SAMPLING-GUIDED HETEROGENEOUS GRAPH NEURAL NETWORK WITH TEMPORAL SMOOTHING FOR SCAL ABLE LONGITUDINAL DATA IMPUTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we propose a novel framework, the Sampling-guided Heterogeneous Graph Neural Network (SHT-GNN), to effectively tackle the challenge of missing data imputation in longitudinal studies. Unlike traditional methods, which often require extensive preprocessing to handle irregular or inconsistent missing data, our approach accommodates arbitrary missing data patterns while maintaining computational efficiency. SHT-GNN models both observations and covariates as distinct node types, connecting observation nodes at successive time points through subject-specific longitudinal subnetworks, while covariate-observation interactions are represented by attributed edges within bipartite graphs. By leveraging subject-wise mini-batch sampling and a multi-layer temporal smoothing mechanism, SHT-GNN efficiently scales to large datasets, while effectively learning node representations and imputing missing data. Extensive experiments on both synthetic and real-world datasets, including the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset, demonstrate that SHT-GNN significantly outperforms existing imputation methods, even with high missing data rates (e.g., 80%). The empirical results highlight SHT-GNN's robust imputation capabilities and superior performance, particularly in the context of complex, large-scale longitudinal data.

028 029

031

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

With the advancement of data collection and analysis techniques, longitudinal data has gained increasing importance across various scientific fields, such as biomedical science, economics, and e-commerce (Pekarčík et al., 2022; Sadowski et al., 2021; Mundo et al., 2021). Observation schedules in longitudinal studies vary widely as some follow regular intervals for each subject while others follow irregular patterns. For example, clinical visits scheduled every two months represent a regular observation schedule, whereas follow-ups at 6 months, 1 year, and 2 years post-treatment is an irregular schedule. Additionally, observation schedules may be either consistent or inconsistent across different subjects—some studies may impose uniform intervals for all participants, while others allow for variability due to individual health conditions or other circumstances.

041 Imputation of missing values poses a major challenge in longitudinal data analysis, both from 042 theoretical and practical standpoints (Daniels & Hogan, 2008; Little & Rubin, 2019). In fact, in 043 complex longitudinal data analysis, such as those in neuroimaging or electronic health record studies, 044 patient follow-up data is collected over time, often resulting in missing values across variables (cross-sectional missingness) and across time (longitudinal missingness). Ideally, as illustrated in Figure 1, each subject would have observations taken at consistent intervals, with an equal number of 046 observations per subject, and each observation would contain fully observed covariates and response 047 variable. However, when missing data arises at specific time points, this misalignment of observations 048 across time can occur. Furthermore, the absence of entire observations at certain time steps can transform regular observation schedules into irregular ones and consistent schedules into inconsistent ones, significantly complicating longitudinal data analysis. 051

Given the frequent occurrence of missing measurements in studies with extensive data collection,
 addressing the missingness of both covariates and response variables is crucial for accurate predictions of the target response variable. For instance, in Alzheimer's disease research, the prediction of key



Figure 1: The ideal format of longitudinal data without missing (left); The regularly and consistently observed longitudinal data with missing in covariate variables and response variable (middle). And the irregular and inconsistent observation schedule in longitudinal data due to missing data (right).

biomarkers often relies on incomplete datasets to assess disease progression. Despite methodological advancements, managing missing data in longitudinal studies presents several ongoing challenges. Firstly, how can we effectively model longitudinal data that is irregularly and inconsistently observed? Secondly, how can we design a cohesive forward pass process for data imputation that accommodates varied observation schedules across subjects? Thirdly, how can we accurately predict target response values in the presence of missing covariate data and seamlessly integrate this process into model training? Fourthly, in the context of large-scale longitudinal studies, how can we ensure the scalability of the missing data imputation method?

In this paper, we address the key challenges by introducing the Sampling-guided Heterogeneous
 Graph Neural Network (SHT-GNN). Unlike traditional imputation methods, our GNN leverages
 scalable modeling of complex data structures to effectively learn from irregular and inconsistent
 longitudinal observations. Our SHT-GNN incorporates three key innovations:

- 1. We handle longitudinal data by sharing trainable parameters across sampled graphs constructed from subject-wise mini-batches. This sampling-guided process ensures scalability for massive longitudinal datasets.
- 2. Subject-wise longitudinal subnetworks are constructed by connecting adjacent observations with directed edges, while covariate values are transformed into attributed edges, linking observation nodes with covariate nodes. This structure allows SHT-GNN to flexibly model longitudinal data under arbitrary missing data conditions.
 - 3. We introduce a temporal smoothing mechanism across observations for each subject using a multi-layer longitudinal subnetwork. The novel MADGap statistic controls the smoothness within the subnetwork, balancing temporal smoothing with the specificity of observation node representations.

We conducted extensive experiments on both real data and synthetic data. The results demonstrate that SHT-GNN consistently achieves state-of-the-art performance across different temporal characteristics and performs exceptionally well even under high missing rates for both covariates and response variable. Furthermore, we applied the proposed SHT-GNN model to the Alzheimer's Disease Neuroimaging Initiative (ADNI) dataset to predict the critical biomarker $A\beta 42/40$. SHT-GNN's strong performance in predicting $A\beta 42/40$ is confirmed through validation with ground truth data and in downstream analyses. Additionally, we conduct extensive ablation studies, demonstrating the longitudinal network's ability to perform temporal smoothing through multi-layer longitudinal subnetworks and the critical role of MADGap in guaranteeing variance between observations.

097 098

099

061

062

063 064

065

066

067

068

069

071

077

078

079

081

082

084

2 RELATED WORK

2.1 STATISTICAL IMPUTATION METHODS

Statistical methods for imputing missing values in longitudinal data are often derived from multiple
 imputation techniques originally developed for cross-sectional data. For example, the 3D-MICE
 method extends the MICE (Multivariate Imputation by Chained Equations) framework to account
 for both cross-sectional and longitudinal dependencies in data (Luo et al., 2018). Some methods
 like trajectory mean and last observation carried forward (LOCF) leverage the time-series nature
 for missing data imputation (Lane, 2008). However, these basic models often fail to capture the
 complex, nonlinear spatio-temporal dependencies inherent in longitudinal data and lack the scalability

108 to manage both continuous and discrete covariates simultaneously. The multi-directional multivariate 109 time series (MD-MTS) method seeks to address these limitations by integrating temporal and cross-110 sectional covariates into a unified imputation framework through extensive feature engineering (Xu 111 et al., 2020). A notable limitation of MD-MTS is its reliance on manually crafted features can be 112 labor-intensive and prone to errors. To overcome such challenges, the time-aware dual-cross-visit (TA-DualCV) method leverages both longitudinal dependencies and covariate correlations using Gibbs 113 sampling for imputation (Gao et al., 2022). Nonetheless, when applied to large-scale longitudinal 114 datasets with numerous repeated observations, those methods involving chained equations and Gibbs 115 sampling can become computationally expensive. 116

117 118

2.2 MACHINE LEARNING IMPUTATION METHODS

119 Machine learning methods for imputing missing values in longitudinal data have primarily utilized 120 recurrent neural networks (RNNs) and their advanced variants, such as long short-term memory 121 (LSTM) networks and gated recurrent units (GRUs). BRITS, for example, leverages bidirectional 122 LSTMs to capture both longitudinal and cross-sectional dependencies by utilizing past and future 123 trends (Cao et al., 2018). Similarly, NAOMI employs a recursive divide-and-conquer approach with 124 bidirectional RNNs to handle missing data imputation (Woillez et al., 2019). CATSI enhances this 125 by combining bidirectional LSTMs with a context vector to account for patient-specific temporal dynamics (Kazijevs & Samad, 2023). A significant limitation of all RNN-based methods, however, 126 is their reliance on consistent time step lengths across different subjects, often requiring padding 127 or truncation to achieve this uniformity. When applied to large datasets, where the number of 128 observations can vary greatly between subjects, padding and truncation may lead to excessive and 129 inefficient computations due to the presence of numerous invalid time steps. Beyond RNN-based 130 approaches, the GP-VAE model integrates deep variational autoencoders with Gaussian processes to 131 address missing data in time series. However, the computational complexity of Gaussian processes, 132 which scales at $O(N^3)$, introduces significant computational bottlenecks when processing large-scale 133 longitudinal data (Fortuin et al., 2020).

134 135

136

2.3 GRAPH-BASED IMPUTATION METHODS

137 Numerous advanced graph-based methods have been developed for data imputation tasks (Zhang & Chen, 2019; Li et al., 2021; Zhang et al., 2023). However, these methods face limitations when 138 handling mixed discrete and continuous features and often fail to capture the temporal dependencies 139 inherent in longitudinal data. For example, RAINDROP employs a graph-guided network to handle 140 irregularly observed time series with varying intervals, effectively addressing missing data imputation 141 and outperforming other methods in downstream prediction tasks (Zhang et al., 2021). Similarly, 142 TGNN4I, which integrates GNNs with gated recurrent units, has shown strong performance in 143 imputing missing data for longitudinal graph data (Oskarsson et al., 2023). Despite their strengths, 144 both RAINDROP and TGNN4I rely on predefined node connections within graph datasets, limiting 145 their applicability to longitudinal tabular data where such connections between subjects do not 146 inherently exist. Additionally, GRAPE demonstrated success using a bipartite graph for feature 147 imputation, but its edge size increases linearly with the number of observations, posing scalability 148 challenges (You et al., 2020). IGRM (Zhong et al., 2023), an extension of GRAPE, constructs a friend network among sample nodes to capture similarities and improve imputation. However, IGRM 149 requires continuous updating of a fully connected graph among all samples, making it computationally 150 expensive and lacking scalability for large-scale longitudinal data. 151

152 153

154

3 PROBLEM FORMULATION

In the longitudinal data structure depicted in Figure 1, the observed variables are divided into two categories: covariates and the response variable. Assume there are *n* subjects, with the *k*-th subject having n_k observations measured. The total number of observations across all subjects is $N = \sum_{k=1}^{n} n_k$. Typically, the set of observed covariates varies across observations in longitudinal studies. For simplicity, we denote the complete set of covariates that can potentially be observed as $X = \{x_1, \ldots, x_p\}$, where *p* represents the total number of covariates. Any missing values within this set are treated as missing data. The covariate data for all observations can thus be represented as a matrix $\mathbf{D} = (D_{il}) \in \mathbb{R}^{N \times p}$, where $l \in \{1, \ldots, p\}$ and $i = \sum_{k=1}^{j-1} n_k + m$ indexes the *m*-th



Figure 2: The flow chart for Sampling-guided Heterogeneous Graph Neural Network.

176 observation of the j-th subject. Similarly, the response variable for all observations is denoted as 177 $\mathbf{Y} = (Y_i) \in \mathbb{R}^{N \times 1}$. Given the presence of missing data, we introduce a mask matrix for the covariate data $\mathbf{M}^O = (M_{il}^O) \in \{0, 1\}^{N \times p}$, where D_{il} is observed if $M_{il}^O = 1$. Likewise, the mask matrix for the response variable is denoted by $\mathbf{M}^Y = (M_i^Y) \in \{0, 1\}^{N \times 1}$, where the value of Y_i is observed 178 179 if $M_i^Y = 1$. Within this irregular or inconsistent longitudinal data structure as depicted in Figure 181 1, the goal is to predict the response variable Y_i for all i where $M_i^Y = 0$, leveraging the available 182 observation information despite missingness.

4 METHOD

175

183

185

186 We propose SHT-GNN to address the challenges of missing data in irregular or inconsistent longi-187 tudinal data mentioned in Section 1. SHT-GNN models the observations and covariates as nodes, 188 connecting them with attributed edges that represent the observed covariate values. Additionally, 189 it employs specially designed longitudinal subnetworks to perform temporal smoothing among 190 observations within the same subject. Initially, observations in longitudinal data can be represented 191 as a graph $\mathcal{G}_L(\mathcal{V}_O, \mathcal{E}_{OO})$ with subject-wise longitudinal subnetworks, where the observation node 192 set $\mathcal{V}_O = \{u_1, u_2, \dots, u_N\}$ represents all observations. The observation nodes that belong to the same subject form their own independent longitudinal subnetwork. In each longitudinal subnetwork, 193 observation nodes at adjacent time points are connected by directed edges. The directed edge set 194 in \mathcal{G}_L is denoted as $\mathcal{E}_{OO} = \{(u_i, u_{i'}, e_{i \to i'}) | u_i, u_{i'} \in \mathcal{V}_O, S(u_i) = S(u_{i'}), u_i \prec u_{i'}\}$, where S(u)195 denotes the subject that observation u belong to, and $u_i \prec u_{i'}$ represents u_i is the direct predecessor 196 of $u_{i'}$ in terms of time. 197

In irregularly and inconsistently observed longitudinal observation data, each observation may have measured different covariates. SHT-GNN establishes a covariate node set $\mathcal{V}_C = \{v_1, v_2, \dots, v_p\}$, 199 and construct edges between observation nodes and covariate nodes to represent the observed covariate values. The longitudinal data matrix $\mathbf{D} \in \mathbb{R}^{N \times p}$ and the missing indicator matrix $\mathbf{M}^O \in$ 200 201 $\{0,1\}^{N \times p}$ can be represented as a bipartite graph $\mathcal{G}_B = (\mathcal{V}, \mathcal{E}_{OC})$. The undirected edge set $\mathcal{E}_{OC} = \mathcal{E}_{OC}$ 202 $\{(u_i, v_l, e_{u_i v_l}) | u_i \in \mathcal{V}_O, v_l \in \mathcal{V}_C, M_{il}^O = 1\}$, where the edge feature $e_{u_i v_l}$ takes the value of the 203 corresponding feature $e_{u_i v_l} = D_{il}$. Then a longitudinal dataset can be represented as the union of one 204 undirected bipartite graph and multiple subject-wise longitudinal subnetworks: $\mathcal{G}(\mathcal{V}, \mathcal{E}) = \mathcal{G}_B \cup \mathcal{G}_L =$ 205 $\mathcal{G}(\mathcal{V}_O \cup \mathcal{V}_C, \mathcal{E}_{OO} \cup \mathcal{E}_{OC})$. After constructing SHT-GNN, the task of response variable prediction 206 under missing covariate imputation can be represented as learning the mapping: $Y_i = [f(\mathcal{G})]_i$ by 207 minimizing the difference between Y_i and \hat{Y}_i , for all *i* where $M_i^Y = 1$. 208

- 209
- LEARNING IN SAMPLING-GUIDED HETEROGENEOUS GRAPH NEURAL NETWORK 4.1
- 210 211 212
- 4.1.1 SUBJECT-WISE MINI-BATCH SAMPLING

It is evident that the spatial complexity of graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ increases with the number of subjects. When 213 working with longitudinal data, it is common to encounter a large number of repeated observations or 214 subjects. When dealing with millions or even tens of millions of observations, directly training a huge 215 graph composed of all subjects is impractical and inefficient (You et al., 2020; Zhong et al., 2023). SHT-GNN performs subject-wise mini-batch sampling at the beginning of each training phase. The corresponding graph is then constructed based on the sampled subjects for the current training phase. In SHT-GNN, the graphs associated with each training phase share the same trainable parameters.

219 Specifically, assume there are n subjects, represented by the subject set $S = \{1, 2, ..., n\}$, where 220 the k-th subject corresponds to an observation set O_k . Across all observation sets, there are a 221 total of N observations. At the beginning of each training phase, s subjects are randomly sampled 222 from n subjects to form a subject-wise mini-batch $S_0 = \text{Sample}(S, s)$, yielding the corresponding 223 observation-wish mini-batch $O_s = \bigcup_{k \in S_0} O_k$. The observations in O_s are then employed to construct 224 the graph $\mathcal{G}_s = (\mathcal{V}_s, \mathcal{E}_s)$, where $\mathcal{V}_s = \mathcal{V}_{O_s} \cup \mathcal{V}_C$, with $\mathcal{V}_{O_s} = \{u_i | i \in O_s\}$. The edge set \mathcal{E}_s is 225 constructed within \mathcal{V}_s according to the definitions provided above. Assuming the trainable parameters 226 in SHT-GNN are denoted by θ , and the loss function under input graph \mathcal{G}_{input} is $L(\mathcal{G}_{input}, \theta)$. The 227 parameter learning process on the sampled graph in each sampling phase can be expressed as:

234

243

249

252 253 254

255 256

257

264 265

269

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta L_{\mathcal{S}_0}(\mathcal{G}_s, \theta^{(t)}), \tag{1}$$

where $L_{S_0}(\mathcal{G}_s, \theta^{(t)})$ represents the loss function obtained from the forward pass in graph \mathcal{G}_s , and η_t represents the learning rate. Here the forward pass in \mathcal{G}_s begins with inductive learning within a multi-layer bipartite graph to derive representations of observation and covariate nodes. Afterward, temporal smoothing is conducted within the longitudinal subnetwork of each subject.

4.1.2 INDUCTIVE LEARNING IN MULTI-LAYER BIPARTITE GRAPH

In SHT-GNN, all the information from the observed data is derived from the attributed edges connecting observation nodes and covariate nodes. The representations of all nodes need to be inductively learned from these edges. Inspired by GraphSAGE, which is a variant of GNNs known for its strong inductive learning capabilities. SHT-GNN modifies GraphSAGE architecture by introducing edge embeddings in the bipartite graph. At the *l*-th layer in \mathcal{G}_B , the message generation function takes the concatenation of the embedding of the source node $\mathbf{h}_v^{(l-1)}$ and the edge embedding $\mathbf{e}_{uv}^{(l-1)}$ as input:

$$\mathbf{m}_{v}^{(l)} = \operatorname{AGG}_{l}\left(\sigma\left(\mathbf{P}^{(l)} \cdot \operatorname{CONCAT}(\mathbf{h}_{u}^{(l-1)}, \mathbf{e}_{uv}^{(l-1)})\right) \mid \forall u \in \mathcal{N}(v, \mathcal{E}_{s})\right) \quad \forall v \in \mathcal{V}_{s},$$
(2)

where AGG_l is message aggregation function for all node, σ is the non-linearity activation function, P^(l) is the trainable weight matrix and \mathcal{N} is the node neighborhood function. Subsequently, all nodes update their embeddings based on the concatenation of aggregated messages and their local representations:

$$\mathbf{h}_{v}^{(l)} = \sigma \left(\mathbf{Q}^{(l)} \cdot \text{CONCAT}(\mathbf{h}_{v}^{(l-1)}, \mathbf{m}_{v}^{(l)}) \right) \qquad \forall v \in \mathcal{V}_{s},$$
(3)

where $\mathbf{Q}^{(l)}$ is the trainable weight matrix. Subsequently, the edge embeddings are then updated based on the updated node embeddings at both ends of each edge:

$$\mathbf{e}_{uv}^{(l)} = \sigma \left(\mathbf{W}^{(l)} \cdot \text{CONCAT}(\mathbf{e}_{uv}^{(l-1)}, \mathbf{h}_{u}^{(l-1)}, \mathbf{h}_{v}^{(l-1)}) \right) \qquad \forall \mathbf{e}_{uv} \in \mathcal{E}_{s}, \tag{4}$$

where $\mathbf{W}^{(l)}$ is the trainable weight matrix.

4.1.3 TEMPORAL SMOOTHING IN MULTI-LAYER LONGITUDINAL SUBNETWORKS

Temporal smoothing is a crucial technique in the imputation of longitudinal data because it leverages the temporal correlation within the observations. After *L* layers of inductive learning in the bipartite graph, SHT-GNN innovatively performs temporal smoothing through *K* layers of message passing and representation updates within the subject-wise longitudinal subnetworks. At the *k*-th layer of each subject-wise longitudinal subnetwork, the message passing function for observations u_i to $u_{i'}$ takes the source node embedding $\mathbf{h}_{u_i}^{(L+k)}$ and the edge weight $w_{u_i \to u_{i'}}^{(L+k)}$ as input:

$$\mathbf{m}_{u_{i} \to u_{i'}}^{(L+k)} = \mathbf{h}_{u_{i}}^{(L+k)} \cdot w_{u_{i} \to u_{i'}}^{(L+k)}, \text{ where } S(u_{i}) = S(u_{i'}), u_{i} \prec u_{i'}, \forall u_{i}, u_{i'} \in \mathcal{V}_{O_{s}},$$
(5)

in which S(u) denotes the subject that observation u belong to and $u_i \prec u_{i'}$ represents that u_i is the direct predecessor of $u_{i'}$ in terms of time. For each pair of message passing described in (5), the target observation node $u_{i'}$ updates its representation as follows:

$$\mathbf{h}_{u_i}^{(L+k+1)} = \sigma \left(\mathbf{U}^{(L+k)} \cdot \text{CONCAT}(\mathbf{h}_{u_{i'}}^{(L+k)}, \ \mathbf{m}_{u_i \to u_{i'}}^{(L+k)}) \right), \tag{6}$$



Figure 3: The variation of information in the representation of observation nodes during multi-layer message passing and representation updates for temporal smoothing within longitudinal subnetworks.

where \mathbf{U}^{L+k} is the trainable parameter matrix for embedding updates. In (5), the weight for edges connecting observation nodes needs to be calculated and updated during training. The edge weight is computed as follows:

$$w_{u_{i}\to u_{i'}}^{(L+k)} = D_{u_{i}\to u_{i'}} \cdot J_{u_{i}\to u_{i'}} \cdot \operatorname{Cos}\left(\mathbf{h}_{u_{i}}^{(L+k)}, \mathbf{h}_{u_{i'}}^{(L+k)}\right).$$
(7)

Here, $D_{u_i \to u_{i'}}$ represents the time decay weight in SHT-GNN. When facing irregular observation schedules in longitudinal data, the time interval between observations u_i and $u_{i'}$ is not fixed. We employ exponential decay functions to compute the time decay weights in (7). We define

$$D_{u_i \to u_{i'}} = \exp\left(-\frac{|T(u_i) - T(u_{i'})|}{\Delta_{max}}\right) \quad (8) \qquad \text{and} \quad J_{u_i \to u_{i'}} = \gamma - \frac{|\mathbb{A} \cap \mathbb{B}|}{|\mathbb{B}|}. \tag{9}$$

290 In (8), T(u) represents the time step associated with observation u, and Δ_{max} represents the longest time interval between adjacent observations for the current subject. We further introduce the Jaccard 291 distance between the sets of observed covariates in the weight calculation for message passing. 292 Suppose the covariate sets observed by \mathbf{h}_{u_i} and \mathbf{h}_{u_i} are denoted as \mathbb{A} and \mathbb{B} , respectively. As shown 293 in (9), a larger $J_{u_i \rightarrow u_{i'}}$ implies that observation u_i contains more covariates that are not observed in observation $u_{i'}$, and γ is a hyperparameter constant for preventing the weights from becoming 295 zero. Therefore, \mathbf{h}_{u_i} need to borrow more information from \mathbf{h}_{u_i} . Beyond the time decay weight and 296 Jaccard distance, the strength of message passing for a specific observation $h_{u,i}$ from its immediate 297 predecessor \mathbf{h}_{u_i} is intuitively determined by the similarity between the currently observed covariates 298 of \mathbf{h}_{u_i} and $\mathbf{h}_{u_{i'}}$. SHT-GNN then introduces the cosine similarity between \mathbf{h}_{u_i} and $\mathbf{h}_{u_{i'}}$ in the edge 299 weight as discribed in (7).

300 Intuitions on why longitudinal subnetworks work. In the SHT-GNN, multi-layer message passing 301 and representation updates within longitudinal subnetworks enable observations from the same subject 302 to leverage information from previous time steps. As shown in Figure 3, in a two-layer longitudinal 303 subnetwork, assume there are time-ordered observations $u_{t_1}, u_{t_2}, u_{t_3} and u_{t_4}$. The message passing 304 and representation update in the first layer allow u_{t_m} to directly draw information from $u_{t_{m-1}}$. 305 Subsequently, while the temporal smoothing in the second layer enable u_{t_m} to indirectly capture information from $u_{t_{m-2}}$ via the representation of $u_{t_{m-1}}$. This iterative process can be extended with 306 307 additional layers, allowing observation nodes to integrate information from progressively earlier time steps, thereby complete temporal smoothing in subject-wise longitudinal subnetworks. 308

309 310

314 315

316

276

277

278

281 282 283

287 288 289

4.1.4 COVARIATE IMPUTATION IN SHT-GNN

After conducting L layers of inductive learning on the bipartite graph and K layers of temporal smoothing within longitudinal subnetworks, edge-level predictions are made at the (L + K)-th layer:

$$\widehat{D}_{il} = \mathbf{O}_{\text{impute}} \left(\text{CONCAT}(\mathbf{h}_{u_i}^{(L+K)}, \mathbf{h}_{v_l}^{(L+K)}) \right) \quad \forall u_i \in V_{O_s}, \ v_l \in V_F,$$

where O_{impute} is a multilayer perceptron (MLP). Here, \widehat{D}_{il} represents the imputation output for the *l*-th covariate for the *i*-th observation.

317 318 319

320

321 322 323 4.1.5 **RESPONSE VARIABLE PREDICTION**

Finally, we complete the prediction of the response variable based on the imputed covariates:

$$\widehat{Y}_{i} = \mathbf{O}_{\text{predict}} \left(\text{CONCAT}(\widehat{D}_{i1}, \widehat{D}_{i2}, ..., \widehat{D}_{ip}) \right) \quad \forall u_{i} \in \mathcal{V}_{O_{s}}$$

where $\mathbf{O}_{\text{predict}}$ is a multilayer perceptron and \widehat{Y}_i represents the predicted response.

4.1.6 Loss function for SHT-GNN

In both the response variable prediction and covariate imputation tasks, the loss function takes the following form:

328

337 338

339

344

345

 $Loss = MSE - \lambda \cdot MADGap.$

Here, MSE = $(\sum_{i=1}^{N} M_i^Y)^{-1} \sum_{i=1}^{N} M_i^Y \cdot (\hat{Y}_i - Y_i)^2$, where M_i^Y is the missing indicator. MADGap (Mean Average Distance Gap) is a statistical measure used to quantify the degree of over-smoothing in GNNs (Chen et al., 2020). A large MADGap value indicates that the node receives more useful information than noise. In SHT-GNN, the multi-layer longitudinal subnetwork between observation nodes represents the process of temporal smoothing. However, multi-layer message passing can lead to over-smoothing, causing the embedding representations of different observations within the same subject to become overly similar. SHT-GNN addresses this by maximizing MADGap to mitigate over-smoothing in the GNN. MADGap is defined for individual nodes as follows:

$$MADGap = MAD_{remote} - MAD_{neighbour}$$

In SHT-GNN, for the k-th subject's subnetwork with n_k time ordered observations $\{u_1, u_2, ..., u_{n_k}\}$, MADGap is calculated as:

$$\mathsf{MADGap}_{k} = \frac{1}{n_{k}} \sum_{m=1}^{n_{k}} \mathbf{1}_{m>2} \cdot \left(\frac{1}{m-2} \sum_{m'=1}^{m-2} \mathsf{Cos}(\mathbf{h}_{u_{m'}}, \mathbf{h}_{u_{m}}) - \mathsf{Cos}(\mathbf{h}_{u_{m}}, \mathbf{h}_{u_{m-1}}) \right).$$

Here $\frac{1}{m-2} \sum_{m'=1}^{m-2} \operatorname{Cos}(\mathbf{h}_{u_m}, \mathbf{h}_{u_{m'}})$ denotes the similarity between the representations of the *m*-th observation and its past ancestor observation nodes in the longitudinal subnetwork, where $\operatorname{Cos}(\mathbf{h}_{u_m}, \mathbf{h}_{u_{m'}})$ represents the cosine similarity between \mathbf{h}_{u_m} and $\mathbf{h}_{u_{m'}}$. Then for all *n* subjects:

$$MADGap = \frac{1}{n} \sum_{k=1}^{n} \left[\frac{1}{n_k} \sum_{m=1}^{n_k} \mathbf{1}_{m>2} \cdot \left(\frac{1}{m-2} \sum_{m'=1}^{m-2} Cos(\mathbf{h}_{u_{m'}}, \mathbf{h}_{u_m}) - Cos(\mathbf{h}_{u_m}, \mathbf{h}_{u_{m-1}}) \right) \right]$$

350 351

352 353

354

5 EXPERIMENT

5.1 BASELINES

We consider eight baseline methods as follows. 1. Mean (Huque et al., 2018): imputes missing values 355 using the covariate-wise mean. 2. Copy-mean Last Observation Carried Forward (LOCF) (Jahangiri 356 et al., 2023): first imputes using the LOCF, then refines the results based on the population's mean 357 trajectory. 3. Multivariate Imputation by Chained Equations (MICE)(Van Buuren & Groothuis-358 Oudshoorn, 2011): employs multiple regressions to model each missing value conditioned on other 359 observed covariate values. 4. 3D-MICE (Kazijevs & Samad, 2023): combines MICE and Gaussian 360 processes for longitudinal data imputation. 5. GRAPE: handles feature imputation through graph 361 representation learning. (You et al., 2020) 6. CASTI (Yin et al., 2020): handles missing data in 362 longitudinal data by employing bidirectional LSTM and MLPs. 7. IGRM (Zhong et al., 2023): enhances feature imputation by leveraging the similarity between observations. 8. Transformer (Zeng et al., 2023): handles missing data imputation in time series with self-attention mechanism. 9. GP-364 VAE (Fortuin et al., 2020): enhances data imputation with Gaussian Process Variational Autoencoder. 365 10. CTA (Wi et al., 2024): enhances feature imputation via continuous-time Autoencoders. 366

367 368

5.2 EXPERIMENT ON SYNTHETIC DATA SIMULATED FROM REAL DATA

369 **Real data introduction.** To comprehensively evaluate the performance of various methods on 370 longitudinal data with diverse temporal characteristics, we first conduct experiments using synthetic 371 data simulated from real-world datasets. We selected the longitudinal behavior modeling (GLOBEM) 372 dataset as the basis for our synthetic data. The GLOBEM dataset is an extensive, multi-year collection 373 derived from mobile and wearable sensing technologies. It contains data from 497 unique subjects, 374 with over 50,000 observations and more than 2,000 covariates, including phone usage, Bluetooth 375 scans, physical activity, and sleep patterns. For simulating synthetic datasets, all covariate values are sampled directly from the original data. Since the GLOBEM dataset does not include a predefined 376 response variable, we simulate response values based on the observed covariates to construct complete 377 synthetic datasets.

378 **Response Variable simulation based on real data.** In line with common assumptions in longitudinal 379 data studies, when simulating the response values for the *m*-th observation of the *k*-th subject, assume the response value Y_i , where $i = \sum_{k=1}^{j-1} n_k + m$, lies within the *t*-th temporal smoothing window. For simplicity of expression, denote it as Y_{km}^t , which is assumed to follow a normal distribution 380 381 382 $N(\mu_k^t, \sigma_k^{t^2})$. Then μ_k^t and $\sigma_k^{t^2}$ respectively represent the mean and variance of the response variable for the t-th temporal smoothing window for the k-th subject. In practice, suppose the k-th subject 384 has n_k observations within the time step set $T = \{1, \ldots, n_k\}$. We first divide T into W temporal smoothing windows. For the m-th observation, which occurs within the t-th temporal smoothing 386 window, the response value Y_{km}^t is modeled as follows: $Y_{km}^t = \mu_k^t + \epsilon_{km}$, where μ_k^t represents the mean response value for the k-th subject within the t-th temporal smoothing window, and ϵ_{km} 387 denotes the random fluctuations for the m-th observation. 388

389 Response variable simulation execution. In our experiment, we run multiple random trials under a 390 set of fixed parameters, reporting the average performance of all methods across these trials. In each 391 random trial, we first randomly select p dimensions from the 2000 covariates in the GLOBEM dataset, 392 denoted as $X = \{x_1, x_2, ..., x_p\}$. We then simulate the response values for each observation based on a specified model $f(X) = f(x_1, x_2, ..., x_p)$, which incorporates both linear and nonlinear elements. 393 The details about two different specified models f(X) are provided in Appendix A.1. According to the longitudinal effect model $Y_{km}^t = \mu_k^t + \epsilon_{km}$ described above, we first impute the mean value of 394 395 the response variable for observations belonging to the same subject and within the same temporal smoothing window. Next, random fluctuations ϵ_{km} are sampled from the distribution $N(0,\epsilon)$ and 397 incorporated into the formula $Y_{km}^t = \mu_k^t + \epsilon_{km}$ to simulate the final response values. 398

Experimental Procedure. After obtaining complete synthetic datasets, we first split all subjects 399 into a training set and a test set according to ratio r. By extracting the observations of all subjects, 400 we obtain the covariate matrix $\mathbf{D}^{\text{train}} \in \mathbb{R}^{N_{\text{train}} \times p}$ and response vector $\mathbf{Y}^{\text{train}} \in \mathbb{R}^{N_{\text{train}}}$ for training, as 401 well as $\mathbf{D}^{\text{test}} \in \mathbb{R}^{N_{\text{test}} \times p}$ and $\mathbf{Y}^{\text{test}} \in \mathbb{R}^{N_{\text{test}}}$ for testing. We respectively generate missing indicator 402 matrices $\mathbf{M}^{\text{train}}$ and \mathbf{M}^{test} for $\mathbf{D}^{\text{train}}$ and \mathbf{D}^{test} according to the missing ratio r_X , along with the missing 403 indicator vector $\mathbf{V}^{\text{train}}$ for $\mathbf{Y}^{\text{train}}$ based on r_Y . For all methods, we employ $\{D_{il}^{\text{train}} | M_{il}^{\text{train}} = 1\}$ and 404 $\{Y_i^{\text{train}}|V_i^{\text{train}}=1\}$ as input for training. In the testing phase, we use $\{D_{il}^{\text{test}}|M_{il}^{\text{test}}=1\}$ as input to predict all missing response values $\{\hat{Y}_i^{\text{test}}|V_i^{\text{test}}=0\}$. Finally, we evaluate the response variable using root mean square error (RMSE) between Y_i^{test} and \hat{Y}_i^{test} for all i under $V_i^{\text{test}}=0$. 405 406 407

408 SHT-GNN configurations. We train SHT-GNN for 20 sampling phases with a sampling size of 200. 409 For each sampled graph, we run 1500 training epochs using the Adam optimizer with a learning rate 410 of 0.001 (Kingma & Ba, 2024). We employ a three-layer bipartite graph and two-layer longitudinal 411 subnetworks for all subjects. We use the ReLU activation function as the non-linear activation 412 function. The dimensions of both node embeddings and edge embeddings are set to 32. The message aggregation function AGG_l is implemented as a mean pooling function MEAN(\cdot). Both **O**_{immute} 413 and $\mathbf{O}_{predict}$ are implemented as multi-layer perceptrons (MLP) with 32 hidden units. The λ in loss 414 function is set to 0.001. 415

416 Baseline implementation. For Dicision Tree, GRAPE and IGRM, we directly conduct response 417 prediction under missing covariate matrix. For other baselines, as no end-to-end response prediction 418 approach is available, we first perform covariate imputation using the baselines, followed by utilizing 419 a Multilayer Perceptron (MLP) as the prediction model. To ensure a fair comparison, we apply the 420 same dimensional settings for these MLPs as those used in SHT-GNN.

421 Results in medium-dimensional covariates and moderate missing ratios. Following the experi-422 mental procedure described earlier, we set the covariate dimension p to 50. For the window size w and 423 the variance of random fluctuation $\sigma(\epsilon)$, we use the parameter combinations $\{w = 3, \sigma(\epsilon) = 0.1\}$, $\{w = 5, \sigma(\epsilon) = 0.1\}$, and $\{w = 7, \sigma(\epsilon) = 0.2\}$. For each configuration, we first split the subjects 424 with a test ratio of r = 0.2, then apply two missing ratio settings to the covariate matrix and response 425 vector: $\{r_X = 0.3, r_Y = 0.3\}$ and $\{r_X = 0.3, r_Y = 0.5\}$. We run All methods for 5 random 426 trials per setting, and report the average RMSE along with its standard deviation for the response 427 prediction on the test set. As shown in Table 1, SHT-GNN outperforms all baselines across all settings, 428 achieving an average reduction of 17.5% in prediction RMSE compared to the best baseline. Across 429 all methods, performance noticeably decline as the temporal smoothing window size and variance 430 increase, particularly for LOCF, Decision Tree, and IGRM. This indicates that the various temporal 431 smoothing characteristics in longitudinal data pose significant challenges for these methods. By

Missing ratio	$r_X = 0.3, r_Y = 0.3$			$r_X = 0.3, r_Y = 0.5$		
Window size Variance	w = 3 $\sigma = 0.1$	w = 5 $\sigma = 0.15$	w = 7 $\sigma = 0.2$	w = 3 $\sigma = 0.1$	w = 5 $\sigma = 0.15$	w = 7 $\sigma = 0.2$
Mean	$0.693 {\pm} 0.009$	$0.826 {\pm} 0.011$	$0.936 {\pm} 0.011$	$0.820 {\pm} 0.012$	$0.833 {\pm} 0.012$	1.051 ± 0.016
LOCF	$0.786 {\pm} 0.021$	$0.813 {\pm} 0.017$	$0.903 {\pm} 0.028$	$0.772 {\pm} 0.024$	$0.787 {\pm} 0.018$	$0.920{\pm}0.021$
MICE	$0.724{\pm}0.051$	$0.851 {\pm} 0.042$	$0.978 {\pm} 0.038$	$0.825 {\pm} 0.051$	$0.863 {\pm} 0.052$	$0.920{\pm}0.041$
3D MICE	$0.689 {\pm} 0.031$	$0.753 {\pm} 0.037$	$0.847 {\pm} 0.029$	$0.741 {\pm} 0.035$	$0.785 {\pm} 0.037$	$0.883{\pm}0.021$
GRAPE	$0.671 {\pm} 0.013$	$0.786 {\pm} 0.020$	$0.865 {\pm} 0.034$	$0.765 {\pm} 0.013$	$0.799 {\pm} 0.027$	$0.935 {\pm} 0.037$
CATSI	$0.701 {\pm} 0.034$	$0.732 {\pm} 0.026$	$0.832 {\pm} 0.023$	$0.749 {\pm} 0.047$	$0.748 {\pm} 0.038$	$0.885 {\pm} 0.027$
IGRM	$0.682 {\pm} 0.015$	$0.768 {\pm} 0.011$	$0.874 {\pm} 0.013$	$0.786 {\pm} 0.034$	$0.831 {\pm} 0.020$	$0.928 {\pm} 0.012$
GP-VAE	$0.733 {\pm} 0.011$	$0.793 {\pm} 0.018$	$0.851 {\pm} 0.021$	$0.731 {\pm} 0.023$	$0.769 {\pm} 0.025$	$0.933 {\pm} 0.030$
Transformer	0.611±0.022	0.691±0.023	0.791±0.013	0.678±0.013	0.718±0.019	0.873±0.021
CTA	0.581±0.010	0.678 ± 0.010	0.778±0.019	0.653±0.015	0.673±0.013	0.835±0.019
Our Method	$0.552{\pm}0.011$	$0.650 {\pm} 0.014$	$0.759{\pm}0.018$	$0.623{\pm}0.013$	$0.653{\pm}0.018$	$0.818 {\pm} 0.020$

Table 1: Performance comparison with different methods under varying temporal smoothing window
 sizes and covariate missing ratios. All RMSE values are 0.1 of the actual values.

integrating inductive learning with temporal smoothing, SHT-GNN demonstrate stable and superior performance in these scenarios.

Results in high-dimensional covariates and high missing ratios. We also compare the performance of different methods under higher covariate dimensions and missing ratios. As shown in Table 4 (Appendix A.3), SHT-GNN consistently outperforms all baselines across different settings, achieving an average reduction of 18% in prediction RMSE compared to the best baseline.

454 455 456

457

448

449

450 451

452

453

5.3 ALZHEIMER'S DISEASE NEUROIMAGING INITIATIVE STUDY DATASET

ADNI dataset introduction. We apply SHT-GNN to the real data from Alzheimer's Disease Neuroimaging Initiative (ADNI) study. We propose our model to predict the CSF biomarker Amyloid beta 42/40 ($A\beta$ 42/40), which has been demonstrated as a crucial biomaker in ADNI study. ADNI dataset containing 1,153 subjects and 10,033 observations. The covariate matrix has 83 dimensions, with a missing ratios of 0.32 for covariate matrix and a missing ratio of 0.83 for $A\beta$ 42/40. More details about ADNI dataset can be found in Appendix A.5,

464
 465
 466
 466
 467
 468
 468
 469
 460
 460
 460
 460
 460
 461
 461
 462
 463
 464
 464
 465
 464
 465
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466

Validation on observed response 467 value. We use RMSE to validate the 468 prediction accuracy. As shown in Ta-469 ble 2, our proposed SHT-GNN signif-470 icantly outperforms other methods in 471 predicting $A\beta 42/40$, with an mean 472 15.5% improvement in accuracy com-473 pared to the best baseline. 474

Validation on diagnostic labels. In 475 the previous section, we evaluated the 476 $A\beta 42/40$ prediction results based on 477 the observed ground truth. However, 478 most observations in ADNI do not 479 have recorded $A\beta 42/40$ values, mak-480 ing it impossible to validate the pre-481 dictions for these observations using 482 RMSE. In the ADNI study, the ultiTable 2: Performance comparison in ADNI dataset. For all methods, we conduct 5-fold cross-validation and report the mean values and standard deviations of the results.

Method	RMSE	AUC	Accuracy
Mean	$0.112 {\pm} 0.002$	$0.671 {\pm} 0.009$	$0.682{\pm}0.008$
LOCF	$0.108 {\pm} 0.003$	$0.706 {\pm} 0.008$	$0.717 {\pm} 0.015$
MICE	$0.109 {\pm} 0.004$	$0.717 {\pm} 0.021$	$0.701 {\pm} 0.021$
3D MICE	$0.103 {\pm} 0.004$	$0.731 {\pm} 0.017$	$0.721{\pm}0.018$
GRAPE	$0.106 {\pm} 0.002$	$0.724{\pm}0.010$	$0.714{\pm}0.011$
CATSI	$0.104{\pm}0.003$	$0.713 {\pm} 0.017$	$0.708 {\pm} 0.013$
IGRM	$0.105 {\pm} 0.002$	$0.721 {\pm} 0.013$	$0.713 {\pm} 0.009$
Transformer	0.110 ± 0.005	0.735 ± 0.021	0.721±0.020
GP-VAE	0.112 ± 0.007	0.738 ± 0.011	0.719±0.017
CTA	0.101±0.003	0.751±0.013	0.721±0.016
Our Method	$\textbf{0.090}{\pm}\textbf{0.001}$	$0.829{\pm}0.012$	$0.782{\pm}0.011$

mate goal is to predict disease progression, so we validate the predictions by employing them to classify the diagnostic labels. Each observation is associated with a binary diagnostic label: either AD or no-AD. We combine the predicted $A\beta 42/40$ values with fully observed demographic features and build a logistic model for classification. The classification performance can reflect the performance

486 of $A\beta 42/40$ predictions. As shown in Table 2, using SHT-GNN-predicted $A\beta 42/40$ combined with 487 basic individual features yields an AUC of 82.92% for the no-AD vs. AD classification. 488

5.4 EXPERIMENTS ON MORE REAL DATASETS

We compared the performance of our method on more longitudinal data and time series, including 492 the AirQuality, Electricity, Energy and PhysioNet ICU datasets. The results in Appendix A.7 shows SHT-GNN's efficiency in borrowing information across features and observations, which 494 significantly enhances its performance, particularly in the context of irregular longitudinal data imputation. However, our method does not show a significant advantage in one-dimensional time 495 series imputation. Regarding performance under different missing mechanisms, apart from the 496 self-masked MCAR, our method still performs excellently under MAR conditions in ADNI and PhysioNet ICU. Additionally, we provide a theoretical understanding in Appendix A.9 that explains 498 the intuition behind our method and its applicability under MCAR and MAR. 499

500

497

489

490 491

493

501

5.5 ABLATION STUDY ON THE LAYER COUNT OF LONGITUDINAL SUBNETWORKS

502 gitudinal subnetworks allow observa-504 tion nodes to borrow information from 505 previous time steps. To validate that 506 longitudinal subnetworks can borrow 507 longer-term historical information by 508 stacking layers, we applied SHT-GNN 509 with one, two and three layers' longitudinal subnetworks under data simu-510 lated under temporal smoothing win-511 dow sizes w = 3, w = 5, and w = 7. 512

As shown in Figure 3, multi-layer lon- Table 3: Performance comparison of longitudinal subnetworks with different layer count across different window sizes. All RMSE values are 0.1 of the actual values.

Window size	w = 3	w = 5	w = 7
None	$0.673 {\pm} 0.019$	$0.757 {\pm} 0.015$	$0.933{\pm}0.018$
One-Layer	$0.621 {\pm} 0.012$	$0.739{\pm}0.015$	$0.889 {\pm} 0.017$
Two-Layer	$0.593 {\pm} 0.015$	$0.687 {\pm} 0.019$	$0.827 {\pm} 0.020$
Three-Layer	$0.552{\pm}0.011$	$0.659 {\pm} 0.014$	$0.763 {\pm} 0.018$

In all cases, we set the subject test ratio r to 0.2, along with the missing rate of both the covariate 513 matrix and the response vector to 0.3. As shown in Table 3, as w increases, SHT-GNN with three-layer 514 subnetworks outperform the one-layer and two-layer versions by a larger margin. This suggests that 515 additional layers in longitudinal subnetworks enhance performance in data with longer-span tem-516 poral smoothing, which demonstrating the multi-layer longitudinal networks designed for temporal 517 smoothing are effective. Additionally, we conduct an ablation study for the MADGap component in 518 SHT-GNN. The results and more details can be found in Appendix A.4.

519 520 521

5.6 SCALABILITY OF SHT-GNN

522 In longitudinal data, a large number of repeated observations often lead to massive data size, making 523 scalability a critical concern for missing data imputation methods. Specifically, we compare the 524 scalability of SHT-GNN with two cutting-edge GNN-based imputation methods: GRAPE and IGRM. 525 For experimental simplicity, we conduct imputation on datasets with varying observation sizes of $N \times 10$ (where N = 50, 500, 1000, and 5000). In SHT-GNN, we fix the subsample size at 500 subjects, 526 resulting in a total of 5000 observations. For all three methods, we report the memory consumption 527 and forward computation time per epoch during training. The results show that both GRAPE and 528 IGRM exhibit significant increases in memory usage and computation time as observation size grows. 529 However, by fixing the sampled subject batch size, SHT-GNN achieves consistent memory usage and 530 computation time per epoch, regardless of the observation size. 531

532

CONCLUSION 6

533 534

535 In this paper, we present SHT-GNN, a scalable and accurate framework for longitudinal data impu-536 tation. Our method combines a sampling-guided training policy with inductive learning, temporal 537 smoothing, and a custom-designed loss function to address the challenges posed by irregular and inconsistent longitudinal missing data. Compared to state-of-the-art imputation techniques, SHT-GNN 538 consistently improves response prediction across both synthetic and real-world datasets. Extensive experiments confirm the efficacy of our multi-layer longitudinal subnetwork for temporal smoothing.

⁵⁴⁰ 7 REPRODUCIBILITY STATEMENT

In this paper, we used the GLOBEM and ADNI datasets for our experiments. For the GLOBEM dataset, access can be requested at https://physionet.org/content/globem/, which is divided into four groups of subjects. In this study, we use the original data without any preprocessing. For the ADNI dataset, access can be requested at https://ida.loni.usc.edu/. The ADNI study is divided into different phases, and we select those subjects with complete basic genetic observations in the ADNI 1, ADNI 2, and ADNI GO phases. In the supplementary materials, we provide the complete code for the proposed SHT-GNN. The original codes for other baselines can be found through the links to the referenced papers in the Section. In this study, all models are trained on a Windows 10 64-bit OS (version 19045) with 32GB of RAM (AMD Ryzen 7 4800H CPU @ 2.9GHz) and 4 NVIDIA GeForce RTX 2060 GPUs with Max-Q Design.

References

- Wei Cao, Dong Wang, Jian Li, Hao Zhou, Lei Li, and Yitan Li. Brits: Bidirectional recurrent imputation for time series. *Advances in neural information processing systems*, 31, 2018.
- ⁵⁵⁷ Deli Chen, Yankai Lin, Wei Li, Peng Li, Jie Zhou, and Xu Sun. Measuring and relieving the over ⁵⁵⁸ smoothing problem for graph neural networks from the topological view. In *Proceedings of the* ⁵⁵⁹ AAAI conference on artificial intelligence, volume 34, pp. 3438–3445, 2020.
- 561 Mark Collier, Alfredo Nazabal, and Christopher KI Williams. Vaes in the presence of missing data. *arXiv preprint arXiv:2006.05301*, 2020.
 - Michael J Daniels and Joseph W Hogan. *Missing data in longitudinal studies: Strategies for Bayesian modeling and sensitivity analysis.* chapman and hall/CRC, 2008.
- Vincent Fortuin, Dmitry Baranchuk, Gunnar Rätsch, and Stephan Mandt. Gp-vae: Deep probabilistic
 time series imputation. In *International conference on artificial intelligence and statistics*, pp. 1651–1661. PMLR, 2020.
- Ge Gao, Farzaneh Khoshnevisan, and Min Chi. Reconstructing missing ehrs using time-aware within and cross-visit information for septic shock early prediction. In 2022 IEEE 10th International Conference on Healthcare Informatics (ICHI), pp. 151–162. IEEE, 2022.
- 573 Md Hamidul Huque, John B Carlin, Julie A Simpson, and Katherine J Lee. A comparison of multiple
 574 imputation methods for missing data in longitudinal studies. *BMC medical research methodology*, 18:1–16, 2018.
- 577 Mina Jahangiri, Anoshirvan Kazemnejad, Keith S Goldfeld, Maryam S Daneshpour, Shayan
 578 Mostafaei, Davood Khalili, Mohammad Reza Moghadas, and Mahdi Akbarzadeh. A wide range
 579 of missing imputation approaches in longitudinal data: a simulation study and real data analysis.
 580 BMC Medical Research Methodology, 23(1):161, 2023.
 - Maksims Kazijevs and Manar D Samad. Deep imputation of missing values in time series health data: A review with benchmarking. *Journal of biomedical informatics*, pp. 104440, 2023.
- 584 DP Kingma and J Ba. Adam: A method for stochastic optimization." arxiv, jan. 29, 2017. *Accessed: Feb*, 6, 2024.
- Peter Lane. Handling drop-out in longitudinal clinical trials: a comparison of the locf and mmrm approaches. *Pharmaceutical Statistics: The Journal of Applied Statistics in the Pharmaceutical Industry*, 7(2):93–106, 2008.
- Jia Li, Jiajin Li, Yang Liu, Jianwei Yu, Yueting Li, and Hong Cheng. Deconvolutional networks on
 graph data. Advances in Neural Information Processing Systems, 34:21019–21030, 2021.
- 593 Roderick JA Little and Donald B Rubin. *Statistical analysis with missing data*, volume 793. John Wiley & Sons, 2019.

- 594 Yuan Luo, Peter Szolovits, Anand S Dighe, and Jason M Baron. 3d-mice: integration of cross-595 sectional and longitudinal imputation for multi-analyte longitudinal clinical data. Journal of the 596 American Medical Informatics Association, 25(6):645–653, 2018. 597 Pierre-Alexandre Mattei and Jes Frellsen. Miwae: Deep generative modelling and imputation of 598 incomplete data sets. In International conference on machine learning, pp. 4413–4423. PMLR, 2019. 600 601 Ariel I Mundo, John R Tipton, and Timothy J Muldoon. Using generalized additive models to analyze 602 biomedical non-linear longitudinal data. bioRxiv, pp. 2021-06, 2021. 603 Joel Oskarsson, Per Sidén, and Fredrik Lindsten. Temporal graph neural networks for irregular data. 604 In International Conference on Artificial Intelligence and Statistics, pp. 4515–4531. PMLR, 2023. 605 Marek Pekarčík, Júlia Durčová, and Jozef Glova. Intangible ict and their importance within global 607 value chains: An empirical analysis based on longitudinal data regression. *Mathematics*, 10(7): 608 1198, 2022. 609 610 Anand Rajaraman. Mining of massive datasets, 2011. 611 Adam Sadowski, Karolina Lewandowska-Gwarda, Renata Pisarek-Bartoszewska, and Per Engelseth. 612 A longitudinal study of e-commerce diversity in europe. Electronic Commerce Research, 21(1): 613 169-194, 2021. 614 615 Stef Van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equations 616 in r. Journal of statistical software, 45:1–67, 2011. 617 Hyowon Wi, Yehjin Shin, and Noseong Park. Continuous-time autoencoders for regular and irregular 618 time series imputation. In Proceedings of the 17th ACM International Conference on Web Search 619 and Data Mining, pp. 826-835, 2024. 620 621 J Woillez, JA Abad, R Abuter, E Aller Carpentier, J Alonso, L Andolfato, P Barriga, J-P Berger, J-L 622 Beuzit, H Bonnet, et al. Naomi: the adaptive optics system of the auxiliary telescopes of the vlti. 623 Astronomy & Astrophysics, 629:A41, 2019. 624 Xiao Xu, Xiaoshuang Liu, Yanni Kang, Xian Xu, Junmei Wang, Yuyao Sun, Quanhe Chen, Xiaoyu 625 Jia, Xinyue Ma, Xiaoyan Meng, et al. A multi-directional approach for missing value estimation in 626 multivariate time series clinical data. Journal of Healthcare Informatics Research, 4(4):365–382, 627 2020. 628 629 Kejing Yin, Liaoliao Feng, and William K Cheung. Context-aware time series imputation for 630 multi-analyte clinical data. Journal of Healthcare Informatics Research, 4:411–426, 2020. 631 Jiaxuan You, Xiaobai Ma, Yi Ding, Mykel J Kochenderfer, and Jure Leskovec. Handling missing 632 data with graph representation learning. Advances in Neural Information Processing Systems, 33: 633 19075-19087, 2020. 634 635 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series 636 forecasting? In Proceedings of the AAAI conference on artificial intelligence, volume 37, pp. 637 11121-11128, 2023. 638 Muhan Zhang and Yixin Chen. Inductive matrix completion based on graph neural networks. arXiv 639 preprint arXiv:1904.12058, 2019. 640 641 Weiqi Zhang, Guanlve Li, Jianheng Tang, Jia Li, and Fugee Tsung. Missing data imputation with 642 graph laplacian pyramid network. arXiv preprint arXiv:2304.04474, 2023. 643 Xiang Zhang, Marko Zeman, Theodoros Tsiligkaridis, and Marinka Zitnik. Graph-guided network 644 for irregularly sampled multivariate time series. arXiv preprint arXiv:2110.05357, 2021. 645 646 Jiajun Zhong, Ning Gui, and Weiwei Ye. Data imputation with iterative graph reconstruction. In 647 Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 11399–11407, 2023.
 - 12

A APPENDIX

A.1 RESPONSE VECTOR SIMULATION

In the experiment involving medium-dimensional covariates and moderate missing ratios, we simulate response variables from a 50-dimensional covariate matrix for each observation. Specifically, we use a response simulation model that includes both linear and nonlinear components. The default model for simulating the response variables is as follows:

$$y = 0.25 \cdot x_1 + 2 \cdot \left(\frac{\log(x_2 + 10)}{25}\right)^2 - 0.4 \cdot x_3 - 0.15 \cdot \left(x_4 + 5 \cdot e^{-5(1.5 - \log(x_5))^2/2}\right) - 0.25 \cdot \log(x_6 + 1) + 0.4 \cdot x_7 + 0.021 \cdot \sin(x_8) + 0.04 \cdot \sqrt{x_9} + 0.1 \cdot e^{x_{10}} + 0.05 \cdot \log(x_{11} + 1) + 0.02 \cdot \tan(x_{12}) + 0.015 \cdot \cos(x_{13}) + 0.07 \cdot \log(x_{14} + 1) + 3.5 \cdot \sqrt{x_{15}} + \epsilon$$

 Each x_i is randomly selected without repetition from all the covariates. Before simulating the response values based on the observed covariates, all covariate values are normalized using the MinMax Scaler (Rajaraman, 2011). The term ϵ represents random noise following a normal distribution N(0, 0.175). In the experiment involving high-dimensional covariates and moderate missing ratios, we simulate response variables from a 100-dimensional covariate matrix for each observation. We use a response simulation model with higher-order inputs and more complex expressions, as shown below:

$$\begin{array}{ll} \begin{array}{ll} & y=0.3\cdot\sqrt{x_{1}}-0.4\cdot x_{2}^{2}+0.15\cdot \log (x_{3}+10^{-6})+0.2\cdot \exp (0.5\cdot x_{4})-0.1\cdot x_{5}\\ & +0.05\cdot \sin (2\pi\cdot x_{6})+0.25\cdot \log 1p(x_{7})-0.1\cdot \cos (2\pi\cdot x_{8})+0.35\cdot \tan (\operatorname{clip}(x_{9},-0.5,0.5)) \\ & +0.05\cdot \operatorname{arcsin}(\operatorname{clip}(x_{10},-1,1))+0.2\cdot x_{11}^{3}-0.3\cdot\sqrt{x_{12}}+0.4\cdot \frac{\log (x_{13}+1)}{10}\\ & +0.15\cdot \sin (2\pi\cdot x_{14})-0.1\cdot \log 1p(x_{15})+0.1\cdot \exp (x_{16})-0.05\cdot \log (x_{17}+1) \\ & +0.2\cdot x_{18}^{2}+0.3\cdot \cos (x_{19})-0.07\cdot \tan (x_{20})+0.05\cdot \frac{\sin (x_{21})}{x_{22}+1}\\ & +0.25\cdot \log 1p(x_{23})+0.15\cdot \operatorname{arcsin}(\operatorname{clip}(x_{24},-1,1))+0.1\cdot x_{25}^{3}-0.05\cdot \sqrt{x_{26}}\\ & +0.07\cdot \log (x_{27}+1)+0.2\cdot \frac{\tan (x_{28})}{1+x_{29}^{2}}-0.1\cdot \exp (x_{30})\\ & +0.3\cdot \log (x_{31}+10)+0.25\cdot x_{32}+\epsilon. \end{array}$$

Similarly, each x_i is randomly selected without repetition from all the covariates. Before simulating the response values based on the observed covariates, all covariate values are normalized using the MinMax Scaler. The term ϵ represents random noise following a normal distribution N(0, 0.2).

A.2 DETAILS FOR BASELINES

Mean: For each covariate with missing, we fill in the missing values using the mean of the observed values for that covariate across all observations.

Copy-mean LOCF: Following the process described in (Jahangiri et al., 2023), missing values are initially imputed using the LOCF (Last Observation Carried Forward) method within each subject to provide an approximation. Next, the population's mean trajectory is used to further refine these imputed values.

MICE: As outlined in (Van Buuren & Groothuis-Oudshoorn, 2011), MICE (Multiple Imputation
 by Chained Equations) performs multiple imputations by modeling each missing value conditioned
 on the non-missing values in the data. A maximum of 20 iterations is used during the imputation
 process.

3D-MICE: Following the procedure described in (Luo et al., 2018), MICE is configured to perform cross-sectional imputation with a maximum of 20 iterations. And Gaussian Process Regression is applied longitudinally to the time-indexed data for each feature. The Gaussian Process uses an RBF

kernel combined with a constant kernel (Kernel = $C(1.0) \times RBF(1.0)$), and the predictions from the Gaussian Process are averaged with the MICE-imputed values to capture both temporal and cross-sectional patterns.

GRAPE: Following the setup in (You et al., 2020), GRAPE is trained for 20,000 epochs using the Adam optimizer with a learning rate of 0.001. We employ two GNN layers with 16 hidden units and ReLU activation. The AGG_l function is implemented as a mean pooling function MEAN(\cdot). Both the edge imputation and response prediction neural networks are implemented as linear layers.

CATSI: Following the setup in (Yin et al., 2020), CATSI is trained for 3,000 epochs using the Adam optimizer with a learning rate of 0.001. We employ the default settings for MLPs and LSTM models in CASTI to impute covariate values.

713IGRM: Following the default settings described in (Zhong et al., 2023), IGRM employs three714GraphSAGE layers with 64 hidden units for bipartite Graph Representation Learning (GRL), and one715GraphSAGE layer for friend network GRL. The Adam optimizer with a learning rate of 0.001 and716ReLU activation function is used. For initializing the friend network, observation nodes are randomly717connected with |U| edges to form the initial network structure, which is updated every 100 epochs718during bipartite graph training.

719 720

729

730

A.3 RESULTS IN HIGH-DIMENSIONAL COVARIATES AND HIGH MISSING RATIOS

721 To further compare the performance of different methods under higher covariate dimensions and 722 higher missing ratios, we set the covariate dimension p to 100. The same settings for temporal 723 smoothing window size w and the variance of random fluctuation $\sigma(\epsilon)$ were applied as before. 724 Additionally, two higher missing ratio settings were employed: $\{r_X = 0.5, r_Y = 0.5\}$ and $\{r_X = 0.5, r_Y = 0.5\}$ 725 $0.5, r_Y = 0.7$. All methods were run for 5 random trials per setting, and the average RMSE 726 of response prediction on the test set are recorded. As shown in Table 4, SHT-GNN consistently 727 outperforms all baselines across all settings, achieving an average reduction of 18% in prediction RMSE compared to the best baseline. 728

Table 4: Performance comparison with different methods under varying temporal smoothing windows and covariate missing ratios. All RMSE values are 0.1 of the actual values.

Missing ratio	$r_X = 0.5, r_Y = 0.5$			$r_X = 0.5, r_Y = 0.7$			
Window size Variance	w = 3 $\sigma = 0.1$	$\begin{array}{c} w=5\\ \sigma=0.15 \end{array}$	w = 7 $\sigma = 0.2$	w = 3 $\sigma = 0.1$	$\begin{array}{c} w=5\\ \sigma=0.15 \end{array}$	w = 7 $\sigma = 0.2$	
Mean	$0.803 {\pm} 0.015$	$0.907 {\pm} 0.022$	$1.054{\pm}0.023$	$0.908 {\pm} 0.019$	$0.893 {\pm} 0.020$	$0.975 {\pm} 0.027$	
LOCF	$0.751 {\pm} 0.029$	$0.813 {\pm} 0.035$	$0.938 {\pm} 0.035$	$0.825 {\pm} 0.023$	$0.823 {\pm} 0.019$	$0.937 {\pm} 0.019$	
MICE	$0.797 {\pm} 0.043$	$0.848 {\pm} 0.041$	$1.014 {\pm} 0.059$	$0.837 {\pm} 0.036$	$0.893 {\pm} 0.041$	$0.974 {\pm} 0.065$	
3D MICE	$0.745 {\pm} 0.021$	$0.773 {\pm} 0.047$	$0.948 {\pm} 0.045$	$0.795 {\pm} 0.036$	$0.793 {\pm} 0.046$	$0.902{\pm}0.046$	
DT	$0.763 {\pm} 0.017$	$0.798 {\pm} 0.018$	$0.981 {\pm} 0.020$	$0.781 {\pm} 0.019$	$0.794{\pm}0.027$	$0.891 {\pm} 0.031$	
GRAPE	$0.724 {\pm} 0.026$	$0.793 {\pm} 0.015$	$0.952{\pm}0.021$	$0.809 {\pm} 0.021$	$0.745 {\pm} 0.021$	$0.944{\pm}0.021$	
CATSI	$0.831 {\pm} 0.029$	$0.725 {\pm} 0.021$	$0.945 {\pm} 0.041$	$0.791 {\pm} 0.031$	$0.755 {\pm} 0.041$	$0.903{\pm}0.031$	
IGRM	$0.795 {\pm} 0.010$	$0.831 {\pm} 0.013$	$0.904{\pm}0.014$	$0.785 {\pm} 0.012$	$0.815 {\pm} 0.019$	$0.895 {\pm} 0.021$	
Our Method	$0.632{\pm}0.010$	$0.672 {\pm} 0.017$	$0.821{\pm}0.014$	$0.658 {\pm} 0.021$	$0.641 {\pm} 0.013$	$0.819 {\pm} 0.020$	

743 744

745 746

A.4 ABLATION STUDY FOR MADGAP IN SHT-GNN

747 In SHT-GNN, the multi-layer longitudinal subnetworks are designed for temporal smoothing. How-748 ever, the degree of smoothing among observations for the same subject may vary across different 749 longitudinal studies. The loss function incorporates MADGap to promote greater variance amonng 750 observation representations for a given subject, enabling the model to effectively trade off between 751 temporal smoothing and representation diversity. To evaluate the impact of MADGap, we test 752 SHT-GNN with and without MADGap under different temporal smoothing window sizes. As shown 753 in Table 5, the results shows that incorporating MADGap enhances SHT-GNN's performance by an average of 6%. Moreover, as the degree of temporal smoothing increases, the design incorporating 754 MADGap exhibited a more substantial advantage over the one without it. This highlights MADGap's 755 capacity to help SHT-GNN capture the unique characteristics of each observation during imputation.

59		SHT-GNN wit	hout MADGap	SHT-GNN w	ith MADGap	
60		$r_{\rm v} = 0.3$	$r_{\rm V} = 0.5$	$r_{\rm v} = 0.3$	$r_{\rm V} = 0.5$	Enhancement
51		$r_X = 0.0$ $r_Y = 0.3$	$r_X = 0.5$ $r_Y = 0.5$	$r_X = 0.3$ $r_Y = 0.3$	$r_X = 0.5$ $r_Y = 0.5$	Linuncement
3	Window size 3	0.591±0.011	0.691±0.013	$0.552{\pm}0.011$	$0.632 {\pm} 0.010$	9.68%
.Л	Window size 5	$0.687 {\pm} 0.012$	$0.718 {\pm} 0.011$	$0.651 {\pm} 0.014$	$0.672 {\pm} 0.017$	5.57%
65	Window size 7	$0.801 {\pm} 0.011$	$0.865 {\pm} 0.011$	$0.769 {\pm} 0.018$	$0.821 {\pm} 0.014$	4.91%

Table 5: Performance comparison of SHT-GNN with and without MADGap across different temporal
 smoothing windows and missing ratios. All RMSE values are 0.1 of the actual values.

A.5 ADNI DATASET INTRODUCTION

769 We apply SHT-GNN to the real data from Alzheimer's Disease Neuroimaging Initiative (ADNI) 770 study. ADNI is a multi-centre longitudinal neuroimaging study with the aim of developing effective 771 treatments that can slow or halt the progression of Alzheimer's Disease (AD). The ADNI participants 772 were followed prospectively, with follow-up time points at 3 months, 6 months, then every 6 773 months until up to 156 months. The ADNI study includes a wide range of clinical data such as 774 cognitive assessments, magnetic resonance imaging (MRI) and cerebrospinal fluid (CSF) biomarkers. 775 Numerous studies show that CSF biomarkers are strong indicators of AD progression, but collecting CSF requires invasive procedures like lumbar puncture, leading to high missing data rates. We 776 propose the SHT-GNN model to predict the CSF biomarker Amyloid beta 42/40 ($A\beta 42/40$), which 777 has been a key biomaker. ADNI dataset containing 1,153 subjects and 10,033 observations. The 778 covariate matrix has 83 dimensions, with a missing ratios of 0.32 for covariate matrix and a missing 779 ratio of 0.83 for $A\beta 42/40$.

781 782

791

792 793

794

766 767

768

A.6 SHT-GNN CONFIGURATIONS FOR ADNI DATASET

We train SHT-GNN for 10 sampling phases with a sampling size of 200. For each sampled graph, we run 1500 training epochs using the Adam optimizer with a learning rate of 0.001. We employ a three-layer bipartite graph and two-layer longitudinal subnetworks for all subjects. We use the ReLU activation function as the non-linear activation function. The dimensions of both node embeddings and edge embeddings are set to 32. The message aggregation function AGG_l is implemented as a mean pooling function MEAN(·). Both **O**_{impute} and **O**_{predict} are implemented as multi-layer perceptrons (MLP) with 32 hidden units. The λ in loss function is set to 0.001.

A.7 EXPERIMENTS ON MORE REAL DATASETS

Table 6: Comparison of methods on different datasets

	AirQuality	Electricity	Energy	PhysioNet-2012
Transformer	0.220 ± 0.019	$0.889 {\pm} 0.071$	0.313±0.018	$0.190 {\pm} 0.019$
GP-VAE	$0.287 {\pm} 0.010$	$0.963 {\pm} 0.056$	$0.401 {\pm} 0.025$	$0.398 {\pm} 0.020$
CTA	$0.196{\pm}0.012$	$0.767 {\pm} 0.042$	$0.205 {\pm} 0.019$	$0.192{\pm}0.016$
SAITS	$0.201 {\pm} 0.009$	$0.894{\pm}0.051$	$0.301 {\pm} 0.012$	$0.190 {\pm} 0.014$
MICE	$0.310 {\pm} 0.023$	$1.319 {\pm} 0.051$	$0.371 {\pm} 0.010$	$0.223 {\pm} 0.021$
3D MICE	$0.293 {\pm} 0.009$	$1.083 {\pm} 0.051$	$0.341 {\pm} 0.011$	$0.209 {\pm} 0.015$
GRAPE	$0.267 {\pm} 0.013$	$0.891 {\pm} 0.029$	$0.251 {\pm} 0.009$	$0.203 {\pm} 0.005$
CATSI	$0.236 {\pm} 0.019$	$0.849 {\pm} 0.071$	$0.201 {\pm} 0.023$	$0.206 {\pm} 0.013$
IGRM	$0.242 {\pm} 0.010$	$0.867 {\pm} 0.045$	$0.231 {\pm} 0.013$	$0.193 {\pm} 0.014$
Our method	$0.212{\pm}0.010$	$0.834{\pm}0.025$	$0.183{\pm}0.011$	$0.187{\pm}0.009$

806 807

804 805

On the AirQuality and Electricity datasets, the SHT-GNN method demonstrates inferior performance
 compared to RNN and VAE-based approaches. This indicates that the SHT-GNN model is not well suited for long-term, single-object time series, which fall outside its intended application scenario.

However, SHT-GNN achieves state-of-the-art performance on the Energy dataset. Unlike the AirQuality and Electricity datasets, the Energy dataset involves time series with multiple subjects and multidimensional features. SHT-GNN leverages its ability to effectively borrow information across features, which significantly enhances its performance, particularly in the context of irregular longitudinal data imputation.

It is worth emphasizing that it can be observed that our method demonstrates a significant advantage on the PhysioNet dataset, which is also clinical longitudinal follow-up data like the ADNI dataset.



Figure 4: A comparison of the performance of missing data imputation across all methods on additional datasets.

A.8 SCALABILITY OF SHT-GNN



Figure 5: The comparison of scalability under different observation sizes across GNN-based methods.

A.9 THEORETICAL UNDERSTANDINGS AND INSIGHTS OF SHT-GNN

From an optimization perspective, the missing data imputation process in our method is conceptually similar to Variational Autoencoders (VAEs), where the goal of the reconstruction step is to minimize the reconstruction error as part of the Evidence Lower Bound (ELBO). As is widely recognized, the training objective of a standard VAE for missing data imputation is expressed as (Collier et al., 2020):

$$X_{obs} \xrightarrow{\text{Encode}} Z \xrightarrow{\text{Decode}} \hat{X}_{obs}$$

864 865 866

867 868

870

871

872

873

874 875

876

877 878 879

883

884

885

Maximize:
$$\log p(x_{obs}) = \int \log q(z|x_{obs}) \cdot \frac{p(z, x_{obs})}{q(z|x_{obs})} dz \ge \int q(z|x_{obs}) \log \frac{p(x_{obs}, z)}{q(z|x_{obs})} dz$$

That is to maximize: $E_{q(z|x_{obs})} \log p(x_{obs}|z) - D_{KL}[q(z|x_{obs}) \parallel p(z)]$

where $E_{q(z|x)} \log p(x|z)$ represents the reconstruction loss, and $D_{KL}[q(z|x) \parallel p(z)]$ is the regularization term. When maximizing $\log p(x_{obs})$, it is guaranteed that the estimated results for missing data will be consistent in both MCAR and MAR scenarios, a point that has been emphasized in many studies (Mattei & Frellsen, 2019; Collier et al., 2020).

In our proposed SHT-GNN, we are also theoretically optimizing $\log p(x_{obs})$. Specifically, the calculation and training process can be described as follows:

$$X_{\rm obs}, Z_O^{\rm init}, Z_F^{\rm init} \xrightarrow{\rm Message Passing, \, Embedding \, Update} \mathcal{G} X_O^L, Z_F^L \xrightarrow{\rm Edge-wise \, Prediction \, as \, Missing \, Data \, Imputation} \hat{X}_{\rm ob}$$

where X_{obs} denotes the observed values, Z_O^{init} and Z_F^{init} represent the initial embedding matrices of all observation and feature nodes, respectively. Z_O^L and Z_F^L denote the embedding matrices of all observation and feature nodes after L layers of forward computation in SHT-GNN. Subsequently, the training objective in SHT-GNN is expressed as:

$$\text{Maximize: } \log p(X_{obs}) = \int \log q(Z_O^L, Z_F^L | X_{obs}) \cdot \frac{p(Z_O^L, Z_F^L, X_{obs})}{q(Z_O^L, Z_F^L | X_{obs})} dZ \ge \int q(Z_O^L, Z_F^L | X_{obs}) \log \frac{p(X_{obs}, Z_O^L, Z_F^L)}{q(Z_O^L, Z_F^L | X_{obs})} dZ$$

889 890 891

892 893

894

888

That is to maximize : $E_{q(Z_O^L, Z_F^L | X_{obs})} \log p(X_{obs} | Z_O^L, Z_F^L) - D_{\mathrm{KL}}[p(Z_O^L, Z_F^L | X_{obs}) \parallel p(Z_O^L, Z_F^L)]$

where $E_{q(Z_O, Z_F | X_{obs})} \log p(X_{obs} | Z_O^L, Z_F^L)$ represents the reconstruction loss.

Here, $E_{q(Z_O, Z_F | X_{obs})} \log p(X_{obs} | Z_O^L, Z_F^L)$ is the joint distribution over all observations, which differs from the $E_{q(z|x_{obs})} \log p(x_{obs}|z)$ in VAE. Previously, in the case of VAE, their expectation is calculated on the sample level, typically using the MSE over the observed values of all samples to approximate the reconstruction loss. The specific form is:

899 900 Loss = $\frac{1}{N} \sum_{i=1}^{N} \sum_{j=1}^{p} m_{ij} \cdot (\hat{X}_{ij} - X_{ij})^2$

⁹⁰¹ where m_{ij} is the missing indicator for the *j*-th feature in the *i*-th observation.

In contrast, in SHT-GNN, our reconstruction loss is in the form of a joint distribution, and it is not possible to estimate it by averaging over the samples. This is why we use edge dropout trick (You et al., 2020) in missing data imputation, because directly using the loss over all observed edges will not provide an effective estimate of $E_{q(Z_O, Z_F|X_{obs})} \log p(X_{obs}|Z_O^L, Z_F^L)$. Specifically, we estimate the overall reconstruction loss by randomly calculating the loss on some edges in each batch of different heterogeneous graphs.

In principle, the SHT-GNN and VAE-based methods share conceptual similarities. Furthermore, the reconstruction loss used in both methods, along with the maximized target $p(X_{obs})$, theoretically indicates that our method is also capable of handling missing data in the MAR scenario, in the same way as VAE-based methods.

913

914

915

916

917