

# CTD-MTSI: CNN–TRANSFORMER DIFFUSION FOR PROBABILISTIC MULTIVARIATE TIME SERIES IMPUTATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent diffusion-based multivariate time series imputation (MTSI) approaches provide strong probabilistic modeling but rely on attention mechanisms for learning both temporal dynamics and cross-variable dependencies. This work introduces CTD-MTSI, a conditional CNN–transformer diffusion architecture that decouples temporal modeling and cross-variable interaction. Depthwise temporal convolutions explicitly capture local dynamics, while attention inter-variable dependencies. Contextual information is injected through mask-aware embeddings and feature-wise modulation during denoising. Experiments on benchmark datasets show that CTD-MTSI consistently outperforms the state-of-the-art diffusion baseline CSDI in both reconstruction accuracy and uncertainty quantification. These results highlight the effectiveness of structured convolutional inductive bias within conditional diffusion models for time series imputation.

**Track:** Research

## 1 INTRODUCTION AND RELATED WORK

Multivariate time series (MTS) data appear in many real-world systems such as healthcare monitoring, finance, traffic, and climate science (Yi et al., 2016; Tan et al., 2013; Nelwamondo et al., 2007; Hudak et al., 2008; van Buuren & Groothuis-Oudshoorn, 2011). Each variable (or sensor) follows its own temporal dynamics, while exhibiting inter-variable correlations. Modeling these temporal and cross-variable dependencies is critical for forecasting, decision making and understanding system dynamics. Multivariate time series imputation (MTSI) aims to recover missing values by exploiting both temporal continuity and relationships among variables.

Early deep learning approaches for MTSI focused on deterministic point estimation. Recurrent models propagate imputed values through bidirectional recurrent structures, capturing temporal dynamics (Che et al., 2016; Cao et al., 2018). Attention-based architectures further improved global dependency modeling by leveraging self-attention and have rapidly gained dominance in time series tasks due to their scalability and strong long-range modeling capacity (Suo et al., 2020; Bansal et al., 2021; Du et al., 2023). However, these methods produce single deterministic reconstructions, which may underestimate uncertainty. To address uncertainty estimation, variational approaches incorporate probabilistic inference (Mulyadi et al., 2022). Gaussian process-based methods (Bonilla et al., 2007; Fortuin et al., 2020) provide principled uncertainty modeling but scale poorly with high-dimensional data. Although these methods improve uncertainty calibration, modeling highly nonlinear dependencies remains challenging.

Diffusion models have recently become the state-of-the-art for probabilistic MTSI, as they formulate imputation as a reverse denoising process conditioned on observed entries. By learning the conditional score function, these models generate multiple plausible imputations and achieve strong performance. Architecturally, most of diffusion-based methods build upon a DiffWave-style residual backbone (Kong et al., 2021), adapted by introducing separate attention mechanisms applied sequentially along temporal and feature dimensions within the same residual diffusion block. However, conditioning is typically implemented via token concatenation, requiring attention layers to implicitly separate observed context from noisy targets. Also, noisy temporal representations may be propagated across variables before sufficient local denoising occurs. Convolutional architectures

054 have received comparatively less recognition than attention-based models in MTSI due to perceived  
 055 inferior modeling capability. Yet, convolution provides strong local inductive bias, motivating this  
 056 work, which investigates whether convolution can be more effectively integrated into probabilistic  
 057 imputation.

058 To this end, CTD-MTSI introduces a conditional CNN–transformer diffusion architecture for prob-  
 059 abilistic MTSI. Contextual information (observed entries) is injected directly into intermediate rep-  
 060 resentations via mask-aware patch embeddings and a feature-wise linear modulation mechanism. A  
 061 hierarchical backbone combines depthwise temporal convolutions with variable attention, provid-  
 062 ing multi-scale temporal modeling. Such structural decoupling simplifies optimization compared to  
 063 jointly modeling temporal and cross-variable dependencies within the same attention-based residual  
 064 block.

## 066 2 METHODOLOGY

068 We consider the problem of probabilistic multivariate time series imputation. Let  $\mathbf{X}_0 \in \mathbb{R}^{V \times L}$  de-  
 069 note a multivariate time series with  $V$  variables observed over  $L$  time steps. Due to missing entries,  
 070 the full signal is only partially observed. We define a binary observation mask  $\mathbf{M} \in \{0, 1\}^{V \times L}$ ,  
 071 where  $M_{v,\ell} = 1$  indicates that the value  $\mathbf{X}_{v,\ell}$  is observed and  $M_{v,\ell} = 0$  otherwise. The partially  
 072 observed data is then given by the element-wise product  $\mathbf{X}_0^{\text{ob}} = \mathbf{M} \odot \mathbf{X}$ , where  $\odot$  denotes Hadamard  
 073 (element-wise) multiplication.

### 075 2.1 CONDITIONAL DENOISING DIFFUSION PROBABILISTIC MODEL

076 To perform probabilistic imputation, we model the conditional distribution  $p(\mathbf{X}_0^{\text{mi}} | \mathbf{X}_0^{\text{ob}}, \mathbf{M})$ , using  
 077 a conditional denoising diffusion probabilistic model (Ho et al., 2020), which allows to generate  
 078 multiple plausible imputations and quantify uncertainty. The diffusion process is applied only to the  
 079 missing part, while  $\mathbf{X}_0^{\text{ob}}$  remains fixed and serves as contextual information throughout the reverse  
 080 process.

081 **Forward process.** We define a Markov chain that gradually transforms the missing values into  
 082 Gaussian noise:

$$084 q(\mathbf{X}_t^{\text{mi}} | \mathbf{X}_{t-1}^{\text{mi}}) = \mathcal{N}(\mathbf{X}_t^{\text{mi}}; \sqrt{\alpha_t} \mathbf{X}_{t-1}^{\text{mi}}, (1 - \alpha_t) \mathbf{I}), \quad (1)$$

085 where  $t$  is the noise level,  $\{\alpha_t\}_{t=1}^T$  is a predefined variance schedule and  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ . Then  
 086  $\mathbf{X}_t^{\text{mi}}$  can be expressed as  $\mathbf{X}_t^{\text{mi}} = \sqrt{\bar{\alpha}_t} \mathbf{X}_0^{\text{mi}} + \sqrt{1 - \bar{\alpha}_t} \epsilon$ , where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Thus, at an arbitrary  
 087 timestep  $t$ , the noisy missing variables  $\mathbf{X}_t^{\text{mi}}$  can be obtained directly from the clean  $\mathbf{X}_0^{\text{mi}}$  by injecting  
 088 Gaussian noise.

089 **Conditional reverse process.** The generative model is defined as a reverse-time Markov chain  
 090 conditioned on the observed context,

$$092 p_{\theta}(\mathbf{X}_{0:T}^{\text{mi}} | \mathbf{X}_0^{\text{ob}}, \mathbf{M}) = p(\mathbf{X}_T^{\text{mi}}) \prod_{t=1}^T p_{\theta}(\mathbf{X}_{t-1}^{\text{mi}} | \mathbf{X}_t^{\text{mi}}, \mathbf{X}_0^{\text{ob}}, \mathbf{M}), \quad \mathbf{X}_T^{\text{mi}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (2)$$

094 where  $\mathbf{X}_T^{\text{mi}}$  is initialized as pure Gaussian noise and progressively refined toward a plausible impu-  
 095 tation. Each reverse transition is modeled as a Gaussian distribution, defined as

$$097 p_{\theta}(\mathbf{X}_{t-1}^{\text{mi}} | \mathbf{X}_t^{\text{mi}}, \mathbf{X}_0^{\text{ob}}, \mathbf{M}) = \mathcal{N}(\mathbf{X}_{t-1}^{\text{mi}}; \mu_{\theta}(\mathbf{X}_t^{\text{mi}}, t | \mathbf{X}_0^{\text{ob}}, \mathbf{M}), \sigma_t^2 \mathbf{I}). \quad (3)$$

098 Following the standard DDPM parameterization, the mean is expressed through a conditional noise  
 099 predictor  $\epsilon_{\theta}(\mathbf{X}_t^{\text{mi}}, t | \mathbf{X}_0^{\text{ob}}, \mathbf{M})$ , which corresponds to predicting and removing the injected noise at  
 100 each diffusion step. The model is trained using the standard denoising objective. At inference time,  
 101 we iteratively sample from  $t = T$  down to  $t = 0$ , obtaining multiple plausible imputations  $\mathbf{X}_0^{\text{mi}}$  that  
 102 are consistent with  $\mathbf{X}_0^{\text{ob}}$  and respect the observed mask. Full derivations and sampling details are  
 103 provided in Appendix A.2.4.

### 105 2.2 NOISE PREDICTOR ARCHITECTURE

106 We parameterize the conditional noise predictor  $\epsilon_{\theta}(\mathbf{X}_t^{\text{mi}}, t | \mathbf{X}_0^{\text{ob}}, \mathbf{M})$  using a CNN-transformer  
 107 backbone. The design combines (i) convolutional patch embeddings and depthwise temporal mixing

with (ii) multi-head attention blocks to capture cross-variable dependencies. A schematic overview of the full backbone is provided in Appendix A.4.

**Patch and mask-aware embedding.** Both the observed context and the noisy missing target are mapped into a shared latent feature space using a convolutional patch stem. The stem consists of a 1D convolution applied along the temporal dimension with kernel size  $p$  (patch size) and stride  $s$  (patch stride). This operation aggregates short temporal windows into latent tokens while reducing the effective sequence length. To make the embedding mask-aware, the binary mask is concatenated as an additional input channel, allowing the network to distinguish observed from missing values. Let  $\mathcal{E}(\cdot)$  denote the stem embedding function. The embedded representations are computed as  $\mathbf{Z}^{\text{ob}} = \mathcal{E}(\mathbf{X}_0^{\text{ob}}, \mathbf{M})$  and  $\mathbf{Z}_t^{\text{mi}} = \mathcal{E}(\mathbf{X}_t^{\text{mi}}, \mathbf{1} - \mathbf{M})$ , where  $\mathbf{Z}^{\text{ob}}, \mathbf{Z}_t^{\text{mi}} \in \mathbb{R}^{B \times V \times H \times T'}$ ,  $B$  denotes the batch size,  $V$  the number of variables,  $H$  the number of latent feature channels (embedding dimension), and  $T' = \lfloor (L - p)/s \rfloor + 1$  the patchified temporal length.

**Fuse conditioning (FiLM).** To inject the observed context, a feature-wise linear modulation (FiLM) mechanism (Perez et al., 2018) is applied at the stem resolution. Per-variable and per-timestep modulation parameters are predicted from the context embedding  $\mathbf{Z}^{\text{ob}}$ , while a separate modulation term is derived from a timestep embedding of the diffusion step  $t$ . Specifically,  $(\gamma^c, \beta^c) = \text{PWConv}(\mathbf{Z}^{\text{ob}})$  and  $(\gamma^t, \beta^t) = \text{MLP}(\text{Emb}(t))$ , which are combined as  $\gamma = \gamma^c + \gamma^t$  and  $\beta = \beta^c + \beta^t$ . The fused representation is then obtained as

$$\mathbf{Z}_t = (1 + \gamma) \odot \mathbf{Z}_t^{\text{mi}} + \beta. \quad (4)$$

Here, PWConv denotes a  $1 \times 1$  convolution applied per variable without pooling, and  $(\gamma^t, \beta^t)$  are broadcast across variables and time. This modulation preserves the backbone feature structure while conditioning the denoiser without directly concatenating noisy and observed channels at the stem.

**Hierarchical CNN–transformer modules.** The fused representation  $\mathbf{Z}_t$  is processed by a hierarchy of backbone levels, each composed of multiple CNN–transformer blocks operating at a fixed temporal resolution. Between levels, temporal downsampling is performed using a strided 1D convolution, progressively enlarging the effective receptive field and enabling multi-scale temporal modeling.

Within each CNN–transformer block, (i) variable attention and (ii) a Conv-FFN for local refinement are alternated sequentially. Variable attention captures cross-variable dependencies, while temporal context is encoded using depthwise temporal convolutions inside the  $Q/K/V$  projections, following modern convolutional designs for efficient long-sequence modeling (Donghao & Xue, 2024; Park et al., 2026). For an intermediate representation  $\mathbf{H} \in \mathbb{R}^{B \times V \times H \times T'}$ , multi-head variable attention is computed as

$$\text{Attn}(\mathbf{H}) = \text{Proj} \left( \text{Softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{T'}} \right) \mathbf{V} \right), \quad \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{B \times H \times V \times T'}, \quad (5)$$

where  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$  are produced via shared depthwise temporal Conv1D mixers and Proj denotes a Conv  $1 \times 1$  projection. Each block is further conditioned on the diffusion timestep through a sinusoidal timestep embedding followed by an MLP, which generates adaptive normalization (Peebles & Xie, 2023) scale and gating parameters that modulate the residual updates.

**Prediction head.** Finally, the latent features from the last backbone level are mapped to the noise prediction through a reconstruction head. A pixel-shuffle style temporal upsampling is first applied to restore the original temporal resolution, followed by a linear projection to the output channel corresponding to the predicted noise. Consistent with the backbone blocks, the prediction head is also conditioned on the diffusion timestep via adaptive layer normalization (AdaLN), whose scale and shift parameters are generated from the timestep embedding.

## 3 EXPERIMENTS

### 3.1 SETUP

The proposed method is evaluated on two commonly used multivariate time series benchmark datasets: ETTh1 and ETTh2 (Zhou et al., 2021). As a baseline, we compare against the state-of-the-art diffusion-based imputation model CSDI (Tashiro et al., 2021). Training follows a self-supervised masking scheme where 30% of observed entries are randomly masked and treated as

162 targets for reconstruction. To ensure a fair comparison, CSDI is trained under the same settings, in-  
 163 cluding identical data splits and random seeds. Performance is measured using mean absolute error  
 164 (MAE), root mean squared error (RMSE), and continuous ranked probability score (CRPS). Further  
 165 implementation details are provided in Appendix A.

### 167 3.2 RESULTS

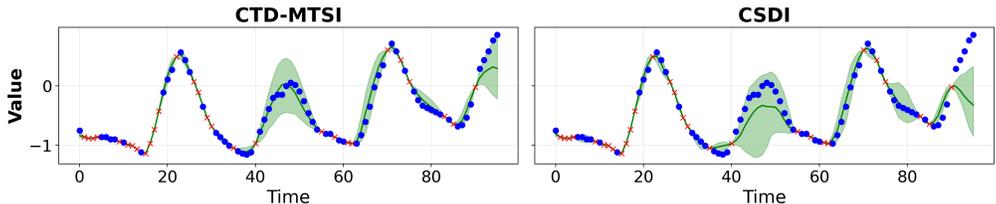
169 Table 1 reports imputation performance under varying point wise missing ratios. Across both  
 170 datasets, CTD-MTSI consistently matches or improves upon CSDI in all metrics under most miss-  
 171 ingness levels. Improvements are particularly consistent in the moderate regime, where temporal  
 172 continuity alone becomes insufficient and cross-variable structure has stronger influence.

173 Table 1: Imputation performance under point missingness. Best results are marked in **red**.

Dataset	Missing	CTD-MTSI (Ours)			CSDI		
		MAE ↓	RMSE ↓	CRPS ↓	MAE ↓	RMSE ↓	CRPS ↓
ETTh1	0.1	<b>0.2310</b>	<b>0.3614</b>	<b>0.0361</b>	0.2383	0.3996	0.0373
	0.3	<b>0.3079</b>	<b>0.6205</b>	<b>0.0473</b>	0.3211	0.6347	0.0495
	0.5	<b>0.4292</b>	<b>0.9079</b>	<b>0.0665</b>	0.4499	0.9496	0.0694
	0.7	0.6974	1.5555	0.1092	<b>0.6712</b>	<b>1.4393</b>	<b>0.1037</b>
ETTh2	0.1	<b>0.3184</b>	<b>0.5465</b>	<b>0.0139</b>	0.3239	0.5749	0.0142
	0.3	<b>0.4166</b>	<b>0.7183</b>	<b>0.0183</b>	0.4248	0.7340	0.0187
	0.5	<b>0.5577</b>	<b>0.9429</b>	<b>0.0246</b>	0.5844	0.9904	0.0258
	0.7	<b>0.8882</b>	<b>1.4597</b>	<b>0.0394</b>	0.8977	1.5010	0.0396

186 While both methods rely on conditional diffusion, CTD-MTSI exhibits greater performance and  
 187 robustness. The consistent gains under moderate missingness (30%–50%) regimes suggest that its  
 188 decoupled temporal–variable modeling provides more stable information propagation during denois-  
 189 ing. CRPS improvements further indicate improved probabilistic calibration. Lower CRPS values  
 190 suggest that CTD-MTSI not only reduces reconstruction error but also produces closer predictive  
 191 distributions.

192 Figure 1 provides a qualitative visualization under 70% missingness on ETTh2. The predicted me-  
 193 dian closely tracks ground truth trajectories, while the uncertainty bands widen appropriately in  
 194 highly ambiguous regions. This behavior supports the quantitative findings and illustrates stable  
 195 uncertainty estimation even under severe sparsity.



204 Figure 1: Visualization of probabilistic imputation results on ETTh2 with 70% random missing-  
 205 ness. Median and 95% CIs shown in green; observed points in red; targets in blue.

## 208 4 CONCLUSION

210 This work introduces CTD-MTSI, a conditional CNN–transformer diffusion architecture for prob-  
 211 abilistic multivariate time series imputation. By structurally decoupling temporal modeling from  
 212 cross-variable interaction, the proposed design provides stable denoising behavior under varying  
 213 levels of sparsity. Experimental results demonstrate competitive and often superior performance  
 214 compared to the state-of-the-art diffusion-based baseline. These findings suggest that convolution,  
 215 when carefully integrated within diffusion models, remains highly effective for time series imputa-  
 tion.

## 216 REFERENCES

- 217  
218 Parikshit Bansal, Prathamesh Deshpande, and Sunita Sarawagi. Missing value imputation on multi-  
219 dimensional time series. *Proc. VLDB Endow.*, 14(11):2533–2545, 2021. doi: 10.14778/3476249.  
220 3476300.
- 221 Edwin V Bonilla, Kian Chai, and Christopher Williams. Multi-task gaussian process prediction. In  
222 *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007.  
223
- 224 Wei Cao, Dong Wang, Jian Li, Hao Zhou, Yitan Li, and Lei Li. Brits: bidirectional recurrent im-  
225 putation for time series. In *Proceedings of the 32nd International Conference on Neural Information*  
226 *Processing Systems*, pp. 6776–6786, 2018.
- 227 Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent  
228 neural networks for multivariate time series with missing values, 2016.
- 229 Luo Donghao and Wang Xue. ModernTCN: A modern pure convolution structure for general time  
230 series analysis. In *The Twelfth International Conference on Learning Representations*, 2024.
- 231 Wenjie Du, David Côté, and Yan Liu. Saits: Self-attention-based imputation for time series. *Expert*  
232 *Systems with Applications*, 219:119619, 2023. ISSN 0957-4174. doi: [https://doi.org/10.1016/j.](https://doi.org/10.1016/j.eswa.2023.119619)  
233 [eswa.2023.119619](https://doi.org/10.1016/j.eswa.2023.119619).
- 234 Vincent Fortuin, Dmitry Baranchuk, Gunnar Raetsch, and Stephan Mandt. Gp-vae: Deep proba-  
235 bilistic time series imputation. In *Proceedings of the Twenty Third International Conference on*  
236 *Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*,  
237 pp. 1651–1661. PMLR, 2020.
- 238 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances*  
239 *in Neural Information Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc.,  
240 2020.
- 241 Andrew T. Hudak, Nicholas L. Crookston, Jeffrey S. Evans, David E. Hall, and Michael J.  
242 Falkowski. Nearest neighbor imputation of species-level, plot-scale forest structure attributes  
243 from lidar data. *Remote Sensing of Environment*, 112(5):2232–2245, 2008. doi: [https://doi.org/](https://doi.org/10.1016/j.rse.2007.10.009)  
244 [10.1016/j.rse.2007.10.009](https://doi.org/10.1016/j.rse.2007.10.009).
- 245 Zhifeng Kong, Wei Ping, Jiayi Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile  
246 diffusion model for audio synthesis. In *International Conference on Learning Representations*,  
247 2021.
- 248 Ahmad Wisnu Mulyadi, Eunji Jun, and Heung-Il Suk. Uncertainty-aware variational-recurrent im-  
249 putation network for clinical time series. *IEEE Transactions on Cybernetics*, 52(9):9684–9694,  
250 2022. doi: 10.1109/TCYB.2021.3053599.
- 251 Fulufhelo V. Nelwamondo, Shakir Mohamed, and Tshilidzi Marwala. Missing data: A comparison  
252 of neural network and expectation maximization techniques. *Current Science*, pp. 1514–1521,  
253 2007.
- 254 Dongik Park, Hyunwoo Ryu, Suahn Bae, Keondo Park, and Hyung-Sin Kim. T1: One-to-one  
255 channel-head binding for multivariate time-series imputation. In *The Fourteenth International*  
256 *Conference on Learning Representations*, 2026.
- 257 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
258 *the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4195–4205, October  
259 2023.
- 260 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron Courville. Film: visual  
261 reasoning with a general conditioning layer. In *Proceedings of the Thirty-Second AAAI Confer-*  
262 *ence on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Con-*  
263 *ference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, 2018.

- 270 Qiuling Suo, Weida Zhong, Guangxu Xun, Jianhui Sun, Changyou Chen, and Aidong Zhang. Glima:  
271 Global and local time series imputation with multi-directional attention learning. In *2020 IEEE In-*  
272 *ternational Conference on Big Data (Big Data)*, pp. 798–807, 2020. doi: 10.1109/BigData50022.  
273 2020.9378408.
- 274
- 275 Huachun Tan, Guangdong Feng, Jianshuai Feng, Wuhong Wang, Yu-Jin Zhang, and Feng Li. A  
276 tensor-based method for missing traffic data completion. *Transportation Research Part C: Emerg-*  
277 *ing Technologies*, 28:15–27, 2013. doi: <https://doi.org/10.1016/j.trc.2012.12.007>.
- 278
- 279 Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. Csd: conditional score-based  
280 diffusion models for probabilistic time series imputation. In *Proceedings of the 35th International*  
281 *Conference on Neural Information Processing Systems*, 2021.
- 282
- 283 Stef van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equa-  
284 tions in r. *Journal of Statistical Software*, 45(3):1–67, 2011. doi: 10.18637/jss.v045.i03.
- 285
- 286 Xiuwen Yi, Yu Zheng, Junbo Zhang, and Tianrui Li. St-mvl: filling missing values in geo-sensory  
287 time series data. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial*  
288 *Intelligence*, pp. 2704–2710, 2016.
- 289
- 290 Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai  
291 Zhang. Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting.  
292 *Proceedings of the AAAI Conference on Artificial Intelligence*, 35:11106–11115, 2021. doi:  
293 10.1609/aaai.v35i12.17325.

## 294 A IMPLEMENTATION DETAILS

### 295 A.1 DATASET DETAILS

#### 296 A.1.1 DATASET DESCRIPTION

297 Experiments are conducted on two multivariate time series datasets from the Electricity Transformer  
298 Temperature (ETT) benchmark (Zhou et al., 2021). The ETT datasets (ETTh1 and ETTh2) contain  
299 electricity transformer measurements recorded at an hourly frequency. Each dataset consists of 7  
300 correlated variables representing oil temperature and load-related features. Table 2 summarizes the  
301 main statistics of the datasets.

302 Table 2: Dataset descriptions.

303 Dataset	304 Variables ( $V$ )	305 Train	306 Valid	307 Test	308 Frequency
309 ETTh1, ETTh2	310 7	311 8,545	312 2,785	313 2,785	314 Hourly

### 315 A.2 EXPERIMENT DETAILS

#### 316 A.2.1 MODEL CONFIGURATION

317 All experiments use a fixed input sequence length of 96. A one-dimensional convolutional embed-  
318 ding layer with kernel size 2 and stride 1 projects the input into a 128-dimensional channel space.  
319 The architecture consists of four CNN-transformer blocks organized into two hierarchical stages.  
320 The first stage contains two blocks using dual-scale depthwise convolutions with kernel sizes 71  
321 and 5, followed by a downsampling layer with kernel size 2 and stride 2. The second stage con-  
322 tains two additional CNN-transformer blocks operating on the downsampled representation, with  
323 kernel sizes 31 and 5. This hierarchical design enables multi-scale temporal modeling, capturing  
both long-range and short-range dependencies. The feed-forward expansion ratio is set to 1.0. The  
configuration remains identical across both datasets.

### 324 A.2.2 EXPERIMENTAL DESIGN

325  
326 All experiments are conducted using three random seeds: 2025, 2026, and 2027. Training is per-  
327 formed on a NVIDIA RTX 6000 GPU. To evaluate generalization under different observability lev-  
328 els, masks are simulated following random point wise missingness. During training, a fixed masking  
329 ratio of 30% is applied using point wise random masking. The trained model is then evaluated under  
330 different test masking ratios (0.1, 0.3, 0.5, 0.7). This setup evaluates whether a model trained under  
331 a single missing ratio can generalize to different levels of observability without retraining.

### 332 A.2.3 EVALUATION METRICS

333  
334 Let  $V$  denote the number of variables and  $L$  the sequence length. Let  $\mathcal{M}$  denote the set of artifi-  
335 cially masked positions used for evaluation. Ground-truth values are denoted by  $y_{v,\ell}$  and imputed  
336 values by  $\hat{x}_{v,\ell}$ , where  $v \in \{1, \dots, V\}$  and  $\ell \in \{1, \dots, L\}$ . Mean absolute error (MAE), root mean  
337 squared error (RMSE), and continuous ranked probability score (CRPS) are employed. All metrics  
338 are computed only on artificially masked positions to ensure consistent and fair evaluation across  
339 models.

340  
341 **MAE** The mean absolute error is defined as

$$342 \text{MAE} = \frac{1}{|\mathcal{M}|} \sum_{(v,\ell) \in \mathcal{M}} |\hat{x}_{v,\ell} - y_{v,\ell}|. \quad (6)$$

343  
344  
345  
346 MAE measures the average magnitude of the imputation error. Each deviation contributes linearly  
347 to the final score, making MAE directly interpretable as the typical reconstruction error.

348  
349 **RMSE** The root mean squared error is defined as

$$350 \text{RMSE} = \sqrt{\frac{1}{|\mathcal{M}|} \sum_{(v,\ell) \in \mathcal{M}} (\hat{x}_{v,\ell} - y_{v,\ell})^2}. \quad (7)$$

351  
352  
353  
354 RMSE computes the square root of the average squared error. By squaring deviations before aver-  
355 aging, larger errors are penalized more heavily. As a result, RMSE is more sensitive to outliers.

356  
357 **CRPS** For probabilistic predictions with cumulative distribution function  $F_{v,\ell}$  and ground truth  
358  $y_{v,\ell}$ , CRPS is defined as

$$359 \text{CRPS} = \frac{1}{|\mathcal{M}|} \sum_{(v,\ell) \in \mathcal{M}} \int_{-\infty}^{\infty} (F_{v,\ell}(z) - \mathbb{I}(z \geq y_{v,\ell}))^2 dz. \quad (8)$$

360  
361  
362  
363  
364 CRPS evaluates the entire predictive distribution rather than a single point estimate. It measures  
365 how close the predicted cumulative distribution is to the empirical distribution concentrated at the  
366 ground-truth value, thus capturing both accuracy and uncertainty calibration.

### 367 A.2.4 TRAINING PROCEDURE

368  
369 Following the conditional diffusion formulation described in Section 2.1. Let  $X_0 \in \mathbb{R}^{V \times L}$  denote  
370 a multivariate time series with  $V$  variables and sequence length  $L$ . Define the binary observation  
371 mask  $M \in \{0, 1\}^{V \times L}$ , where  $M_{v,\ell} = 1$  indicates observed values and  $M_{v,\ell} = 0$  indicates missing  
372 values. The observed and missing components are denoted as

$$373 X_0^{\text{ob}} = M \odot X_0, \quad X_0^{\text{mi}} = (1 - M) \odot X_0.$$

374  
375  
376 During training, a self-supervised masking strategy is adopted: a subset of originally observed en-  
377 tries is randomly masked and treated as missing targets. The diffusion process is applied only to the  
missing part  $X_0^{\text{mi}}$ , while  $X_0^{\text{ob}}$  remains fixed and serves as conditioning context.

At diffusion step  $t$ , Gaussian noise  $\epsilon \sim \mathcal{N}(0, I)$  is injected into the missing component

$$X_t^{\text{mi}} = \sqrt{\bar{\alpha}_t} X_0^{\text{mi}} + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad (9)$$

where  $\{\alpha_t\}_{t=1}^T$  is a predefined variance schedule and  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ . The model is trained to predict the injected noise using a conditional noise predictor  $\epsilon_\theta(X_t^{\text{mi}}, t | X_0^{\text{ob}}, M)$ . The denoising objective is minimized. Since diffusion is applied only to the missing component, the loss is computed exclusively on the missing positions

$$\mathcal{L}(\theta) = \mathbb{E} \left[ \frac{1}{\sum_{v,\ell} (1 - M_{v,\ell})} \sum_{v=1}^V \sum_{\ell=1}^L (1 - M_{v,\ell}) (\epsilon_{v,\ell} - \epsilon_\theta(\cdot)_{v,\ell})^2 \right]. \quad (10)$$

The self-supervised training procedure of CTD-MTISI is summarized in Algorithm 1.

---

**Algorithm 1** Training of CTD-MTISI
 

---

- 1: **Input:** training data distribution  $q(X_0)$ , number of iterations  $N_{\text{iter}}$ , noise schedule  $\{\alpha_t\}_{t=1}^T$  with  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$
- 2: **Output:** trained conditional noise predictor  $\epsilon_\theta$
- 3: **for**  $i = 1$  **to**  $N_{\text{iter}}$  **do**
- 4:   Sample  $t \sim \text{Uniform}(\{1, \dots, T\})$ ,  $X_0 \sim q(X_0)$
- 5:   Sample point-wise random masking on observed entries by updating  $M$  (e.g., set 30% of ones in  $M$  to zero)
- 6:   Compute conditional context and missing targets:

$$X_0^{\text{ob}} \leftarrow M \odot X_0, \quad X_0^{\text{mi}} \leftarrow (1 - M) \odot X_0$$

- 7:   Sample Gaussian noise  $\epsilon \sim \mathcal{N}(0, I)$  with the same shape as  $X_0^{\text{mi}}$
- 8:   Apply forward noising to the missing part:  $X_t^{\text{mi}} \leftarrow \sqrt{\bar{\alpha}_t} X_0^{\text{mi}} + \sqrt{1 - \bar{\alpha}_t} \epsilon$
- 9:   Predict noise using CTD-MTISI:  $\hat{\epsilon} \leftarrow \epsilon_\theta(X_t^{\text{mi}}, t | X_0^{\text{ob}}, M)$
- 10:   Take a gradient step minimizing the denoising objective (computed on missing entries):

$$\nabla_\theta \left( \frac{1}{\sum_{v,\ell} (1 - M_{v,\ell})} \sum_{v=1}^V \sum_{\ell=1}^L (1 - M_{v,\ell}) (\epsilon_{v,\ell} - \hat{\epsilon}_{v,\ell})^2 \right)$$

11: **end for**

---

The model is optimized using Adam with an initial learning rate of  $10^{-3}$  and a batch size of 16. Training runs for 200 epochs. A multi-step learning rate scheduler is employed, reducing the learning rate by 90% at 75% and 90% of the total training epochs.

---

**Algorithm 2** Imputation (Conditional Reverse Sampling) with CTD-MTISI
 

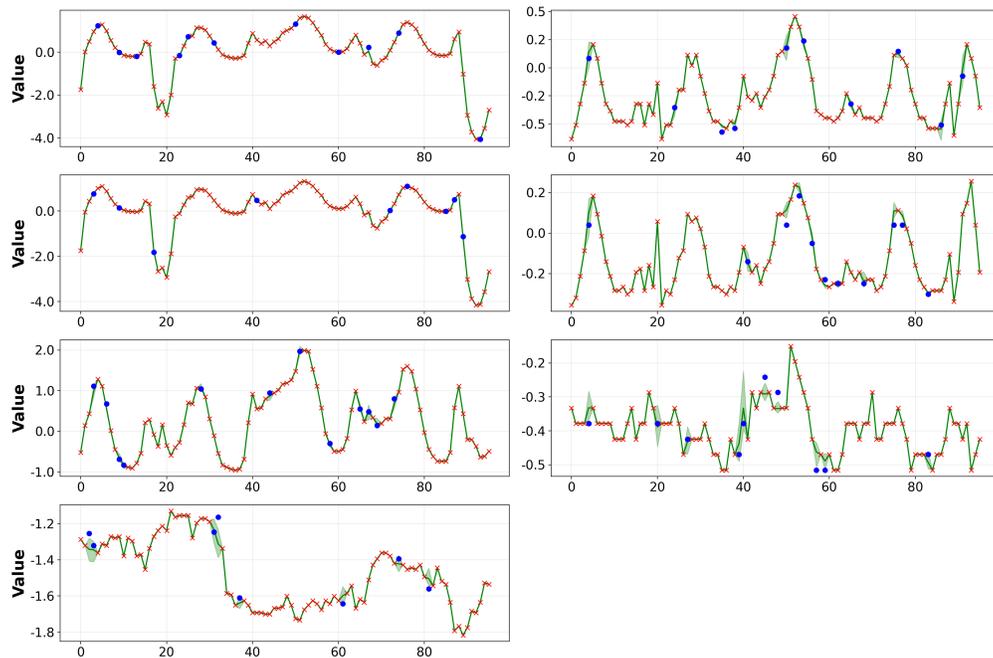
---

- 1: **Input:** partially observed sample  $X_0$ , mask  $M \in \{0, 1\}^{V \times L}$ , trained  $\epsilon_\theta$
  - 2: **Output:** imputed missing values  $\hat{X}_0^{\text{mi}}$
  - 3: Compute observed context:  $X_0^{\text{ob}} \leftarrow M \odot X_0$
  - 4: Initialize missing part with Gaussian noise:  $X_T^{\text{mi}} \sim \mathcal{N}(0, I)$
  - 5: **for**  $t = T$  **down to** 1 **do**
  - 6:   Predict noise:  $\hat{\epsilon} \leftarrow \epsilon_\theta(X_t^{\text{mi}}, t | X_0^{\text{ob}}, M)$
  - 7:   Compute DDPM mean:  $\mu_\theta \leftarrow \frac{1}{\sqrt{\alpha_t}} \left( X_t^{\text{mi}} - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon} \right)$
  - 8:   Sample reverse step:  $z \sim \mathcal{N}(0, I)$ ,  $X_{t-1}^{\text{mi}} \leftarrow \mu_\theta + \sigma_t z$
  - 9: **end for**
  - 10: Set  $\hat{X}_0^{\text{mi}} \leftarrow X_0^{\text{mi}}$  and return the full imputed series:  $\hat{X}_0 \leftarrow X_0^{\text{ob}} + \hat{X}_0^{\text{mi}}$
- 

### A.3 ADDITIONAL QUALITATIVE RESULTS

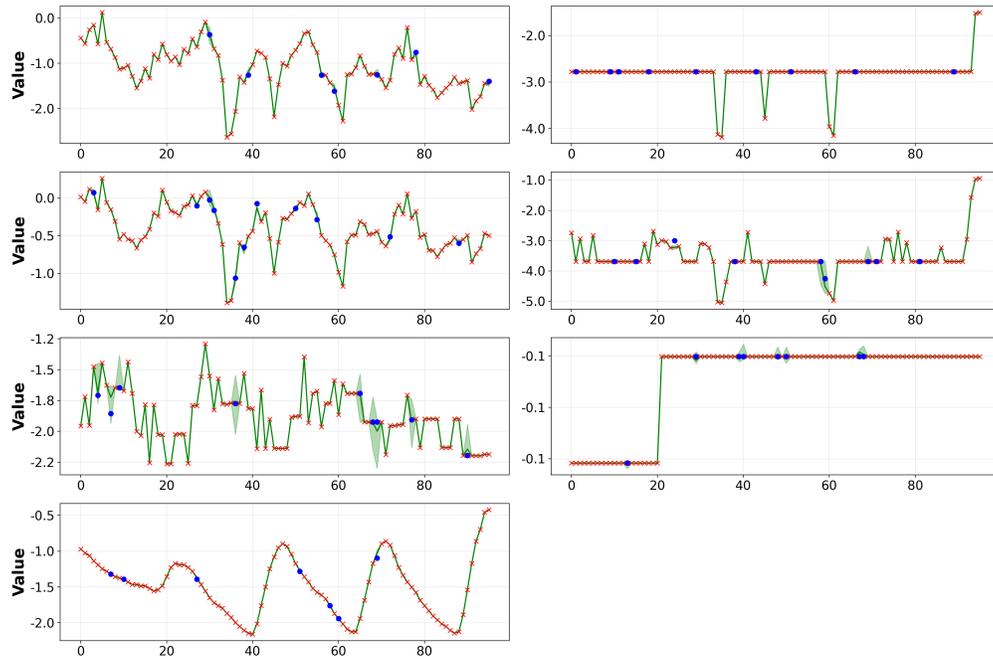
We provide additional qualitative imputation visualizations under varying point wise missing ratios (10%, 30%, 50%, and 70%) for both ETTh1 and ETTh2. These figures illustrate median predictions

432 and 95% confidence intervals, highlighting reconstruction accuracy and uncertainty calibration un-  
433 der increasing sparsity.  
434



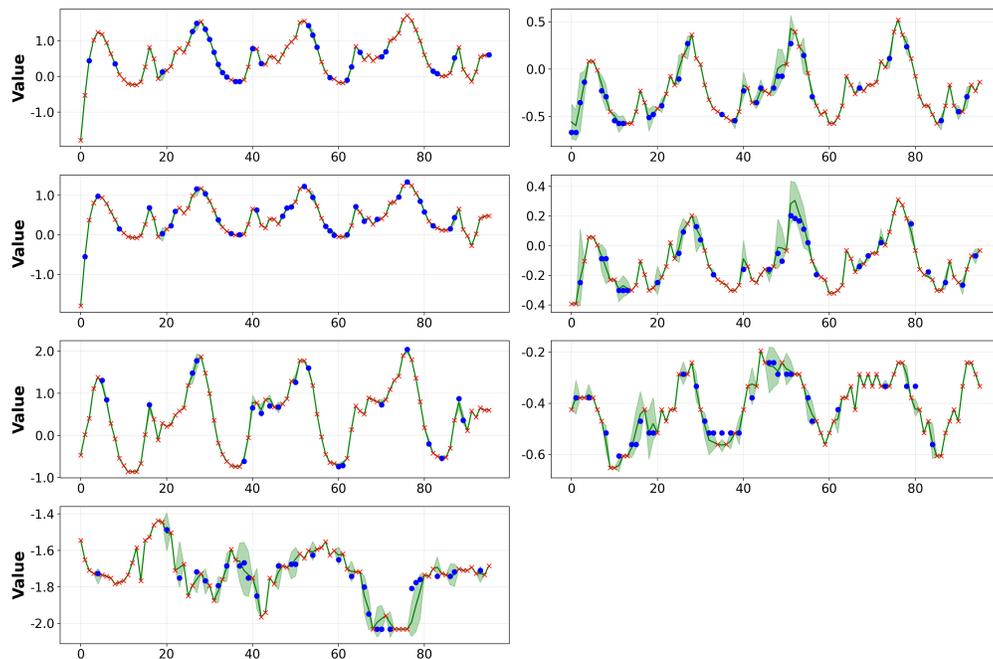
457 Figure 2: Visualization of probabilistic imputations results on ETTh1 (10% missingness). The  
458 results is for a time series sample with all 7 features. The median and 95% CIs are shown in green,  
459 observed points in red and targets in blue.  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508



509 Figure 3: Visualization of probabilistic imputations results on ETTh2 (10% missingness). The  
510 results is for a time series sample with all 7 features. The median and 95% CIs are shown in green,  
511 observed points in red and targets in blue.

512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535



536 Figure 4: Visualization of probabilistic imputations results on ETTh1 (30% missingness). The  
537 results is for a time series sample with all 7 features. The median and 95% CIs are shown in green,  
538 observed points in red and targets in blue.

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

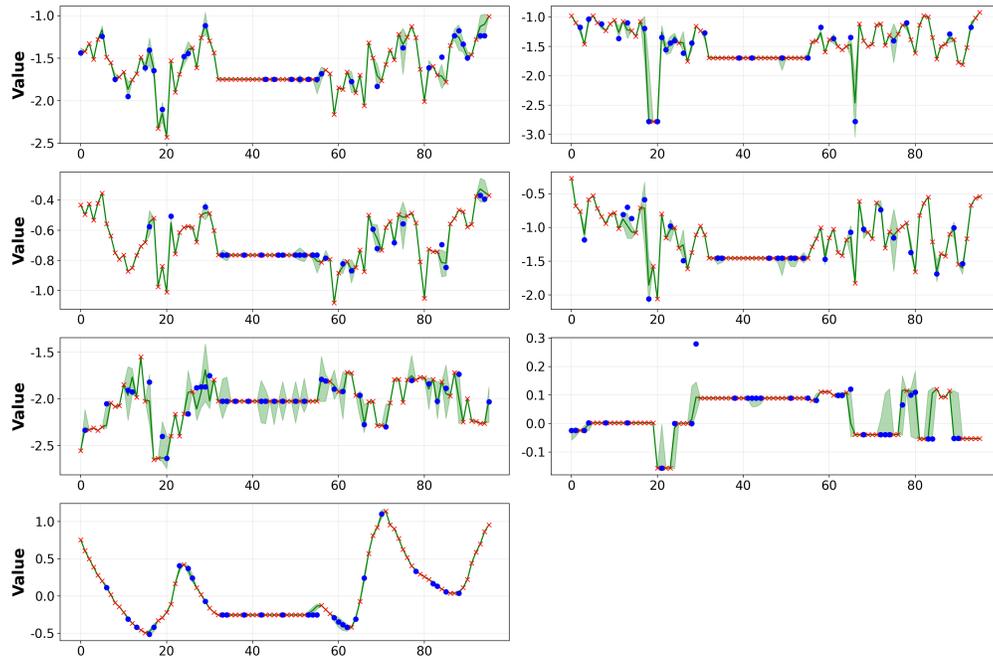


Figure 5: Visualization of probabilistic imputations results on ETTh2 (30% missingness). The results is for a time series sample with all 7 features. The median and 95% CIs are shown in green, observed points in red and targets in blue.

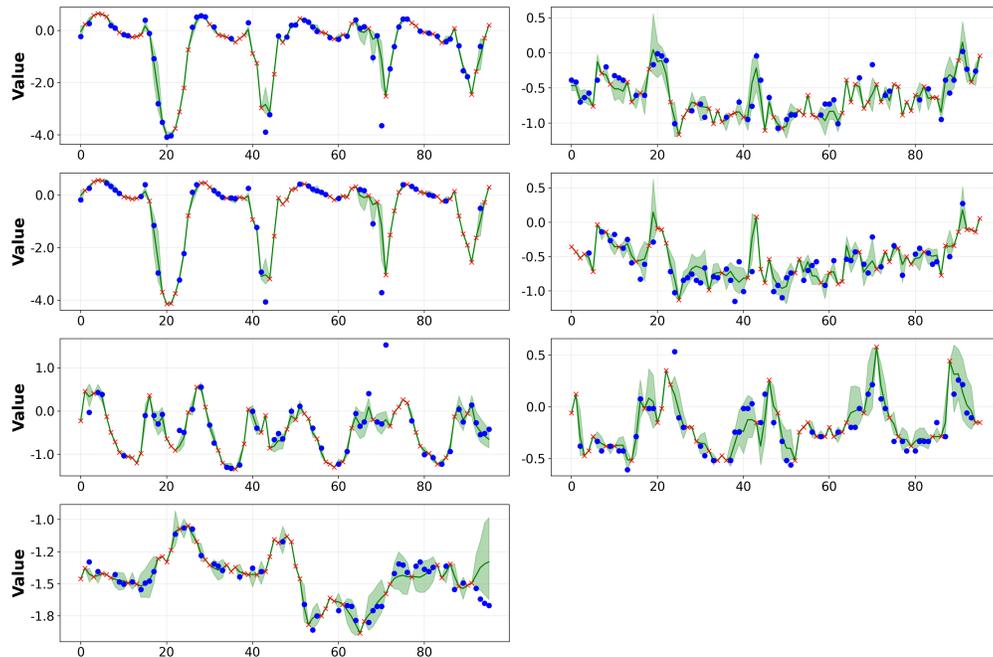


Figure 6: Visualization of probabilistic imputations results on ETTh1 (50% missingness). The results is for a time series sample with all 7 features. The median and 95% CIs are shown in green, observed points in red and targets in blue.

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

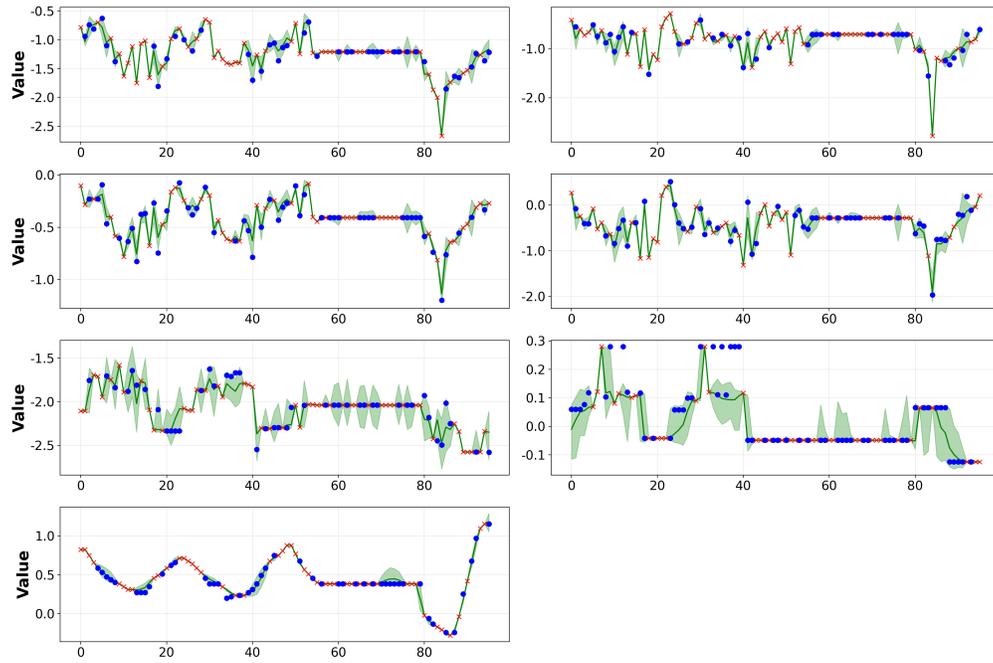


Figure 7: Visualization of probabilistic imputations results on ETTh2 (50% missingness). The results is for a time series sample with all 7 features. The median and 95% CIs are shown in green, observed points in red and targets in blue.

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

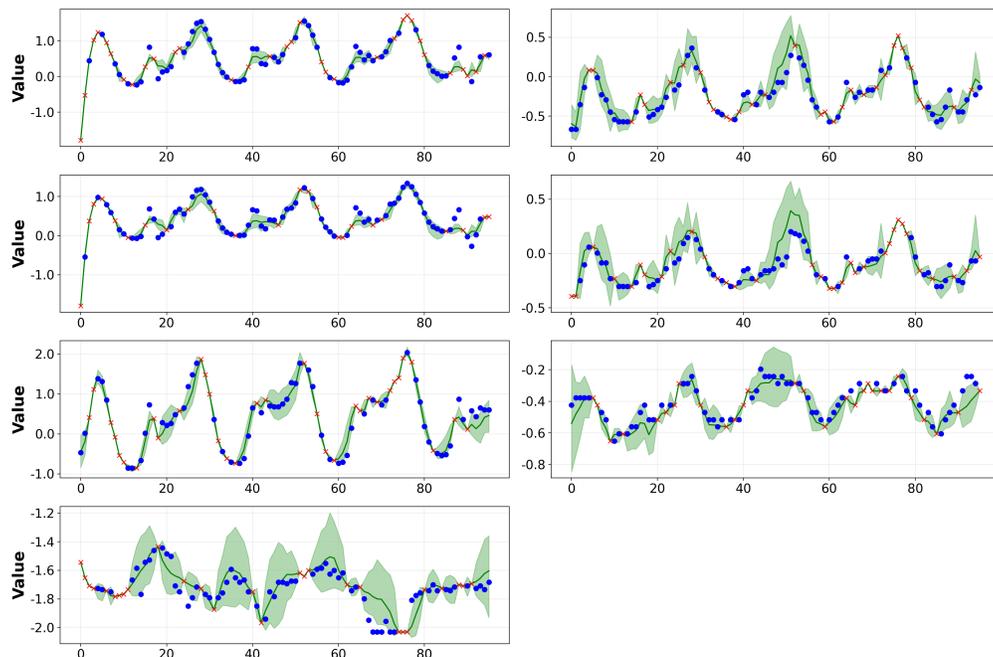


Figure 8: Visualization of probabilistic imputations results on ETTh1 (70% missingness). The results is for a time series sample with all 7 features. The median and 95% CIs are shown in green, observed points in red and targets in blue.

648  
 649  
 650  
 651  
 652  
 653  
 654  
 655  
 656  
 657  
 658  
 659  
 660  
 661  
 662  
 663  
 664  
 665  
 666  
 667  
 668  
 669  
 670  
 671  
 672  
 673  
 674  
 675  
 676  
 677  
 678  
 679  
 680  
 681  
 682  
 683  
 684  
 685  
 686  
 687  
 688  
 689  
 690  
 691  
 692  
 693  
 694  
 695  
 696  
 697  
 698  
 699  
 700  
 701

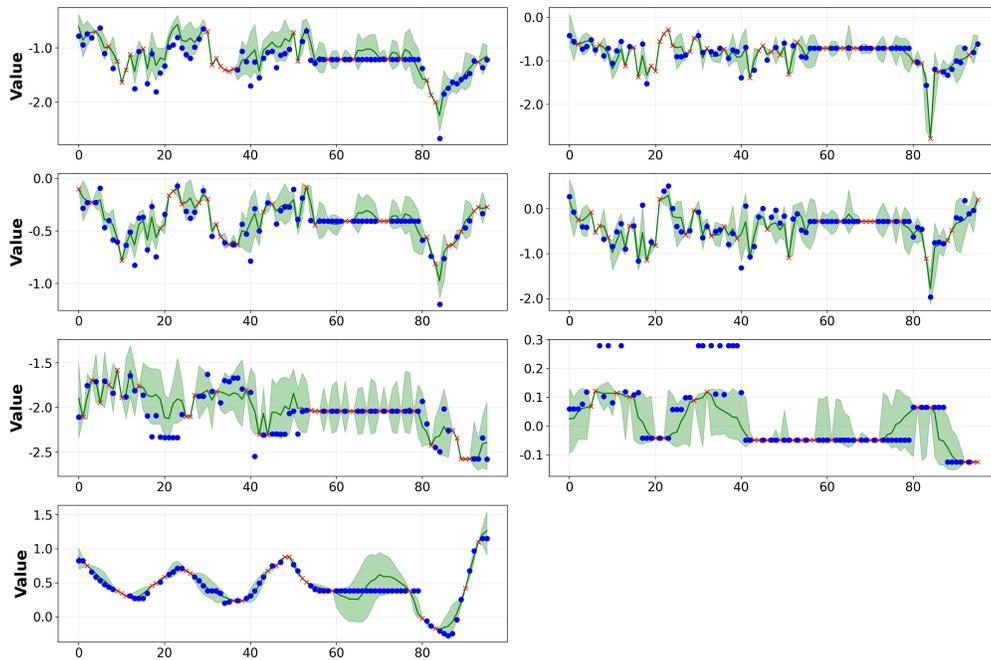


Figure 9: Visualization of probabilistic imputations results on ETTh2 (70% missingness). The results is for a time series sample with all 7 features. The median and 95% CIs are shown in green, observed points in red and targets in blue.

## A.4 NOISE PREDICTOR ARCHITECTURE DIAGRAM

Figure 10 illustrates the overall architecture of the CTD-MTSI noise predictor. The model receives the noisy data  $X_t^{\text{mi}}$ , the observed data  $X_0^{\text{ob}}$ , and the diffusion timestep  $t$  as inputs. The backbone consists of (i) a mask-aware convolutional patch embedding stem, (ii) FiLM-based conditioning, (iii) hierarchical CNN–transformer modules combining depthwise temporal convolutions and variable attention, and (iv) a reconstruction head that outputs the predicted noise  $\epsilon_\theta$ .

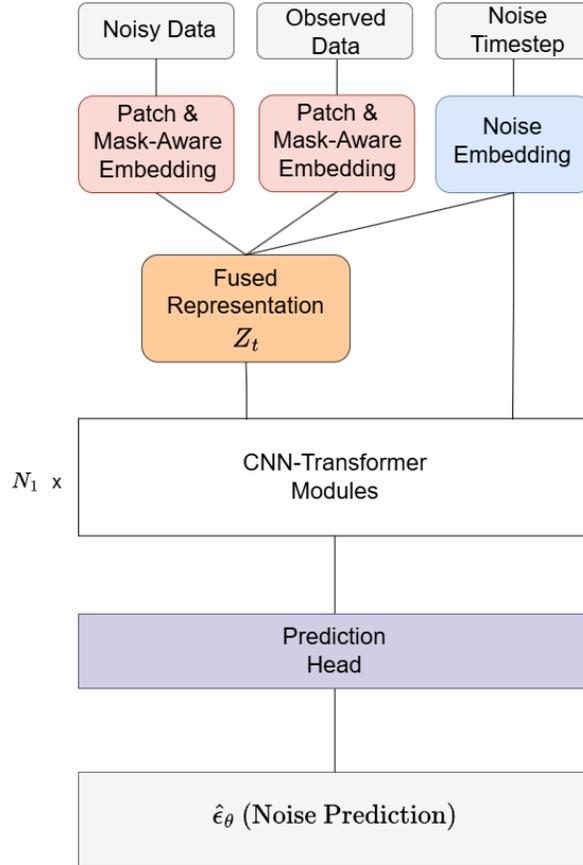


Figure 10: Diagram of the CTD-MTSI noise predictor backbone.

To provide further detail, Figure 11 illustrates the internal structure of a single CNN–transformer module used in the hierarchical backbone, combine with temporal downsampling. The module alternates depthwise temporal mixing and variable attention, followed by a convolutional feed-forward network, with noise-conditioned normalization applied at each residual branch.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

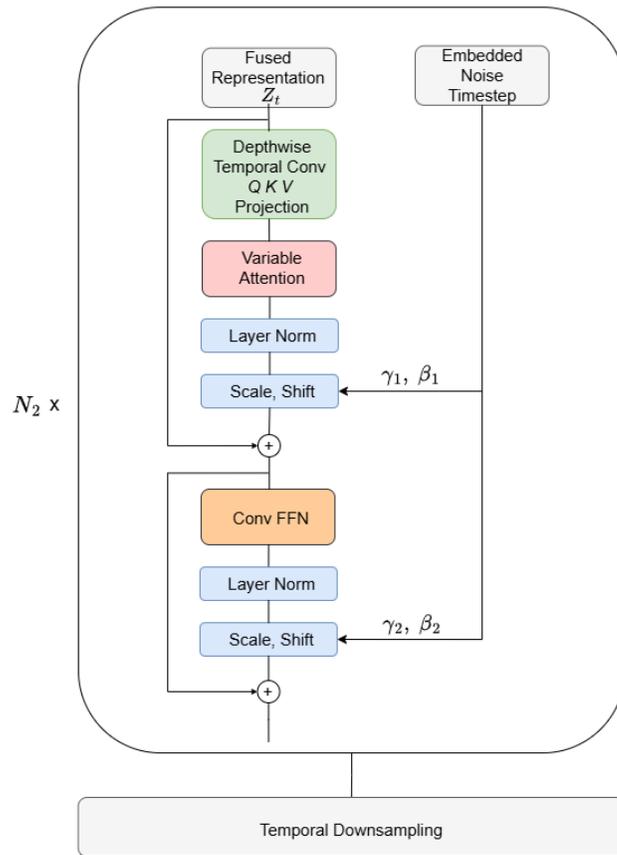


Figure 11: Diagram of the CNN-transformer module.