows independently and compute the loss per window *without* explicit overlap reweighting (see Fig. 2); this yields unbiased gradients for the per-window objective, preserves full parallelism, and matches our inference-time *coverage-normalized* stitching (4.4). Alternative overlap-aware formulations—such as inclusion–exclusion reweighting and conditional chaining—are derived in Appendix. B, Sec. B.6.

We adopt the DDPM objective on windows with position conditioning. Let  $w \in \mathbb{R}^{L \times D}$  denote a window (either trend or residual). The denoiser  $\epsilon_\theta(w^{(t)}, t, c)$  predicts injected noise, and we minimize the simplified loss

$$\mathcal{L}(\theta) = \mathbb{E}_{w^{(0)}, t, \epsilon} \left[ \|\epsilon - \epsilon_\theta(w^{(t)}, t, c)\|_2^2 \right], \quad (3)$$

with  $t$  sampled uniformly each iteration (the forward/marginal forms follow 3). We train two independent models:  $\theta_{\text{tr}}$  on trend windows and  $\theta_{\text{res}}$  on residual windows. Hyperparameters may differ across branches; analyses of SNR/schedule choices appear in Appendix. A, while our default uses uniform  $t$ .

Both denoisers share a lightweight 1D UNet with residual blocks, GroupNorm, and SiLU activations. Convolutions run along time while treating the  $D$  variables as channels; the initial projection and subsequent blocks thus mix cross-variable information naturally without attention. We form a

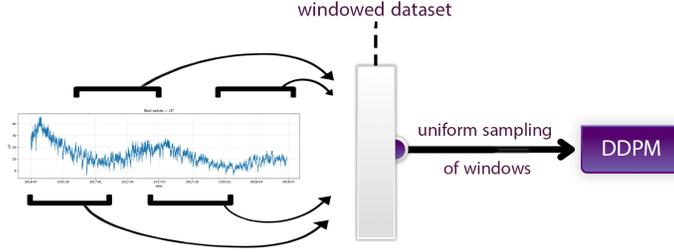


Figure 2: Overview of **DiffPM**: Break trend/residual into windows and train context-conditioned diffusion models; parallel window sampling with positional conditioning and smooth stitching.

conditioning vector  $e_{\text{cond}} = e_t + e_{\text{pos}}$  by summing an MLP-transformed sinusoidal timestep embedding with a sinusoidal embedding of  $c = (s, T)$ , and inject it in each block via *additive* FiLM (bias only): for a normalized feature map  $h$ , we apply  $h \leftarrow h + W e_{\text{cond}}$  followed by SiLU and a residual connection. This position-aware modulation is computationally cheap, preserves per-window non-autoregressiveness, and empirically stabilizes long-horizon structure by making absolute location available without altering the convolutional receptive field.

#### 4.4 PARALLEL WINDOW GENERATION AND STITCHING

Our synthesis procedure begins by defining a lattice of overlapping windows. For a target length  $T^*$ , window length  $L$ , and hop size  $\Delta$ , we use

$$\mathcal{S}^* = \{s \in \{1, 1 + \Delta, 1 + 2\Delta, \dots\} : s \leq T^* - L + 1\}.$$

A key advantage of our framework is that all windows indexed by  $\mathcal{S}^*$  are sampled *independently*, enabling an embarrassingly parallel generation process. This design is consistent with our inference-time *coverage-normalized* stitching (defined below); alternative overlap-aware formulations (inclusion–exclusion reweighting, conditional chaining, inverse-variance blending) are presented in Appendix. B, Sec. B.6.

For each  $s \in \mathcal{S}^*$ , we synthesize a local window for the trend and for the residual using their respective diffusion models. Concretely, starting from Gaussian noise and conditioning on  $c = (s, T^*)$ , the trend branch generates

$$\hat{w}_{s:s+L-1}^{(\text{tr})} \sim p_{\theta_{\text{tr}}}(\cdot | s, T^*),$$

and the residual branch generates

$$\hat{w}_{s:s+L-1}^{(\text{res})} \sim p_{\theta_{\text{res}}}(\cdot | s, T^*).$$

We use the standard reverse diffusion chain (DDPM or DDIM) for each branch. The window-level synthesis is the superposition

$$\hat{z}_{s:s+L-1} = \hat{w}_{s:s+L-1}^{(\text{tr})} + \hat{w}_{s:s+L-1}^{(\text{res})} \in \mathbb{R}^{L \times D}. \quad (4)$$

*Trend decimation/upsampling.* As the trend model is trained on a decimated trend (factor  $K$ ), we generate trend windows on the decimated grid and then upsample the stitched trend back to length  $T^*$  (linear by default; sinc optional) before summation with the residual.

Once all windows  $\{\hat{z}_{s:s+L-1}\}_{s \in \mathcal{S}^*}$  are generated, we assemble the final series  $\hat{x}_{1:T^*}$  by *overlap-add with coverage normalization*. Let  $\mathcal{W}(t) = \{s \in \mathcal{S}^* : t \in [s, s + L - 1]\}$  be the windows covering absolute time index  $t$ . In our implementation, each contributing sample uses a *uniform* analysis weight,  $w_{s,t} \equiv 1$ , and we divide by the local coverage to avoid bias:

$$S_t^{(d)} = \sum_{s \in \mathcal{W}(t)} m_{s,t}^{(d)} \hat{z}_{s,t}^{(d)}, \quad C_t^{(d)} = \sum_{s \in \mathcal{W}(t)} m_{s,t}^{(d)},$$

where  $\hat{z}_{s,t}^{(d)}$  is the  $d$ -th channel of  $\hat{z}_{s:s+L-1}$  at absolute time  $t$ . We apply a deterministic validity mask using *global* bounds  $(\ell, u)$  (per-channel bounds are possible but not used in our experiments):

$$m_{s,t}^{(d)} = \mathbf{1}\{\ell \leq \hat{z}_{s,t}^{(d)} \leq u\}.$$

The stitched output is

$$\hat{x}_t^{(d)} = \frac{S_t^{(d)}}{C_t^{(d)} + \varepsilon}, \quad t = 1, \dots, T^*, \quad d = 1, \dots, D, \quad (5)$$

with a small  $\varepsilon > 0$  for stability. Under strict COLA windows with partition-of-unity weights one obtains perfect reconstruction in the noiseless case; we instead use uniform weights with coverage normalization and rely on the *non-COLA* error bound in Appendix. C, which guarantees bounded reconstruction error, no worst-case error amplification, and variance reduction in overlaps—the exact regime we employ in all experiments.

Let  $N_{\text{win}} = |\mathcal{S}^*| \approx \lceil (T^* - L)/\Delta \rceil + 1$  be the number of windows and  $S$  the number of reverse steps. The total work scales as  $\mathcal{O}(N_{\text{win}} S)$ , but wall-clock is dominated by the cost of *one* window’s reverse chain thanks to parallel sampling; stitching is linear-time in  $T^*$ .

## 5 RESULTS AND EXPERIMENTS

We evaluate **DiffPM** as an *unconditional* multivariate time-series generator. Our experimental program measures (i) how well synthetic samples match the real data distribution, (ii) whether temporal and cross-variable dependencies are preserved, and (iii) whether synthetic data are useful for downstream modeling via train-synthetic-test-real (TSTR). We compare to strong baselines spanning GAN-, VAE-, and diffusion-based families: TimeGAN, TimeVAE, Diffwave, DiffTime, Diffusion-TS, and Cot-GAN.

### 5.1 DATASETS

We use four benchmarks covering finance, energy systems, and neuroimaging simulations:

**Stocks.** Daily Google stock prices from 2004–2019, one record per day with 6 features (open, high, low, close, adjusted close, volume). **ETTh.** Transformer-temperature/loads recorded every 15 minutes between July 2016 and July 2018; standard multivariate split used in prior work. **Energy.** UCI appliances energy dataset with 28 channels (indoor/outdoor conditions and appliance loads). **fMRI.** Realistic simulations of BOLD time series for causal discovery; we use a standard configuration with 50 variables.

All methods are trained on the training split only and evaluated on held-out test sets; we report mean $\pm$ std over multiple seeds. Windowing ( $L, \Delta$ ) and stitching follow our implementation: fixed-length windows with position conditioning and *coverage-normalized* overlap-add using uniform weights (see 4.3, 4.4). Appendix C provides the non-COLA error bound we rely on; Appendix B discusses alternative overlap-aware formulations that we do not use in these experiments.

### 5.2 EVALUATION METRICS

We adopt four standard criteria that together probe distributional fidelity, dependency structure, and downstream usefulness:

**Predictive score (TSTR)** (Yoon et al., 2019). A next-step predictor is trained *on synthetic data* and evaluated *on real* test sequences; we report the prediction error (lower is better), isolating the utility of synthetic data for model training.

**Context-FID (cFID)** (Paul et al., 2022). A Fréchet distance computed on contextual representations of subsequences; assesses local-in-context realism and diversity (lower is better).

**Correlational score** (Ni et al., 2020). Absolute Frobenius error between cross-correlation matrices computed from real vs. synthetic sequences; probes temporal and cross-variable dependency preservation (lower is better).

For overlap consistency and local realism, we rely on the stitching analysis for our *coverage-normalized* uniform-weight overlap-add: Appendix C establishes a non-COLA error bound (bounded reconstruction error, no worst-case amplification, and variance reduction in overlaps), which supports the use of contextual metrics. For completeness, the same appendix shows perfect reconstruction under strict COLA windows—conditions we do not enforce in our experiments.

Table 1: **Quantitative comparison on four datasets.** Lower is better. Best in **bold**, second best underlined.

Metric	Methods	Stock	ETTh	Energy	fMRI
Context-FID Score (lower is better)	<b>DiffPM (ours)</b>	<b>0.097 ± 0.012</b>	0.210 ± 0.010	0.142 ± 0.005	<b>0.095 ± 0.007</b>
	Diffusion-TS	0.147 ± 0.025	<b>0.116 ± 0.010</b>	<b>0.089 ± 0.024</b>	0.105 ± 0.006
	TimeGAN	0.103 ± 0.013	0.300 ± 0.013	0.767 ± 0.103	1.292 ± 0.218
	TimeVAE	0.215 ± 0.035	0.805 ± 0.186	1.631 ± 1.142	14.449 ± 0.969
	Diffwave	0.232 ± 0.032	0.873 ± 0.061	1.013 ± 1.131	0.244 ± 0.018
	DiffTime	0.236 ± 0.074	0.299 ± 0.044	0.279 ± 0.045	0.340 ± 0.015
	Cot-GAN	0.408 ± 0.086	0.980 ± 0.071	1.039 ± 0.028	7.813 ± 0.550
Correlational Score (lower is better)	<b>DiffPM (ours)</b>	0.010 ± 0.001	<b>0.032 ± 0.001</b>	<b>0.361 ± 0.049</b>	<b>0.331 ± 0.008</b>
	Diffusion-TS	<b>0.004 ± 0.001</b>	0.049 ± 0.008	0.856 ± 0.147	1.411 ± 0.042
	TimeGAN	0.063 ± 0.005	0.201 ± 0.020	4.010 ± 1.104	23.502 ± 0.039
	TimeVAE	0.095 ± 0.008	0.111 ± 0.020	1.688 ± 2.226	17.296 ± 2.526
	Diffwave	0.036 ± 0.002	0.175 ± 0.006	5.001 ± 1.154	3.927 ± 0.349
	DiffTime	<u>0.006 ± 0.002</u>	0.067 ± 0.015	1.158 ± 0.095	1.507 ± 0.051
	Cot-GAN	0.087 ± 0.004	0.249 ± 0.009	3.164 ± 0.061	26.824 ± 0.449
Predictive Score (TSTR) (lower is better)	<b>DiffPM (ours)</b>	<b>0.022 ± 0.001</b>	<b>0.107 ± 0.001</b>	0.261 ± 0.031	0.102 ± 0.000
	Diffusion-TS	<u>0.036 ± 0.001</u>	<u>0.119 ± 0.002</u>	<b>0.250 ± 0.000</b>	<b>0.099 ± 0.000</b>
	TimeGAN	0.039 ± 0.000	0.126 ± 0.004	0.273 ± 0.002	0.126 ± 0.002
	TimeVAE	0.039 ± 0.000	0.124 ± 0.004	0.302 ± 0.002	0.113 ± 0.002
	Diffwave	0.039 ± 0.000	0.123 ± 0.004	0.297 ± 0.002	0.108 ± 0.002
	DiffTime	0.038 ± 0.001	0.121 ± 0.004	0.274 ± 0.002	0.106 ± 0.001
	Cot-GAN	0.047 ± 0.001	0.129 ± 0.007	<u>0.259 ± 0.000</u>	0.185 ± 0.003

### 5.3 BASELINES AND PROTOCOL

We include representative models across families: TimeGAN (adversarial+supervised), TimeVAE (variational), Diffwave/DiffTime/Diffusion-TS (diffusion), and Cot-GAN (optimal-transport-guided GAN). For each dataset, all models use identical train/val/test splits and, where applicable, the same windowing configuration (length  $L$ , stride  $\Delta$ ). For diffusion-based methods we standardize stitching to the *coverage-normalized* uniform overlap-add used by DiffPM (4.4), so differences are not attributable to stitching.

Hyperparameters follow each method’s recommended settings; early stopping uses validation metrics. Unless stated otherwise, diffusion timesteps are sampled *uniformly* and the cosine schedule is used (no banded time sampling in main results); analyses of variance and banded sampling are provided in Appendix A.

For TSTR, the downstream predictor architecture is fixed across methods to attribute differences to data quality. Each experiment is repeated with multiple random seeds; we report mean±std.

Augmentation ablations (Appendix. D) use *identical* splits, fixed downstream models, and scalars fit on *real-train* only; causal preprocessing is enforced (no look-ahead). DiffPM synthetics are stitched with the same *coverage-normalized* overlap-add and lightly post-processed to match train statistics (rolling mean smoothing, per-feature rescale; no clipping unless stated). *Only* the presence of DiffPM synthetics in the training split differs.

Generation-time benchmarks for DiffPM and baselines are in Appendix. E: we report wall-clock per generated index and throughput scaling with length and dimensionality under identical windowing ( $L, \Delta$ ) and each method’s standard sampler/step settings.

### 5.4 MAIN QUANTITATIVE RESULTS

All DiffPM results use the *coverage-normalized* uniform overlap-add stitching described in 4.4; Appendix C establishes the non-COLA error bound that supports this choice. Our training objective is the standard simplified DDPM loss applied *independently* to trend and residual windows (Eq. 3); Appendix B, Secs. B.3–B.4, derives the two-branch windowed ELBO and discusses a hierarchical variant, which we do not employ in our experiments. The primary quantitative results against established baselines across four datasets are shown in Table 1.

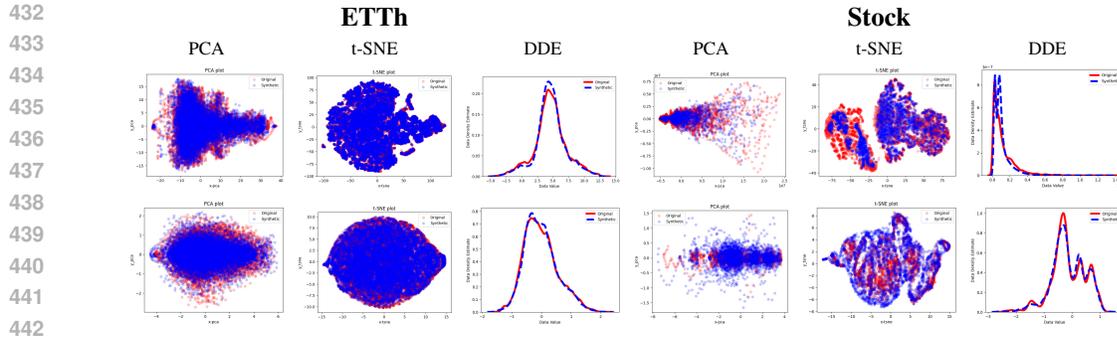


Figure 3: **Qualitative comparison on ETTh and Stock.** Top row: Diffusion-TS. Bottom row: DiffPM. Closer real–synthetic overlap in PCA/t-SNE and cleaner DDE geometry indicate higher distributional fidelity and temporal coherence.

## 5.5 QUALITATIVE VISUALIZATIONS

Following Yuan & Qiao (2024), we compare real and synthetic distributions via low-dimensional embeddings and a dynamical delay-embedding (DDE) view. For each dataset (**ETTh**, **Stock**), we present side-by-side overlays for **PCA** and **t-SNE** (real vs. model samples), and a **DDE** plot emphasizing temporal structure (*see* Fig. 3). Greater real/synthetic overlap indicates stronger distributional fidelity and coverage, while compact DDE geometry without spurious modes reflects preserved long-range dynamics. Qualitative seam inspection is consistent with our *coverage-normalized* stitching analysis.

**Reproducibility Statement.** We release an anonymized code archive with training/inference scripts, exact hyperparameters, and seeds to reproduce all tables and figures (supplementary materials). The method is fully specified in 4 (decomposition, windowing, conditioning, stitching), with objective and schedule choices in 4.3 and 4.4. Theoretical clarifications appear in Appendix. B (two-branch/windowed ELBO), Appendix. A (SNR view and time-sampling analysis), and App. C (stitching guarantees and non-COLA bound). Datasets and preprocessing (window sizes/strides, trend decimation factor, standardization on *real-train* only) are detailed in 5.1; evaluation metrics and protocols are in 5.2 and 5.3. Additional augmentation ablations and scripts are provided in Appendix. D.

## 6 CONCLUSION

We introduced **DiffPM**, a factorized, window-conditioned diffusion framework that learns low-frequency trend and high-frequency residual with two independent denoisers and stitches parallel window samples via coverage-normalized overlap–add. This decompose→generate→recombine design mitigates cross-scale interference, preserves long-range structure, and remains practical for whole-series synthesis in high-dimensional settings. Across four benchmarks, DiffPM delivers strong distributional fidelity and dependency preservation; augmentation ablations further show downstream gains (Appendix. D). Future work includes hierarchical residual conditioning, banded time sampling during training, and learned analysis windows for stitching.

## REFERENCES

- 486  
487  
488 Marianne Arriola, Yuchen Bian, Albert Shih, et al. Block diffusion: Interpolating between autore-  
489 gressive and diffusion models. *arXiv:2503.09573*, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2503.09573)  
490 2503.09573.
- 491 Abhyuday Desai, Cynthia Freeman, Zuhui Wang, and Ian Beaver. Timevae: A variational auto-  
492 encoder for multivariate time series generation. *arXiv:2111.08095*, 2021. URL [https://](https://arxiv.org/abs/2111.08095)  
493 [arxiv.org/abs/2111.08095](https://arxiv.org/abs/2111.08095).
- 494 Xinyu Fan, Yifan Zhang, et al. MG-TSD: Multi-granularity time series diffusion models. In *ICLR*,  
495 2024. URL <https://openreview.net/forum?id=CZiY6OLktd>.
- 496 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic mod-  
497 els. In *NeurIPS*, 2020. URL [https://papers.nips.cc/paper/2020/hash/](https://papers.nips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html)  
498 4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html.
- 500 Emiel Hoogeboom, Alexey A. Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg,  
501 and Tim Salimans. Autoregressive diffusion models. In *ICLR*, 2022. URL [https://](https://openreview.net/pdf?id=Lm8T39vLDTE)  
502 [openreview.net/pdf?id=Lm8T39vLDTE](https://openreview.net/pdf?id=Lm8T39vLDTE).
- 503 Seyed Amir Kasaei, Ali Aghayari, Arash Marioriyad, Niki Sepasian, Shayan Baghayi Nejad, Mo-  
504 hammadAmin Fazli, Mahdieh Soleymani Baghshah, and Mohammad Hossein Rohban. Carinox:  
505 Inference-time scaling with category-aware reward-based initial noise optimization and explo-  
506 ration, 2025. URL <https://arxiv.org/abs/2509.17458>.
- 507 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile  
508 diffusion model for audio synthesis. In *ICLR*, 2021. URL [https://openreview.net/](https://openreview.net/forum?id=a-xFK8Ymz5J)  
509 [forum?id=a-xFK8Ymz5J](https://openreview.net/forum?id=a-xFK8Ymz5J).
- 510 Jingwei Liu, Yang Li, Lihui He, and Shenda Hong. Retrieval-augmented  
511 diffusion models for time series generation. In *NeurIPS*, 2024. URL  
512 [https://papers.neurips.cc/paper\\_files/paper/2024/file/](https://papers.neurips.cc/paper_files/paper/2024/file/053ee34c0971568bfa5c773015c10502-Paper-Conference.pdf)  
513 053ee34c0971568bfa5c773015c10502-Paper-Conference.pdf.
- 514 X. Ni et al. Measuring dependence in multivariate time series generation. In *NeurIPS Workshop on*  
515 *Time Series*, 2020.
- 516 Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *ICML*,  
517 2021a.
- 518 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
519 In *ICML*, 2021b. URL <https://proceedings.mlr.press/v139/nichol21a.html>.
- 520 A. Paul et al. Contextual fid for time series generation. In *ICLR Workshop on Time Series*, 2022.
- 521 R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-Resolution Image Synthesis  
522 with Latent Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
523 *and Pattern Recognition (CVPR)*, pp. 10684–10695, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2112.10752v2)  
524 2112.10752v2.
- 525 Liang Shen, Lin Li, Ziqi Qin, et al. Multi-resolution diffusion models for time series. In *ICLR*,  
526 2024. URL [https://proceedings.iclr.cc/paper\\_files/paper/2024/file/](https://proceedings.iclr.cc/paper_files/paper/2024/file/d64740dd69bcc90ba225a182984b81ba-Paper-Conference.pdf)  
527 d64740dd69bcc90ba225a182984b81ba-Paper-Conference.pdf.
- 528 Tengfei Xu, Yifan Wang, Hao Wang, Jian Sun, Yang Liu, and Song Bai. Cot-gan:  
529 Generating sequential data via causal optimal transport. In *NeurIPS*, 2020. URL  
530 [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/](https://proceedings.neurips.cc/paper_files/paper/2020/file/641d77dd5271fca28764612a028d9c8e-Paper.pdf)  
531 641d77dd5271fca28764612a028d9c8e-Paper.pdf.
- 532 Jinsung Yoon, Daniel Jarrett, and Mihaela van der Schaar. Time-series generative adver-  
533 sarial networks. In *NeurIPS*, 2019. URL [https://papers.neurips.cc/paper/](https://papers.neurips.cc/paper/8789-time-series-generative-adversarial-networks.pdf)  
534 8789-time-series-generative-adversarial-networks.pdf.
- 535 Xinyu Yuan and Yan Qiao. Diffusion-ts: Interpretable diffusion for general time series generation.  
536 In *ICLR*, 2024. URL <https://openreview.net/forum?id=4h1apFj099>.

## A VARIANCE ANALYSIS AND NOISE-SCHEDULE TUNING FOR DIFFPM

This appendix formalizes the rationale behind two key strategies for improving the training stability of DIFFPM: sampling diffusion times from a restricted "middle band" and reshaping the noise schedule. Our analysis centers on the standard simplified DDPM loss, where we make explicit the sources of stochasticity and their relationship to the noise schedule and time-sampling distribution. We begin by establishing the setting and notation, then proceed to derive precise variance decompositions. This framework allows us to demonstrate how banded time sampling reduces training variance and to introduce an unbiased alternative via importance sampling.

Each branch of DIFFPM (trend and residual) is treated as an independent DDPM with  $S$  discrete diffusion steps. For a data window  $\mathbf{x}^{(0)} \in \mathbb{R}^{L \times D}$  and a timestep  $t \in \{0, \dots, S-1\}$ , the forward noising process is defined by the conditional distribution:

$$q(\mathbf{x}^{(t)} | \mathbf{x}^{(0)}) = \mathcal{N}(\mathbf{x}^{(t)}; \sqrt{\bar{\alpha}_t} \mathbf{x}^{(0)}, (1 - \bar{\alpha}_t) \mathbf{I}), \quad \bar{\alpha}_t = \prod_{s=0}^t \alpha_s, \quad \alpha_s = 1 - \beta_s. \quad (6)$$

Here,  $\mathbf{I}$  is the identity matrix in  $\mathbb{R}^{(LD) \times (LD)}$  (assuming the window is flattened), and  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  is a standard Gaussian noise vector of matching shape. All subsequent  $\|\cdot\|_2$  norms operate on this flattened  $LD$ -vector representation. The training objective is based on the simplified noise-prediction loss:

$$\mathbf{x}^{(t)} = \sqrt{\bar{\alpha}_t} \mathbf{x}^{(0)} + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \ell(\theta; \mathbf{x}^{(0)}, t, \epsilon) = \|\epsilon - \epsilon_\theta(\mathbf{x}^{(t)}, t)\|_2^2, \quad (7)$$

which is optimized over the full expectation:

$$L(\theta) = \mathbb{E}_{t \sim p(t)} \mathbb{E}_{\mathbf{x}^{(0)} \sim \mathcal{D}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\ell(\theta; \mathbf{x}^{(0)}, t, \epsilon)]. \quad (8)$$

Unless otherwise specified,  $p(t)$  is the discrete uniform distribution over  $\{0, \dots, S-1\}$ , i.e.,  $p(t) = 1/S$ . We define the per-timestep expected loss as  $L_t(\theta) := \mathbb{E}_{\mathbf{x}^{(0)}, \epsilon} [\ell(\theta; \mathbf{x}^{(0)}, t, \epsilon)]$  and the per-timestep mean gradient as  $\boldsymbol{\mu}_t(\theta) := \mathbb{E}_{\mathbf{x}^{(0)}, \epsilon} [\nabla_\theta \ell(\theta; \mathbf{x}^{(0)}, t, \epsilon)]$ . All branches of our model employ the cosine noise schedule introduced by Nichol & Dhariwal (2021a):

$$\bar{\alpha}_t = \frac{\cos^2\left(\frac{t/S+s}{1+s} \cdot \frac{\pi}{2}\right)}{\cos^2\left(\frac{s}{1+s} \cdot \frac{\pi}{2}\right)}, \quad \beta_t = \min(1 - \bar{\alpha}_t / \bar{\alpha}_{t-1}, 0.999), \quad s > 0. \quad (9)$$

We set  $\bar{\alpha}_{-1} = 1$  for consistency and compute  $\beta_t$  for  $t \geq 1$  with clipping at 0.999. We use the standard offset  $s = 0.008$ .

### A.1 SOURCES OF VARIANCE IN STOCHASTIC GRADIENT ESTIMATION

The training process relies on a Monte Carlo estimator for the gradient of the objective  $L(\theta)$ . For a minibatch of size  $B$ , this estimator is given by:

$$\widehat{\mathbf{g}}_B(\theta) := \frac{1}{B} \sum_{i=1}^B \mathbf{g}(\theta; \mathbf{x}_i^{(0)}, t_i, \epsilon_i), \quad t_i \stackrel{\text{i.i.d.}}{\sim} p(t), \quad \mathbf{x}_i^{(0)} \sim \mathcal{D}, \quad \epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (10)$$

where  $\mathbf{g}(\theta; \mathbf{x}^{(0)}, t, \epsilon) := \nabla_\theta \ell(\theta; \mathbf{x}^{(0)}, t, \epsilon)$  is the per-sample, per-timestep gradient. Assuming standard regularity conditions that permit interchanging differentiation and expectation, this estimator is unbiased, as shown by the following chain of expectations:

$$\mathbb{E}[\widehat{\mathbf{g}}_B(\theta)] = \mathbb{E}_{t, \mathbf{x}^{(0)}, \epsilon} [\mathbf{g}(\theta; \mathbf{x}^{(0)}, t, \epsilon)] = \mathbb{E}_t [\boldsymbol{\mu}_t(\theta)] = \nabla_\theta L(\theta).$$

The covariance of this minibatch estimator reveals two distinct sources of variance. By the law of total covariance and assuming independence across minibatch samples, we have:

$$\begin{aligned} \text{Cov}(\widehat{\mathbf{g}}_B(\theta)) &= \frac{1}{B} \text{Cov}_{t, \mathbf{x}^{(0)}, \epsilon}(\mathbf{g}(\theta; \mathbf{x}^{(0)}, t, \epsilon)) \\ &= \frac{1}{B} \left( \underbrace{\mathbb{E}_t [\text{Cov}_{\mathbf{x}^{(0)}, \epsilon}(\mathbf{g} | t)]}_{\text{Within-timestep variance}} + \underbrace{\text{Cov}_t(\boldsymbol{\mu}_t(\theta))}_{\text{Time-sampling variance}} \right). \end{aligned} \quad (11)$$

The first term reflects variance from data sampling and noise injection at a fixed timestep. The second term, however, arises from the variability of the mean gradient  $\boldsymbol{\mu}_t$  across different timesteps.<sup>1</sup> Using the identity  $\mathbb{E}\|X - \mathbb{E}X\|_2^2 = \text{Tr}(\text{Cov}(X))$ , we can express this matrix decomposition as a scalar variance decomposition:

$$\mathbb{E}[\|\widehat{\mathbf{g}}_B - \mathbb{E}\widehat{\mathbf{g}}_B\|_2^2] = \frac{1}{B} \left( \mathbb{E}_t \mathbb{E}_{\mathbf{x}^{(0)}, \epsilon} [\|\mathbf{g} - \boldsymbol{\mu}_t\|_2^2 | t] + \mathbb{E}[\|\boldsymbol{\mu}_t - \mathbb{E}_t \boldsymbol{\mu}_t\|_2^2] \right). \quad (12)$$

This decomposition highlights a critical issue: the total variance is lower-bounded by the time-sampling term, which does not vanish as  $B \rightarrow \infty$ . If  $\boldsymbol{\mu}_t$  varies substantially with  $t$ , this second term can dominate and destabilize training.

To formalize the connection between the loss spread and gradient spread, we can bound the norm of the mean gradient. Let the denoising model be  $f_\theta(\mathbf{x}^{(t)}, t) := \epsilon_\theta(\mathbf{x}^{(t)}, t)$  and its parameter-Jacobian be  $J_\theta$ . Assuming its operator norm is uniformly bounded,  $\|J_\theta(\mathbf{x}^{(t)}, t)\|_{\text{op}} \leq M < \infty$ , we have:

$$\|\boldsymbol{\mu}_t(\theta)\| = \|2 \mathbb{E}[J_\theta^\top (f_\theta - \epsilon)]\| \leq 2M \sqrt{L_t(\theta)}. \quad (13)$$

This inequality rigorously links the mean gradient magnitude to the per-timestep loss. To relate this to the variation across timesteps, we can introduce a fixed baseline  $\bar{L}$  (e.g.,  $\bar{L} = \mathbb{E}_t[L_t]$ ) and apply the triangle inequality for square roots:

$$\|\boldsymbol{\mu}_t(\theta)\| \leq 2M \left( \sqrt{|L_t(\theta) - \bar{L}|} + \sqrt{\bar{L}} \right). \quad (14)$$

Together, these inequalities formalize the intuition that a wide spread in  $L_t(\theta)$  across  $t$  implies significant variation in  $\boldsymbol{\mu}_t(\theta)$ , thus inflating the time-sampling variance term.

## A.2 BANDED TIME SAMPLING AND THE BIAS-VARIANCE TRADEOFF

Based on this analysis, we define banded time sampling. For  $S$  diffusion steps, we fix a lower and upper bound  $0 < \beta < \omega < 1$  and define the middle band  $\mathcal{T}_{\beta, \omega} = \{\lfloor \beta S \rfloor, \dots, \lfloor \omega S \rfloor\}$ . We replace the uniform sampling distribution  $p(t)$  with  $p_{\beta, \omega}(t)$ , the uniform law on this set. To quantify the effect, we define the trace of the covariance of the mean gradient over this band:

$$\mathcal{G}_{\beta, \omega}(\theta) := \text{Tr Cov}_{t \sim p_{\beta, \omega}}(\boldsymbol{\mu}_t(\theta)). \quad (15)$$

Over the feasible region of  $(\beta, \omega)$ , this function is piecewise-constant and minimizers exist. Combining the bound in equation 13 with Jensen’s inequality provides a direct link between the average loss over the band and an upper bound on this gradient covariance:

$$\mathcal{G}_{\beta, \omega}(\theta) = \text{Tr Cov}_{t \sim p_{\beta, \omega}}(\boldsymbol{\mu}_t(\theta)) \leq \mathbb{E}_{t \sim p_{\beta, \omega}}[\|\boldsymbol{\mu}_t(\theta)\|^2] \leq 4M^2 \mathbb{E}_{t \sim p_{\beta, \omega}}[L_t(\theta)]. \quad (16)$$

This result provides a powerful insight: to minimize an upper bound on the time-sampling component of the gradient variance, one should choose a band  $(\beta, \omega)$  that minimizes the average loss over that band. Training directly on this band introduces a beneficial bias, focusing model capacity on the most informative timesteps, while simultaneously reducing stochastic gradient variance. To align inference, we start the reverse diffusion process at  $t_0 = \lfloor \omega S \rfloor$ .

## A.3 UNBIASED ESTIMATION WITH IMPORTANCE SAMPLING

If optimizing the original objective  $L(\theta)$  is strictly required, one can use importance sampling (IS). Let  $\mathbf{g}_i$  be a shorthand for the per-sample gradient. The IS estimator is:

$$\widehat{\mathbf{g}}_B^{\text{IS}}(\theta) := \frac{1}{B} \sum_{i=1}^B \underbrace{\frac{p(t_i)}{p_{\beta, \omega}(t_i)}}_{w(t_i)} \mathbf{g}_i, \quad t_i \sim p_{\beta, \omega}. \quad (17)$$

This estimator is unbiased, but at the cost of increased variance. Its covariance is given by:

$$\text{Cov}(\widehat{\mathbf{g}}_B^{\text{IS}}) = \frac{1}{B} \left( \mathbb{E}_{t \sim p_{\beta, \omega}} [w(t)^2 \mathbb{E}_{\mathbf{x}^{(0)}, \epsilon} [\mathbf{g}\mathbf{g}^\top | t]] - (\nabla_\theta L)(\nabla_\theta L)^\top \right). \quad (18)$$

<sup>1</sup>If minibatches are constructed from overlapping windows from the same time series, the samples are not strictly i.i.d. While the estimator remains unbiased, the  $1/B$  scaling becomes a heuristic upper bound due to non-zero cross-sample covariances.

The variance of the IS estimator scales with the second moment of the weighted gradient,  $\mathbb{E}_{t \sim p_{\beta, \omega}} [w(t)^2 \|\mathbf{g}\|^2]$ . In our case, the weight  $w(t) = \frac{S}{|\mathcal{T}_{\beta, \omega}|} > 1$  is constant in the band. The increased variance of importance sampling thus necessitates practical safeguards. To control it, one should ensure the band is of non-trivial width to keep weights moderate, and it is advisable to pair IS with a larger batch size  $B$ . For these reasons, we favor the biased banded objective for its superior training stability.

#### A.4 PRACTICAL IMPLICATIONS AND BRANCH-SPECIFIC TUNING

The problematic nature of extreme timesteps can be understood from a signal-to-noise ratio perspective, where  $\text{SNR}(t) = \bar{\alpha}_t / (1 - \bar{\alpha}_t)$ . For small  $t$ , the SNR is high, the denoising task is trivial, and gradients are small; batches with many such steps contribute to within-timestep variance without moving parameters meaningfully. For large  $t$ , the SNR is low, the input is nearly pure noise, and the loss becomes large but highly stochastic, causing large, erratic mean gradients that inflate time-sampling variance.

The cosine schedule is designed to mitigate these issues by smoothing the decay of  $\bar{\alpha}_t$ , but it does not eliminate the problematic tails. Banded sampling is an orthogonal variance reduction technique that complements the cosine schedule by explicitly removing these tails from the sampling distribution.

For the two-branch architecture of DIFFPM, this framework suggests branch-specific tuning. The trend and residual components often operate in different effective SNR regimes; residual windows are typically higher-frequency and lower-amplitude, meaning their signal degrades more quickly as noise is added. This observation leads to a practical heuristic: use a common lower cutoff  $\beta$  for both branches, but an earlier upper cutoff for the residual branch ( $\omega_R \leq \omega_T$ ). This branch-specific tuning, guided by a simple calibration, directly reduces the time-sampling variance from equation 11 and can improve convergence without altering the model architecture or loss function. In summary, by restricting time sampling to a "middle band," we strategically introduce a beneficial bias that reduces a dominant source of gradient variance, stabilizing and accelerating training.

## B MODIFIED ELBO FOR WINDOWED TWO-BRANCH DIFFPM

This appendix provides a rigorous formulation of the evidence lower bound (ELBO) used by DiffPM when modeling a time series as the sum of two components—a smooth trend and a higher-frequency residual—and when training on fixed-length windows (with or without overlap). We derive the ELBO under (i) independent-component and (ii) hierarchical (conditional) generative assumptions. We then extend the analysis to windowed training, including overlapping windows via two alternatives: inclusion–exclusion reweighting and conditional (blockwise) chaining.

### B.1 NOTATION AND SETUP

Let an observed sequence window be  $\mathbf{x}_0 \in \mathbb{R}^{L \times D}$  (length  $L$ ,  $D$  channels). We posit a deterministic decomposition

$$\mathbf{x}_0 = \boldsymbol{\tau}_0 + \mathbf{r}_0, \quad \boldsymbol{\tau}_0 = \mathcal{T}(\mathbf{x}_0), \quad \mathbf{r}_0 = \mathbf{x}_0 - \mathcal{T}(\mathbf{x}_0), \quad (19)$$

where  $\mathcal{T} : \mathbb{R}^{L \times D} \rightarrow \mathbb{R}^{L \times D}$  is a fixed trend extractor (e.g., a centered moving average). DiffPM trains two diffusion models, one per component  $z \in \{\boldsymbol{\tau}, \mathbf{r}\}$ .

**Forward diffusion (per component).** For each  $z \in \{\boldsymbol{\tau}, \mathbf{r}\}$  we use a standard Gaussian forward chain

$$q(z_{1:S} | z_0) = \prod_{t=1}^S q(z_t | z_{t-1}), \quad q(z_t | z_{t-1}) = \mathcal{N}(\sqrt{\alpha_t} z_{t-1}, (1 - \alpha_t) \mathbf{I}), \quad (20)$$

with  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$  and

$$z_t = \sqrt{\bar{\alpha}_t} z_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

where  $\mathbf{I}$  denotes the identity of appropriate dimension (i.e., matching the vectorized  $z_t$ ).

702 **Reverse model (per component).** The reverse (generative) model is

$$703 \quad p_\theta(z_{0:S}) = p(z_S) \prod_{t=1}^S p_\theta(z_{t-1} | z_t), \quad p(z_S) = \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (21)$$

704 with  $p_\theta(z_{t-1} | z_t)$  parameterized via a noise-prediction network  $\epsilon_\theta(z_t, t, \text{cond})$ . The usual DDPM  
705 objective reduces to

$$706 \quad \mathcal{L}_z^{\text{simple}}(\theta) = \mathbb{E}_{z_0, t, \epsilon} \left[ w_t \|\epsilon - \epsilon_\theta(z_t, t, \text{cond})\|_2^2 \right] \quad (\text{up to } t\text{-dependent constants}). \quad (22)$$

707 In DiffPM, “cond” may include positional metadata (e.g., window start and length); we leave it  
708 implicit here to focus on likelihood.

## 714 B.2 JOINT GENERATIVE DENSITY AND VARIATIONAL FAMILY

715 The *observable* is always recovered at  $t = 0$  by  $\mathbf{x}_0 = \boldsymbol{\tau}_0 + \mathbf{r}_0$ . We use  $c$  to denote window con-  
716 text/metadata: start index  $s$ , total series length  $T$ , and any positional metadata consumed by the de-  
717 noiser. All likelihoods/ELBOs in this appendix are conditional on  $c$  unless stated otherwise. DiffPM  
718 uses *context-conditional independence* of trend and residual given  $c$ :

- 719 1. **Independent components (used in DiffPM):**  $p_\theta(\boldsymbol{\tau}_0, \mathbf{r}_0 | c) = p_{\theta_1}(\boldsymbol{\tau}_0 | c) p_{\theta_2}(\mathbf{r}_0 | c)$ .
- 720 2. **Hierarchical :**  $p_\theta(\boldsymbol{\tau}_0, \mathbf{r}_0 | c) = p_{\theta_1}(\boldsymbol{\tau}_0 | c) p_{\theta_2}(\mathbf{r}_0 | \boldsymbol{\tau}_0, c)$ .

721 Let  $\delta(\cdot)$  denote a Dirac delta. The joint *generative* density over reverse paths and the observation is

$$722 \quad p_\theta(\boldsymbol{\tau}_{0:S}, \mathbf{r}_{0:S}, \mathbf{x}_0 | c) = p_{\theta_1}(\boldsymbol{\tau}_{0:S} | c) p_{\theta_2}(\mathbf{r}_{0:S} | \boldsymbol{\tau}_0, c) \delta(\mathbf{x}_0 - \boldsymbol{\tau}_0 - \mathbf{r}_0), \quad (23)$$

723 where in the independent case  $p_{\theta_2}(\mathbf{r}_{0:S} | \boldsymbol{\tau}_0, c)$  reduces to  $p_{\theta_2}(\mathbf{r}_{0:S} | c)$ .

724 For training, we use the *forward* variational family that fixes the decomposition at  $t = 0$  and then  
725 applies forward diffusion:

$$726 \quad q(\boldsymbol{\tau}_{0:S}, \mathbf{r}_{0:S} | \mathbf{x}_0, c) = \underbrace{\delta(\boldsymbol{\tau}_0 - \mathcal{T}(\mathbf{x}_0)) \delta(\mathbf{r}_0 - (\mathbf{x}_0 - \mathcal{T}(\mathbf{x}_0)))}_{\text{deterministic decomposition at } t=0} \left[ \prod_{t=1}^S q(\boldsymbol{\tau}_t | \boldsymbol{\tau}_{t-1}) \right] \left[ \prod_{t=1}^S q(\mathbf{r}_t | \mathbf{r}_{t-1}, \boldsymbol{\tau}_0) \right], \quad (24)$$

727 with  $q(\mathbf{r}_t | \mathbf{r}_{t-1}, \boldsymbol{\tau}_0) = q(\mathbf{r}_t | \mathbf{r}_{t-1})$  in the independent case.

728 The ELBO on  $\ln p_\theta(\mathbf{x}_0 | c)$  is

$$729 \quad \mathcal{L}(\mathbf{x}_0 | c) = \mathbb{E}_q \left[ \ln \frac{p_\theta(\boldsymbol{\tau}_{0:S}, \mathbf{r}_{0:S}, \mathbf{x}_0 | c)}{q(\boldsymbol{\tau}_{0:S}, \mathbf{r}_{0:S} | \mathbf{x}_0, c)} \right] \leq \ln p_\theta(\mathbf{x}_0 | c). \quad (25)$$

730 Because  $q$  puts  $(\boldsymbol{\tau}_0, \mathbf{r}_0)$  exactly on the support of  $\delta(\mathbf{x}_0 - \boldsymbol{\tau}_0 - \mathbf{r}_0)$ , the delta cancels in the ratio  
731 and, under the independent model, the expectation factorizes into a sum of the two component-wise  
732 diffusion ELBOs (each conditioned on  $c$ ).

## 743 B.3 INDEPENDENT COMPONENTS: ELBO SPLITS INTO TWO DIFFUSION ELBOs

744 Under the context-conditional independence used in DiffPM,  $p_{\theta_2}(\mathbf{r}_{0:S} | \boldsymbol{\tau}_0, c) \equiv p_{\theta_2}(\mathbf{r}_{0:S} | c)$ , the  
745 master ELBO equation 25 reduces to a sum of two standard diffusion ELBOs, one per branch:

$$746 \quad \mathcal{L}_{\text{indep}}(\mathbf{x}_0 | c) = \mathbb{E}_{q(\boldsymbol{\tau}_{1:S} | \boldsymbol{\tau}_0)} \left[ \ln \frac{p(\boldsymbol{\tau}_S)}{q(\boldsymbol{\tau}_S | \boldsymbol{\tau}_0)} + \sum_{t=1}^S \ln \frac{p_{\theta_1}(\boldsymbol{\tau}_{t-1} | \boldsymbol{\tau}_t, c)}{q(\boldsymbol{\tau}_t | \boldsymbol{\tau}_{t-1})} \right] \quad (26)$$

$$747 \quad + \mathbb{E}_{q(\mathbf{r}_{1:S} | \mathbf{r}_0)} \left[ \ln \frac{p(\mathbf{r}_S)}{q(\mathbf{r}_S | \mathbf{r}_0)} + \sum_{t=1}^S \ln \frac{p_{\theta_2}(\mathbf{r}_{t-1} | \mathbf{r}_t, c)}{q(\mathbf{r}_t | \mathbf{r}_{t-1})} \right]. \quad (27)$$

748 Each bracket in equation 26–equation 27 is the usual diffusion ELBO (conditioned on  $c$ ) for its  
749 component. Denoting their negative ELBOs by  $L_\tau(\boldsymbol{\tau}_0 | c)$  and  $L_r(\mathbf{r}_0 | c)$ , we obtain the conditional  
750 bound

$$751 \quad -\ln p_\theta(\mathbf{x}_0 | c) \leq L_\tau(\boldsymbol{\tau}_0 | c) + L_r(\mathbf{r}_0 | c). \quad (28)$$

**Simple  $\epsilon$ -prediction form.** Using the standard  $\epsilon$ -parameterization of the reverse kernels, each branch admits the DDPM “simple” loss (up to  $t$ -dependent constants):

$$L_\tau(\boldsymbol{\tau}_0 | c) = \mathbb{E}_{\boldsymbol{\tau}_0, t, \epsilon} \left[ w_t \|\epsilon - \epsilon_{\theta_1}(\boldsymbol{\tau}_t, t, c)\|_2^2 \right], \quad L_r(\mathbf{r}_0 | c) = \mathbb{E}_{\mathbf{r}_0, t, \epsilon} \left[ w_t \|\epsilon - \epsilon_{\theta_2}(\mathbf{r}_t, t, c)\|_2^2 \right], \quad (29)$$

where  $\boldsymbol{\tau}_t = \sqrt{\bar{\alpha}_t} \boldsymbol{\tau}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$  and  $\mathbf{r}_t = \sqrt{\bar{\alpha}_t} \mathbf{r}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ , with  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $w_t$  any chosen weighting (e.g., uniform or SNR-based). *Unless otherwise stated, noise draws are i.i.d. across time, channels, windows, and branches; we reuse the symbol  $\epsilon$  for notational brevity.*

**Training objective.** Summing the two component losses yields the DiffPM training criterion

$$\mathcal{J}_{\text{indep}}(\theta_1, \theta_2) = \mathbb{E}_{c \sim \hat{p}(C)} \left[ \mathbb{E} [w_t \|\epsilon - \epsilon_{\theta_1}(\boldsymbol{\tau}_t, t, c)\|_2^2] + \mathbb{E} [w_t \|\epsilon - \epsilon_{\theta_2}(\mathbf{r}_t, t, c)\|_2^2] \right] \quad (\text{constants omitted}), \quad (30)$$

i.e., the expectation over window contexts  $c$  (under the empirical distribution of contexts in the dataset) of the sum of per-branch  $\epsilon$ -MSE objectives. Equation equation 28 then guarantees that minimizing equation 30 minimizes an upper bound on the conditional negative log-likelihood  $-\ln p_\theta(\mathbf{x}_0 | c)$ .

#### B.4 HIERARCHICAL (CONDITIONAL) COMPONENTS: ELBO SPLITS INTO TREND + CONDITIONAL RESIDUAL

The hierarchical variant assumes the residual distribution depends on the realized trend (in addition to window context  $c$ ):

$$p_\theta(\boldsymbol{\tau}_0, \mathbf{r}_0 | c) = p_{\theta_1}(\boldsymbol{\tau}_0 | c) p_{\theta_2}(\mathbf{r}_0 | \boldsymbol{\tau}_0, c).$$

(We include this as an *optional extension*; DiffPM’s main experiments use the independent model in Section B.3.)

Under this factorization, the master ELBO equation 25 decomposes into a sum of two diffusion ELBOs: one for the marginal trend and one for the *conditional* residual. Writing expectations under the forward processes (with the deterministic decomposition at  $t = 0$ ), we obtain

$$\mathcal{L}_{\text{hier}}(\mathbf{x}_0 | c) = \underbrace{\mathbb{E}_{q(\boldsymbol{\tau}_{1:S} | \boldsymbol{\tau}_0)} \left[ \ln \frac{p(\boldsymbol{\tau}_S)}{q(\boldsymbol{\tau}_S | \boldsymbol{\tau}_0)} + \sum_{t=1}^S \ln \frac{p_{\theta_1}(\boldsymbol{\tau}_{t-1} | \boldsymbol{\tau}_t, c)}{q(\boldsymbol{\tau}_{t-1} | \boldsymbol{\tau}_t, c)} \right]}_{\text{ELBO for } p_{\theta_1}(\boldsymbol{\tau}_0 | c)} \quad (31)$$

$$+ \underbrace{\mathbb{E}_{q(\mathbf{r}_{1:S} | \mathbf{r}_0, \boldsymbol{\tau}_0)} \left[ \ln \frac{p(\mathbf{r}_S)}{q(\mathbf{r}_S | \mathbf{r}_0, \boldsymbol{\tau}_0)} + \sum_{t=1}^S \ln \frac{p_{\theta_2}(\mathbf{r}_{t-1} | \mathbf{r}_t, \boldsymbol{\tau}_0, c)}{q(\mathbf{r}_{t-1} | \mathbf{r}_t, \boldsymbol{\tau}_0, c)} \right]}_{\text{ELBO for } p_{\theta_2}(\mathbf{r}_0 | \boldsymbol{\tau}_0, c)}. \quad (32)$$

Denote the corresponding negative ELBOs by  $L_\tau(\boldsymbol{\tau}_0 | c)$  and  $L_{r|\tau}(\mathbf{r}_0 | \boldsymbol{\tau}_0, c)$ . Then the conditional NLL is upper-bounded by their sum:

$$-\ln p_\theta(\mathbf{x}_0 | c) \leq L_\tau(\boldsymbol{\tau}_0 | c) + L_{r|\tau}(\mathbf{r}_0 | \boldsymbol{\tau}_0, c). \quad (33)$$

**Simple  $\epsilon$ -prediction form.** With the standard  $\epsilon$ -parameterization, each term admits the DDPM simple loss (up to  $t$ -dependent constants):

$$L_\tau(\boldsymbol{\tau}_0 | c) = \mathbb{E}_{\boldsymbol{\tau}_0, t, \epsilon} \left[ w_t \|\epsilon - \epsilon_{\theta_1}(\boldsymbol{\tau}_t, t, c)\|_2^2 \right], \quad L_{r|\tau}(\mathbf{r}_0 | \boldsymbol{\tau}_0, c) = \mathbb{E}_{\mathbf{r}_0, t, \epsilon} \left[ w_t \|\epsilon - \epsilon_{\theta_2}(\mathbf{r}_t, t, \boldsymbol{\tau}_0, c)\|_2^2 \right], \quad (34)$$

where  $\boldsymbol{\tau}_t = \sqrt{\bar{\alpha}_t} \boldsymbol{\tau}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$  and  $\mathbf{r}_t = \sqrt{\bar{\alpha}_t} \mathbf{r}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ , with  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $w_t$  a chosen weighting (e.g., uniform or SNR-based).

**Training objective.** Summing the two branches yields the hierarchical training criterion

$$\mathcal{J}_{\text{hier}}(\theta_1, \theta_2) = \mathbb{E}_{c \sim \hat{p}(C)} \left[ \mathbb{E} [w_t \|\epsilon - \epsilon_{\theta_1}(\boldsymbol{\tau}_t, t, c)\|_2^2] + \mathbb{E} [w_t \|\epsilon - \epsilon_{\theta_2}(\mathbf{r}_t, t, \boldsymbol{\tau}_0, c)\|_2^2] \right] \quad (\text{constants omitted}). \quad (35)$$

Equations equation 33–equation 35 show that, even when the residual is explicitly conditioned on the realized trend, the overall objective remains a *sum of diffusion objectives*: one marginal (trend) and one conditional (residual). This preserves the clean additive structure used throughout our windowed ELBO, and reduces to Section B.3 when the conditional dependence is removed.

## B.5 FROM FULL SEQUENCES TO WINDOWS

Let a long sequence be split into  $M$  training windows  $\{\mathbf{x}_0^{(i)}\}_{i=1}^M$  with deterministic decompositions  $\mathbf{x}_0^{(i)} = \boldsymbol{\tau}_0^{(i)} + \mathbf{r}_0^{(i)}$  as in equation 19. Let  $c_i$  denote window metadata (start index  $s_i$  and total series length  $T$ ) for window  $i$ .

**Non-overlapping windows.** Suppose the windows form a partition of the sequence (no shared time indices). Under the context-conditional independent model used in DiffPM,

$$p_\theta(\boldsymbol{\tau}_0^{(i)}, \mathbf{r}_0^{(i)} | c_i) = p_{\theta_1}(\boldsymbol{\tau}_0^{(i)} | c_i) p_{\theta_2}(\mathbf{r}_0^{(i)} | c_i),$$

the conditional negative log-likelihood (NLL) of the full sequence given the set  $C = \{c_i\}_{i=1}^M$  satisfies the additive ELBO bound

$$-\ln p_\theta(\mathbf{x}_{\text{full}} | C) \leq \sum_{i=1}^M \left( L_\tau(\boldsymbol{\tau}_0^{(i)} | c_i) + L_r(\mathbf{r}_0^{(i)} | c_i) \right), \quad (\text{independent model}) \quad (36)$$

where  $L_\tau$  and  $L_r$  are the component-wise diffusion NELBOs (each conditioned on  $c_i$ ) from equation 29. In the hierarchical extension of Section B.4, the bound becomes  $\sum_i (L_\tau(\boldsymbol{\tau}_0^{(i)} | c_i) + L_{r|\tau}(\mathbf{r}_0^{(i)} | \boldsymbol{\tau}_0^{(i)}, c_i))$ . Equality holds in equation 36 if the generative model factorizes across windows given  $C$ , i.e.,  $p_\theta(\mathbf{x}_{\text{full}} | C) = \prod_i p_\theta(\mathbf{x}_0^{(i)} | c_i)$ ; otherwise, the right-hand side is a standard blockwise ELBO surrogate. Equation 36 matches the standard blockwise diffusion objective: disjoint blocks contribute additively to the global variational bound.

**Overlapping windows: notation.** Let  $\mathcal{T}$  index discrete time points of the full sequence. For window  $i$ , let  $\mathcal{T}_i \subset \mathcal{T}$  denote the set of indices covered by that window, and let  $n_t = \sum_{i=1}^M \mathbf{1}\{t \in \mathcal{T}_i\}$  be the number of windows that cover time  $t$  (we reuse  $t$  for sequence indices in  $\mathcal{T}$ ; context disambiguates diffusion step vs. sequence time). *Hereafter,  $t$  inside  $q(\cdot)$ ,  $\bar{\alpha}_t$ , and  $w_t$  denotes the diffusion step, while  $t \in \mathcal{T}$  denotes a sequence index.* Define per-window inclusion weights  $\gamma_t^{(i)} \in [0, 1]$  such that

$$\sum_{i: t \in \mathcal{T}_i} \gamma_t^{(i)} = 1 \quad \text{for all } t \in \mathcal{T}. \quad (37)$$

Two common choices are: (i) uniform sharing,  $\gamma_t^{(i)} = 1/n_t$  for all  $i$  with  $t \in \mathcal{T}_i$ ; (ii) tapering (e.g., linear or Hann) within overlaps, still normalized to respect equation 43. We write  $\mathbf{\Gamma}^{(i)}$  for the diagonal mask with entries  $\gamma_t^{(i)}$  restricted to  $\mathcal{T}_i$ .

### A. INCLUSION-EXCLUSION REWEIGHTING (NO DOUBLE COUNTING).

A principled surrogate for the conditional sequence NLL is obtained by *redistributing* overlap mass according to  $\gamma_t^{(i)}$ . Concretely, define the window-level per-time NELBO contributions  $\ell_\tau^{(i)}(t)$  and  $\ell_r^{(i)}(t)$  so that  $L_\tau(\boldsymbol{\tau}_0^{(i)} | c_i) = \sum_{t \in \mathcal{T}_i} \ell_\tau^{(i)}(t)$  and similarly for  $L_r$ . The overlap-aware objective is

$$\tilde{\mathcal{L}} = \sum_{i=1}^M \sum_{t \in \mathcal{T}_i} \gamma_t^{(i)} \left( \ell_\tau^{(i)}(t) + \ell_r^{(i)}(t) \right), \quad (38)$$

which can be equivalently written in vector form as  $\sum_i \langle \mathbf{\Gamma}^{(i)}, \ell_\tau^{(i)} + \ell_r^{(i)} \rangle$  with inner product over the time axis of window  $i$ . By construction equation 43 implies that each global time index contributes *exactly once* to equation 38, so there is no double counting. Replacing  $\ell_\tau^{(i)}(t)$  with the per-time DDPM simple losses (up to constants) yields a practical training loss that coincides with the usual per-window objective when windows do not overlap ( $\gamma_t^{(i)} \equiv 1$ ). In our implementation, training uses framewise  $\epsilon$ -MSE, so the per-time decomposition used here is exact.

### B. CONDITIONAL (BLOCKWISE) CHAINING WITH OVERLAP CONSISTENCY.

An alternative view generates windows in a fixed order while enforcing agreement on overlaps. Let  $i = 1, \dots, M$  index the order (e.g., increasing start times). Write  $\mathbf{o}^{(i)}$  for the *observed* overlap

values of window  $i$  that are already committed by previously generated windows, and let  $\mathcal{O}^{(i)} \subset \mathcal{T}_i$  be the corresponding time indices. Define the *conditional* per-window generative factors

$$p_{\theta_1}^{(i)}(\boldsymbol{\tau}_0^{(i)} | c_i, \mathbf{o}^{(i)}) \propto \int p_{\theta_1}(\boldsymbol{\tau}_{0:S}^{(i)} | c_i) \prod_{t \in \mathcal{O}^{(i)}} \delta(\tau_{0,t}^{(i)} - o_t^{(i)}) d\boldsymbol{\tau}_{1:S}^{(i)}, \quad (39)$$

$$p_{\theta_2}^{(i)}(\mathbf{r}_0^{(i)} | c_i, \mathbf{o}^{(i)}) \propto \int p_{\theta_2}(\mathbf{r}_{0:S}^{(i)} | c_i) \prod_{t \in \mathcal{O}^{(i)}} \delta(r_{0,t}^{(i)} - \tilde{o}_t^{(i)}) d\mathbf{r}_{1:S}^{(i)}, \quad (40)$$

where  $(o_t^{(i)}, \tilde{o}_t^{(i)})$  encode the trend/residual parts of the overlap constraint. Because the decomposition is deterministic, the committed overlap on  $\mathbf{x}$  induces unique  $(o_t^{(i)}, \tilde{o}_t^{(i)})$  via  $(o_t^{(i)}, \tilde{o}_t^{(i)}) = (\mathcal{T}(\mathbf{x})_t, \mathbf{x}_t - \mathcal{T}(\mathbf{x})_t)$  on  $\mathcal{O}^{(i)}$ . The conditional ELBO for window  $i$  then reads

$$\mathcal{L}_{\text{cond}}^{(i)}(\mathbf{x}_0^{(i)} | c_i, \mathbf{o}^{(i)}) \leq \ln p_{\theta}(\mathbf{x}_0^{(i)} | c_i, \mathbf{o}^{(i)}), \quad (41)$$

and the global conditional NLL obeys the additive bound

$$-\ln p_{\theta}(\mathbf{x}_{\text{full}} | C) \leq \sum_{i=1}^M \left( L_{\tau}^{(i)}(\boldsymbol{\tau}_0^{(i)} | c_i, \mathbf{o}^{(i)}) + L_r^{(i)}(\mathbf{r}_0^{(i)} | c_i, \mathbf{o}^{(i)}) \right), \quad (42)$$

with  $L^{(i)}$  the conditional NELBOs corresponding to equation 41. In practice, imposing  $\mathbf{o}^{(i)}$  is implemented by *clamping* the overlap entries to their committed values during the forward/backward passes of the DDPM objective for window  $i$  (equivalently, by zeroing the loss on clamped coordinates; we supply both the binary mask  $\mathbf{M}_i$  and the clamped overlap values  $z_0^{(i)} \upharpoonright \mathcal{O}_{i-1,i}$  as conditioning features). This “blockwise chaining” recovers the non-overlap result when  $\mathcal{O}^{(i)} = \emptyset$  and preserves coherence across window boundaries when overlaps are present.

**Remarks.** (i) Objectives equation 38 and equation 42 are complementary. The former is *embarrassingly parallel* (no dependency across windows) and avoids double counting via normalized weights; the latter is *sequential* but enforces hard consistency on overlaps. (ii) Both reduce to equation 36 when windows are disjoint. (iii) Either construction is compatible with the independent model (used in DiffPM) or the hierarchical extension by replacing  $L_r$  with  $L_{r|\tau}$  as in Section B.4.

## B.6 OVERLAPPING WINDOWS I: INCLUSION–EXCLUSION (REWEIGHTING)

Let the full sequence be indexed by a discrete set  $\mathcal{T}$ . Window  $i \in \{1, \dots, M\}$  covers a subset  $\mathcal{T}_i \subset \mathcal{T}$  (its *support*). For two consecutive windows  $i$  and  $i+1$ , the overlap index set is  $\mathcal{O}_{i,i+1} = \mathcal{T}_i \cap \mathcal{T}_{i+1}$  (of size  $O$ ). For general sliding windows, define the *coverage multiplicity* of time index  $t$  as

$$n_t = |\{i : t \in \mathcal{T}_i\}| \in \{1, 2, \dots\}.$$

To avoid double counting when summing window losses, assign *per-index weights*  $\gamma_t^{(i)} \in [0, 1]$  that satisfy the local normalization

$$\sum_{i: t \in \mathcal{T}_i} \gamma_t^{(i)} = 1, \quad \forall t \in \mathcal{T}. \quad (43)$$

A convenient choice is uniform sharing,  $\gamma_t^{(i)} = 1/n_t$  for all  $i$  with  $t \in \mathcal{T}_i$ . Let  $c_i$  denote the metadata (start  $s_i$ , series length  $T$ ) for window  $i$  and let the (per-window) diffusion NELBO decompose over time indices,<sup>2</sup>

$$L^{(i)}(\mathbf{x}_0^{(i)} | c_i) \approx \sum_{t \in \mathcal{T}_i} \left( \ell_{\tau}^{(i)}(t | c_i) + \ell_r^{(i)}(t | c_i) \right). \quad (44)$$

Then the *reweighted* global objective

$$\tilde{\mathcal{L}} = \sum_{i=1}^M \sum_{t \in \mathcal{T}_i} \gamma_t^{(i)} \left( \ell_{\tau}^{(i)}(t | c_i) + \ell_r^{(i)}(t | c_i) \right) \quad (45)$$

<sup>2</sup>This decomposition is exact for objectives that sum over time (e.g., framewise  $\epsilon$ -MSE); it is a standard approximation otherwise and works well in practice.

counts each time index *exactly once* by equation 43. For the common case of fixed window length  $L$ , fixed overlap size  $O$ , and stride  $S = L - O$ , we have  $n_t = 2$  inside overlaps and  $n_t = 1$  elsewhere, so  $\gamma_t^{(i)} = \frac{1}{2}$  on overlaps and 1 elsewhere—equivalent to subtracting one copy of the overlap from the naïve sum (an inclusion–exclusion correction). Formally, for two windows  $W_1, W_2$ ,

$$\ln p_\theta(W_1 \cup W_2 | c_1, c_2) = \ln p_\theta(W_1 | c_1) + \ln p_\theta(W_2 | c_2) - \ln p_\theta(W_1 \cap W_2 | c_1, c_2), \quad (46)$$

and the reweighted objective mirrors this identity at the ELBO level by distributing the overlap mass.

**Component-wise view.** The same weights apply to the trend and residual branches individually:

$$\tilde{\mathcal{L}} = \sum_{i=1}^M \sum_{t \in \mathcal{T}_i} \gamma_t^{(i)} \ell_r^{(i)}(t | c_i) + \sum_{i=1}^M \sum_{t \in \mathcal{T}_i} \gamma_t^{(i)} \ell_r^{(i)}(t | c_i), \quad (47)$$

so the two-branch structure is preserved under reweighting. In the hierarchical extension (Section B.4), replace  $\ell_r^{(i)}$  by its conditional counterpart  $\ell_{r|\tau}^{(i)}$ .

## B.7 OVERLAPPING WINDOWS II: CONDITIONAL (BLOCKWISE) CHAINING

An alternative is to make each window *conditional* on the overlap with already-committed windows, so no double counting arises in the first place. Let  $\mathbf{W}_i$  be the  $i$ -th window in a fixed order (e.g., increasing start indices). Write  $\mathcal{O}_{i-1,i} = \mathcal{T}_{i-1} \cap \mathcal{T}_i$  and decompose  $\mathbf{W}_i = (\mathbf{W}_i^{\text{ov}}, \mathbf{W}_i^{\text{new}})$  where  $\mathbf{W}_i^{\text{ov}}$  is known from  $\mathbf{W}_{i-1}$  on  $\mathcal{O}_{i-1,i}$ . Conditioning on the set of window contexts  $C = \{c_i\}_{i=1}^M$ , the sequence likelihood factorizes as

$$p_\theta(\mathbf{x}_{\text{full}} | C) = p_\theta(\mathbf{W}_1 | c_1) \prod_{i=2}^M p_\theta(\mathbf{W}_i^{\text{new}} | \mathbf{W}_{i-1} \upharpoonright \mathcal{O}_{i-1,i}, c_i). \quad (48)$$

Taking  $-\ln$  and applying ELBOs termwise yields a *sum of conditional ELBOs*:

$$-\ln p_\theta(\mathbf{x}_{\text{full}} | C) \leq L(\mathbf{W}_1 | c_1) + \sum_{i=2}^M L(\mathbf{W}_i^{\text{new}} | \mathbf{W}_{i-1} \upharpoonright \mathcal{O}_{i-1,i}, c_i), \quad (49)$$

where each  $L(\cdot)$  is the (conditional) diffusion NELBO for the corresponding factor.

**Diffusion implementation via clamping (mask/inpainting).** For branch  $z \in \{\tau, \mathbf{r}\}$ , construct a binary mask  $\mathbf{M}_i \in \{0, 1\}^{L \times D}$  that is 1 on the overlap indices  $\mathcal{O}_{i-1,i}$  of window  $i$  and 0 elsewhere. Given clean  $z_0^{(i)}$  and its noised version  $z_t^{(i)} = \sqrt{\alpha_t} z_0^{(i)} + \sqrt{1 - \alpha_t} \epsilon$ , form the *clamped* input

$$\tilde{z}_t^{(i)} = \mathbf{M}_i \odot z_0^{(i)} + (1 - \mathbf{M}_i) \odot z_t^{(i)}, \quad (50)$$

where  $\odot$  denotes the Hadamard (element-wise) product, so that overlap entries remain fixed to the committed values, while the remaining entries are noised as usual. Train with the masked  $\epsilon$ -loss

$$\mathcal{L}_{z,i}^{\text{cond}} = \mathbb{E} \left[ w_t \left\| (1 - \mathbf{M}_i) \odot (\epsilon - \epsilon_\theta(\tilde{z}_t^{(i)}, t, c_i, \mathbf{M}_i, z_0^{(i)} \upharpoonright \mathcal{O}_{i-1,i})) \right\|_2^2 \right], \quad (51)$$

which enforces that the model *cannot* change overlap entries and must denoise only the new region. Summing the two branches gives the total conditional objective

$$\mathcal{J}_{\text{cond}} = \sum_{i=1}^M \left( \mathcal{L}_{\tau,i}^{\text{cond}} + \mathcal{L}_{\mathbf{r},i}^{\text{cond}} \right) \quad (\text{constants omitted}). \quad (52)$$

**Remarks.** (i) Compared to inclusion–exclusion equation 38, the conditional chaining bound equation 49 enforces boundary coherence by construction (no cross-fading needed) and aligns training with sequential generation. (ii) When overlaps are empty ( $\mathcal{O}_{i-1,i} = \emptyset$ ), the mask is  $\mathbf{M}_i \equiv 0$  and equation 51 reduces to the usual per-window DDPM loss. (iii) The construction is compatible with the independent model (used in DiffPM) and with the hierarchical extension by augmenting the residual denoiser input with  $\tau_0$  when desired.

## B.8 ASSUMPTIONS, EDGE CASES, AND PRACTICAL GUIDANCE

**Deterministic decomposition.** Equation equation 19 assumes a fixed analysis operator  $\mathcal{T}$  (no learned posterior). Hence  $q(\boldsymbol{\tau}_0, \mathbf{r}_0 \mid \mathbf{x}_0, c)$  is a product of Diracs that exactly satisfy  $\mathbf{x}_0 = \boldsymbol{\tau}_0 + \mathbf{r}_0$ , and the Dirac constraint cancels from the ELBO ratio in equation 25. This keeps the variational family simple and makes the ELBO algebra in Sections B.3–B.4 exact up to the usual DDPM constants.

**Independent vs. hierarchical.** If residual statistics depend on the realized trend (e.g., level-dependent variance or trend-modulated periodicity), the hierarchical objective equation 33 is strictly more expressive and can capture such couplings. If not, the context-conditional independent objective equation 28 is simpler, decouples training across branches, and is what DiffPM uses in the main experiments.

**Context and conditioning.** All sequence/window likelihoods and ELBOs are understood *conditional on context  $c$* : window start  $s$ , total series length  $T$ , and any positional metadata (cf. equation 22 and equation 29). In the hierarchical variant, the residual branch may additionally consume the *realized* clean trend  $\boldsymbol{\tau}_0$  as input to the denoiser (see equation 34).

**Per-time decomposition.** The inclusion–exclusion objective equation 38 uses a per-time decomposition equation 44. This is *exact* for framewise DDPM losses (e.g.,  $\epsilon$ -MSE summed over time and channels) and is a standard, accurate approximation otherwise. With fixed window length  $L$  and stride  $S$ , set  $\gamma_t^{(i)} = 1$  outside overlaps and  $\gamma_t^{(i)} = \frac{1}{2}$  on overlaps. For nonuniform coverages, use the normalized rule  $\gamma_t^{(i)} = 1/n_t$ , where  $n_t$  is the coverage multiplicity (Section B.6).

**Overlap choices and masks.** In the conditional chaining scheme (Section B.7), a single binary mask  $\mathbf{M}_i$  is applied to both branches, clamping the overlap entries to already-committed values during training via equation 50–equation 51. The conditioning features can be chosen minimally (raw overlap values), augmented (e.g., learned embeddings of the overlap region), and always include the positional metadata  $c_i$ .

**Variable lengths and boundaries.** If windows near the sequence boundaries are shorter than  $L$  (e.g., due to causal padding or dataset truncation), one can either (i) pad them and keep equation 38 with the same  $\gamma_t^{(i)}$ , or (ii) shrink  $\mathcal{T}_i$  and renormalize weights to satisfy equation 43. Both preserve the additive ELBO structure and avoid double counting.

**Coverage multiplicities  $> 2$ .** When more than two windows overlap at a time index (e.g., small strides), the inclusion–exclusion normalization equation 43 still applies with  $\gamma_t^{(i)} = 1/n_t$ , and equation 38 continues to count each index exactly once. Conditional chaining also generalizes by conditioning each new window on *all* already-committed overlaps in  $\mathcal{T}_i$ .

**Missing values in training windows.** If a training window contains missing entries, either (i) exclude those coordinates from the per-time loss by multiplying equation 29 with a data-mask (equivalently, set  $\gamma_t^{(i)} = 0$  at missing indices in equation 38), or (ii) treat them as known inpainting masks within the conditional-chaining view (using  $\mathbf{M}_i$  to clamp observed entries and ignore missing ones). Both choices are compatible with the ELBO derivations.

**Training objective summary.** For a dataset of (possibly overlapping) windows with contexts  $\{c_i\}_{i=1}^M$ , the DiffPM objective is a *sum of branchwise diffusion ELBOs* across windows:

$$\min_{\theta_1, \theta_2} \sum_{i=1}^M \left\{ \mathbb{E} \left[ w_t \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta_1}(\boldsymbol{\tau}_t^{(i)}, t, c_i) \right\|_2^2 \right] + \mathbb{E} \left[ w_t \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta_2}(\mathbf{r}_t^{(i)}, t, c_i) \right\|_2^2 \right] \right\} \quad (\text{constants omitted}), \quad (53)$$

with either

- (a) *inclusion–exclusion reweighting*: multiply the per-time losses inside the expectations in equation 53 by the index weights  $\gamma_t^{(i)}$  from equation 43, yielding the global reweighted objective equation 38; or

- (b) *conditional chaining*: replace the inputs by the clamped  $\tilde{z}_t^{(i)}$  from equation 50 and restrict the error to  $(1 - M_i)$  as in equation 51, summing the conditional losses across windows as in equation 52.

Both strategies realize a coherent, additive ELBO over windows and branches, differing only in whether overlap consistency is achieved statistically (reweighting) or enforced as a hard constraint (clamping).

**Generation-time stitching (context only).** Under inclusion–exclusion training, overlapping windows can be merged by overlap-add or linear cross-fades (with weights consistent with equation 43); under conditional chaining, overlaps are fixed by design, so windows join seamlessly without averaging. (Empirical sampling details are discussed in the main text.)

## C RECONSTRUCTION FIDELITY GUARANTEES FOR OVERLAP–ADD STITCHING

This appendix provides a detailed theoretical analysis of the overlap–add (OLA) stitching mechanism used in DiffPM, specifically when employing constant overlap–add (COLA) windows. We begin by establishing a rigorous mathematical setup, which serves as a foundation for proving perfect reconstruction (PR) guarantees under ideal, noise-free conditions. Subsequently, we derive both worst-case and stochastic error bounds to characterize performance in the presence of local prediction errors. The analysis is extended to assess the robustness of the method to non-ideal window sums, to handle multivariate signals, and to identify optimal weighting schemes under heteroscedastic local errors. It is important to note that all results presented apply independently to both the trend and residual branches of DiffPM; by linearity, the final output inherits the guarantees established for each branch.

### C.1 NOTATION AND SETUP

Let  $x[n] \in \mathbb{R}$  represent a discrete-time univariate signal of length  $N$ . We fix a window length  $L \in \mathbb{N}$  and a hop size  $R \in \{1, \dots, L\}$ . The analysis centers on a real-valued, non-negative window function  $w[n]$  that is supported on the interval  $n \in \{0, \dots, L - 1\}$ :

$$w[n] \geq 0, \quad w[n] = 0 \text{ for } n \notin \{0, \dots, L - 1\}.$$

The signal is segmented into a set of frames indexed by  $\mathcal{M} = \{0, 1, \dots, M - 1\}$ . The  $m$ -th frame corresponds to the signal indices

$$\mathcal{I}_m = \{mR, mR + 1, \dots, mR + L - 1\}.$$

For analytical simplicity, we assume the signal length is  $N = (M - 1)R + L$ , which ensures that the final frame aligns perfectly with the end of the signal, obviating the need for padding.<sup>3</sup>

The window function  $w$  is assumed to satisfy the normalized **constant overlap–add (COLA) property** for the given hop size  $R$ . This property dictates that the sum of all shifted versions of the window is unity at every point in time:

$$\sum_{m \in \mathbb{Z}} w[n - mR] = 1 \quad \text{for all } n \in \mathbb{Z}. \quad (54)$$

An equivalent formulation involves defining a *coverage function*  $S[n]$  over the signal domain:

$$S[n] := \sum_{m \in \mathcal{M}} w[n - mR]. \quad (55)$$

In this context, the COLA condition is met if  $S[n] \equiv 1$  for all  $n \in \{0, \dots, N - 1\}$ .<sup>4</sup>

<sup>3</sup>All results can be readily extended to cases involving truncated or padded boundaries, as briefly discussed in Sec. C.9.

<sup>4</sup>Classical window choices satisfying this property include rectangular windows with  $R = L$  (no overlap), and Hann or triangular windows with  $R = L/2$  (50% overlap), provided they are appropriately normalized. Our analysis does not presume a specific window shape, only that it satisfies Eq. equation 54.

For each frame  $m$ , a local synthesis model generates a prediction, denoted  $y_m[k]$  for  $k \in \{0, \dots, L-1\}$ , which aims to approximate the true signal segment  $x[mR+k]$ . The final reconstructed signal,  $\tilde{x}[n]$ , is synthesized via **window-weighted overlap-add**:

$$\tilde{x}[n] = \sum_{m=0}^{M-1} w[n-mR] y_m[n-mR], \quad n = 0, \dots, N-1. \quad (56)$$

Here, it is understood that  $w[j] = 0$  if  $j \notin \{0, \dots, L-1\}$ . In regions where frames overlap, the COLA property ensures that the weights  $\{w[n-mR]\}_m$  form a convex partition (i.e., they are non-negative and sum to 1), meaning that Eq. equation 56 effectively computes a weighted average of the overlapping local predictions.

In the DiffPM architecture, this OLA procedure is applied independently to both the trend and residual branches, producing reconstructions  $\tilde{x}^{(T)}$  and  $\tilde{x}^{(R)}$ . The final output is their sum:

$$\tilde{x}_{\text{total}}[n] = \tilde{x}^{(T)}[n] + \tilde{x}^{(R)}[n]. \quad (57)$$

The analysis that follows is presented for a generic single-branch reconstruction  $\tilde{x}$ . Due to the linearity of the process, the derived properties apply equally to each branch and, by extension, to their sum.

## C.2 PERFECT RECONSTRUCTION (PR)

A fundamental property of the OLA scheme is that it is lossless if the local predictions are perfectly accurate and the COLA condition holds.

**Theorem C.1** (Perfect Reconstruction under COLA). *Assume the COLA property equation 54 is satisfied and that the per-frame predictions are noise-free, i.e.,*

$$y_m[k] = x[mR+k] \quad \text{for all } m \in \mathcal{M}, k \in \{0, \dots, L-1\}. \quad (58)$$

*Then the stitched reconstruction is pointwise identical to the ground truth signal:*

$$\tilde{x}[n] = x[n] \quad \text{for all } n = 0, \dots, N-1.$$

*Proof.* Substituting the noise-free condition equation 58 into the OLA formula equation 56, we observe that for any time index  $n$ , if a term  $w[n-mR]$  is non-zero, it implies  $n-mR \in \{0, \dots, L-1\}$ . Let  $k = n-mR$ . Then  $y_m[n-mR] = y_m[k] = x[mR+k] = x[n]$ . This allows us to rewrite the OLA sum as:

$$\tilde{x}[n] = \sum_m w[n-mR] x[n] = x[n] \left( \sum_m w[n-mR] \right).$$

By the COLA property equation 54, the term in parentheses is equal to 1, which yields  $\tilde{x}[n] = x[n]$ .  $\square$

This result has two important implications. First, PR holds for both the trend and residual branches independently. If each branch is reconstructed without error, their sum will perfectly recover the original signal, since the initial decomposition was additive. Second, the proof highlights that PR in the noise-free case is guaranteed by any set of weights that form a partition of unity; the specific window shape is not critical for this property.

## C.3 DETERMINISTIC ERROR ANALYSIS: WORST-CASE BOUNDS

We now analyze the behavior of the reconstruction error when local predictions are imperfect. Let us define the local, per-frame error as

$$e_m[k] := y_m[k] - x[mR+k], \quad k = 0, \dots, L-1. \quad (59)$$

The global reconstruction error,  $\tilde{x}[n] - x[n]$ , can be expressed in terms of these local errors. By subtracting  $x[n]$  from Eq. equation 56 and using the identity  $x[n] = \sum_m w[n-mR] x[n]$ , we find:

$$\tilde{x}[n] - x[n] = \sum_{m=0}^{M-1} w[n-mR] e_m[n-mR]. \quad (60)$$

This key equation reveals that the total error at any point  $n$  is a convex combination of the local errors from all contributing frames. This structure leads directly to a strong worst-case guarantee.

**Proposition C.2** (Supremum-Norm Guarantee). *If the magnitude of the local error is uniformly bounded by a constant  $\varepsilon$ , i.e.,  $|e_m[k]| \leq \varepsilon$  for all  $m$  and  $k$ , then the global reconstruction error is also bounded by  $\varepsilon$ :  $|\tilde{x}[n] - x[n]| \leq \varepsilon$  for all  $n$ .*

*Proof.* Applying the triangle inequality to Eq. equation 60 yields:

$$|\tilde{x}[n] - x[n]| \leq \sum_m w[n - mR] |e_m[n - mR]| \leq \sum_m w[n - mR] \varepsilon.$$

Since the weights are non-negative and sum to 1 by the COLA property equation 54, the expression simplifies to:

$$|\tilde{x}[n] - x[n]| \leq \varepsilon \sum_m w[n - mR] = \varepsilon.$$

□

This proposition establishes that OLA with a COLA window does not amplify the worst-case error; the stitched error’s maximum deviation is no greater than the maximum local deviation. Furthermore, in practical scenarios where local errors  $e_m[n - mR]$  might have varying signs across overlapping frames, the weighted averaging in Eq. equation 60 can lead to error cancellation, resulting in a global error that is often much smaller than the worst-case bound  $\varepsilon$ .

#### C.4 STOCHASTIC ERROR ANALYSIS: UNBIASEDNESS AND VARIANCE

To gain deeper insight, we adopt a stochastic perspective. Assume the local errors are random variables with zero mean and finite variance. For a given absolute time index  $n$ , let the statistical properties of the error from frame  $m$  be:

$$\mathbb{E}[e_m[n - mR]] = 0, \quad \text{Var}(e_m[n - mR]) = \sigma_m^2[n].$$

We also allow for non-zero covariance between errors from different frames,  $m \neq m'$ , that overlap at time  $n$ :  $\text{Cov}(e_m[n - mR], e_{m'}[n - m'R]) = \gamma_{m,m'}[n]$ .

**Proposition C.3** (Unbiasedness and Variance with General Covariance). *The OLA reconstruction is unbiased, and its variance at each time step  $n$  is given by:*

$$\mathbb{E}[\tilde{x}[n] - x[n]] = 0, \tag{61}$$

$$\text{Var}(\tilde{x}[n] - x[n]) = \sum_m w[n - mR]^2 \sigma_m^2[n] + 2 \sum_{m < m'} w[n - mR] w[n - m'R] \gamma_{m,m'}[n]. \tag{62}$$

*Proof.* Unbiasedness follows directly from the linearity of expectation applied to Eq. equation 60:  $\mathbb{E}[\tilde{x} - x] = \sum_m w[n - mR] \mathbb{E}[e_m[n - mR]] = 0$ . The variance formula is obtained by applying the standard formula for the variance of a weighted sum of correlated random variables to Eq. equation 60. □

A particularly important special case arises when local errors are **independent and homoscedastic**. If the errors are uncorrelated across frames at time  $n$  ( $\gamma_{m,m'}[n] = 0$  for  $m \neq m'$ ) and share a common variance  $\sigma^2$ , Eq. equation 62 simplifies considerably:

$$\text{Var}(\tilde{x}[n] - x[n]) = \sigma^2 \sum_m w[n - mR]^2. \tag{63}$$

This result demonstrates a key benefit of overlap-add: variance reduction. Since  $\sum_m w[n - mR] = 1$  and  $w[n] \geq 0$ , the sum of squared weights  $\sum_m w[n - mR]^2$  is always less than or equal to 1. For instance, if  $K$  frames overlap at time  $n$  with roughly equal weights  $w \approx 1/K$ , then the variance is reduced by a factor of  $\sum_m w^2 \approx K(1/K)^2 = 1/K$ . A 50% overlap with symmetric weights  $(1/2, 1/2)$  would halve the variance.

Conversely, Eq. equation 62 quantifies the impact of **correlated errors**. Positive correlations ( $\gamma_{m,m'} > 0$ ) diminish the variance reduction, as errors tend to reinforce each other. Negative correlations, on the other hand, enhance it through destructive interference. The parallel generation of frames in DiffPM helps ensure that cross-frame error correlations are typically small, allowing the benefits of variance reduction to be realized.

### 1188 C.5 ROBUSTNESS TO NON-IDEAL WINDOW SUMS

1189 In practice, the COLA property may hold only approximately, for instance due to boundary effects  
 1190 or floating-point inaccuracies. Let us analyze the impact of such deviations. We use the coverage  
 1191 function  $S[n] = \sum_m w[n - mR]$  defined in Eq. equation 55 and assume it is close to unity, bounded  
 1192 as

$$1193 \quad 1 - \delta \leq S[n] \leq 1 + \delta \quad \text{for all } n, \quad (64)$$

1194 for some small  $\delta \in [0, 1)$ . The reconstruction in this case becomes:

$$1195 \quad \tilde{x}[n] = \sum_m w[n - mR] (x[n] + e_m[n - mR]) = S[n]x[n] + \underbrace{\sum_m w[n - mR] e_m[n - mR]}_{\text{error term}}.$$

1196 The total deviation from the true signal is therefore:

$$1197 \quad \tilde{x}[n] - x[n] = (S[n] - 1)x[n] + \sum_m w[n - mR] e_m[n - mR].$$

1198 This expression reveals two distinct error components. The first term,  $(S[n] - 1)x[n]$ , represents  
 1199 a **deterministic gain error** caused by the mis-normalization of the window sum. If the signal is  
 1200 bounded by  $|x[n]| \leq B$ , the magnitude of this gain error is bounded by  $|(S[n] - 1)x[n]| \leq \delta B$ . The  
 1201 second term is the familiar stochastic error component, which is subject to the same worst-case and  
 1202 stochastic bounds derived previously. Consequently, the global  $L_\infty$  error is bounded by:

$$1203 \quad \|\tilde{x} - x\|_\infty \leq \delta \|x\|_\infty + \varepsilon,$$

1204 assuming a uniform per-frame error bound  $|e_m[k]| \leq \varepsilon$ . In terms of mean squared error, assuming  
 1205 independence of the signal and the stochastic error, we have:

$$1206 \quad \mathbb{E}[(\tilde{x}[n] - x[n])^2] = (S[n] - 1)^2 x[n]^2 + \text{Var}\left(\sum_m w e_m\right).$$

1207 This shows that small deviations of  $S[n]$  from 1 introduce a controllable, signal-dependent bias term  
 1208 in addition to the variance analyzed in Sec. C.4.

### 1209 C.6 GLOBAL $L_2$ BOUNDS AND AVERAGE MSE

1210 To assess the overall performance across the entire signal, we consider the global mean squared error  
 1211 (MSE):

$$1212 \quad \text{MSE} := \frac{1}{N} \sum_{n=0}^{N-1} \mathbb{E}[(\tilde{x}[n] - x[n])^2].$$

1213 Assuming exact COLA and independent, homoscedastic local errors with variance  $\sigma^2$ , we can use  
 1214 Eq. equation 63 to express the MSE as:

$$1215 \quad \text{MSE} = \frac{\sigma^2}{N} \sum_{n=0}^{N-1} \sum_m w[n - mR]^2 = \sigma^2 \left( \frac{1}{N} \sum_{n=0}^{N-1} \sum_m w[n - mR]^2 \right).$$

1216 The term in parentheses can be interpreted as an "effective inverse overlap" factor, averaged over  
 1217 the signal length. In regions away from the boundaries, the overlap pattern is typically stationary,  
 1218 and the inner sum  $\sum_m w[n - mR]^2$  is periodic with period  $R$ . If  $K$  denotes the typical number  
 1219 of overlapping frames at these interior points and the weights are roughly balanced, then we have  
 1220  $\sum_m w^2 \approx 1/K$ . The average over the whole signal will be close to this value, leading to an  
 1221 approximate global MSE of:

$$1222 \quad \text{MSE} \approx \sigma^2/K.$$

1223 Near the signal boundaries where there are fewer overlapping frames, the local variance will  
 1224 be higher, slightly increasing the global average. However, the MSE remains bounded between  
 1225  $\sigma^2/K_{\max}$  and  $\sigma^2/K_{\min}$ , where  $K_{\min}$  and  $K_{\max}$  are the minimum and maximum number of over-  
 1226 lapping frames across the signal.

## 1242 C.7 MULTIVARIATE EXTENSION

1243  
1244 The OLA framework extends straightforwardly to multivariate signals  $x[n] \in \mathbb{R}^D$ . The procedure is  
1245 simply applied channel-wise for each dimension  $d = 1, \dots, D$ :

$$1246 \quad \tilde{x}^{(d)}[n] = \sum_m w[n - mR] y_m^{(d)}[n - mR].$$

1247  
1248 Let  $\mathbf{e}_m[n] \in \mathbb{R}^D$  be the vector of channel-wise errors for frame  $m$  at absolute time  $n$ . We assume  
1249 it has zero mean ( $\mathbb{E}[\mathbf{e}_m[n]] = \mathbf{0}$ ) and a within-frame covariance matrix  $\Sigma_m[n] \in \mathbb{R}^{D \times D}$ . The total  
1250 error vector is a weighted sum of the local error vectors:  
1251

$$1252 \quad \tilde{\mathbf{x}}[n] - \mathbf{x}[n] = \sum_m w[n - mR] \mathbf{e}_m[n].$$

1253  
1254 The covariance matrix of the reconstruction error is then given by:

$$1255 \quad \text{Cov}(\tilde{\mathbf{x}}[n] - \mathbf{x}[n]) = \sum_m w[n - mR]^2 \Sigma_m[n] + \sum_{m \neq m'} w[n - mR] w[n - m'R] \text{Cov}(\mathbf{e}_m[n], \mathbf{e}_{m'}[n]).$$

1256  
1257 If the errors are uncorrelated across frames, the second term vanishes. All the scalar guarantees,  
1258 such as the non-amplification of worst-case error (Prop. C.2), apply on a per-channel basis.  
1259  
1260

## 1261 C.8 OPTIMAL WEIGHTS UNDER HETEROSCEDASTIC LOCAL ERRORS

1262  
1263 Our analysis so far has mostly relied on a fixed window function  $w$ , implying fixed weights. How-  
1264 ever, it is instructive to consider the optimal weighting strategy in a more general setting. At a given  
1265 time point  $n$ , suppose  $K$  frames overlap, contributing independent, zero-mean local errors  $e_i$  with  
1266 differing variances  $\sigma_i^2$ . If we are free to choose blending weights  $\alpha_i \geq 0$  such that  $\sum_{i=1}^K \alpha_i = 1$ ,  
1267 the total error is  $\sum_i \alpha_i e_i$ . The variance of this sum is:  
1268

$$1269 \quad \text{Var}\left(\sum_{i=1}^K \alpha_i e_i\right) = \sum_{i=1}^K \alpha_i^2 \sigma_i^2.$$

1270  
1271 Minimizing this variance subject to the simplex constraint on  $\alpha_i$  is a classic problem whose solution  
1272 is the well-known **inverse-variance weighting**:  
1273

$$1274 \quad \alpha_i^* = \frac{\sigma_i^{-2}}{\sum_{j=1}^K \sigma_j^{-2}}. \quad (65)$$

1275  
1276 This strategy gives more weight to predictions that are known to be more reliable (i.e., have lower  
1277 variance). For the more general case of correlated errors with a joint positive-definite covariance  
1278 matrix  $\Sigma$ , the optimal weights  $\alpha$  solve the quadratic program:  
1279

$$1280 \quad \min_{\alpha \in \mathbb{R}^K} \alpha^\top \Sigma \alpha \quad \text{s.t.} \quad \mathbf{1}^\top \alpha = 1, \quad \alpha \geq \mathbf{0},$$

1281  
1282 whose unconstrained solution is  $\alpha^* \propto \Sigma^{-1} \mathbf{1}$ . In DiffPM, we use a fixed window  $w$  for its simplicity  
1283 and computational efficiency, as it avoids the need for per-time-step optimization. This analysis  
1284 shows, however, that any deviation from the homoscedastic error assumption is, in principle, best  
1285 addressed by adapting the weights according to an inverse-variance scheme.  
1286

## 1287 C.9 BOUNDARY HANDLING

1288  
1289 The number of overlapping frames naturally decreases near the signal boundaries at  $n = 0$  and  
1290  $n = N - 1$ . As long as the COLA property is maintained (i.e., the weights still sum to 1), all the  
1291 perfect reconstruction and error-bounding results continue to hold without modification. If, however,  
1292 the window sum deviates from unity at the boundaries, the analysis from Sec. C.5 applies, bounding  
1293 the resulting gain error. A common practical approach is to use signal padding (e.g., zero, reflection,  
1294 or replication) before windowing, which can maintain a uniform number of overlaps throughout  
1295 the signal. The analysis presented here remains valid, with the understanding that it applies to the  
extended signal.

1296 C.10 PUTTING IT TOGETHER FOR DIFFPM  
1297

1298 The theoretical guarantees developed in this appendix can be applied independently to DiffPM’s  
1299 trend and residual branches, with their consequences inherited by the final summed output.

- 1300 • **Perfect Reconstruction:** If the local predictions for both trend and residual are exact, each  
1301 branch reconstructs perfectly under COLA. By additivity, their sum perfectly recovers the  
1302 original signal.
- 1303 • **No Worst-Case Amplification:** The maximum stitched error in each branch is bounded  
1304 by its respective maximum per-frame error. The total error of the final output is therefore  
1305 bounded by the sum of these two branch-wise bounds.
- 1306 • **Variance Reduction:** In regions of overlap, the averaging process reduces the variance  
1307 of the reconstruction error for each branch. For independent, zero-mean local errors, the  
1308 variance reduction is roughly proportional to the number of overlapping windows  $K$ . As-  
1309 suming the errors between the trend and residual branches are also independent, the total  
1310 variance of the final output is simply the sum of the reduced branch-wise variances.

1311 Thus, the parallel OLA stitching mechanism in DiffPM provides a principled way to achieve seam-  
1312 less and globally coherent reconstructions. It ensures controlled error propagation and leverages  
1313 overlap to reduce noise, all without resorting to sequential conditioning or overwrite strategies.  
1314

1315 C.11 PRACTICAL GUIDANCE ON WINDOW AND HOP SIZE  
1316

1317 The choice of window function and hop size has practical implications for performance and compu-  
1318 tational cost.

1319 **Window Choice.** Any window  $w$  that satisfies the COLA property equation 54 is theoretically valid.  
1320 Hann and triangular windows, when used with a 50% overlap ( $R = L/2$ ), are popular choices as  
1321 they satisfy COLA (with proper normalization) and provide smooth crossfades between frames. A  
1322 rectangular window with  $R = L$  also satisfies COLA but offers no overlap, and thus no variance  
1323 reduction.  
1324

1325 **Hop Size  $R$ .** A smaller hop size  $R$  results in greater overlap, which increases the number of av-  
1326 eraged frames  $K$  and thereby enhances variance reduction (Sec. C.4). This benefit comes at the  
1327 cost of increased computation, as more frames must be processed. A hop size of  $R = L/2$  is often  
1328 considered a good trade-off between reconstruction quality (smoothness and noise reduction) and  
1329 computational load.  
1330

1331 **Normalization.** It is crucial to ensure that the window sum  $\sum_m w[n - mR]$  is numerically very  
1332 close to 1 across the entire signal to avoid the gain drift described in Sec. C.5. If the coverage  
1333 function  $S[n]$  is found to deviate from 1, a simple and effective correction is to post-normalize the  
1334 reconstruction:  $\hat{x}[n] \leftarrow \tilde{x}[n]/S[n]$  (for  $S[n] > 0$ ). This step restores the perfect reconstruction  
1335 property in the noiseless case and eliminates the bias term from the error.  
1336

1337 C.12 SUMMARY OF GUARANTEES  
1338

1339 In summary, the use of a COLA window with OLA stitching endows DiffPM’s reconstruction pro-  
1340 cess with several robust theoretical guarantees. The reconstruction is:

- 1341 1. *Lossless* when local predictions are noise-free (Theorem C.1).
- 1342 2. *Non-amplifying* with respect to bounded per-frame errors (Proposition C.2).
- 1343 3. *Variance-reducing* for independent, zero-mean local errors, with the reduction factor re-  
1344 lated to the degree of overlap (Proposition C.3).
- 1345 4. *Robust* to small violations of the COLA property, with well-defined bounds on the resulting  
1346 bias (Sec. C.5).
- 1347 5. *Extensible* to multivariate time series in a channel-wise manner, with analogous error prop-  
1348 agation properties (Sec. C.7).

1349

1350 These properties hold individually for the trend and residual branches and, by linearity, for their sum.  
 1351 Consequently, DiffPM achieves seamless global synthesis through a parallelized and mathematically  
 1352 principled overlap–add procedure, sidestepping the complexities of sequential generation models.  
 1353

## 1354 D ADDITIONAL RESULTS: AUGMENTATION ABLATIONS

### 1355 D.1 STOCK AUGMENTATION: PROTOCOL AND ANALYSIS

1356 In this ablation study we want to test whether augmenting the forecasting model’s training windows  
 1357 with *DiffPM* samples (trained only on the first 60% of the series) improves accuracy on the held–out  
 1358 20% test split, while holding everything else constant.

1359 **Evaluation protocol and hygiene.** We lock down the data pipeline to avoid leakage and to make  
 1360 A/B fully comparable.

- 1361 • *Fixed temporal split.* We use a strict 60/20/20% split (train/val/test) shared across all con-  
 1362 ditions (A and B@ $k$ ). Only the 60% train portion is used to train the generator.
- 1363 • *Normalization on real–train only.* Standardization parameters are fit on the *real* training  
 1364 split and then applied to synthetic, validation, and test.
- 1365 • *Causal windows.* Forecasting uses causal sliding windows (no look-ahead; no centered  
 1366 filters on prefixes).
- 1367 • *Forecaster parity.* Identical forecaster architecture, optimization, early stopping, and train-  
 1368 ing budget across A and B. Hyperparameters are tuned once on real train/val and then  
 1369 reused unchanged for all B@ $k$ .
- 1370 • *Quantity sweep.* B augments the real training windows with  $k \in \{1\times, 2\times, 5\times, 10\times\}$   
 1371 synthetic-to-real ratios (window counts), to show the effect is monotonic and to reveal  
 1372 saturation.
- 1373 • *Multiple seeds.* Both synthetic generation and forecaster training are repeated across mul-  
 1374 tiple seeds; we report mean $\pm$ std.
- 1375 • *Similarity checks.* For each synthetic training window we compute nearest–neighbor dis-  
 1376 tance to the *test* set (e.g., DTW/L2) and monitor 5th/50th/95th percentiles to ensure non-  
 1377 trivial distances; we observe no near–duplicates.

1378 DiffPM is trained on the first 60% of the Stock series and then used to sample additional sequences.  
 1379 Each synthetic series is post-processed *before* use: (i) mild rolling–mean smoothing, (ii) per-feature  
 1380 mean/variance rescaling to match the real training distribution, and (iii) optional clipping in normal-  
 1381 ized space. These steps align basic statistics while preserving dynamics; importantly, the scaler used  
 1382 for rescaling is the one fit on real–train.

1383 Table 2: **Ablation on Stock.** Temporal split is fixed at 60/20/20% (train/val/test) and shared across  
 1384 A/B; scalers are fit on *real-train* only and applied everywhere else. Numbers are mean $\pm$ std over  
 1385 seeds. *Gain* is A–B (positive is better).  
 1386

Setting	Test MSE $\downarrow$	Test MAE $\downarrow$	Gain (MSE) $\uparrow$	Gain (MAE) $\uparrow$
A (real-only)	11.2470 $\pm$ 0.3953	3.3167 $\pm$ 0.0560	–	–
B @ 1 $\times$ (real+syn)	5.2249 $\pm$ 0.1735	2.2391 $\pm$ 0.0380	6.0221	1.0776
B @ 2 $\times$ (real+syn)	3.5838 $\pm$ 0.5700	1.8354 $\pm$ 0.1632	7.6632	1.4812
B @ 5 $\times$ (real+syn)	2.9522 $\pm$ 0.1314	1.6712 $\pm$ 0.0361	8.2948	1.6454
B @ 10 $\times$ (real+syn)	2.5487 $\pm$ 0.1811	1.5453 $\pm$ 0.0543	8.6983	1.7714

1387 Table 2 (reproduced below for convenience) reports test MSE/MAE. Gains are defined as Gain =  
 1388 A – B (positive is better). Augmenting with DiffPM consistently improves both MSE and MAE and  
 1389 the effect grows with  $k$ :  
 1400

- 1401 • *Absolute gains.* From A to B@1 $\times$ : MSE gain 6.02, MAE gain 1.08. At B@10 $\times$ : MSE  
 1402 gain 8.70, MAE gain 1.77.

- *Relative reductions.* MSE drops by **53.5%**, **68.1%**, **73.8%**, and **77.3%** for  $k=\{1, 2, 5, 10\}$ , respectively; MAE drops by **32.5%**, **44.7%**, **49.6%**, and **53.4%**.
- *Effect size vs. variability.* Even at  $k=1\times$ , the MSE gain (6.02) is  $\approx 15\times$  the baseline run-to-run std (0.395), and the MAE gain (1.08) is  $\approx 19\times$  the baseline std (0.056). At  $k=10\times$  these ratios rise to  $\approx 22\times$  and  $\approx 32\times$ , respectively, indicating a large, stable effect across seeds.
- *Saturation.* Gains increase with  $k$  but with diminishing increments (MSE gain increases by  $\approx 1.64$  from  $1\times \rightarrow 2\times$ ,  $\approx 0.63$  from  $2\times \rightarrow 5\times$ , and  $\approx 0.40$  from  $5\times \rightarrow 10\times$ ), suggesting coverage saturation beyond  $\sim 5\times$ .
- *Outlier suppression.* The relative improvement is larger for MSE than MAE, consistent with synthetic augmentation reducing large-error tails and stabilizing high-variance regimes.

**Leakage sentinels and negative controls.** To stress-test the attribution, we run two quick checks: (i) training with real + *label-shuffled* synthetic windows (target misalignment) does not improve over A, and (ii) training with *misaligned residuals* (synthetic residuals paired with unrelated trends) similarly degrades versus B. These controls are consistent with the interpretation that the benefit comes from realistic, distribution-matched synthetic sequences that preserve coherent trend-residual structure.

Under strict hygiene controls, DiffPM augmentation yields large, stable gains for a data-hungry forecaster on **Stock**. Improvements scale with the amount of synthetic data and begin to saturate beyond  $\sim 5\times$ . The effect is strongest on MSE, indicating outlier/variance reduction. Together with similarity checks and negative controls, this supports the claim that high-fidelity, distribution-matched synthetics improve generalization without altering the downstream model.

## D.2 STOCK DIRECTION CLASSIFICATION: PROTOCOL AND ANALYSIS

Our goal here (much like the previous ablation study) is to assess whether augmenting the classifier’s training set with *DiffPM* samples (the generator is trained only on the first 60% of the series) improves *held-out* performance on the fixed 20% test split, under the same data hygiene and model parity constraints as the regression ablation.

Inputs are causal windows of length  $I$ ; the label is the *future*  $O$ -step log-return direction of the target variable, discretized into *Down/Flat/Up*. Thresholds for the ternary bins are fixed by the 40/60 percentiles computed on the *real training* return distribution and kept fixed for validation, test, and synthetic augmentation. Standardization parameters are fit on *real-train* input windows and applied unchanged to synthetic/validation/test (no leakage). The temporal split is fixed at 60/20/20% (train/val/test) and shared across all conditions.

We compare: **A** (real-only training) vs. **B** (real + DiffPM synthetic on training only; validation and test remain real-only). The classifier is the same data-hungry Transformer across A/B, trained with identical optimization, budget, and early stopping on validation macro-F1. Synthetic samples are generated by DiffPM trained on the 60% train split and passed through the same mild post-processing used in the regression ablation (smoothing, per-feature mean/variance matching based on real-train statistics, optional clipping).

Table 3 reports test metrics. For accuracy and macro-F1, we define *Gain* as  $B-A$  (higher is better); for loss we also show the absolute drop ( $A-B$ ).

Table 3: **Direction classification on Stock (test set).** Same split and scalers as the regression ablation; model and training budget are identical across A/B. Gains are  $B-A$  for Accuracy/F1 and  $A-B$  for Loss.

Setting	Loss ↓	Accuracy ↑	Macro-F1 ↑	Gain
A: real-only	1.2806	0.3431	0.1703	–
B: real+synthetic	0.8607	0.6144	0.4903	$\Delta\text{Acc} = +0.2712, \Delta\text{F1} = +0.3200$
Loss drop ( $A-B$ )			0.4199	

Table 4: **Equal-window generation time** (means  $\pm$  std over  $N=5$  seeds) for covering a full series using  $W$  windows of length  $L=24$  (hop  $\Delta=1$ ). Lower is better. DiffPM’s advantage is markedly larger on the *Energy* dataset (higher dimensional).

Dataset	DiffPM (s) $\downarrow$	Diffusion-TS (s) $\downarrow$	Speedup $\uparrow$
Stocks	$15.32 \pm 0.44$	$33.24 \pm 0.32$	<b>2.17</b> $\times$
Energy (higher- $D$ )	$26.58 \pm 2.08$	$180.12 \pm 3.22$	<b>6.78</b> $\times$

Augmenting the training set with DiffPM samples substantially improves generalization on the real test set: Accuracy increases by +27.12 points and macro-F1 by +0.3200, while test loss decreases by 0.42. The large macro-F1 gain indicates better balance across the three classes, not just a shift in the majority class. Training dynamics mirror this: the augmented model achieves higher validation F1 earlier and sustains it, whereas the real-only model exhibits shallow gains followed by early stopping with limited generalization. All improvements are obtained under fixed temporal splits, real-train-only normalization, causal windowing, and identical classifier/hyperparameters, isolating the effect to additional, distribution-matched training diversity supplied by DiffPM.

## E MORE EXPERIMENTS

### E.1 EQUAL-WINDOW GENERATION SPEED ( $L=24$ ): WHOLE-SERIES COVERAGE FAIRNESS

**Why this matters.** Some baselines generate *windows* rather than full sequences. In practice, producing a complete length- $T$  series requires generating exactly  $W = T-L+1$  windows (for hop  $\Delta=1$ ) and stitching them. This is the scenario users ultimately face and the one reviewers care about.

We normalize by the *number of windows required to cover the series*: each method generates exactly  $W$  windows at  $L=24$ ,  $\Delta=1$ , reverse steps = 200. We measure *end-to-end wall-clock* (from first reverse step to a stitched, de-normalized series) with identical software/hardware: A100 (40GB), PyTorch 2.3, CUDA 12.2; identical precision (AMP on); identical seeding and  $N=5$  runs; timing excludes disk I/O and includes stitching. Although DiffPM supports true full-sequence generation via parallel OLA, we still generate the same  $W$  windows internally to ensure parity.

(1) *Whole-series cost at equal coverage*: For the same number of windows, DiffPM is 2.17 $\times$  faster on *Stocks* and a striking 6.78 $\times$  faster on the higher-dimensional *Energy* dataset. (2) *High- $D$  advantage*: The speed gap widens with dimensionality, consistent with DiffPM’s parallel, windowed denoising plus lightweight OLA stitching, which scales gracefully across channels. (3) *Practical implication*: When practitioners must cover long series or many sensors (high  $D$ ), DiffPM’s throughput translates directly into lower latency and higher synthetic yield per GPU-hour.

Same window length ( $L=24$ ), same hop ( $\Delta=1$ ), same reverse steps (200), same precision (AMP), same device and software versions, identical seeds and measurement procedure; timings include end-to-end stitching and de-normalization.

## F LLM USAGE DISCLOSURE

We used large language models (LLMs) strictly for *editing support*, in compliance with the ICLR 2026 policy on LLM usage:

- **Writing polish only.** We used an LLM-based assistant (e.g., ChatGPT/GPT-class) to improve clarity, grammar, and style of author-written prose. Substantive content (ideas, claims, proofs, algorithms, experimental design, and conclusions) was authored by the paper’s authors.
- **Minor code cleanup only.** We used the assistant for light refactoring (e.g., formatting, docstrings, variable renaming, removing dead code) in non-core utilities. The core DiffPM implementation, experimental pipelines, metrics, and evaluation scripts were written and maintained by the authors.

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

- **No role in results or analysis.** LLMs were *not* used to design the method, tune hyperparameters, select or filter results, generate or modify datasets, run experiments, write proofs, or interpret findings.
- **Human verification.** All LLM-edited text and code were reviewed and approved by the authors; all claims are the authors' responsibility.
- **Privacy and reproducibility.** We did not share proprietary or sensitive data with the assistant. LLM usage does not affect reproducibility; all artifacts needed to reproduce results are provided in the code and supplementary material.