

Abstract

TL;DR: Conditional information bottleneck (CIB) enables mitigating strict regularization of information bottleneck while preserving temporal dynamics.

Clinical time series imputation presents a significant challenge because it requires capturing the underlying temporal dynamics from partially observed time series data input. However, direct application of **information bottleneck (IB)** framework to time series data without considering temporal context can lead to a substantial loss of temporal dependencies, which, in turn, can degrade the overall imputation performance and further clinical decisions.

To address such a challenge, we propose a novel **conditional information bottleneck (CIB)** approach for time series imputation. Variational decomposition of CIB motivates us to develop a novel deep learning method that can approximately achieve the proposed CIB objective for time series imputation as a combination of evidence lower bound and novel temporal kernel-enhanced contrastive optimization.

Information Bottleneck (IB) on Imputation

Definition 1. (Imputation) Let \mathbf{X}^o and \mathbf{X}^m be random variables for the partially observed features and missing features of \mathbf{X} , respectively, such that $\mathbf{X} = \mathbf{X}^o \cup \mathbf{X}^m$. Then, we define imputation as an unsupervised IB as follows:

$$\min_{\phi, \theta} I_\phi(\mathbf{Z}; \mathbf{X}^o) - \beta I_\theta(\mathbf{Z}; \mathbf{X}) \quad (1)$$

where $\beta \in \mathbb{R}_+$ is a Lagrangian multiplier, and ϕ and θ correspond to learnable parameters that define probabilistic mappings $q_\phi(\mathbf{Z}|\mathbf{X}^o)$ and $q_\theta(\mathbf{X}|\mathbf{Z})$, respectively.

Goal: finding the distribution of latent representation \mathbf{Z} and the corresponding parameters that preserves the core information for accurately reconstructing the (complete) original input \mathbf{X} while suppressing redundant information from its incomplete observation, \mathbf{X}^o

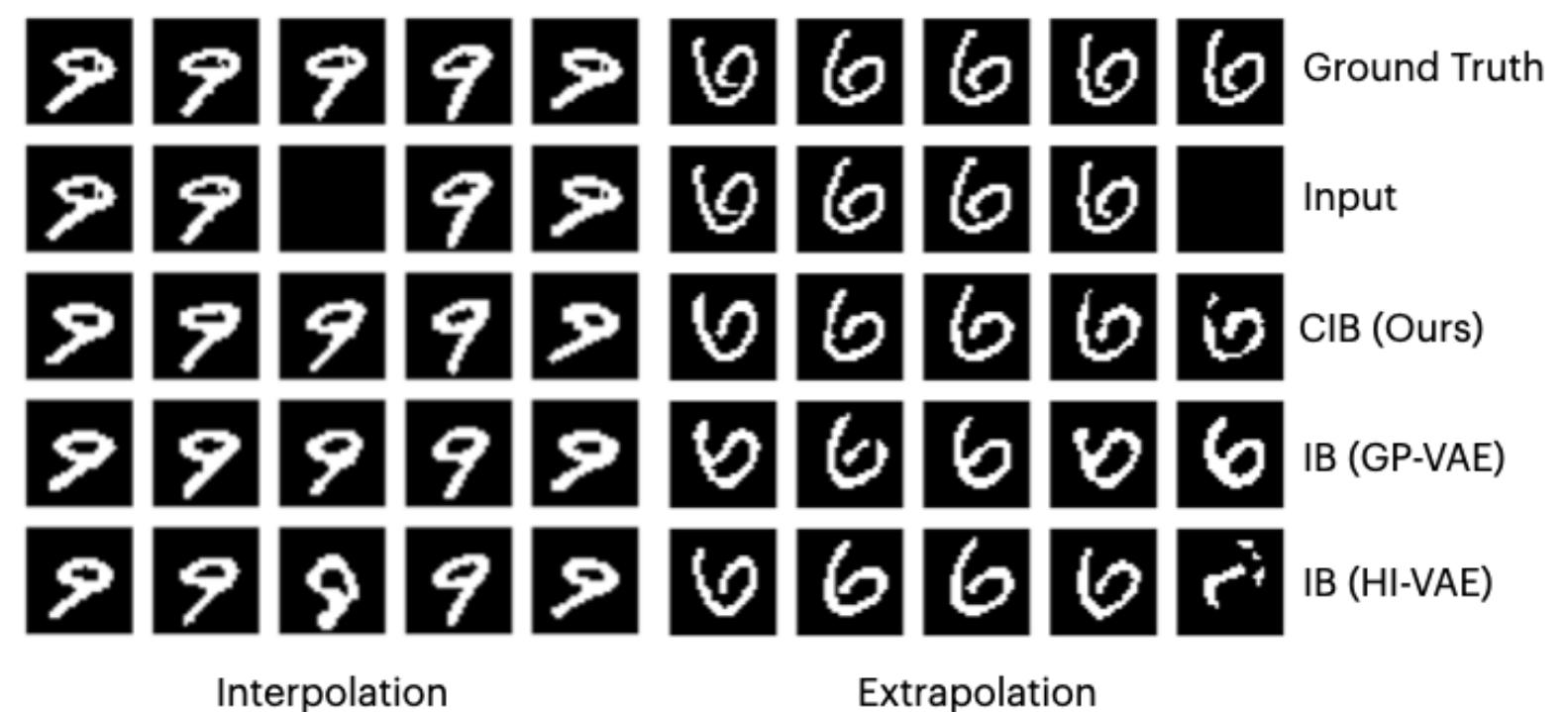


Figure 1B. Motivating experimental results on interpolation (left) and extrapolation (right). Because features in a single time step are completely missing, a model must collect information from other time steps. The conventional IB approach (HI-VAE) shows deteriorating performance in both cases. Another IB approach (GP-VAE) using a Gaussian process prior demonstrates enhanced performance for interpolation but often significantly loses time series characteristics for extrapolation (i.e., the writing style is corrupted). The CIB approach (Ours) exhibits improved imputation performance for both cases.

Unconditional regularization may lead to a significant loss of information regarding the temporal context. Therefore, apply regularization conditioned on remaining time steps, conditional information bottleneck (CIB).

Conditional Information Bottleneck (CIB)

Definition 2. (Time Series Imputation) Let \mathbf{X}_t^o and $\mathbf{X}_{t'}^m$ be random variables for the partially observed features and missing features of \mathbf{X}_t at time step t . Then, given the observed time series input $\mathbf{X}_{1:T}^o$, we define time series imputation at time step t as an unsupervised CIB as follows:

$$\min_{\phi, \theta} \underbrace{I_\phi(\mathbf{Z}_t; \mathbf{X}_t^o | \mathbf{X}_{1:t}^o)}_{\text{Conditional Regularization}} - \underbrace{\beta I_\theta(\mathbf{X}_t; \mathbf{Z}_t)}_{\text{Reconstruction}} \quad (2)$$

where $\mathbf{X}_{1:t}^o$ represents the random variables for the remaining input observations, excluding \mathbf{X}_t^o .

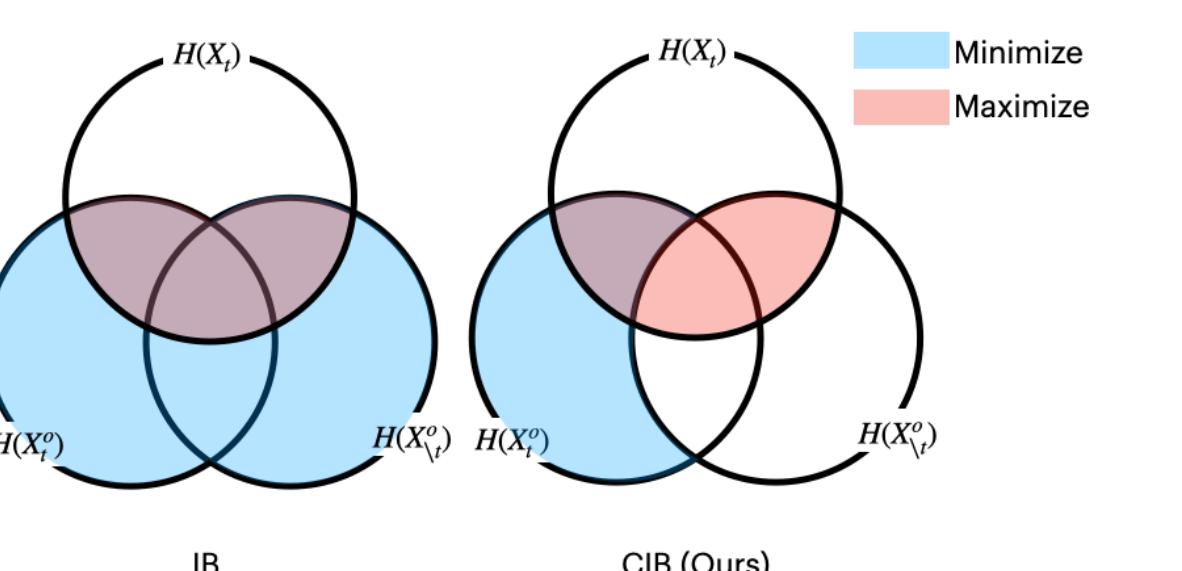


Figure 1A. Conceptual illustration of the IB and CIB principles.

Variational Decomposition of CIB

Maximizing Reconstruction $\min_{\phi, \theta} -I(\mathbf{X}_t; \mathbf{Z}_t)$

$$I_\theta(\mathbf{X}_t; \mathbf{Z}_t) \geq \mathbb{E}_{\mathbf{x}_{1:T}^o \sim p_{\text{data}}} [\mathbb{E}_{\mathbf{z}_t \sim q_\phi(\mathbf{z}_t | \mathbf{x}_{1:T}^o)} [\log p_\theta(\mathbf{x}_t | \mathbf{z}_t)]] \stackrel{\text{def}}{=} -\mathcal{L}_{\phi, \theta}^1 \quad (3)$$

(Voloshynovskiy et al., 2019)

Minimizing Conditional Regularization $\min_{\phi, \theta} I_\phi(\mathbf{Z}_t; \mathbf{X}_t^o | \mathbf{X}_{1:t}^o)$

Chain rule of mutual information gives:

$$\min_{\phi, \theta} I(\mathbf{Z}_t; \mathbf{X}_t^o | \mathbf{X}_{1:t}^o) = \min_{\phi, \theta} I(\mathbf{Z}_t; \mathbf{X}_{1:T}^o) - I(\mathbf{Z}_t; \mathbf{X}_{t'}^o). \quad (4)$$

1. Conciseness term:

$$I(\mathbf{Z}_t; \mathbf{X}_{1:T}^o) \leq \mathbb{E}_{\mathbf{x}_{1:T}^o \sim p_{\text{data}}} [D_{\text{KL}}(q_\phi(\mathbf{z}_t | \mathbf{x}_{1:T}^o) || p(\mathbf{z}_t))] \stackrel{\text{def}}{=} \mathcal{L}_\phi^2 \quad (5)$$

2. Temporal dynamics term:

$$I(\mathbf{Z}_t; \mathbf{X}_{1:T}^o) \geq \mathbb{E}_{\mathbf{x}_{1:T}^o \sim p_{\text{data}}} \left[\log \left(\frac{\sum_{t' \in \{1, \dots, T\} \setminus t} \exp(\mathbf{z}_t^T \tilde{\mathbf{z}}_{t'} / \tau)}{\sum_{\mathbf{x}_{1:T}^o \in \mathcal{X}_{1:T}^o} \sum_{t' \in \{1, \dots, T\}, \mathbf{z}_{t'} \sim q_\phi(\mathbf{z}_{t'} | \mathbf{x}_{1:T}^o)} \exp(\mathbf{z}_t^T \mathbf{z}_{t'}^- / \tau)} \right) \right] \stackrel{\text{def}}{=} -\mathcal{L}_\phi^3 \quad (6)$$

(A. v. d. Oord et al., 2018)

Introducing inductive bias on temporal dynamics:

Temporal kernels with Conditional alignment

$$I(\mathbf{Z}_t; \mathbf{X}_{1:T}^o) \geq \mathbb{E}_{\mathbf{x}_{1:T}^o \sim p_{\text{data}}} \left[\log \left(\frac{\sum_{t' \in \{1, \dots, T\} \setminus t} c_{t,t'} \exp(\mathbf{z}_t^T \tilde{\mathbf{z}}_{t'} / \tau)}{\sum_{\mathbf{x}_{1:T}^o \in \mathcal{X}_{1:T}^o} \sum_{t' \in \{1, \dots, T\}, \mathbf{z}_{t'} \sim q_\phi(\mathbf{z}_{t'} | \mathbf{x}_{1:T}^o)} \exp(\mathbf{z}_t^T \mathbf{z}_{t'}^- / \tau)} \right) \right] \stackrel{\text{def}}{=} -\mathcal{L}_\phi^{3'} \quad (7)$$

$c_{\text{cauchy}}(\tau, \tau') = \sigma^2 (1 + (\tau - \tau')^2 / l^2)^{-1}$

Experiments

Image Sequence Dataset

Methods	HealingMNIST (missing with MNAR pattern)			RotatedMNIST (interpolation & extrapolation)		
	NLL(\downarrow)	MSE(\downarrow)	AUROC(\uparrow)	NLL(\downarrow)	MSE(\downarrow)	AUROC(\uparrow)
No Imp.	-	0.293 \pm 0.000	0.920 \pm 0.000	-	-	0.133 \pm 0.000
Mean Imp.	-	0.168 \pm 0.000	0.938 \pm 0.000	-	-	0.085 \pm 0.000
Forward Imp.	-	0.177 \pm 0.000	0.946 \pm 0.000	-	-	0.080 \pm 0.000
VAE	0.480 \pm 0.002	0.232 \pm 0.000	0.922 \pm 0.000	1.773 \pm 0.127	0.133 \pm 0.000	-
HI-VAE	0.290 \pm 0.001	0.134 \pm 0.003	0.962 \pm 0.001	0.207 \pm 0.007	0.087 \pm 0.001	-
GP-VAE	0.261 \pm 0.001	0.114 \pm 0.002	0.960 \pm 0.002	0.190 \pm 0.001	0.080 \pm 0.004	-
Ours(Uniform)	0.204 \pm 0.002	0.090 \pm 0.001	0.967 \pm 0.001	0.184 \pm 0.001	0.077 \pm 0.001	-
Ours(Cauchy)	0.202 \pm 0.004	0.088 \pm 0.002	0.967 \pm 0.000	0.184 \pm 0.001	0.076 \pm 0.002	-

Table 1. Imputation and prediction performance on image sequence dataset.

Robustness Analysis

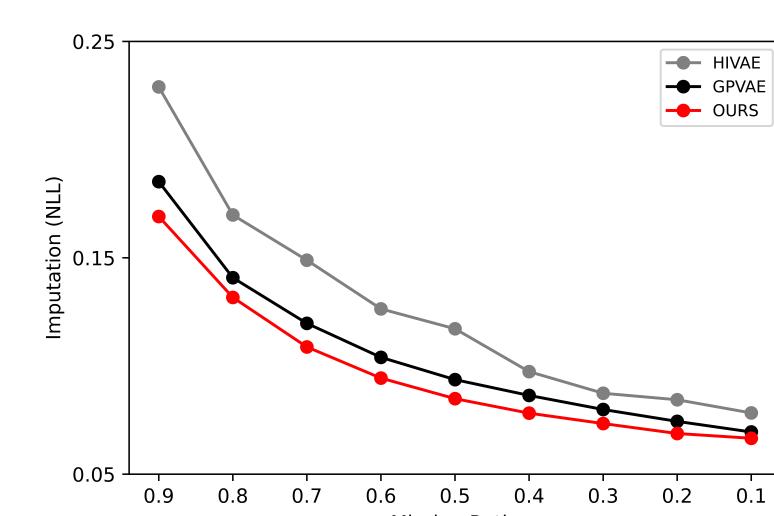
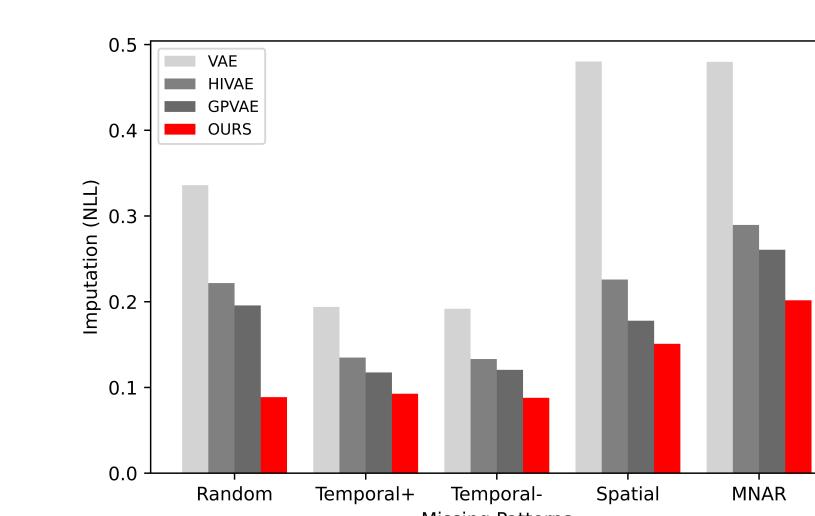


Figure 2. Robustness on missing patterns (Left) and missing ratios (Right). Spatial (i.e., neighboring pixels have correlated missing probabilities), and Temporal+/- (i.e., positive/negative temporal correlation). Missing ratio experiments performed with random missing pattern.

Phisyonet2012

Methods	Phisyonet2012 (mortality prediction)		
	NLL(\downarrow)	MSE(\downarrow)	AUROC(\uparrow)
No Imp.	-	0.962 \pm 0.000	0.692 \pm 0.000
Mean Imp.	-	0.511 \pm 0.000	0.703 \pm 0.000
Forward Imp.	-	0.613 \pm 0.000	0.710 \pm 0.000
BRITS	-	0.529 \pm 0.004	0.700 \pm 0.005
SAITS	-	0.501 \pm 0.024	0.713 \pm 0.007
VAE	1.400 \pm 0.000	0.962 \pm 0.000	0.691 \pm 0.001
HI-VAE	1.345 \pm 0.009	0.852 \pm 0.018	0.696 \pm 0.004
GP-VAE	1.227 \pm 0.007	0.616 \pm 0.013	0.730 \pm 0.006
Ours(Uniform)	1.183 \pm 0.007	0.528 \pm 0.014	0.744 \pm 0.009
Ours(Cauchy)	1.179 \pm 0.006	0.521 \pm 0.012	0.744 \pm 0.009

Table 2. Imputation and prediction performance on the clinical dataset.

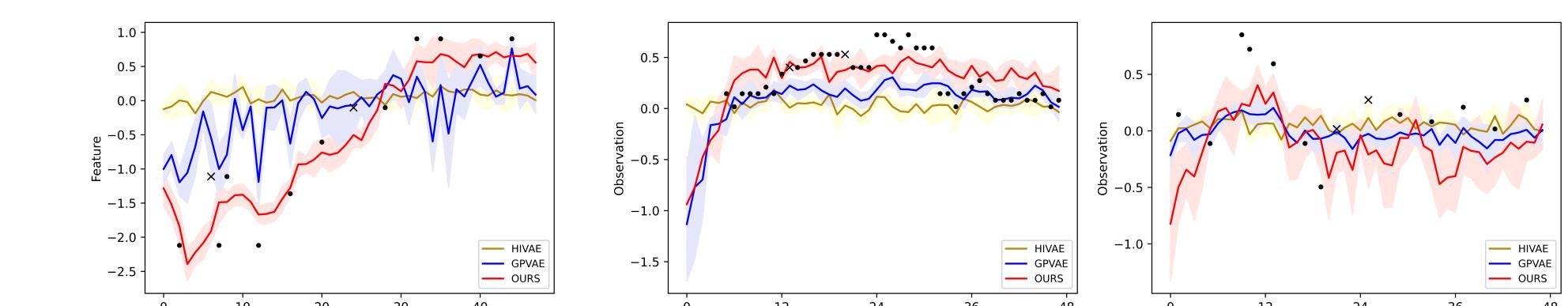


Figure 3. Qualitative results on Physionet2012.

References

- B. Dufumier et al., (2021). "Conditional Alignment and Uniformity for Contrastive Learning with Continuous Proxy Labels." In ArXiv
- V. Fortuin et al., (2020). "GP-VAE: Deep Probabilistic Time Series Imputation" In International conference on artificial intelligence and statistics.
- A. Nazabal et al., (2020). "Handling Incomplete Heterogeneous Data Using VAEs." In Pattern Recognition.
- S. Voloshynovskiy et al., (2019). "Information Bottleneck Through Variational Glasses." In: NeurIPS Workshop on Bayesian Deep Learning.
- A. v. d. Oord et al., (2018). "Representation Learning with Contrastive Predictive Coding." In: ArXiv.

