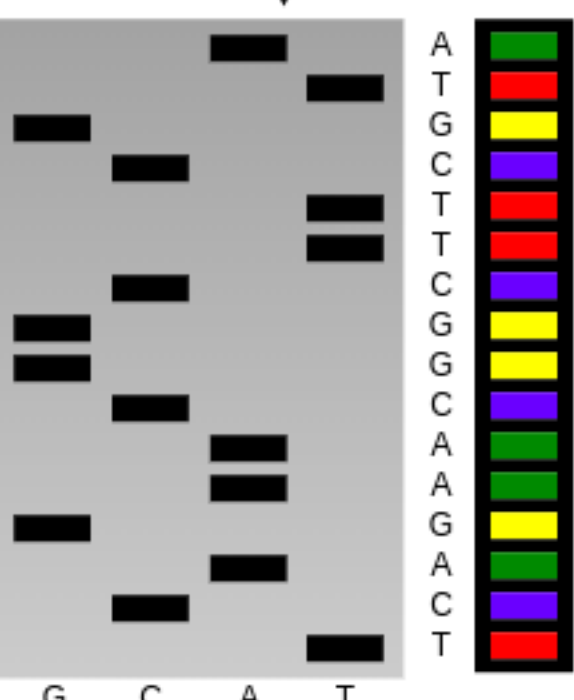
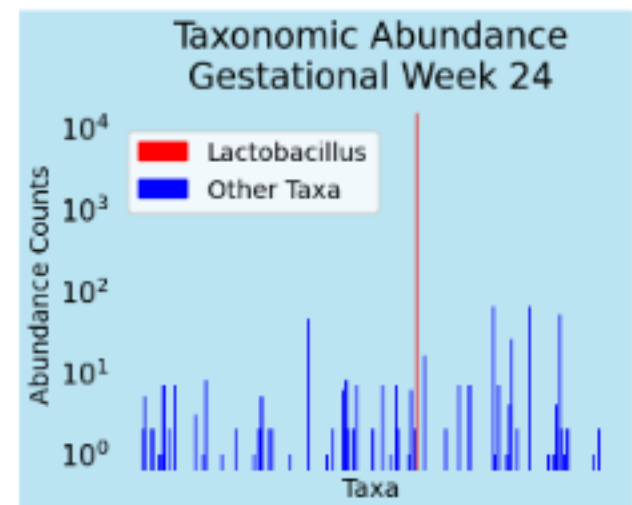
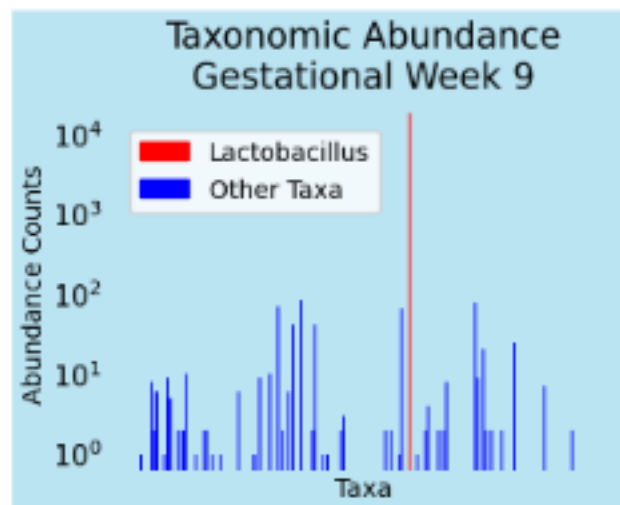




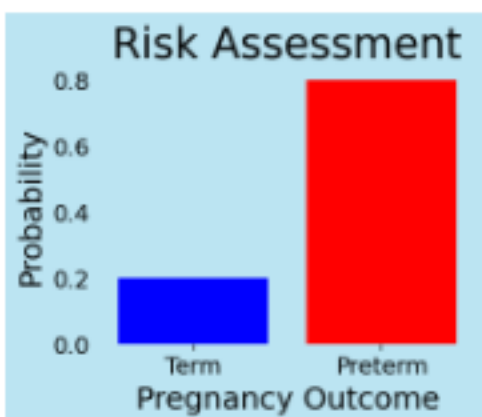
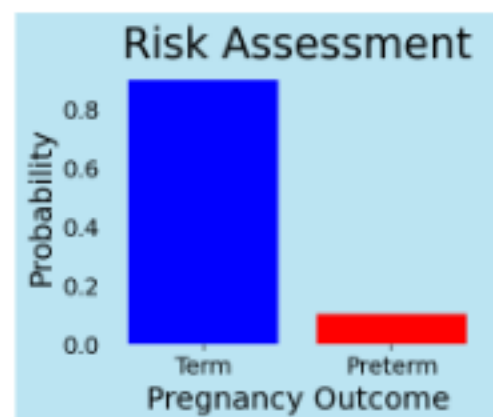
V4 Hypervariable Region



DADA2

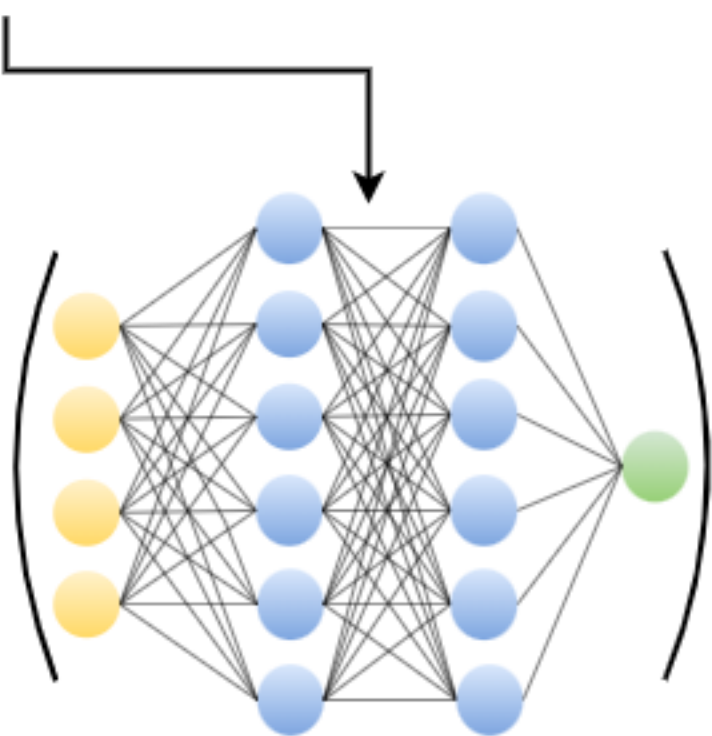


Week-Wise Taxonomic Abundance



Term v/s Preterm Prediction

$$\frac{d}{dt}$$



**Neural Differential
Equations**

1
2
3
4
5
6
7
8
9 **Deep Learning Enables Early-Stage Prediction of Preterm**
10 **Birth Using Vaginal Microbiota**

11
12
13 Kaushik Karambelkar^{1, #}

Mayank Baranwal^{1, 2, *}

14
15 ¹Data and Decision Sciences, Tata Consultancy Services Ltd., Thane, Maharashtra, India – 400607

16 ²Systems and Control Engineering, Indian Institute of Technology, Bombay, Maharashtra, India – 400076

17 #Email (Institutional): kaushik.karambelkar@tcs.com

18 *Corresponding Author; Email (Institutional): baranwal.mayank@tcs.com, mbaranwal@iitb.ac.in

19
20
21
22 **Time-Series Modeling of Vaginal Microbiota**

23
24 **Keywords:** Preterm Birth, Deep Learning, Vaginal Microbiome, Time-Series Modeling

Abstract

Objective: Preterm birth (PTB) is one of the leading issues concerning infant health and is a problem that plagues all parts of the world. Vaginal microbial communities have recently garnered attention in the context of PTB, however, the vaginal microbiome varies greatly from individual to individual, and this variation is more pronounced in racially, ethnically and geographically diverse populations. Additionally, microbial communities have been reported to evolve during the duration of the pregnancy, and capturing such a signature may require higher, more complex modeling paradigms. In this study, we develop a neural controlled differential equations (CDEs) based framework for identifying early PTBs in racially diverse cohorts from irregularly sampled vaginal microbial abundance data.

Methods: We obtained relative abundances of microbial species within vaginal microbiota using 16S rRNA sequences obtained from vaginal swabs at various stages of pregnancy. We employed a recently introduced deep learning paradigm known as “Neural CDEs” to predict PTBs. This method, previously unexplored, analyzes irregularly sampled microbial abundance profiles in a time-series format.

Results: Our framework is able to identify signatures in the temporally evolving vaginal microbiome during trimester 2 and can predict incidences of PTB (mean test set ROC-AUC = 0.81, accuracy = 0.75, f1 score = 0.71) significantly better than traditional ML classifiers, thus enabling effective early-stage PTB risk assessment.

Conclusion and Significance: Our method is able to differentiate between term and preterm outcomes with a substantial accuracy, despite being trained using irregularly sampled microbial abundance profiles, thus overcoming the limitations of traditional time-series modeling methods.

1. Introduction

Preterm births (PTBs) are live births that occur before 37 weeks of pregnancy and are a major public health concern worldwide. It is estimated that about 15 million babies are born pre-term globally each year, putting the global PTB rate at about 11% (Blencowe et al., 2012). PTB is among the leading causes of neonatal mortality and morbidity, especially in low- and middle-income economies (Perin et al., 2022). The global PTB rate is also on a steady rise, thus making it a significant burden (Chawanpaiboon et al., 2019). Approximately 18% of the deaths among children under the age of 5 years happen within the first 28 days of life and can be attributed to complications arising from PTB (Walani, 2020). Additionally, it can lead to long-term health complications such as respiratory illnesses, neurodevelopmental disorders and learning disabilities, arising from the developmental issues associated with PTB (Townsi et al., 2018; Chung et al., 2020). PTB is also not a problem that is specific to underdeveloped or developing countries, although the ill-effects of it may be more pronounced in low-income countries. Incidences of PTB are also found in high-income parts of the world, albeit with a lower frequency, and the rates of PTB have not been on the decline either in most parts of the world (Chawanpaiboon et al., 2019).

The pathophysiology behind PTB is not completely understood yet, although certain risk factors, including but not limited to smoking habits, alcohol intake and reproductive history, have been identified to be associated with increased pre-term delivery risk (Blencowe et al., 2012; Pfänder et al., 2013; Stock et al., 2020). The ill-effects of PTB can be mitigated, and a healthy, full-term gestation outcome may also be achieved if appropriate interventions are administered (Newnham et al., 2014). Their success, however, depends on identifying at-risk subjects at earlier stages of their pregnancy, as these approaches are effective when administered during the earlier stages (Blencowe et al., 2012; Newnham et al., 2014). Current methods for assessing pre-term pregnancy outcomes involve the use of physical and biochemical markers, which are not accurately determinative of a potential incidence of premature culmination of pregnancy in the future (Georgiou et al., 2015).

Many non-pathogenic bacteria, viruses and fungi inhabit various areas of the human body, such as the gut, mouth, and reproductive tracts (Sender et al., 2016), and are collectively referred to as the “microbiome”. Microbiomes are essential for normal functioning of the respective organs, and maintain a symbiotic relationship with the human body and drive key biochemical reactions, (Gordon et al., 1971) and dysregulated microbiomes are often implicated in various diseases. These microbial communities are also present in the reproductive tracts, and have been reported to influence the pregnancy outcome (MacIntyre et al., 2015). There is evidence linking the composition of vaginal microbiomes to risk of PTB, and the abundance levels of

80 specific microbiota, such as various species of the *Lactobacillus* genus, have the potential to be indicative of
81 PTB even at earlier stages of pregnancy (Brown et al., 2019; Romero et al., 2014). Vaginal microbial
82 communities can be categorized into specific Community State Types (CSTs), which are typically
83 characterized by abundances of various *Lactobacillus* species (Romero et al., 2014). CSTs are associated with
84 increased or decreased risk of abnormalities such as Bacterial Vaginosis (BV), Urinary Tract Infections (UTIs)
85 and even PTB (Gudnadottir et al., 2022). Moreover, alpha-diversity indices, such as Shannon and Simpson
86 diversity, which can quantify the diversity of vaginal microbiota, have been harnessed for predicting PTB
87 (DiGiulio et al., 2015; Haque et al., 2017). However, the vaginal microbiome differs considerably from
88 individual to individual, especially across races (Sun et al., 2022; Gupta et al., 2017). Additionally, the
89 microbial abundance may further vary depending on the sequence processing methods used on 16S ribosomal
90 RNA (rRNA) data, which is typically used to estimate taxonomic abundance at various levels of classification
91 (Bharti et al., 2019). Consequently, the success of diversity indices for estimating PTB risk may be specific to
92 certain cohorts, or be influenced by the sequence processing pipeline and consequently, may not translate
93 across cohorts, as is our observation in this study. Machine Learning (ML)-based approaches have also been
94 explored in this context, which leverage features such as abundance of various taxa, phylotype counts, CST of
95 the vaginal microbiome, age, race and more, for PTB risk assessment.

96 The vaginal microbiome evolves as the pregnancy progresses (MacIntyre et al., 2015; Romero et al.,
97 2014), and the numerous changes that it undergoes may contain a signature for identifying PTB risk. Currently,
98 there is a severe lack of approaches that exploit the temporal dynamics of vaginal microbiomes, by looking at
99 it as a time-series problem, for PTB risk assessment, and most current predictive methodologies use static data.
100 Deep learning approaches such as recurrent neural networks (RNNs), have previously been used for modeling
101 the dynamics of gut and other microbiota in various contexts (Baranwal et al., 2022; Fung et al., 2023; Medina
102 et al., 2022) and have found success. To the best of our knowledge, such methods have not been applied to
103 vaginal microbiota, especially in the context of PTB risk assessment, so far. This may be partly attributed to
104 the fact that RNN-based approaches demand data sampled at regular intervals, which is challenging to collect
105 as study subjects are often irregular in clinical visits. With this study, we present a deep learning-based
106 approach, “neural controlled differential equations (CDEs)” that is capable of differentiating between term and
107 preterm births using time-series vaginal microbiome data, which overcomes the dependence on regularly
108 sampled microbial data. We also highlight the limitations of alpha diversity indices and traditional ML methods
109 for PTB prediction in racially and ethnically diverse patient cohorts. We show that modeling the temporal
110 dynamics of microbiota using deep learning methods results in more reliable PTB risk scoring than simple
111 ML-based methods. Our best model, utilizing neural CDEs, outperforms any ML-based PTB prediction
112 approaches so far. On the basis of this study, we show the potential of vaginal microbiota for PTB prediction,
113 and that such approaches can be pushed towards complete clinical viability with further efforts.

115 2. Materials and Methods

117 2.1. Dataset

118 We obtained 16S rRNA sequences collected from human patient samples. The data was sourced from a
119 previously published study by Callahan et al. on refinement of a vaginal microbiome signature of preterm birth
120 (Callahan et al., 2017). The dataset is publicly available under the open access category in the Sequence Read
121 Archive (SRA), BioProject ID PRJNA393472. It consists of 16S rRNA sequence samples, spread across 133
122 racially and ethnically diverse subjects, and sampled at different points of time during the pregnancy.

125 2.2. 16S rRNA Sequence Processing

126 It has been widely established that the hypervariable regions (V1-V9) within 16S rRNA gene can be used for
127 phylogenetic studies and genus or species-level classification in diverse microbial populations (Weisburg et
128 al, 1991). Furthermore, certain hypervariable regions (such as the V4) are semi-conserved and can reliably
129 predict specific taxonomic levels (Yang et al., 2016). The procedure to convert the 16S rRNA sequence data
130 to microbial abundance involves various stages of processing. In the first step, quality control checks are
131 performed and sequencing artifacts, low quality reads, etc. are removed from the read sequences. Secondly,
132 the preprocessed sequences are aligned against a chosen reference database, and a taxonomic class is assigned
133

134
135 to each sequence. The sequences are then grouped into clusters, which represent Operational Taxonomic Units
136 (OTUs), based on sequence similarities. Lastly, the abundances of each OTU in samples are estimated (Schloss
137 et al., 2009; Edgar, 2013; Estaki et al., 2020; Callahan et al., 2016). The microbial abundance obtained depends
138 on the specific processing steps, and variations in processing steps can result in different abundance values
139 (Schloss et al., 2009; Edgar, 2013; Estaki et al., 2020; Callahan et al., 2016).

140 The dataset (PRJNA393472) derived from SRA contains sequences generated after amplifying and
141 sequencing the V4 hypervariable region of the 16S rRNA gene. We used the DADA2 processing pipeline
142 (Callahan et al., 2016) to derive microbial abundance data from the sequence reads. The metadata and
143 taxonomic abundance tables were generated using the SRA cloud (Katz et al., 2021) and abundances were
144 obtained at various levels of taxonomic classification. We retained genus-level abundances for all our analyses
145 since abundances at further levels were captured at a much lower resolution.
146

147 **2.3. Processing Taxonomic Abundance Data**

148
149 We eliminated the samples for which metadata information for certain key fields, viz., gestational age at the
150 time of sample collection, gestational age at delivery, etc., was missing. We eliminated the genera abundance
151 for samples collected during trimester 3 (gestational age > 24 weeks) for our analyses, with the intention of
152 being able to predict instances of preterm delivery sufficiently early. Furthermore, we removed samples
153 collected during or before the 8th week of gestation, as they were present for very few (6 out of 133) subjects.
154 Furthermore, for some of the analyses, we transformed the genera abundance data to sample-wise relative
155 abundance, and filtered out genera with high skewness and high kurtosis, thus removing some of the genera
156 whose abundances contributed to noise. We retained 70% of the subjects (90 out of 133) as the training dataset
157 and the rest were used for validating the approaches. The training and test datasets were kept consistent across
158 all the analyses. The processed taxonomic abundance data and the corresponding metadata are made available
159 in the code repository (see “Code Availability”, Section 5).
160

161 **2.4. Diversity Metrics**

162
163 Alpha diversity metrics have been reported to be potentially indicative of preterm birth, and a highly-diverse
164 vaginal microbiome is correlated with increased risk of preterm delivery (DiGiulio et al., 2015; Hyman et al.,
165 2014; Haque et al., 2017). We computed Shannon, Simpson, Chao1 and Gini alpha diversity indices, as well
166 as Taxonomic Composition Skew (TCS) (Haque et al., 2017), a diversity index specifically tailored for the
167 vaginal microbiome. Unlike other diversity indices, TCS takes into account that vaginal microbiomes are
168 usually dominated by the *Lactobacillus* species and other genera are in the minority. TCS responds in a
169 different manner, to changes in abundances of sparse and dominant taxa, and thus is possibly more suitable for
170 quantifying the diversity in vaginal microbiomes. We checked for statistically significant differences in alpha-
171 diversity index values between the term and preterm classes during various gestational periods using a two-
172 sided, independent *t*-test. The standard diversity metrics were computed using the scikit-bio python library
173 (version 0.5.8).
174

175 **2.5. Traditional ML Approaches**

176
177 We used two ML classifiers: Decision Tree (DT) and Random Forest (RF), to predict term/preterm outcomes.
178 This constituted a secondary baseline for benchmarking the performance of higher, more complex deep
179 learning-based prediction approaches. For each patient subject, the microbial abundance profile closest to the

180 week of delivery and obtained during the period between the 9th and the 24th week of gestation, following the
 181 hypothesis that composition of vaginal microbial communities closer to the period of delivery are better
 182 indicative of preterm delivery risk. The resultant training and test sets contained 93 and 40 samples respectively.
 183

184 2.6. Deep Learning Approaches

185
 186 Machine learning classifiers, such as Support Vector Classifiers (SVCs), as well as tree-based classifiers such
 187 as DT and RF have been explored extensively for preterm birth prediction using vaginal microbiota, most often
 188 in tandem with other features such as physical markers and patient history. However, these classifiers have
 189 largely lacked the capability of making reliable predictions. Surprisingly, deep learning models have hardly
 190 been explored for this particular problem. Given the time-series nature of the data, we focused on deep learning
 191 algorithms for sequential data in this study.
 192

193 2.6.1. Recurrent Neural Networks

194
 195 Recurrent Neural Network (RNN) is a type of neural network designed to handle sequential or time-series
 196 data. The issue with standard RNNs however, is that they have difficulty in learning long-term dependencies
 197 in long sequences, due to the issue of vanishing/exploding gradients (Pascanu et al., 2012). Long Short-Term
 198 Memory (LSTM) is a type of RNN that is capable of learning long sequences, and are possibly more
 199 appropriate for the week-wise taxonomic abundance dataset. LSTM maintains a hidden state, which stores
 200 short-term information, and a cell state, which stores long-term information. The initial hidden and cell states
 201 are generally set to zero vectors. A LSTM cell at each time step updates the hidden and cell states based on
 202 the states at the previous time step and the input data at the current time step.
 203

$$\begin{aligned}
 204 \quad & (h_{t_0}, c_{t_0}) = (\vec{0}, \vec{0}) \\
 205 \quad & z_{t_i} = \begin{cases} x_{t_i} & \text{if } x_{t_i} \text{ is available} \\ \hat{x}_{t_i} & \text{if } x_{t_i} \text{ is unavailable} \end{cases} \\
 206 \quad & (h_{t_i}, c_{t_i}) = \text{LSTM}(z_{t_i}, [h_{t_{i-1}}, c_{t_{i-1}}]), \text{ for all } i \in \{1, 2, \dots, T\} \\
 207 \quad & y = \sigma(W \cdot h_{t_T} + b)
 \end{aligned}$$

208
 209 However, the conventional LSTM system demands a continuous and uniform time-series dataset, i.e.,
 210 the time-steps must represent uniform intervals and the input data should be available for each step. While the
 211 taxonomic abundance data is uniformly sampled (week-wise), data for some subjects is not present for some
 212 weeks. Additionally, the first week for which data is available is different for each subject, and thus, a zero-
 213 vector initialization will not be appropriate for the initial hidden state. To address these issues, we modified
 214 the LSTM network accordingly. Firstly, we initialized the hidden state, h_{t_0} , with a trainable embedding layer,
 215 and the cell state, c_{t_0} , was initialized as a zero vector. Secondly, at each time step, we had the LSTM cell
 216 generate the taxonomic abundance forecast, x_{t_i} , at each time step, and used the forecast whenever the input
 217 data x_{t_i} was not available. The final hidden state, h_{t_T} , was fed to a linear layer with parameters W and b
 218 followed by a sigmoid activation, $\sigma(\cdot)$, to predict the term/pre-term outcome (y). The entire network was
 219 trained end-to-end.

220 Figure 1 outlines the described LSTM network. The LSTM implementation assumes that the input
 221 data is continuously and regularly sampled. However, in our case, data for some intervals may be missing. To
 222 overcome this, we masked the data for missing time steps by using zero vectors to substitute the missing data
 223 points. In parallel, we also used a vector indicating the coordinates of the masked time intervals for each sample,
 224 for which the model used the forecast, i.e., the model-predicted microbial abundance instead of the ground
 225 truth values. Additional method details are provided in the supplementary material (Section 2.1) and the
 226 hyperparameter values are listed in supplementary Table 3. The LSTM model was implemented using the
 227 pytorch library (Paszke et al., 2019) (version 2.0.1, cuda version 11.8).
 228

229 2.6.2. Neural Differential Equations

230

231 A significant limitation associated with clinical data pertains to its irregular sampling of data points, which
 232 presents challenges in constructing effective machine learning models that can effectively harness the inherent
 233 time-series information. The irregularity in data sampling introduces two notable drawbacks: firstly, the size
 234 of input data, contingent upon the number of sampling instances, differs among various subjects; secondly, the
 235 timing of sampling instances is not strictly discrete, thereby restricting the applicability of commonly
 236 employed RNN models that assume uniform intervals between sampled data points. To overcome this, we
 237 leverage a recently introduced class of deep learning models - Neural Ordinary Differential Equations (ODEs)
 238 that combine a neural network with ODEs and allow for continuous interpolation between two randomly
 239 spaced sampling instants.

240 Notably, Neural ODEs exclusively consider the evolution of time-series data commencing at a fixed
 241 time point, denoted as t_0 , which is accompanied by an initial condition represented as x_{t_0} . In the context of
 242 our research, x_{t_0} signifies the initial abundance of genera at this specific time t_0 , which might correspond to
 243 the onset of gestation week 9. Regrettably, the trajectory of microbial abundance varies from subject to subject,
 244 and the initial abundance data for all subjects may not be accessible for subsequent analysis, as the microbial
 245 profiles of each subject were not uniformly sampled at the same time point, namely t_0 . Consequently, while
 246 Neural ODEs excel in interpolation tasks, they cannot be seamlessly integrated into our framework due to the
 247 absence of consistent initial abundance data. Nevertheless, it transpires that addressing this problem,
 248 specifically how to integrate incoming information, has already been thoroughly explored within the realm of
 249 mathematics, particularly in the field of rough analysis, which is dedicated to the examination of CDEs (Lyons,
 250 1994; Lyons et al., 2007). Kidger et al. (Kidger et al., 2020) have introduced a novel framework known as
 251 Neural CDEs, which extends CDEs to Neural ODE models. To put it simply, Neural CDEs can be seen as
 252 continuous-time counterparts of Recurrent Neural Network (RNN) models. These models can be trained
 253 efficiently using a method called “adjoint backpropagation”, which is elaborated on briefly in the
 254 supplementary materials (Section 3.1), and detailed mathematical representation of it can be found in (Chen et
 255 al., 2018). In brief, the Neural CDE model can be summarized through the following sequence of operations:
 256

$$\begin{aligned}
 257 \quad & \text{(Initialization)} \quad h_{t_0} = \psi_{\theta}(t_0, x_{t_0}) \\
 258 \quad & \text{(CDE)} \quad \frac{dh}{dt}(t) = f_{\theta}(h_t) \frac{dX}{dt}(t) \\
 259 \quad & \text{(Result)} \quad y = l_{\theta}(h_{t_T})
 \end{aligned}$$

260
 261 Here, ψ_{θ} and l_{θ} correspond to linear models responsible for transforming the initial taxa abundance (along
 262 with the time-stamp t_0) into the initial hidden state, h_{t_0} , and the final hidden state, h_{t_T} , into the output label,
 263 respectively. The map, ψ_{θ} is used to avoid translational invariance to first sampled time instant. X is the natural
 264 cubic spline with knots at t_0, \dots, t_T such that $X_{t_i} = (x_i, t_i)$. Natural cubic splines allow for smooth interpolation
 265 and minimum regularity for handling certain edge cases. f_{θ} is neural network model depending on parameters,
 266 θ . Due to its dependence on cubic splines, Neural CDEs (Kidger et al., 2020) can be applied to irregularly
 267 sampled time series, even with temporally-scattered initial conditions. Thus, we chose to apply Neural CDEs
 268 to predicting preterm birth using the irregularly-sampled microbial abundance dataset. A comprehensive
 269 mathematical introduction to Neural CDEs is outside the scope of this paper. For those interested in delving

270 deeper into the mathematical details, we recommend consulting (Kidger et al., 2020) for a more thorough
 271 explanation. The torchcde library (version 0.2.5) was used to implement the Neural CDE model in python.
 272 Hyperparameter values for the model are listed in supplementary Table 4.
 273

274 3. Results

275 3.1. Microbial abundance dataset contains racially and ethnically diverse subjects

276
 277 The 16S rRNA sequence data was converted to taxonomic abundance (Methods, Section 2.2), which led to
 278 approximately 290,000 abundance counts, spanning taxonomic counts at various levels of classification, out
 279

280 of which approximately 65,000 corresponded to genus-level, belonging to 2,326 unique samples which
281 collected at various weeks of gestation throughout the pregnancy, spread across 133 subjects of diverse race
282 and ethnicities (Figure 2a), out of which 85 subjects delivered at term and 48 subjects delivered preterm (Figure
283 2 b, d). We aligned the abundance counts subject- and gestational week-wise for ease of interpretation (Figure
284 2c). Abundance counts for multiple samples derived from the same subject collected during the same week of
285 gestation, if any, were replaced by the mean of those counts to ensure consistency, as future analyses were
286 carried out on week-wise data. The distribution of number of genera present in each sample (i.e., number of
287 genera with non-zero abundance in each sample) is visualized in Figure 2e. Microbial abundance profiles
288 corresponding to 43 out of the 133 subjects were reserved as the test set. (refer methods Section 2.3).

289 290 **3.2. Diversity metrics do not reliably identify at-risk PTB subjects**

291
292 Alpha-diversity indices were computed on samples collected during trimester 1 and trimester 2 (gestational
293 weeks 9 to 24, see methods Section 2.3). The visualization of Chao1, Gini, Shannon, Simpson and TCS alpha-
294 diversity indices computed at various weeks of gestation for subjects who delivered at term and preterm is
295 presented in Figure 3 a-e, respectively. The results of two-sided, independent *t*-tests to examine differences in
296 diversity index values across term and preterm groups are presented in Table 1. Although *t*-test reveals a
297 statistically significant difference in the Chao1 diversity index between term and preterm groups during
298 gestational weeks 9-12 ($p = 0.005$) and 13-17 ($p = 0.02$), there are too few samples in the preterm group during
299 these gestational periods to draw reliable conclusions (see Figure 3). No signature that can distinguish
300 term/preterm birth is observable from these metrics. As can be seen from the plot, the alpha-diversity indices
301 are not predictive of a preterm delivery outcome, and perform worse than random classifiers (prediction
302 accuracy on test set < 50%, refer methods Section 2.4).

303

304 305 **3.3. Statistical analyses indicate presence of signatures for PTB prediction in evolving** 306 **microbiomes**

307
308 To reduce the set of features used for the classification task, we removed the genera with high skewness and
309 kurtosis. Skewness and kurtosis were computed on relative abundance of microbial genera during the period
310 of trimester 1 - trimester 2 (gestational weeks 9 to 24). The abundance distributions of a large number of genera
311 have a high positive skewness and kurtosis (see Figure 4 a, b), and may contribute to noise. Thus, we excluded
312 the genera which had highly skewed relative abundances (i.e., skewness > 10) as well as genera with high
313 kurtosis (kurtosis > 10). As a result, only 6 genera were retained, namely, *Lactobacillus*, *Anaerococcus*,
314 *Gardnerella*, *Peptoniphilus*, *Finegoldia* and *Prevotella*. Except for *Finegoldia*, these genera have been identified
315 to be linked to PTB risk previously. *Lactobacillus* is the most dominant genus within the vaginal microbiome,
316 and low counts of *Lactobacillus* have previously been stated to be indicative of increased PTB risk (Bayar et
317 al., 2020; Gudnadottir et al., 2022). Certain species of *Anaerococcus* are found to be associated with increased

318
319 PTB risk (Ansari et al., 2021), however, there are also reports that *Anaerococcus* species may be protective in
320 nature (dos Anjos Borges et al., 2023). There is strong evidence linking high *Gardnerella vaginalis* presence
321 with PTB and bacterial vaginosis, which also increases PTB risk (Nelson et al., 2009; Ng et al., 2023). In some
322 populations, increased counts of certain *Peptoniphilus* species were also found to be associated with high PTB
323 risk (Park et al., 2022a). Similar evidence exists associating
324 high *Prevotella* abundances with increased PTB risk (Freitas et al., 2018; Fettweis et al., 2019; Park et al.,
325 2022b). However, there is insufficient evidence to establish a link between these observations and the changes
326 that vaginal microbiota undergo throughout the duration of the pregnancy.

327 We further re-computed the relative abundance based on the abundance counts of these 6 genera only,
328 and visualized the gestational week-wise abundances and corresponding term/preterm delivery outcomes. The
329 results are presented in Figure 5 a-f. This analysis highlights that the composition of the 6 genera mentioned
330 above, changes significantly during the period from end of trimester 1 to trimester 2 of pregnancy. The trends
331 indicate that low abundance of the Lactobacillus genus during trimester 2 (gestational
332 Gardnerella counts with PTB risk (Figure 5c, $p = 0.0073$), and the effect is more pronounced during gestational
333 weeks 16-18. Most of the samples with increased Prevotella counts during weeks 19-21 belonged to subjects
334 who went on to deliver preterm (Figure 5e, $p \approx 10^{-6}$).

336 **3.4. ML classifiers do not make adequate PTB risk assessment**

337 We tested the performance of two ML classifiers, viz., RF and DT towards prediction of preterm birth. For
338 this, we isolated the latest available microbial abundance profiles of training set subjects as well as the test set
339 subjects, collected during the period between the 9th and the 24th week of gestation, with the rationale that the
340 state of the microbiome closer to the period of delivery should be better predictive of term/preterm delivery.
341 Optimal hyperparameters for both RF and DT were identified by performing a grid search on pre-defined
342 parameter search spaces using 3-fold cross-validation on the training set, and are listed in supplementary Table
343 2. The resultant models were validated on the test set by computing ROC-AUC, accuracy, precision-recall and
344 f1 score. The RF model performed significantly better (Figure 6 b, d) compared to DT (Figure 6 a, c) which
345 performs worse than a random predictor. The detailed results of both classifiers are presented in Table 2
346 Neither of the models, however, make adequately reliable predictions on the test set.

349 **3.5. LSTMs do not decode the temporal dynamics of vaginal microbiomes**

350 LSTM was trained on the week-wise genera abundance data sampled during gestational weeks 9 to 24, after
351 making appropriate adjustments to account for the irregularity in sampling (see Methods, Section 2.6.1, Figure
352 1). When trained on the entire set of genera with non-zero variance in the training set, the model overfits and
353 does not generalize well to making predictions on the test set (training set accuracy = 100%, test set accuracy
354 < 60%). Even when trained on the set of genera with low skewness and kurtosis, the model fails to make
355 sufficiently accurate predictions on the test set (accuracy = 63%), and is outperformed by the RF model
356 described above. We attribute this lack of predictivity to the increased estimations that the model makes to fill
357 the temporal gaps in the data, and it may necessitate availability of additional data samples to be able to make
358 these estimations more accurately, either in terms of more patient subjects or increased density of samples per
359 subject.

362 **3.6. Neural CDEs are capable of achieving PTB prediction with a substantial accuracy**

363 The neural CDE model trained on relative abundances of genera selected by the skewness and kurtosis filtering
364 outperforms all other models described above. The resultant model performs reasonably well on the test set
365 (mean test set ROC-AUC = 0.82, accuracy = 74.5% precision = 0.65, recall = 0.71, f1 score = 0.71). The results
366 are presented in Figure 7 a, b. We then trained the model after shuffling the term/preterm labels for subjects in
367 the training set. As expected, the model behaved similar to a random predictor on the test set (ROC-AUC <
368 0.5). Albeit not on the same validation dataset, our approach describes better results than the best submission
369 in the DREAM challenge for term vs preterm prediction (ROC-AUC = 0.68, accuracy = 67%, sensitivity =
370 0.48, specificity = 0.79) (Golob et al., 2023), despite using microbial abundances only up to the end of the 2nd
371 trimester, i.e., the 24th week of gestation, as opposed to the 32nd week of gestation in the DREAM challenge.

375
376

377 4. Discussion

378
379 Current in vitro or in vivo approaches lack the ability to detect PTB incidences at earlier stages reduces the
380 effectiveness of prophylactic or therapeutic interventions that can be administered to mitigate neonatal health
381 concerns associated with it. Current risk assessment approaches involve physical examinations, and factors
382 such as cervical length may be used to estimate the risk. Additionally, abnormality in levels of biochemical
383 markers such as Pregnancy-Associated Protein A (PAPP-A) (Smith et al., 2002; Gundu et al., 2016),
384 Cervicovaginal Interleukins (Manning et al., 2019; Park et al., 2020), etc., may help detect PTB as well.
385 However, there is a lack of definitive confidence intervals for these physical or biochemical tests. Recent
386 studies have highlighted the utility of diversity of vaginal microbial communities towards PTB prediction, by
387 establishing correlations between alpha-diversity indices associated with abundances of various microbial
388 species (or genera) and incidences of PTB (DiGiulio et al., 2015; Hyman et al., 2014; Haque et al., 2017).
389 However, microbial communities are highly diverse across various individuals, and more so for individuals
390 belonging to different ethnic populations (Sun et al., 2022; Gupta et al., 2017). We have demonstrated above,
391 that in a heterogenous dataset, with microbial profiles derived from ethnically and racially diverse subjects,
392 diversity metrics could not accurately estimate PTB risk. This indicates that previously reported success of
393 diversity indices in identifying subjects at high risk of PTB may be dataset dependent, either with respect to
394 the subject cohort or to the pipeline used in computation of microbial abundance from 16S rRNA sequences,
395 based on our observations on a mixed-race dataset with 16S rRNA sequences transformed using a standardized
396 method.

397 Traditional ML methods, which have been explored previously in the context of predicting PTB using
398 vaginal microbial species abundance, fail to learn an abundance signature associated with PTB in our dataset.
399 While others (Park et al., 2022a) report success with machine learning methods in ethnically homogenous
400 cohorts, we found that the predictive performance did not translate to mixed-race cohorts. As we have
401 demonstrated above, vaginal microbial communities evolve throughout the duration of pregnancy, and
402 abundance levels of certain species or genera may change significantly as the pregnancy progresses. Learning
403 a PTB-associated signature in an evolving microbiome may be out of scope of such models as they are not
404 designed to handle time-series datasets. We explored the utility of LSTM, a type of RNN, which is able to
405 work with sequential datasets. The architecture of RNN-based approaches requires input datasets to be
406 regularly and continuously sampled. As far as human patient subject data is concerned, obtaining such a dataset
407 is a challenge, as study subjects may not be regular or consistent in clinical visits. For this purpose, we filtered
408 the dataset such that a single sample was present across each of the gestational weeks, which constituted the
409 time intervals for LSTM. We suitably modified the LSTM workflow to overcome missing time intervals,
410 however it proved to be incapable of learning any signature associated with PTB.

411 Neural differential equations have recently gained traction with regards to analyzing sequential data.
412 Since it uses differential equations to model the temporal dynamics, it can handle irregularly and/or
413 inconsistently sampled data. Neural CDEs are more efficient than neural ODEs (Kidger et al., 2020), and are
414 even capable of working with partially sampled datasets, although we have not harnessed that in this study.
415 We found considerable success in using Neural CDEs to predict PTB, in spite of working with a dataset sourced
416 from an ethnically varied population (refer Figure 2), and outperformed all other approaches that we tested.
417 To the best of our knowledge, this is the first effort towards modeling the temporal dynamics of vaginal
418 microbial communities using deep learning, and the first instance of applying neural differential equations for
419 a problem of this kind.

420 The DREAM challenge for PTB prediction (Golob et al., 2023), was issued in 2019 with the goal of
421 driving efforts for PTB prediction using the vaginal microbiome. One of the sub-problems for the DREAM
422 challenge consisted of predicting term births (≥ 37 weeks of gestation) and preterm births (< 37 weeks of
423 gestation) using vaginal microbiomes. The dataset for this challenge was derived from 9 different studies, and
424 amounted to 3578 samples collected from 1268 individuals (Golob et al., 2023). The dataset used in this study
425 was also part of the DREAM challenge. We also considered using some of the other datasets in the DREAM
426 challenge while outlining this study, but dropped either due to not being labelled week-wise or due to
427 insufficient week-wise samples per patient, for modeling the temporal dynamics. Most of the top submissions
428 in the challenge used tree-based classifiers. On our test dataset, neural CDEs show better predictivity (mean
429 test set ROC-AUC = 0.82, accuracy = 75%, sensitivity = 0.71, specificity = 0.85) than the best submission in
430 the DREAM challenge on their validation dataset (ROC-AUC = 0.69, accuracy = 67%, sensitivity = 0.48,
431 specificity = 0.79) (Golob et al., 2023), despite using microbial abundances only up to the end of the 2nd

432 trimester, i.e., the 24th week of gestation. The DREAM challenge for PTB prediction used taxonomic
433 abundance data up to the 32nd week of gestation. Our emphasis was on early-stage PTB prediction, and we were
434 able to achieve better predictive performance in spite of restricting the input data till the 2nd trimester.

435 Predictive approaches using data other than microbial communities also exist. For instance, Tarca et
436 al. (Tarca et al., 2021) report the results of the DREAM challenge for PTB prediction using the maternal blood
437 transcriptome and the proteome. The top performing models report better results than what was reported on
438 microbial communities, with a ROC-AUC of 0.76 when proteomics data from weeks 27-33 was used. However,
439 in another sub-challenge where early-stage data (weeks 17-22) was used, the top performing model had a
440 ROC-AUC of 0.62. Obtaining blood transcriptomic or proteomic data may pose difficulties due to the
441 involvement of invasive procedures requiring clinical expertise to perform. On the other hand, microbial
442 abundance data is sourced from vaginal swabs, which can be obtained without invasive procedures, by patient
443 subjects themselves. Several attempts have been made at predicting PTB using biochemical marker (Aung et
444 al., 2019; Leow et al., 2020), however, obtaining such data may require regular clinical visits, and their viability
445 in racially-diverse populations is unknown.

446 Poorer and remote parts of the world may even lack the medical infrastructure or presence of adequate
447 facilities that are required for PTB assessment and prevention. For example, in remote areas in India, there are
448 clinics called Anganwadis, which roughly translates to courtyard shelter. As of 2018, the Ministry of Women
449 and Child Development reports the existence of 1.4 million Anganwadi centres spread out across the country.
450 Anganwadis provide limited healthcare facilities for maternal and infant health and lack the funding and
451 facilities, or even trained medical personnel required to mitigate PTB and its ill-effects. Given the simplicity
452 of obtaining samples from which microbial abundance is derived, reliable approaches for PTB risk assessments
453 developed on microbiota will greatly help such remote clinics.

454 While we have demonstrated the capability of Neural CDEs towards PTB prediction using the vaginal
455 microbiome, further effort can be made for increasing its clinical viability. Firstly, our dataset is limited in size
456 (133 patient subjects), and we believe that larger datasets with better racial and ethnic representation may help
457 learn signatures which take into account the diversity of vaginal microbiomes across individuals/races.
458 Secondly, predicting the extent of preterm birth (extremely preterm, very preterm, moderate to late preterm)
459 is also important as far as administering interventions is concerned, as they may have varying impact on
460 maternal and infant health and may require different strategies. This may be achieved by predicting the
461 gestational week of delivery, or by treating PTB as a multi-class problem with different extents of PTB as the
462 classes, on more high-quality datasets. We strongly believe that vaginal microbial communities may be the
463 key to achieving early-stage PTB prediction, and our findings strongly encourage future efforts for pregnancy
464 microbiome data generation and further refinements in modeling procedures, which may take us closer to
465 achieving full clinical viability.

466 **5. Data & Code Availability**

467 The code and data files for this study have been made available on <https://tinyurl.com/3427p6y4>. The
468 repository contains a readme file which describes the contents of the individual data files along with the
469 demographics of the train and test data, as well as a brief description of the code files.

470 **6. Declarations**

471 **6.1. Competing Interests**

472 The authors declare no competing interests

473 **6.2. Funding**

474 This research did not receive any specific grant from any funding agency in the public, commercial or not-for-
475 profit sector.

476 **6.3. Ethical Approval**

485

486

The data used in this study was derived from a public, non-controlled access dataset from within the Sequence Read Archive (SRA), and hence, ethical approval was not required.

487

488

489

6.4. Author Contributions

490

491

M.B. and K.K. designed the study. K.K. performed the computational analyses along with assistance from M.B., and M.B. and K.K. drafted the final manuscript.

492

493

494

6.5. Acknowledgements

495

496

The authors would like to thank Dr. Mohammed Haque, Dr. Anirban Dutta, and Mr. Nishal Kumar Pinna for their help in processing the 16S rRNA sequences from the SRA cloud, and Mr. Sunil Nagpal and Dr. Anirban Dutta for outlining the state of the art for PTB prediction using diversity metrics. We would also like to thank Dr. Rajgopal Srinivasan for proofreading and helping improve the manuscript.

497

498

499

500

501

Figure 1: Visualization of the LSTM Network. Week 0 represents the earliest week of gestation for which microbial data is available. This is passed through an embedding layer, which generates the initial hidden state. Subsequent hidden states are determined by the previous hidden and cell states, and the microbial data at the respective time step. If the microbial data is not available, the output of the previous hidden state is used instead. The final hidden state is passed through a linear layer with a single output neuron to predict the term/preterm outcome, and the entire network is trained end-to-end.

502

503

504

505

506

507

Figure 2: Distribution of (a) race, (b) number of subjects who delivered at term/pre-term, (c) sample gestational age at the time of collection, (d) gestational age of subject at the time of delivery, and (e) number of unique genera detected in each sample in the microbial abundance dataset

508

509

510

511

Figure 3: (a) Chao1, (b) Gini, (c) Shannon, (d) Simpson and (e) TCS diversity metrics computed on microbial abundance data collected during trimester 1 and trimester 2 of pregnancy. Blue and red points represent samples derived from subjects who delivered at term and pre-term, respectively.

512

513

514

515

Figure 4: Distribution of (a) skewness and (b) kurtosis: computed on relative abundances of genera.

516

517

Figure 5: (a)-(f): Variation in relative abundances for Anaerococcus, Finegoldia, Gardnerella, Lactobacillus, Peptoniphilus and Prevotella, respectively, during gestational weeks 9-24.

518

519

520

Figure 6: Receiver operating characteristic (ROC) curve for the (a) decision tree (DT) and (b) random forest (RF) classifiers on the training and test datasets, and validation metrics (AUC, accuracy, precision, recall and f1 score) computed for the (c) DT and (d) RF classifiers.

521

522

523

524

Figure 7: (a) ROC curves and (b) classification metrics on the training and test datasets for the Neural CDE model.

525

526

Table 1: *p*-values representing significance of differences between diversity metric values across the term and preterm groups for gestational weeks (a) 9-12, (b) 13-16, (c) 17-20 and (d) 21-24.

527

528

Table 2: Comparative performance of various machine learning methods

529

530

531

References

532

533

Ansari, A. et al. (2021), 'Molecular mechanism of microbiota metabolites in preterm birth: Pathological and therapeutic insights', *International Journal of Molecular Sciences* **22**(15), 8145.

534

535

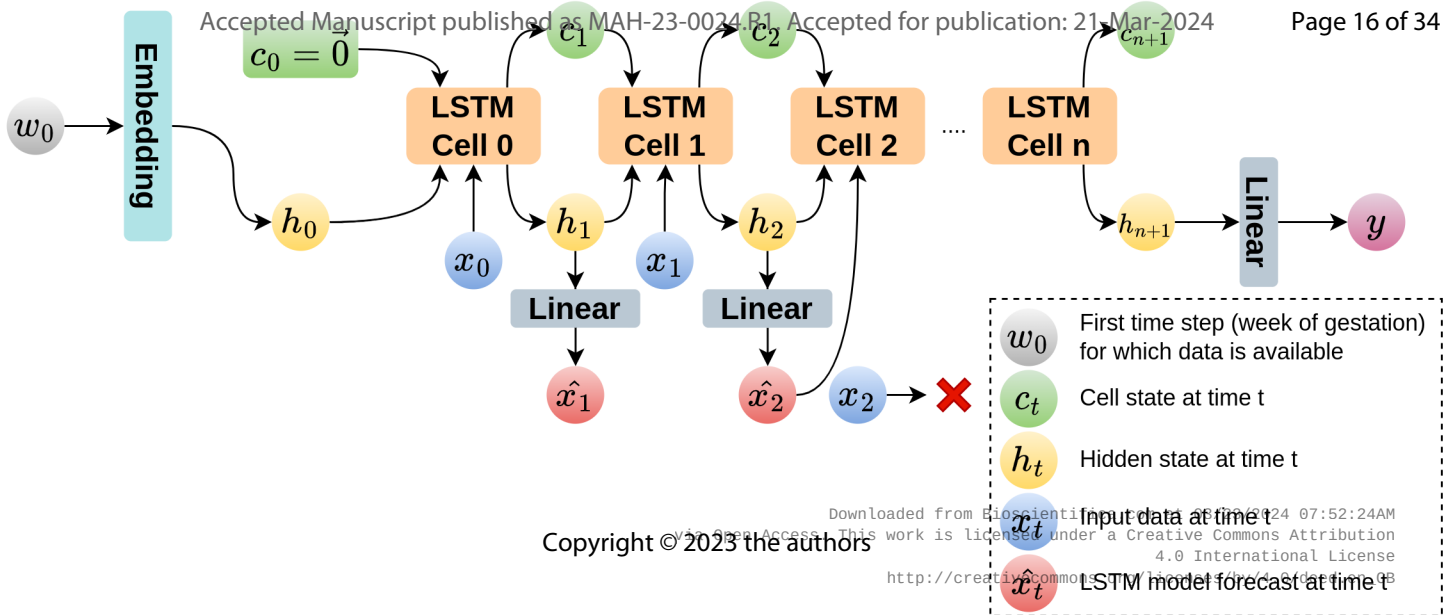
URL: <https://doi.org/10.3390/ijms22158145>

536

- 537 Aung, M. T. et al. (2019), 'Prediction and associations of preterm birth and its subtypes with eicosanoid
538 enzymatic pathways and inflammatory markers', *Scientific Reports* **9**(1).
539 URL: <https://doi.org/10.1038/s41598-019-53448-z>
- 540 Baranwal, M. et al. (2022), 'Recurrent neural networks enable design of multifunctional synthetic human gut
541 microbiome dynamics', *eLife* **11**.
542 URL: <https://doi.org/10.7554/elife.73870>
- 543 Bayar, E. et al. (2020), 'The pregnancy microbiome and preterm birth', *Seminars in Immunopathology* **42**(4),
544 487–499.
545 URL: <https://doi.org/10.1007/s00281-020-00817-w>
- 546 Bharti, R. et al. (2019), 'Current challenges and best-practice protocols for microbiome analysis', *Briefings
547 in Bioinformatics* **22**(1), 178–193.
548 URL: <https://doi.org/10.1093/bib/bbz155>
- 549 Blencowe, H. et al. (2012), 'National, regional, and worldwide estimates of preterm birth rates in the year
550 2010 with time trends since 1990 for selected countries: a systematic analysis and implications', *The
551 Lancet* **379**(9832), 2162–2172.
552 URL: [https://doi.org/10.1016/s0140-6736\(12\)60820-4](https://doi.org/10.1016/s0140-6736(12)60820-4)
- 553 Brown, R. G. et al. (2019), 'Establishment of vaginal microbiota composition in early pregnancy and its
554 association with subsequent preterm prelabor rupture of the fetal membranes', *Translational
555 Research* **207**, 30–43.
556 URL: <https://doi.org/10.1016/j.trsl.2018.12.005>
- 557 Callahan, B. J. et al. (2016), 'DADA2: High-resolution sample inference from illumina amplicon data',
558 *Nature Methods* **13**(7), 581–583.
559 URL: <https://doi.org/10.1038/nmeth.3869>
- 560 Callahan, B. J. et al. (2017), 'Replication and refinement of a vaginal microbial signature of preterm birth in
561 two racially distinct cohorts of US women', *Proceedings of the National Academy of Sciences*
562 **114**(37), 9966– 9971.
563 URL: <https://doi.org/10.1073/pnas.1705899114>
- 564 Chawanpaiboon, S. et al. (2019), 'Global, regional, and national estimates of levels of preterm birth in 2014:
565 a systematic review and modelling analysis', *The Lancet Global Health* **7**(1), e37–e46.
566 URL: [https://doi.org/10.1016/s2214-109x\(18\)30451-0](https://doi.org/10.1016/s2214-109x(18)30451-0)
- 567 Chen et al. (2018), 'Neural ordinary differential equations', *arXiv* .
568 URL: <https://arxiv.org/abs/1806.07366>
- 569 Chung, E. H. et al. (2020), 'Neurodevelopmental outcomes of preterm infants: a recent literature review',
570 *Translational Pediatrics* **9**(S1), S3–S8.
571 URL: <https://doi.org/10.21037/tp.2019.09.10>
- 572 DiGiulio, D. B. et al. (2015), 'Temporal and spatial variation of the human microbiota during pregnancy',
573 *Proceedings of the National Academy of Sciences* **112**(35), 11060–11065.
574 URL: <https://doi.org/10.1073/pnas.1502875112>
- 575 dos Anjos Borges, L. G. et al. (2023), 'Vaginal and neonatal microbiota in pregnant women with preterm
576 premature rupture of membranes and consecutive early onset neonatal sepsis', *BMC Medicine* **21**(1).
577 URL: <https://doi.org/10.1186/s12916-023-02805-x>
- 578 Edgar, R. C. (2013), 'UPARSE: highly accurate OTU sequences from microbial amplicon reads', *Nature
579 Methods* **10**(10), 996–998.
580 URL: <https://doi.org/10.1038/nmeth.2604>
- 581 Estaki, M. et al. (2020), 'QIIME 2 enables comprehensive end-to-end analysis of diverse microbiome data
582 and comparative studies with publicly available data', *Current Protocols in Bioinformatics* **70**(1).
583 URL: <https://doi.org/10.1002/cpbi.100>
- 584 Fettweis, J. M. et al. (2019), 'The vaginal microbiome and preterm birth', *Nature Medicine* **25**(6), 1012–
585 1021.
586 URL: <https://doi.org/10.1038/s41591-019-0450-2>
- 587 Freitas, A. C. et al. (2018), 'Increased richness and diversity of the vaginal microbiota and spontaneous
588 preterm birth', *Microbiome* **6**(1).
589 URL: <https://doi.org/10.1186/s40168-018-0502-8>
- 590 Fung, D. L. X. et al. (2023), 'A self-knowledge distillation-driven CNN-LSTM model for predicting disease
591 outcomes using longitudinal microbiome data', *Bioinformatics Advances* **3**(1).

- 592 URL: <https://doi.org/10.1093/bioadv/vbad059>
- 593 Georgiou, H. M. et al. (2015), 'Predicting preterm labour: Current status and future prospects', *Disease*
- 594 *Markers* 2015, 1–9.
- 595 URL: <https://doi.org/10.1155/2015/435014>
- 596 Golob, J. L. et al. (2023), 'Microbiome preterm birth DREAM challenge: Crowdsourcing machine learning
- 597 approaches to advance preterm birth research', *Preprint*.
- 598 URL: <https://doi.org/10.1101/2023.03.07.23286920>
- 599 Gordon, H. A. et al. (1971), 'The gnotobiotic animal as a tool in the study of host microbial relationships',
- 600 *Bacteriological Reviews* 35(4), 390–429.
- 601 URL: <https://doi.org/10.1128/br.35.4.390-429.1971>
- 602 Gudnadottir, U. et al. (2022), 'The vaginal microbiome and the risk of preterm birth: a systematic review and
- 603 network meta-analysis', *Scientific Reports* 12(1).
- 604 URL: <https://doi.org/10.1038/s41598-022-12007-9>
- 605 Gundu, S. et al. (2016), 'Correlation of first-trimester serum levels of pregnancy-associated plasma protein a
- 606 with small-for-gestational-age neonates and preterm births', *International Journal of Gynecology &*
- 607 *Obstetrics* 133(2), 159–163.
- 608 URL: <https://doi.org/10.1016/j.ijgo.2015.09.022>
- 609 Gupta, V. K. et al. (2017), 'Geography, ethnicity or subsistence-specific variations in human microbiome
- 610 composition and diversity', *Frontiers in Microbiology* 8.
- 611 URL: <https://doi.org/10.3389/fmicb.2017.01162>
- 612 Haque, M. M. et al. (2017), 'First-trimester vaginal microbiome diversity: A potential indicator of preterm
- 613 delivery risk', *Scientific Reports* 7(1).
- 614 URL: <https://doi.org/10.1038/s41598-017-16352-y>
- 615 Hyman, R. W. et al. (2014), 'Diversity of the vaginal microbiome correlates with preterm birth',
- 616 *Reproductive Sciences* 21(1), 32–40.
- 617 URL: <https://doi.org/10.1177/1933719113488838>
- 618 Katz, K. et al. (2021), 'The sequence read archive: a decade more of explosive growth', *Nucleic Acids*
- 619 *Research* 50(D1), D387–D390.
- 620 URL: <https://doi.org/10.1093/nar/gkab1053>
- 621 Kidger, P. et al. (2020), 'Neural controlled differential equations for irregular time series'.
- 622 URL: <https://arxiv.org/abs/2005.08926>
- 623 Leow, S. M. et al. (2020), 'Preterm birth prediction in asymptomatic women at mid-gestation using a panel
- 624 of novel protein biomarkers: the prediction of PreTerm labor (PPeTaL) study', *American Journal of*
- 625 *Obstetrics & Gynecology* MFM 2(2), 100084.
- 626 URL: <https://doi.org/10.1016/j.ajogmf.2019.100084>
- 627 Lyons, T. (1994), 'Differential equations driven by rough signals (i): An extension of an inequality of Ito
- 628 young', *Mathematical Research Letters* 1(4), 451–464.
- 629 Lyons, T. J. et al. (2007), *Differential equations driven by rough paths*, Springer.
- 630 MacIntyre, D. A. et al. (2015), 'The vaginal microbiome during pregnancy and the postpartum period in a
- 631 european population', *Scientific Reports* 5(1).
- 632 URL: <https://doi.org/10.1038/srep08988>
- 633 Manning, R. et al. (2019), 'Predictive value of cervical cytokine, antimicrobial and microflora levels for pre-
- 634 term birth in high-risk women', *Scientific Reports* 9(1).
- 635 URL: <https://doi.org/10.1038/s41598-019-47756-7>
- 636 Medina, R. H. et al. (2022), 'Machine learning and deep learning applications in microbiome research',
- 637 *ISME Communications* 2(1).
- 638 URL: <https://doi.org/10.1038/s43705-022-00182-9>
- 639 Nelson, D. B. et al. (2009), 'Preterm labor and bacterial vaginosis-associated bacteria among urban women',
- 640 *Journal of Perinatal Medicine* 37(2).
- 641 URL: <https://doi.org/10.1515/jpm.2009.026>
- 642 Newnham, J. P. et al. (2014), 'Strategies to prevent preterm birth', *Frontiers in Immunology* 5.
- 643 URL: <https://doi.org/10.3389/fimmu.2014.00584>
- 644 Ng, B. K. et al. (2023), 'Maternal and fetal outcomes of pregnant women with bacterial vaginosis', *Frontiers*
- 645 *in Surgery* 10.
- 646 URL: <https://doi.org/10.3389/fsurg.2023.1084867>

- 647 Park, S. et al. (2020), ‘Cervicovaginal fluid cytokines as predictive markers of preterm birth in symptomatic
648 women’, *Obstetrics & Gynecology Science* **63**(4), 455–463.
649 URL: <https://doi.org/10.5468/ogs.19131>
- 650 Park, S. et al. (2022a), ‘Predicting preterm birth through vaginal microbiota, cervical length, and WBC using
651 a machine learning model’, *Frontiers in Microbiology* **13**.
652 URL: <https://doi.org/10.3389/fmicb.2022.912853>
- 653 Park, S. et al. (2022b), ‘Ureaplasma and prevotella colonization with lactobacillus abundance during
654 pregnancy facilitates term birth’, *Scientific Reports* **12**(1).
655 URL: <https://doi.org/10.1038/s41598-022-13871-1>
- 656 Pascanu et al. (2012), ‘On the difficulty of training recurrent neural networks’, *arXiv* .
657 URL: <https://arxiv.org/abs/1211.5063>
- 658 Paszke, A. et al. (2019), ‘Pytorch: An imperative style, high-performance deep learning library’, *arXiv* .
659 URL: <https://arxiv.org/abs/1912.01703>
- 660 Perin, J. et al. (2022), ‘Global, regional, and national causes of under-5 mortality in 2000–19: an updated
661 systematic analysis with implications for the sustainable development goals’, *The Lancet Child &
662 Adolescent Health* **6**(2), 106–115.
663 URL: [https://doi.org/10.1016/s2352-4642\(21\)00311-4](https://doi.org/10.1016/s2352-4642(21)00311-4)
- 664 Pfänder, M. et al. (2013), ‘Preterm birth and small for gestational age in relation to alcohol consumption
665 during pregnancy: stronger associations among vulnerable women? results from two large western-
666 european studies’, *BMC Pregnancy and Childbirth* **13**(1).
667 URL: <https://doi.org/10.1186/1471-2393-13-49>
- 668 Romero, R. et al. (2014), ‘Erratum to: The composition and stability of the vaginal microbiota of normal
669 pregnant women is different from that of non-pregnant women’, *Microbiome* **2**(1).
670 URL: <https://doi.org/10.1186/2049-2618-2-10>
- 671 Schloss, P. D. et al. (2009), ‘Introducing mothur: Open-source, platform-independent, community-supported
672 software for describing and comparing microbial communities’, *Applied and Environmental
673 Microbiology* **75**(23), 7537–7541.
674 URL: <https://doi.org/10.1128/aem.01541-09>
- 675 Sender, R. et al. (2016), ‘Revised estimates for the number of human and bacteria cells in the body’, *PLOS
676 Biology* **14**(8), e1002533.
677 URL: <https://doi.org/10.1371/journal.pbio.1002533>
- 678 Smith, G. C. S. et al. (2002), ‘Early pregnancy levels of pregnancy-associated plasma protein a and the risk
679 of intrauterine growth restriction, premature birth, preeclampsia, and stillbirth’, *The Journal of
680 Clinical Endocrinology & Metabolism* **87**(4), 1762–1767.
681 URL: <https://doi.org/10.1210/jcem.87.4.8430>
- 682 Stock, S. J. et al. (2020), ‘Maternal smoking and preterm birth: An unresolved health challenge’, *PLOS
683 Medicine* **17**(9), e1003386.
684 URL: <https://doi.org/10.1371/journal.pmed.1003386>
- 685 Sun, S. et al. (2022), ‘Race, the vaginal microbiome, and spontaneous preterm birth’, *mSystems* **7**(3).
686 URL: <https://doi.org/10.1128/msystems.00017-22>
- 687 Tarca, A. L. et al. (2021), ‘Crowdsourcing assessment of maternal blood multi-omics for predicting
688 gestational age and preterm birth’, *Cell Reports Medicine* **2**(6), 100323.
689 URL: <https://doi.org/10.1016/j.xcrm.2021.100323>
- 690 Townsi, N. et al. (2018), ‘The impact of respiratory viruses on lung health after preterm birth’, *European
691 Clinical Respiratory Journal* **5**(1), 1487214.
692 URL: <https://doi.org/10.1080/20018525.2018.1487214>
- 693 Walani, S. R. (2020), ‘Global burden of preterm birth’, *International Journal of Gynecology & Obstetrics
694* **150**(1), 31–33.
695 URL: <https://doi.org/10.1002/ijgo.13195>
- 696 Weisburg, W. G. et al. (1991), ‘16s ribosomal DNA amplification for phylogenetic study’, *Journal of
697 Bacteriology* **173**(2), 697–703.
698 URL: <https://doi.org/10.1128/jb.173.2.697-703.1991>
- 699 Yang, B. et al. (2016), ‘Sensitivity and correlation of hypervariable regions in 16s rRNA genes in
700 phylogenetic analysis’, *BMC Bioinformatics* **17**(1).
701 URL: <https://doi.org/10.1186/s12859-016-0992-y>



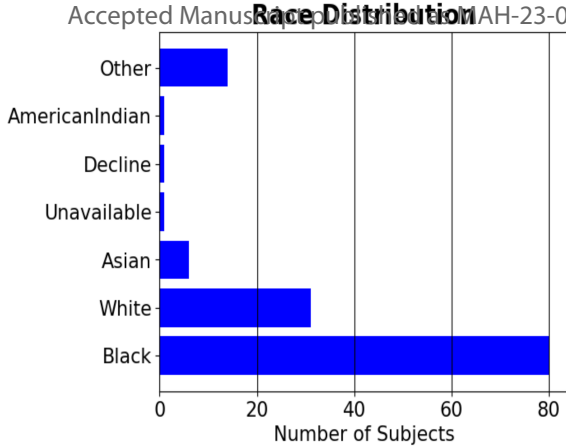
Copyright © 2023 the authors

Downloaded from Bioscientific Data Journal 08/20/2024 07:52:24AM

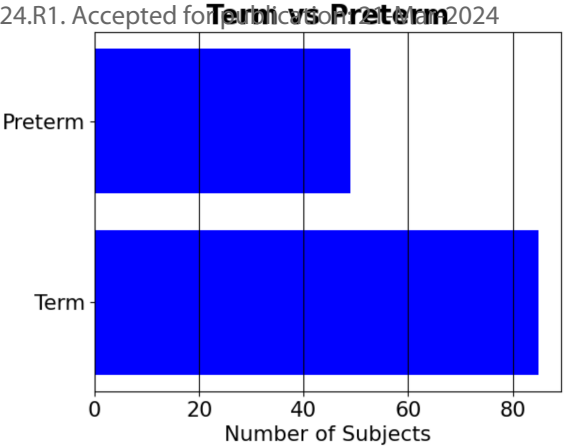
via Open Access. This work is licensed under a Creative Commons Attribution

4.0 International License

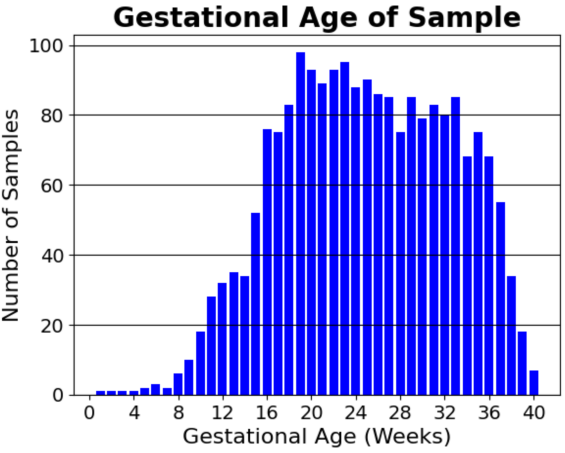
<http://creativecommons.org/licenses/by/4.0/>



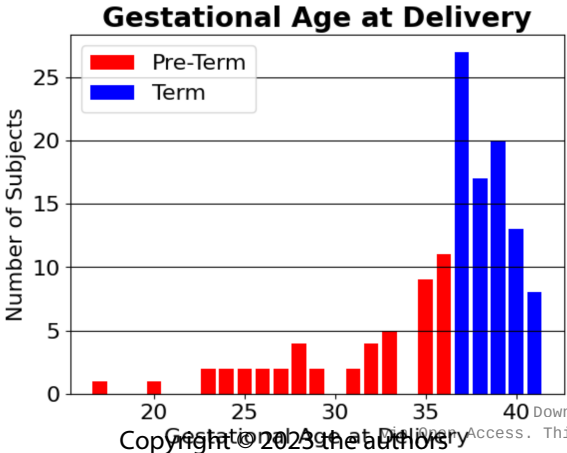
(a)



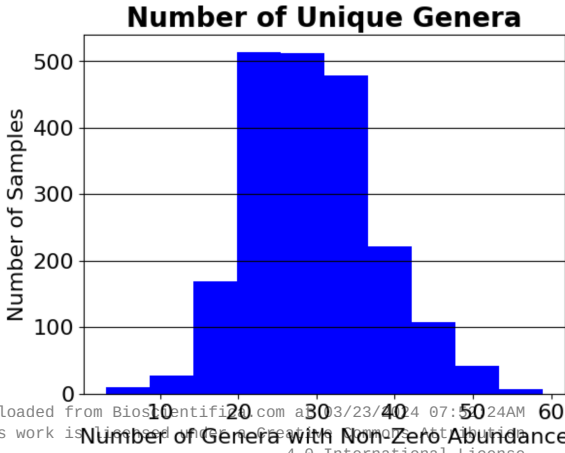
(b)



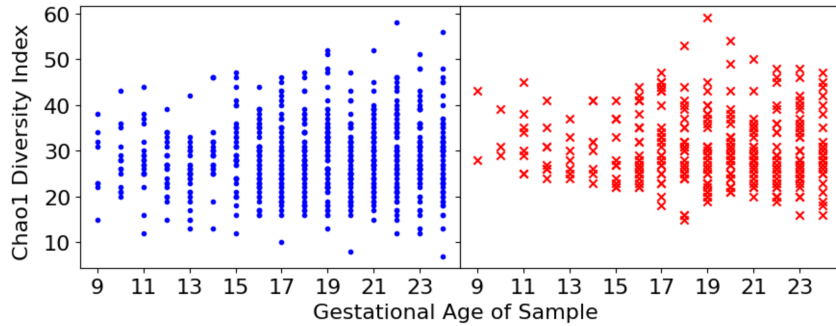
(c)



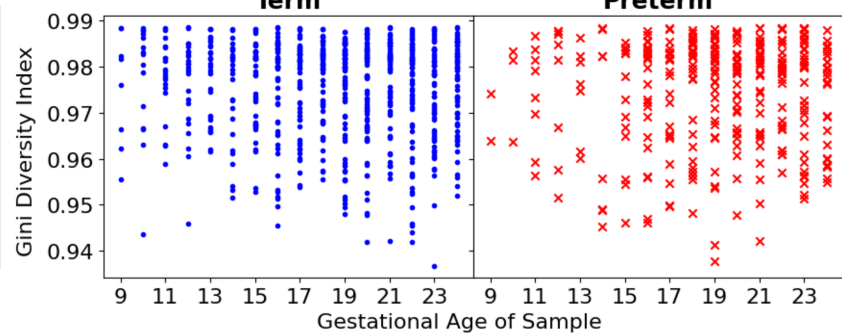
(d)



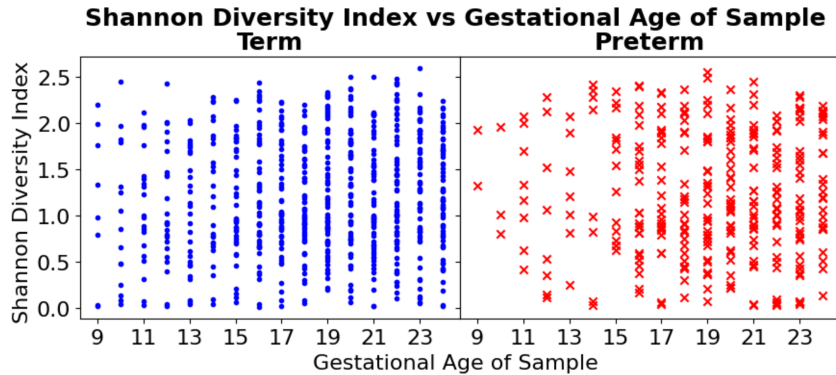
(e)



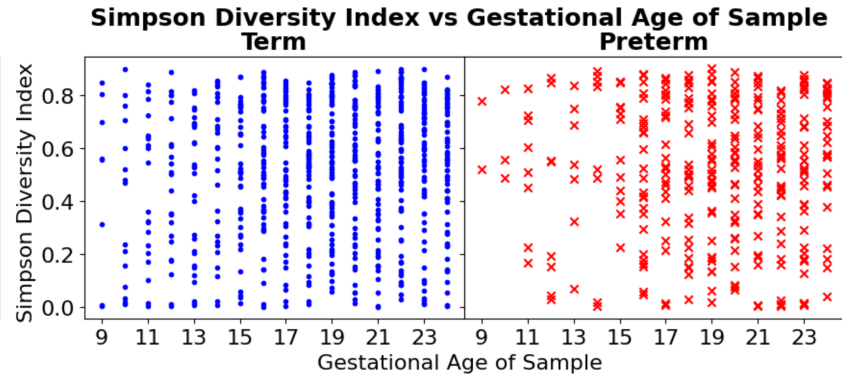
(a)



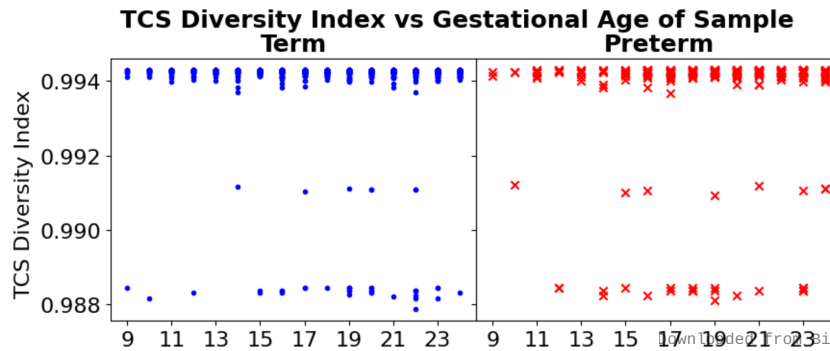
(b)



(c)

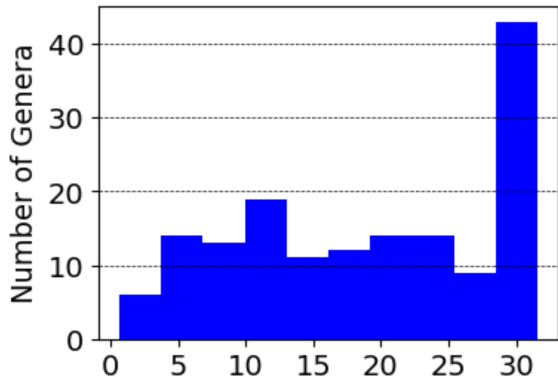


(d)



(e)

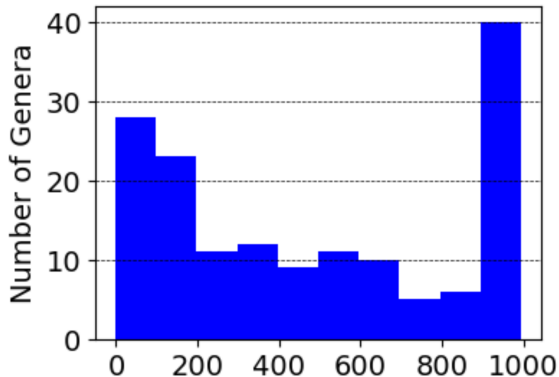
Skewness Distribution



Skewness

(a)

Kurtosis Distribution



Kurtosis

(b)

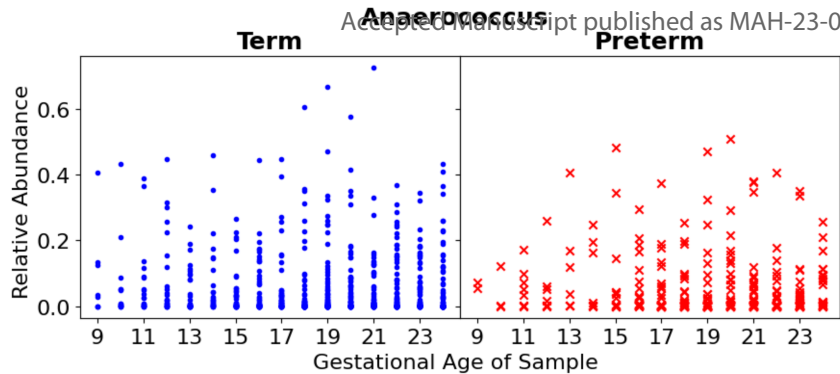
Downloaded from Bioscientifica.com at 03/23/2024 07:52:24AM

Copyright © 2023 the authors

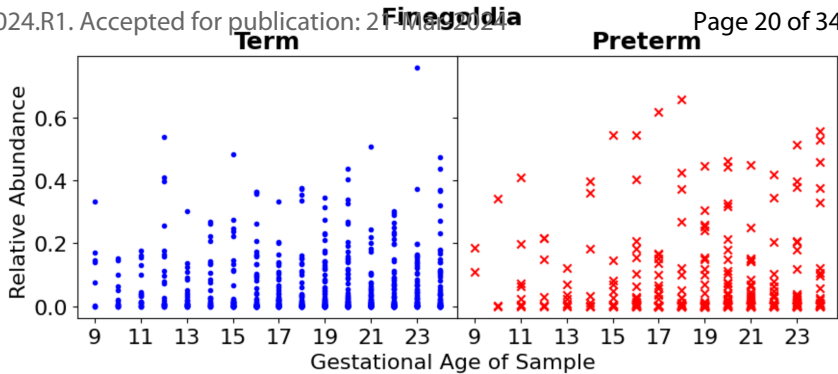
This work is licensed under a Creative Commons Attribution

4.0 International License

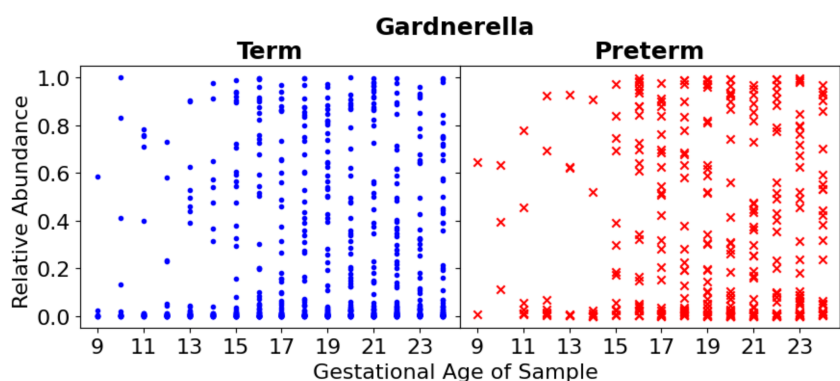
http://creativecommons.org/licenses/by/4.0/deed.en_GB



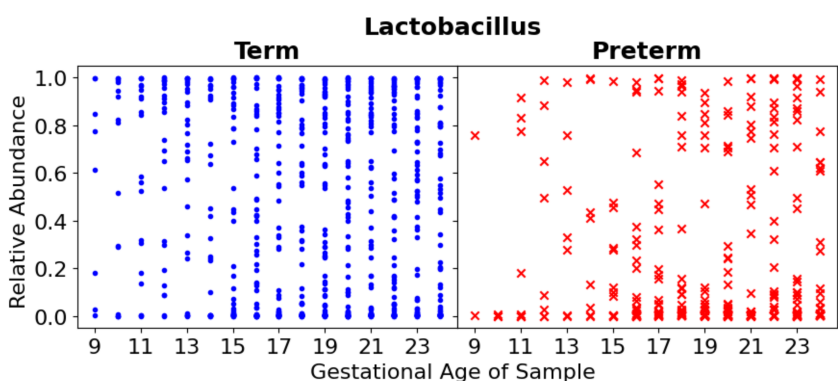
(a)



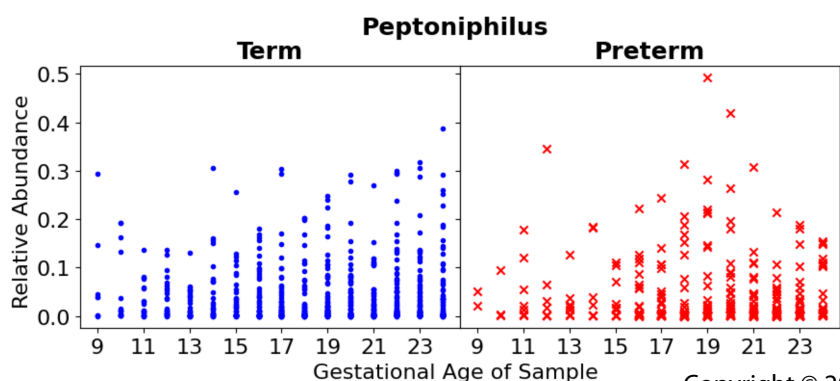
(b)



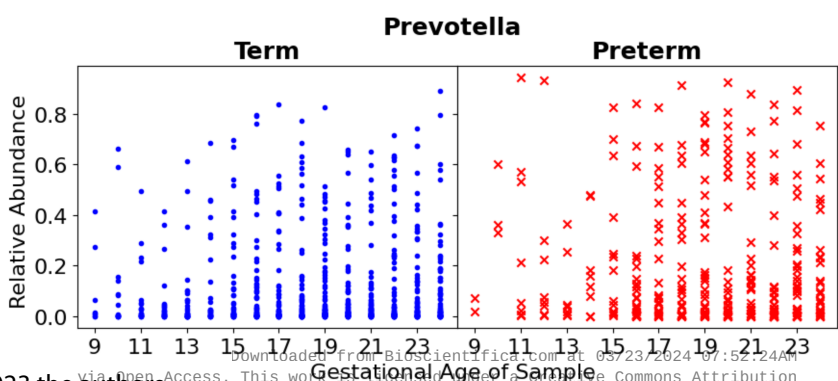
(c)



(d)



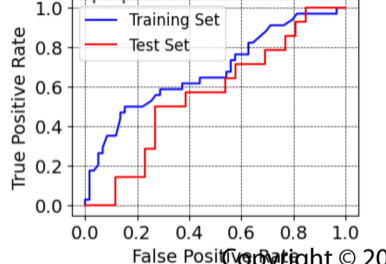
(e)



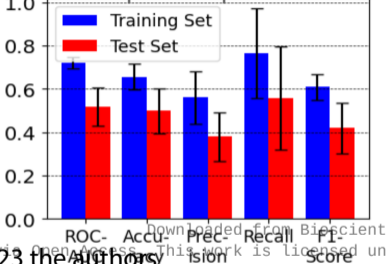
(f)



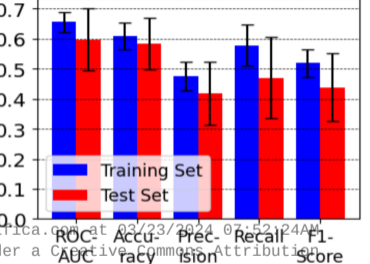
(a)



(b)

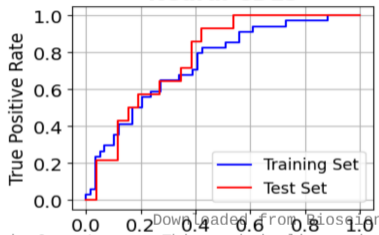


(c)



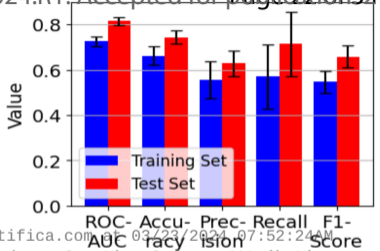
(d)

**ROC Curves:
Neural CDEs**



(a)

Metrics: Neural CDE



(b)

Downloaded from Bioscientifica.com at 03/23/2024 07:52:24AM
via Open Access. This work is licensed under a Creative Commons Attribution
4.0 International License
Copyright © 2023 the authors

http://creativecommons.org/licenses/by/4.0/deed.en_GB

Table 1: *P*-values representing significance of differences between diversity metric values across the term and preterm groups for gestational weeks 9-12, 13-16, 17-20 and 21-24.

Gestational week/Metric	<i>P</i>-value
Gestation weeks 9 – 12	
Shannon	0.8034
Simpson	0.9271
Gini	0.4074
Chao1	0.0049
TCS	0.0931
Gestation weeks 13 – 16	
Shannon	0.7225
Simpson	0.4616
Gini	0.9042
Chao1	0.0190
TCS	0.5035
Gestation weeks 17 – 20	
Shannon	0.6221
Simpson	0.4852
Gini	0.7739
Chao1	0.3352
TCS	0.8353
Gestation weeks 21 – 24	
Shannon	0.3058
Simpson	0.2270
Gini	0.2935
Chao1	0.5482
TCS	0.3751

Method	Metrics				
	ROC-AUC	Accuracy	Precision	Recall	F1-score
Decision Tree	0.46 ± 0.09	0.45 ± 0.12	0.39 ± 0.11	0.58 ± 0.22	0.41 ± 0.14
Random Forest	0.60 ± 0.10	0.59 ± 0.09	0.42 ± 0.10	0.48 ± 0.13	0.44 ± 0.12
Ours (Neural-CDE)	0.82 ± 0.02	0.74 ± 0.04	0.65 ± 0.05	0.71 ± 0.12	0.71 ± 0.03

Table 2: Comparative performance of various machine learning methods

Deep Learning Enables Early-Stage Prediction of Preterm Birth Using Vaginal Microbiota

Supplementary Material

1 Machine Learning Classifiers: Hyperparameters

1.1 Hyperparameter search spaces

Hyperparameters for both, the Random Forest (RF) and Decision Tree (DT) models were tuned by carrying out 3-fold cross-validation on the training set. For each set of hyperparameters, the training set was divided into 3 folds, and classification performance was evaluated on each fold after training on the 2 remaining folds. Accuracy, ROC-AUC, precision, recall and F1-score was computed for each fold. Finally, the optimal hyperparameters were selected based on mean F1-score (i.e., the hyperparameter set with the highest mean F1-score was selected as the optimal set). The search spaces of the hyperparameters are presented in supplementary table 1. The explanation and significance of parameters can be found in the scikit-learn documentation.

1.2 Optimal Hyperparameter Values

The optimal values of the hyperparameters, as identified by grid search, are listed in supplementary table 2.

2. Long Short-Term Memory (LSTM):

2.1 Implementation details

To facilitate the LSTM model to make predictions on a non-continuous and irregularly sampled dataset, we made two modifications to it. Let w_0 be the first time step for which data is available. Firstly, the hidden state is initialized by an embedding layer, which takes w_0 as input, and the cell state is initialized as a zero vector. Since week 9 is the earliest possible time point for which data is used, we denote it as $t = 0$. The LSTM model is run from weeks 9 to 24 ($t = 0$ to $t = 15$). From $t = 0$ to $t = w_0 - 1$, the LSTM cells at each time step do not update the hidden state, and simply act as dummy cells. At $t = w_0$, the hidden state is still the same as what was initialized with the embedding, and the LSTM cell at $t = w_0$ updates the hidden state based on the input data.

The LSTM cell at each time step is followed by a linear layer which generates the “forecast”, or the predicted genera abundance at that time step. As part of our second modification, for the time steps beyond $t = w_0$ for which input data is missing, we use this forecast to update the hidden state. Since the missing sampling intervals are different for each subject, we have to process each subject with its own set of rules to update the model weights. Thus, the batch size is restricted to 1.

2.2 Model Parameters

We used the Adam optimizer implemented within pytorch for training the model. The model parameters for LSTM are listed in supplementary table 3.

3. Neural Controlled Differential Equations (CDEs)

3.1 Training using adjoint backpropagation

Neural differential equation models are solutions to a time-dependent system whose function is a neural network that can be learned. This neural network is typically a linear fully connected network, however, a non-linearity can be introduced by using a tanh or a relu activation. In case of CDEs, there exists a “hidden state” (h), essentially a parametrized version of the input data, and the ODE can be represented as:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{h}} = \mathcal{N}_{\theta}(\mathbf{h}) \frac{\partial \mathcal{L}}{\partial \mathbf{h}}$$

Where \mathcal{N}_{θ} is the neural network with parameters θ . Training the model essentially refers to update these parameters based on the gradient of the cost function (\mathcal{L}) with respect to the parameter vector, scaled by a learning rate (η).

$$\theta_{t+1} = \theta_t + \eta \frac{\partial \mathcal{L}}{\partial \theta}$$

And the gradient of the cost function can be represented as:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \int_0^1 g(\mathbf{h}, \theta) \frac{\partial \mathcal{L}}{\partial \theta}$$

Where g is the loss function (binary cross-entropy in our case). The above expression essentially maps this gradient to a scalar, using which the parameters are updated. This can be further represented as:

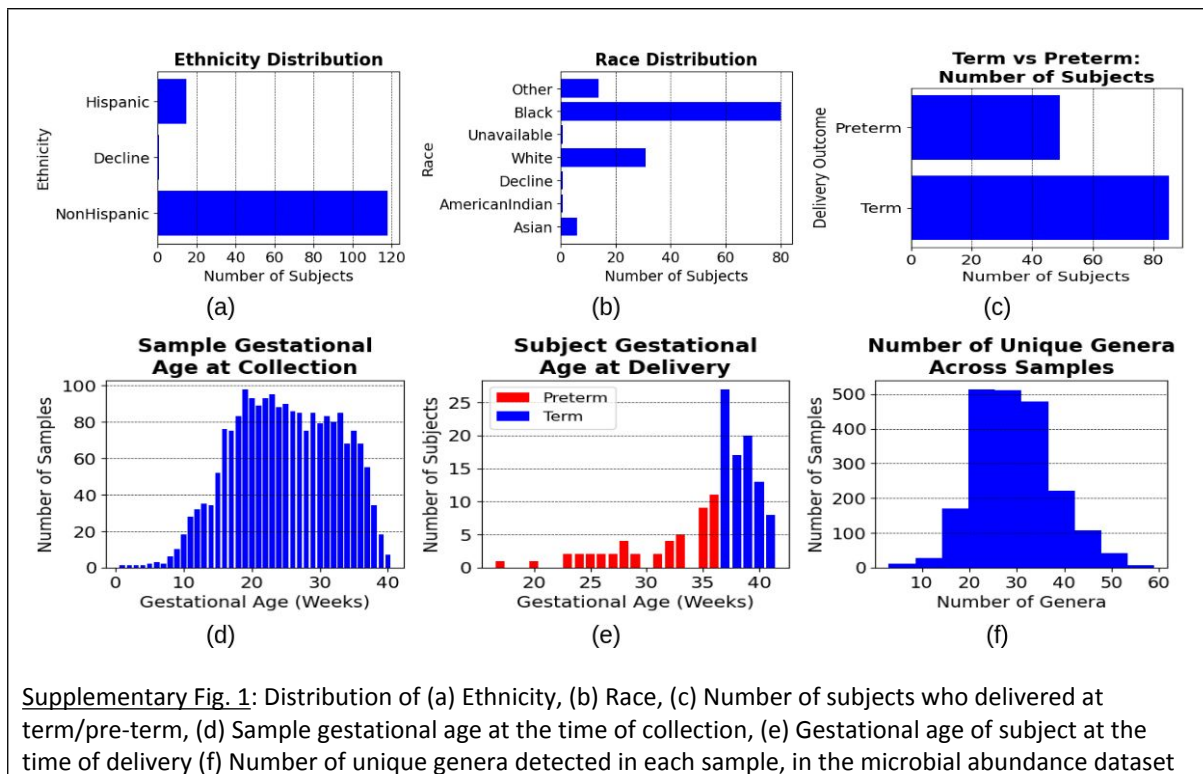
$$\frac{\partial \mathcal{L}}{\partial \theta} = \int_0^1 \frac{\partial \mathcal{L}}{\partial \theta} + \frac{\partial \mathcal{L}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \theta}$$

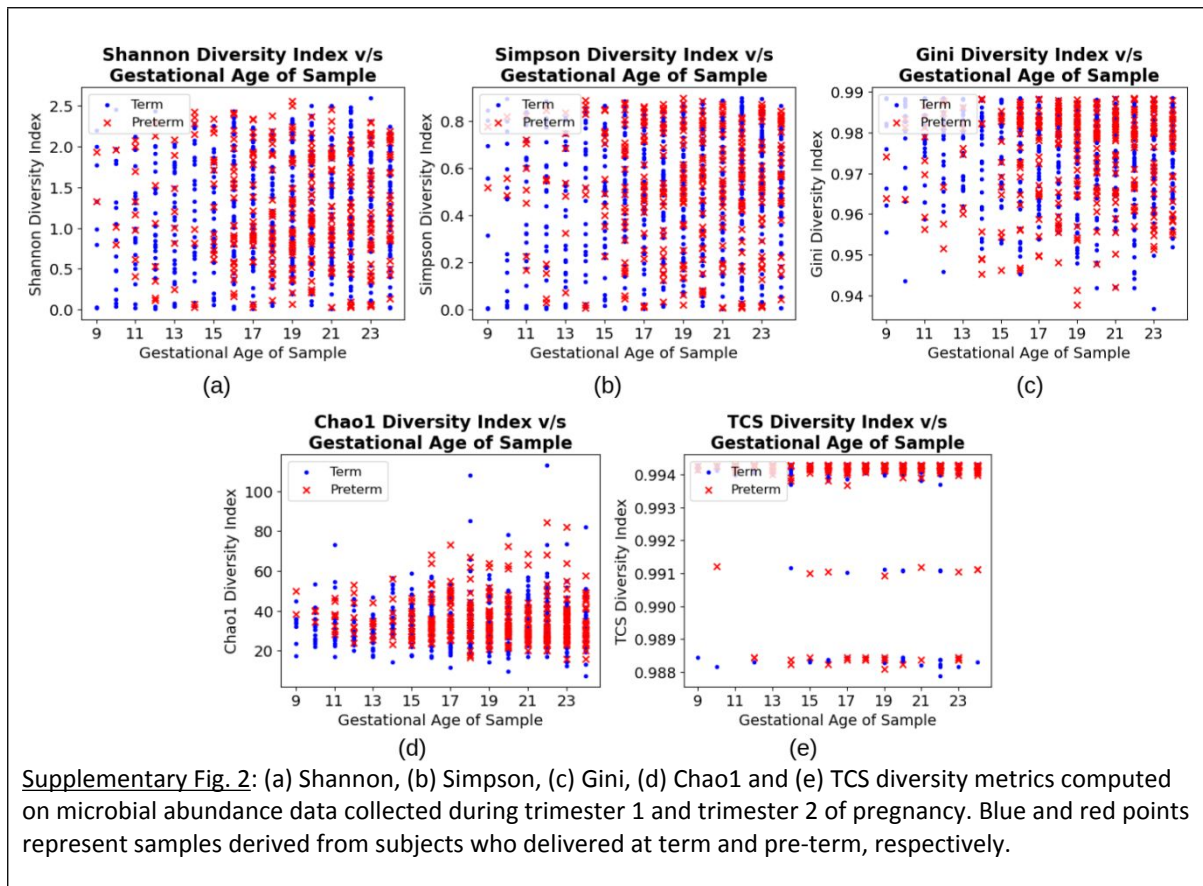
Adjoint sensitivity analysis or the adjoint state method essentially refers to computing the cost function gradient using the above expression in an efficient manner. The exact mathematical derivation for this is comprehensively described in Chen et al. [1].

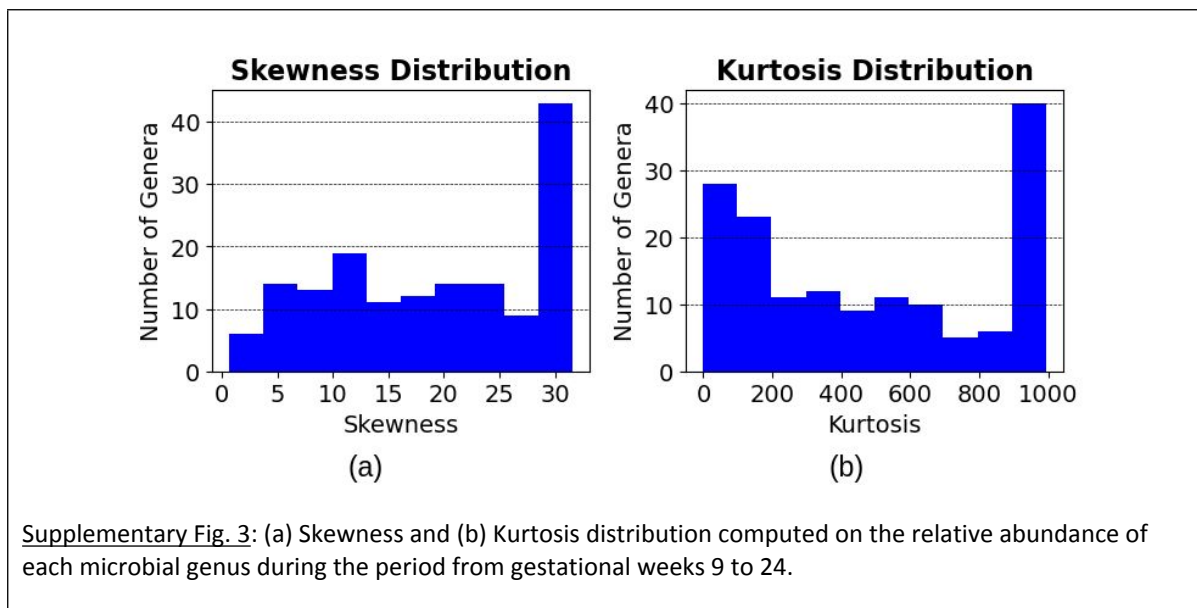
3.2 Model parameters

Neural CDE model parameters are listed in supplementary table 4.

Supplementary Figures







Supplementary Tables

Parameter	Search Spaces
n_estimators (RF only)	50, 100, 150
max_depth	None, 2, 3, 4
min_samples_split	0.25, 0.33, 0.5, 2, 4
max_features	sqrt, log2, 0.25, 0.33, 0.5, 2, 4
max_leaf_nodes	None, 2, 4
min_samples_leaf	0.25, 0.33, 0.5, 2, 4
criterion	entropy, gini, log_loss
splitter (DT only)	best, random
class_weight	balanced, balanced_subsample

Supplementary Table 1: Hyperparameter search spaces for the Random Forest and Decision Tree classifiers. The search spaces for both models are common unless specified otherwise.

Parameter	Optimal Value	
	Random Forest	Decision Tree
n_estimators	50	-
max_depth	4	3
min_samples_split	0.33	0.5
max_features	0.25	0.25
max_leaf_nodes	2	4
min_samples_leaf	0.33	0.2
criterion	entropy	gini
splitter	-	best
class_weight	balanced_subsample	balanced

Supplementary Table 2: Optimal hyperparameter values for the Random Forest and Decision Tree classifiers.

Parameter	Value	Explanation
Input channels	30	Dimension of the input (number of genera in the input dataset)
Hidden channels	128	Dimension of the hidden state
Embedding input dimensions	14	Number of possible inputs to embedding layer (i.e., number of unique weeks between 9 and 24 which serve as first available data sample)
Learning rate	1e-4	Learning rate for updating weights
Learning rate decay	1 (No decay)	Reduction in learning rate across epochs
Loss function	Binary cross-entropy	The function to minimize while optimizing the LSTM model
Number of epochs	35	Number of epochs (iterations) for which the model is trained
Batch size	1 (Fixed)	Number of samples per training batch. Batch size is restricted to 1 due to different time steps for which data is missing, which forces handling each sample individually

Supplementary Table 3: Parameters for the long short-term memory model

Parameter	Value	Explanation
Input channels	6	Dimension of the input (number of genera in the input dataset)
Hidden channels	64	Dimension of the hidden state
Output channels	1 (Fixed)	Dimension of the output (fixed to 1, since we are only outputting a single probability)
Interpolation type	Cubic	Method of interpolation of missing data from nearest observations
Learning rate	1e-4	Learning rate for updating weights
Learning rate decay	0.5	Reduction in learning rate across epochs
Learning rate decay step size	15	Number of epochs after which learning rate decay is applied
Loss function	Binary cross-entropy	The function to minimize while optimizing the LSTM model
Positive class weight	1.5	Weighted penalty to apply to negative class predictions (i.e., encourage predictions of positive class)
Batch size	10	Number of samples per training batch

Supplementary Table 4: Parameters for the Neural CDE model.

References

- [1]. R. T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. Neural ordinary differential equations. ArXIV, 2018.