

1 **Dense phenotyping from electronic health records enables machine-learning-based**  
2 **prediction of preterm birth**

3  
4 Abin Abraham<sup>1,2</sup>, Brian Le<sup>3</sup>, Idit Kosti<sup>3,4</sup>, Peter Straub<sup>1,5</sup>, Digna R. Velez-Edwards<sup>1,6,7</sup>, Lea K. Davis<sup>1,8,9</sup>, J.  
5 M. Newton<sup>7</sup>, Louis J. Muglia<sup>10</sup>, Antonis Rokas<sup>6,11</sup>, Cosmin A. Bejan<sup>6</sup>, Marina Sirota<sup>3,4</sup>, John A.  
6 Capra<sup>1,6,11,12\*</sup>

7  
8 **Affiliations:**

- 9 <sup>1</sup>Vanderbilt Genetics Institute, Vanderbilt University, Nashville, TN 37235, USA.  
10 <sup>2</sup>Vanderbilt University Medical Center, Vanderbilt University, Nashville, TN 37232, USA.  
11 <sup>3</sup>Bakar Computational Health Sciences Institute, University of California, San Francisco, San Francisco, CA, USA.  
12 <sup>4</sup>Department of Pediatrics, University of California, San Francisco, San Francisco, CA, USA.  
13 <sup>5</sup>Division of Genetic Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, TN,  
14 USA.  
15 <sup>6</sup>Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA.  
16 <sup>7</sup>Department of Obstetrics and Gynecology, Vanderbilt University Medical Center, Nashville, TN, USA.  
17 <sup>8</sup>Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA.  
18 <sup>9</sup>Department of Psychiatry and Behavioral Sciences, Division of Genetic Medicine, Vanderbilt University Medical  
19 Center, Nashville, TN, USA.  
20 <sup>10</sup>Burroughs-Wellcome Fund, Research Triangle Park, North Carolina, USA.  
21 <sup>11</sup>Department of Biological Sciences, Vanderbilt University.  
22 <sup>12</sup>Bakar Computational Health Sciences Institute, University of California, San Francisco

23  
24 \*Corresponding author: [tony.capra@ucsf.edu](mailto:tony.capra@ucsf.edu)  
25

26 **Abstract:** Identifying pregnancies at risk for preterm birth, one of the leading causes of worldwide infant  
27 mortality, has the potential to improve prenatal care. However, we lack broadly applicable methods to  
28 accurately predict preterm birth risk. The dense longitudinal information present in electronic health  
29 records (EHRs) is enabling scalable and cost-efficient risk modeling of many diseases, but EHR resources  
30 have been largely untapped in the study of pregnancy. Here, we apply machine learning to diverse data  
31 from EHRs to predict singleton preterm birth. Leveraging a large cohort of 35,282 deliveries, we find that  
32 machine learning models based on billing codes alone can predict preterm birth risk at various gestational  
33 ages (e.g., ROC-AUC=0.75, PR-AUC=0.40 at 28 weeks of gestation) and outperform comparable  
34 models trained using known risk factors (e.g., ROC-AUC=0.65, PR-AUC=0.25 at 28 weeks). Examining  
35 the patterns learned by the model reveals it stratifies deliveries into interpretable groups, including high-  
36 risk preterm birth sub-types enriched for distinct comorbidities. Our machine learning approach also  
37 predicts preterm birth sub-types (spontaneous vs. indicated), mode of delivery, and recurrent preterm  
38 birth. Finally, we demonstrate the portability of our approach by showing that the prediction models  
39 maintain their accuracy on a large, independent cohort (5,978 deliveries) from a different healthcare  
40 system. By leveraging rich phenotypic and genetic features derived from EHRs, we suggest that machine  
41 learning algorithms have great potential to improve medical care during pregnancy.

## 42 43 **Introduction**

44 Preterm birth, occurring before 37 weeks of completed gestation, affects approximately 10% of  
45 pregnancies globally[1–3] and is the leading cause of infant mortality worldwide[4,5]. The causes of  
46 preterm birth are multifactorial, since different biological pathways and environmental exposures can  
47 trigger premature labor[6]. Large epidemiological studies have identified many risk factors, including  
48 multiple gestations[1], cervical anatomic abnormalities[7], and maternal age[8]. Notably, even though a  
49 history of preterm birth [9] is one of the strongest risk factors, the recurrence rate remains low at <  
50 30%[10,11]. Additionally, maternal race is associated with risk for preterm birth; Black women have

51 twice the prevalence compared to white women[1,12]. Preterm births have a heterogenous clinical  
52 presentation and cluster based on maternal, fetal, or placental conditions[3]. These obstetric and systemic  
53 comorbidities (e.g. pre-existing diabetes, cardiovascular disease) can also increase risk for preterm  
54 birth[13,14].

55 Despite our understanding of numerous risk factors, there are no accurate methods to predict preterm  
56 birth. Some biomarkers associate with preterm birth, but their best performance is limited to a subset of  
57 all cases[15,16]. Recently, analysis of maternal cell-free RNA and integrated -omic models have  
58 emerged as promising approaches[17–19], but initial results were based on a small pregnancy cohort and  
59 require further validation. In silico classifiers based on demographic and clinical risk factors have the  
60 advantage of not requiring serology or invasive testing. However, even in large cohorts (>1 million  
61 individuals), demographic- and risk-factor-based models report limited discrimination (AUC=0.63-  
62 0.74)[20–23]. To date, we lack effective screening tools and preventative strategies for prematurity[24].

63 Electronic health records (EHRs) are scalable, readily available, and cost-efficient for disease-risk  
64 modeling[25]. EHRs capture longitudinal data across a broad set of phenotypes with detailed temporal  
65 resolution. EHR data can be combined with socio-demographic factors and family medical history to  
66 comprehensively model disease risk[26–28]. EHRs are also increasingly being augmented by linking  
67 patient records to molecular data, such as DNA and laboratory test results[29]. Since preterm birth has a  
68 substantial heritable risk[30], combining rich phenotypes with genetic risk may lead to better prediction.

69 Machine learning models have shown promise for accurate risk stratification across a variety of clinical  
70 domains[31–33]. However, despite the rapid adoption of machine learning in translational research, a  
71 review of 107 risk prediction studies reported that most models used only few variables, did not consider  
72 longitudinal data, and rarely evaluated model performance across multiple sites[34]. Studies using  
73 machine learning to predict preterm birth have relied on small cohorts, subsets of preterm birth, and are  
74 rarely replicated in external datasets[22,35–37]. Pregnancy research is especially well poised to benefit

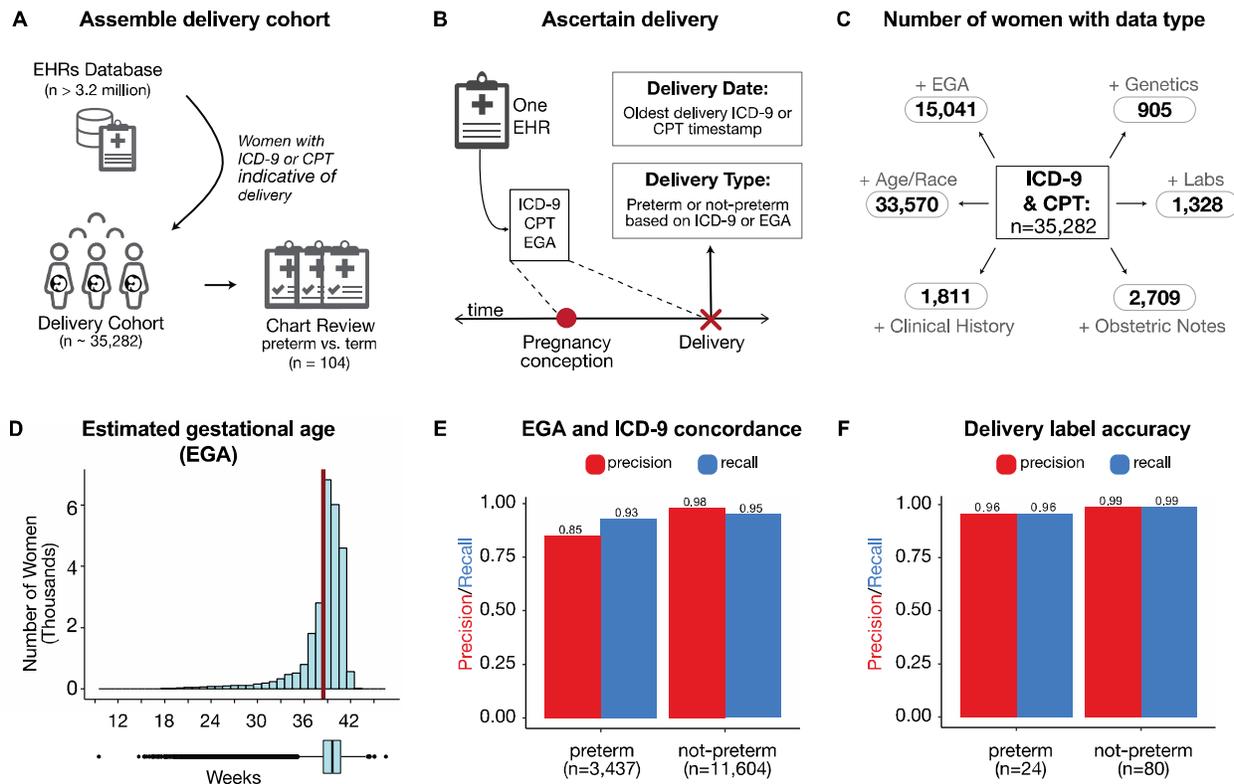
75 from machine learning approaches[26]. Per standard of care during pregnancy, women are carefully  
76 monitored with frequent prenatal visits, medical imaging, and clinical laboratories tests. Compared to  
77 other clinical contexts, pregnancy and the corresponding clinical surveillance occur in a defined  
78 timeframe based on gestational length. Thus, EHRs are well-suited for modeling pregnancy  
79 complications, especially when combined with the well documented outcomes at the end of pregnancy.  
80 In this study, we combine multiple sources of data from EHRs to predict preterm birth using machine  
81 learning. From Vanderbilt's EHR database (>3.2 million records) and linked genetic biobank (>100,000  
82 individuals), we identified a large cohort of women (n=35,282) with documented deliveries. We trained  
83 models (gradient boosted decision trees) that combine demographic factors, clinical history, laboratory  
84 tests, and genetic risk with billing codes (ICD-9 and CPT) to predict preterm birth. We find models  
85 trained on only billing codes show potential for predicting preterm birth. Billing code based models  
86 outperform a similar model using only known clinical risk factors. Across a variety of clinical contexts,  
87 such as second or spontaneous preterm birth, our models maintain accuracy. By investigating the patterns  
88 learned by our models, we identify clusters with distinct preterm birth risk and comorbidity profiles.  
89 Finally, we demonstrate the generalizability of our billing-code-based models on an external, independent  
90 cohort from the University of California, San Francisco (UCSF, n=5,978). Prediction models trained at  
91 Vanderbilt maintain high accuracy in the external cohort with only a modest drop in performance. Our  
92 findings provide a proof-of-concept that machine learning on rich phenotypes in EHRs show promise for  
93 portable, accurate, and non-invasive prediction of preterm birth. The strong predictive performance across  
94 clinical context and preterm birth subtypes argues that machine learning models have the potential to add  
95 value during the management of pregnancy; however, further work is needed before these models can be  
96 applied in clinical settings.

97

98

99

100 **Results**



101  
 102  
 103  
 104 **Figure 1. Definition and attributes of Vanderbilt delivery cohort.** (A) Schematic overview of the assembly of the  
 105 delivery cohort from electronic health records (EHRs). Using billing codes, women with at least one delivery were  
 106 extracted from the EHR database (n=35,282). (B) Delivery date and type were ascertained using ICD-9, CPT, and/or  
 107 estimated gestational age (EGA) from each woman’s EHR (Methods). From this cohort, 104 randomly selected  
 108 EHRs were chart reviewed to validate the preterm birth label for the first recorded delivery. (C) Number of women  
 109 in billing code cohort with estimated gestational age (+EGA), demographics (+Age, self- or third-party reported  
 110 Race), clinical labs (+Labs), clinical obstetric notes (+Obstetric notes), patient clinical history (+Clinical History),  
 111 and genetic data (+Genetics). (D) The EGA distribution at delivery (mean 38.5 weeks (red line); 38.0-40.3 weeks,  
 112 25<sup>th</sup> and 75<sup>th</sup> percentiles). Less than 0.015% (n=49) deliveries have EGA below 20 weeks. (E) The concordance  
 113 between estimated gestational age (EGA) within three days of delivery and ICD-9 based delivery type for the 15,041  
 114 women with sufficient data for both. Precision and recall values were > 93% across labels except for preterm  
 115 precision (85%). (F) Accuracy of delivery type phenotyping. The phenotyping algorithm was evaluated by chart  
 116 review of 104 randomly selected women. The approach has high precision and recall for binary classification of  
 117 ‘preterm’ or ‘not-preterm’.

118

119

120 **Assembling pregnancy cohort and ascertaining delivery type from Vanderbilt EHRs**

121 From the Vanderbilt EHR database (>3.2 million patients), we identified a ‘delivery cohort’ of 35,282

122 women with at least one delivery in the Vanderbilt hospital system (Fig. 1A). In addition to ICD and CPT

123 billing codes, we extracted demographic data, past medical histories, obstetric notes, clinical labs, and  
124 genome-wide genetic data for the delivery cohort. Because billing codes were the most prevalent data in  
125 this cohort (n=35,282), we quantified the pairwise overlap between billing codes and each other data type.  
126 The largest subset included women with billing codes paired with demographic data (n=33,570). The  
127 smallest subset was women with billing codes paired with genetic data (n=905; Fig. 1C). The mean  
128 maternal age at the first delivery in the delivery cohort was 27.3 years (23.0–31.0 years, 25<sup>th</sup> and 75<sup>th</sup>  
129 percentiles, Fig. S1A). The majority of women in the cohort self- or third-party reported as white  
130 (n=21,343), Black (n=6,178), or Hispanic (n=3,979; Fig. S1B). The estimated gestational age (EGA)  
131 distribution had a mean of 38.5 weeks (38.0 to 40.3 weeks, 25<sup>th</sup> to 75<sup>th</sup> percentile; Fig. 1D). The rate of  
132 multiple gestations (e.g. twins, triplets) was (7.6%, n=1,353). Since multiple gestation pregnancies are  
133 more likely to deliver preterm, we developed prediction models using singleton pregnancies unless  
134 otherwise stated.

135 To determine the delivery date and type (preterm vs. not-preterm) at scale across our large cohort, we  
136 developed a phenotyping algorithm using delivery-specific billing codes and estimated gestational age at  
137 delivery. For women with multiple pregnancies, we only considered the earliest delivery. We find that  
138 delivery-specific billing codes that can be used to label preterm births have high concordance (PPV $\geq$ 0.85,  
139 Recall  $\geq$ 0.95) with EGA based delivery labels (Fig. 1E). Our final algorithm combined billing codes and  
140 EGA when available (n=15,041, Fig. 1C). To evaluate the accuracy of the ascertained delivery labels, a  
141 domain expert blinded to the delivery type reviewed clinical notes from 104 EHRs selected at random  
142 from the delivery cohort. The algorithm had high accuracy: precision (positive predictive value) of 96%  
143 and recall (sensitivity) of 96% using the chart reviewed label as the gold standard (Fig. 1F).

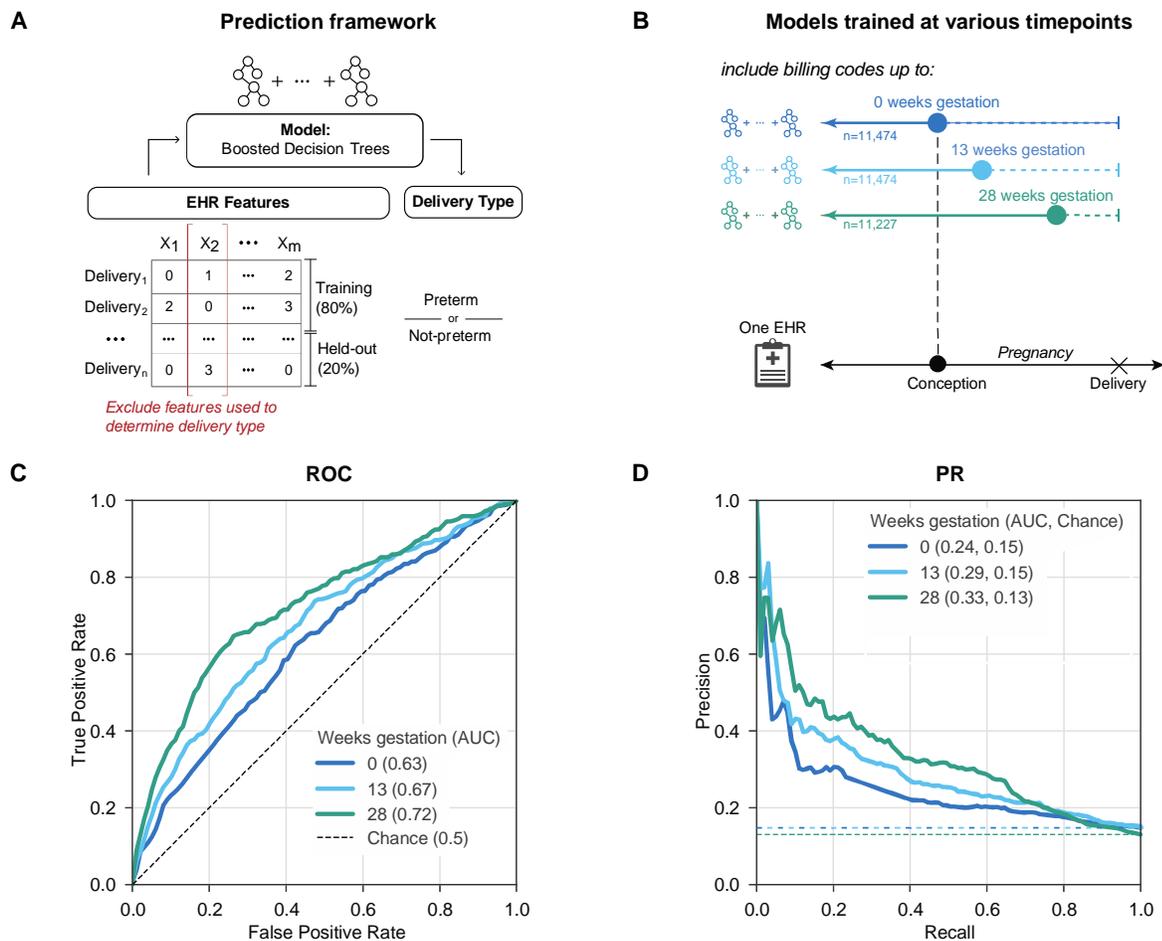
144

#### 145 ***Boosted decision trees using billing codes identify preterm deliveries***

146 Using this richly phenotyped delivery cohort, we evaluated how well the entire clinical phenome, defined  
147 as billing codes (ICD-9 and CPT) before and after delivery, could identify preterm births. With counts of

148 each billing code (excluding those used to ascertain delivery type), we trained gradient boosted decision  
149 trees[38] to classify each mother's first delivery as preterm or not-preterm. Boosted decisions trees are  
150 well-suited for EHR data because they require minimal transformation of the raw data, are robust to  
151 correlated features, and capture non-linear relationships[39]. Moreover, boosted decision trees have been  
152 successfully applied on a variety of clinical tasks[27,40,41].

153 In all evaluations, we held out 20% of the cohort as a test set and used the remaining 80% for training and  
154 validation (Fig. 2A). Boosted decision tree models trained on ICD-9 and CPT codes accurately identified  
155 preterm births (singletons and multiple gestations) with PR-AUC=0.86 (chance=0.22) and ROC-  
156 AUC=0.95 (Fig. S2A, B). While the combined ICD-9 and CPT based model achieved the best  
157 performance, models trained on either ICD-9 or CPT individually also performed well (PR-AUC  $\geq$ 0.82;  
158 chance=0.22, ROC-AUC  $\geq$ 0.93). All three models demonstrated good calibration with low Brier scores  
159 ( $\leq$ 0.092; Fig. S2C). Thus, billing codes across an EHR show potential as a discriminatory feature for  
160 *predicting* preterm birth.



161  
 162 **Figure 2. Machine learning classifiers accurately predict preterm birth using billing codes present before 28**  
 163 **weeks of gestation.** (A) Machine learning framework for training and evaluating all models. We train models  
 164 (boosted decision trees) on 80% of each cohort to predict the delivery as preterm or not-preterm. EHR features used  
 165 to ascertain delivery type are excluded from training. Performance is reported on the held-out cohort consisting of  
 166 20% of deliveries using area under the ROC and precision-recall curves (ROC-AUC, PR-AUC). (B) We trained  
 167 models using billing codes (ICD-9 and CPT) present before each of the following timepoints during pregnancy: 0,  
 168 13, and 28 weeks of gestation. These timepoints were selected to approximate gestational trimesters. Women who  
 169 already delivered were excluded at each timepoint. To facilitate comparison across timepoints, we downsampled  
 170 cohorts available so that the models were trained on a cohort with similar numbers of women (n=11,227 to 11,474).  
 171 (C) The ROC-AUC increased from conception at 0 weeks (0.63, dark blue line) to 28 weeks of gestation (0.72,  
 172 green line) compared to a chance (black dashed line) AUC of 0.5. (D) The model at 28 weeks of gestation achieved  
 173 the highest PR-AUC (0.33). This is an underestimate of the possible performance; the accuracy improves further  
 174 when all women with data available at 28 weeks are considered (Fig. 4B,C). Chance (dashed lines) represents the  
 175 preterm birth prevalence in each cohort.

176

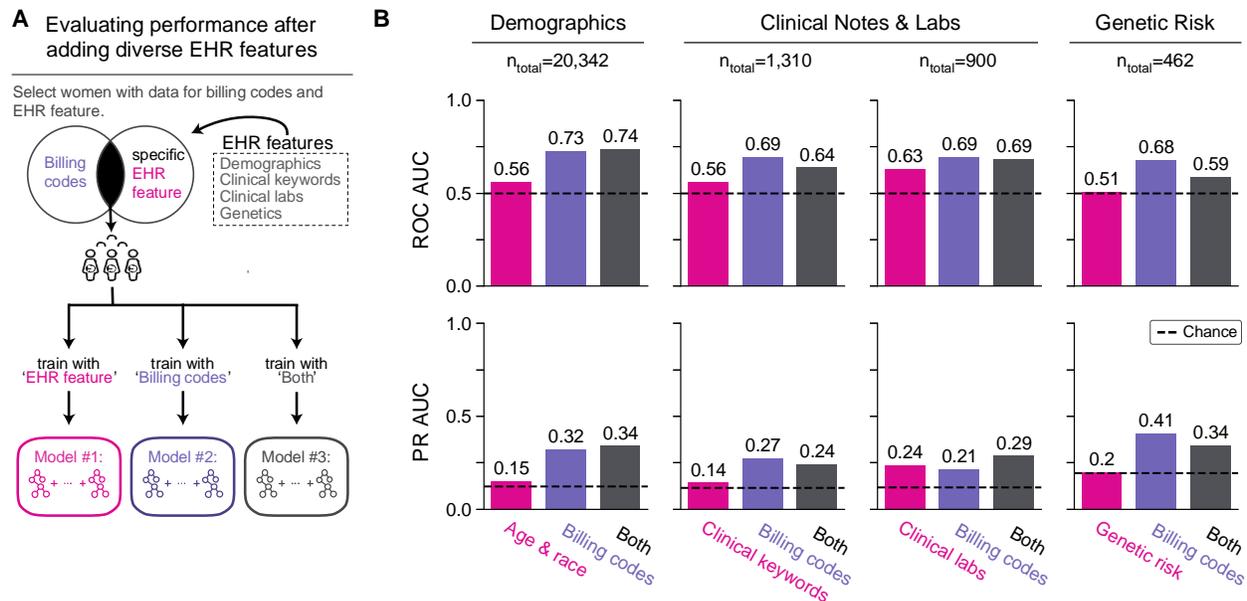
177 **Accurate prediction of preterm birth at 28 weeks of gestation**

178 To evaluate preterm birth prediction in a clinical context, we trained a boosted decision tree model (Fig.

179 2A) on billing codes present before each of the following timepoints: 0, 13, and 28 weeks of gestation

180 (Fig. 2B). These timepoints were selected to approximately reflect pregnancy trimesters. We  
181 downsampled to achieve comparable number of singleton deliveries across each timepoint ( $n=11,227$  to  
182  $11,474$ ) to mitigate sample size as a potential confounder while comparing performance. We only  
183 considered active pregnancies at each timepoint; for example, a delivery at 27 weeks would not be  
184 included in the 28 week model, since the outcome would already be known. The ROC-AUC increased  
185 from conception (0 weeks; 0.63) to the highest performance at 28 weeks (0.72; Fig. 2C). The PR-AUC  
186 (Fig. 2D), which accounts for preterm birth prevalence, is highest at 28 weeks (0.33, chance=0.13).  
187 However, as we show in the next section, this is an underestimate of the ability to predict preterm delivery  
188 at 28 weeks due to the down-sampling of the number of training examples. As expected, when we  
189 included multiple gestations, the model performed even better (PR-AUC=0.42 at 28 weeks, chance=0.14;  
190 Fig. S3). Results were similar when models were trained using billing codes available before different  
191 timepoints from the date of delivery (Fig. S4).

192 To test whether differences in contact with the health system between cases and controls were driving  
193 performance, we trained a classifier based on the total number of codes in an individual's EHR before  
194 delivery to predict preterm birth. This simple classifier failed to discriminate between delivery types with  
195 PR-AUC and ROC-AUC only slightly higher than chance (PR-AUC=0.19, chance=0.19; ROC-  
196 AUC=0.56, chance=0.5, Fig. S5). Therefore, cumulative disease burden or the number of contacts alone  
197 are not informative for predicting preterm birth.



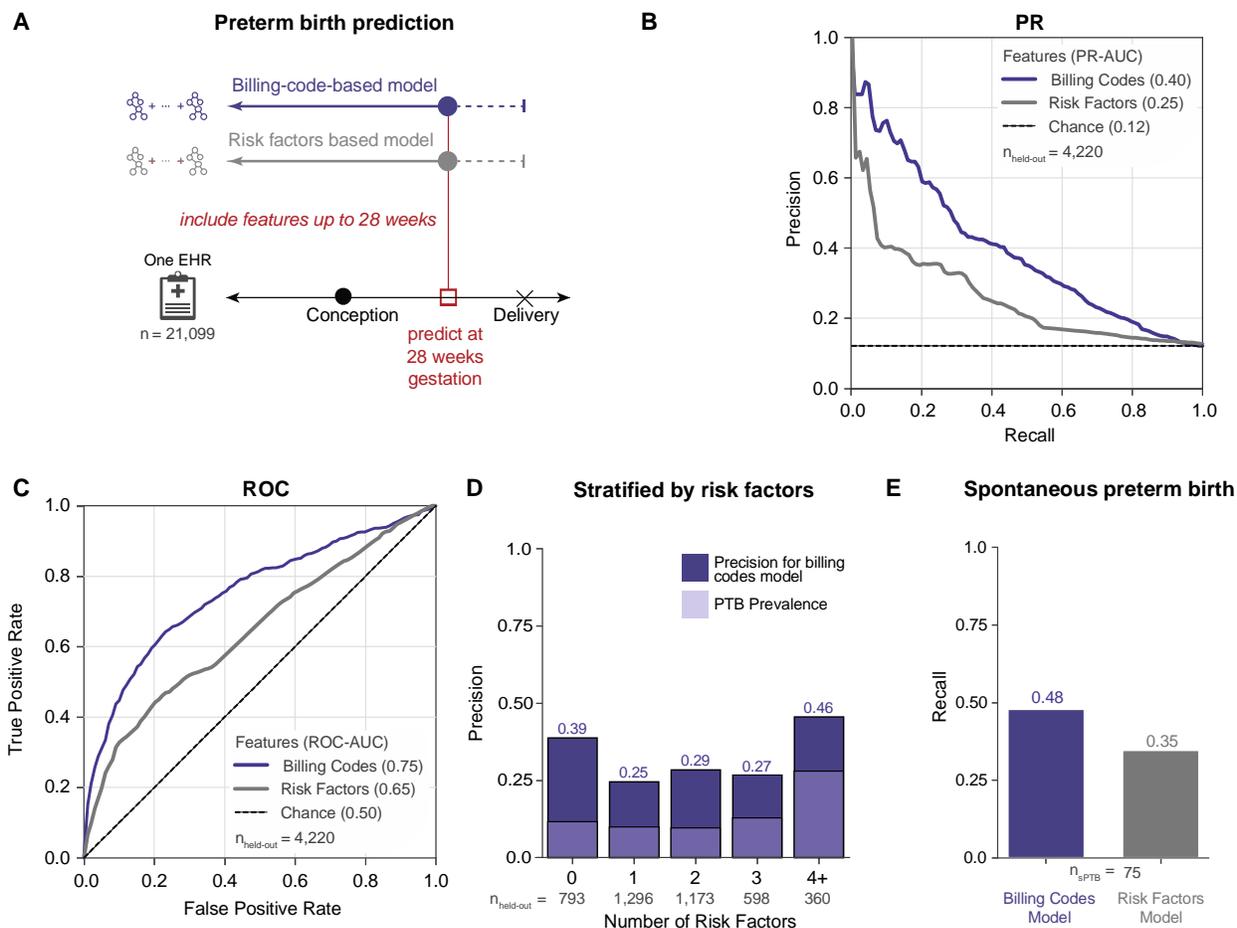
198  
 199 **Figure 3. Combining demographic, clinical, and genetic features does not substantially improve preterm birth**  
 200 **prediction compared to using only billing codes.** (A) Framework for evaluating change in preterm birth prediction  
 201 performance after incorporating diverse types of EHR features with billing codes (ICD-9 and CPT codes). We used  
 202 only features and billing codes occurring before 28 weeks of gestation. EHR features are grouped by sets of  
 203 demographic factors (age and race), clinical keywords (UMLS concept unique identifiers from obstetric notes),  
 204 clinical labs, and genetic risk (polygenic risk score for preterm birth). We compared three models for each feature  
 205 set: 1) using only the feature set being evaluated (pink), 2) using only billing codes ('Billing codes', purple), and 3)  
 206 using the feature set combined with billing codes ('Both', gray). For each feature set, we considered the subset of  
 207 women who had at least one recorded value for the EHR feature and billing codes. All three models for a given EHR  
 208 feature set considered the same pregnancies, but there are differences in the cohorts considered across features set  
 209 due to differences in data availability; n<sub>total</sub> is the number of women (training and held-out) for each feature set. (B)  
 210 Each of the three models (x-axis) and their ROC-AUC and PR-AUC (y-axes) are shown. Each of the additional  
 211 EHR features performed worse than the billing codes only model and did not substantially improve performance  
 212 when combined with the billing codes. Dotted lines represent chance of 0.5 for ROC-AUC and the preterm birth  
 213 prevalence for PR-AUC. Even when including EHR features before and after delivery in this framework revealed  
 214 the same pattern with no substantial improvement in predictive performance compared to the billing codes only  
 215 model (Fig. S6).  
 216  
 217

### 218 *Integrating other EHR features does not improve model performance*

219 In addition to billing codes, EHRs capture aspects of an individual's health through different types of  
 220 structured and unstructured data. We tested whether incorporating additional features from EHRs can  
 221 improve preterm birth prediction. Models were evaluated using data available at 28 weeks of gestation;  
 222 we selected this time point as a tradeoff between being sufficiently early for some potential interventions  
 223 and late enough for sufficient data to be present to enable accurate predictions using billing codes. From  
 224 the EHRs, we extracted sets of features including demographic variables (age, race), clinical keywords

225 from obstetric notes, clinical lab tests ran during the pregnancy, and predicted genetic risk (polygenic risk  
226 score for preterm birth). To measure the performance gain for each feature set, we compared models  
227 trained using: the feature set only, billing codes only, and billing codes combined with the feature set  
228 (Fig. 3A). Within each feature set, the same pregnancies comprised the training and held-out sets for the  
229 three models. However, the number of deliveries (training + held-out sets) varied widely across feature  
230 sets (n=462 to 20,342) due to the differing availability of each feature type.

231 Models using only demographic factors, clinical keywords, and genetic risk had ROC-AUC and PR-AUC  
232 similar to chance (Fig. 3B). Clinical labs had moderate predictive power with ROC-AUC of 0.63 and PR-  
233 AUC of 0.24 (Fig. 3B). Compared to models using only billing codes, adding additional feature sets did  
234 not substantially improve performance (Fig. 3B). We note that some features sets, such as clinical labs  
235 and genetic risk, were evaluated on held-out sets with small numbers of deliveries (180 and 92,  
236 respectively). However, even after increasing the sample size by including women who may have features  
237 either before or after delivery, we did not observe a consistent gain in performance compared to models  
238 trained using only billing codes (Fig. S6).



239  
240  
241 **Figure 4. Billing-code-based model outperforms a model based on clinical risk factors.** (A) We compared the  
242 performance of boosted decision trees trained using either billing codes (ICD-9 and CPT) present before 28 weeks  
243 of gestation (purple) or known clinical risk factors (gray) to predict preterm delivery. Clinical risk factors (Methods)  
244 included self- or third-party reported race (Black, Asian, or Hispanic), age at delivery (> 34 or <18 years old), non-  
245 gestational diabetes, gestational diabetes, sickle cell disease, presence of fetal abnormalities, pre-pregnancy BMI  
246 >35, pre-pregnancy hypertension (>120/80), gestational hypertension, preeclampsia, eclampsia, and cervical  
247 abnormalities. Both models were trained and evaluated on the same cohort of women (n = 21,099). (B) Precision-  
248 recall and (C) ROC curves for model using billing codes (purple line) or clinical risk factors (gray line). Preterm  
249 births are predicted more accurately by models using billing codes at 28 weeks of gestation (PR-AUC = 0.40, ROC-  
250 AUC = 0.75) than using clinical risk factors as features (PR-AUC = 0.25, ROC-AUC = 0.65). For the precision-  
251 recall curves chance performance is determined by the preterm birth prevalence (dashed black line). (D) Billing-  
252 code-based prediction model performance stratified by number of risk factors for an individual. The billing-code-  
253 based model detects more preterm cases and has higher precision (dark purple) across all numbers of risk factors  
254 compared to preterm (PTB) prevalence (light purple). (E) The model using billing codes also performs well at  
255 predicting the subset of spontaneous preterm births in the held-out set (recall = 0.48) compared to risk factors (recall  
256 = 0.35).

### 257 **Models using billing codes outperforms prediction from risk factors**

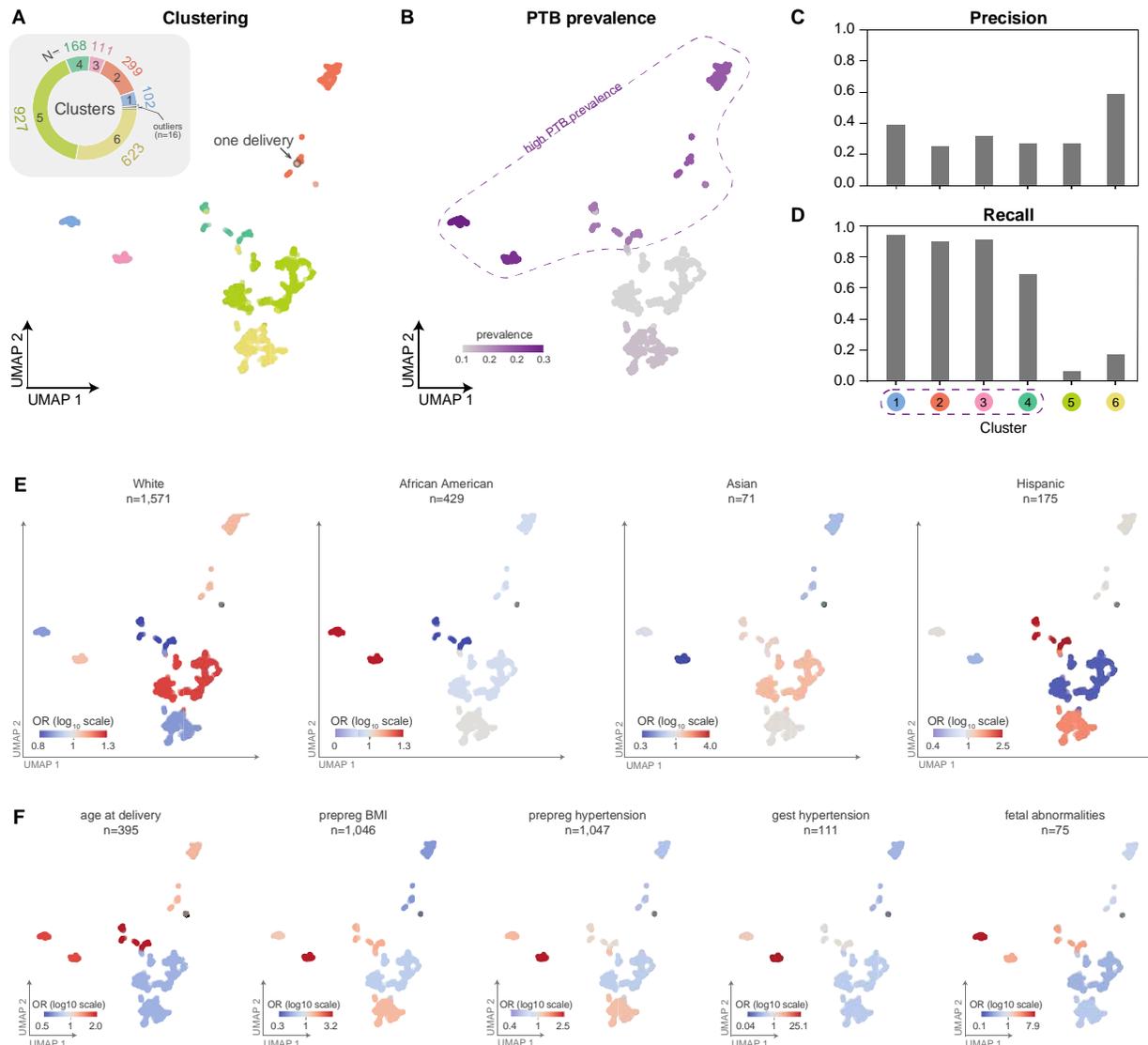
258 Although there are well known risk factors for preterm birth, few validated risk calculators exist and even  
259 fewer are routinely implemented in clinical practice[42]. We evaluated how a prediction model  
260 incorporating only common risk factors associated with moderate to high risk for preterm birth compared

261 to a model using billing codes, which captured a broad range of comorbidities, at 28 weeks of gestation.  
262 We included maternal and fetal risk factors that occurred before and during the pregnancy and across  
263 many organ systems[3,13,23,43], race[20], age at delivery[44–46], pre-gestational and gestational  
264 diabetes[47], sickle cell disease[48], fetal abnormalities[13], pre-pregnancy hypertension, gestational  
265 hypertension (including pre-eclampsia or eclampsia)[1,49], and cervical abnormalities[50] (Methods).  
266 The billing-code-based model significantly outperformed a model trained with clinical risk factors at  
267 predicting preterm birth at 28 weeks of gestation (PR-AUC=0.40 vs. 0.25, ROC-AUC=0.75 vs. 0.65; Fig.  
268 4B, C). The stronger performance of the billing-code-based classifier was true for women across the  
269 spectrum of comorbidity burden; it had higher precision across individuals with different numbers of risk  
270 factors. Performance peaked for individuals with 0 (precision=0.39) and 4+ (precision=0.46) risk factors,  
271 but we did not observe a trend between model performance and increasing number of clinical risk factors  
272 (Fig. 4D). This suggests that machine learning approaches incorporating a comprehensive clinical  
273 phenome can add value to predicting preterm birth.

#### 274 ***Machine learning models can predict spontaneous preterm births***

275 The multifactorial etiologies of preterm birth lead to clinical presentations with different comorbidities  
276 and trajectories. Medically-indicated and idiopathic spontaneous preterm births are distinct in etiologies  
277 and outcomes. Identifying pregnancies that ultimately result in spontaneous preterm deliveries is  
278 particularly valuable, and we anticipated that spontaneous preterm birth would be more challenging to  
279 predict than preterm birth overall. To test this, we identified spontaneous preterm births in the held-out set  
280 at 28 weeks of gestation by excluding women with medically induced labor, a cesarean section delivery,  
281 or PPROM (Methods). We intentionally used a conservative phenotyping strategy that aimed to minimize  
282 false positive spontaneous preterm births to evaluate the model's ability to predict spontaneous preterm  
283 births. The prediction model trained using billing codes up to 28 weeks of gestation classified 48%  
284 (recall) of all spontaneous preterm births (n=75) as preterm; this is significantly higher than the risk factor  
285 only model (recall = 35%; Fig. 4E).

286  
287



288  
289

**Figure 5. Machine-learning-based clustering of deliveries identifies sub-groups with distinct preterm birth prevalence, clinical features, and prediction accuracy.** (A) For the model predicting preterm birth at 28 weeks of gestation using billing codes (ICD-9 and CPT, Figure 4A), we assigned deliveries from the held-out test set ( $n=2,246$ ) to one of six clusters (colors) using density-based clustering (HDBSCAN) on the SHAP feature importance matrix. For visualization of the clusters, we used UMAP to embed the deliveries into a low dimensional space based on the matrix of feature importance values. Inset pie chart displays count of individuals in each cluster. (B) The preterm birth prevalence (colorbar) in each cluster. The algorithm discovered four clusters with high PTB prevalence (enclosed by dashed line). (C) Precision and (D) recall for preterm birth classification within each cluster. (E) The enrichment (odds ratios, colorbar in  $\log_{10}$  scale) of race as derived from EHRs for each cluster (Table S1). (F) The enrichment ( $\log_{10}$  odds ratio) of relevant clinical risk factors in each cluster (Table S2). Risk factors include: age at delivery ( $> 34$  or  $< 18$  years old), pre-pregnancy BMI (prepreg BMI), pre-pregnancy hypertension (prepreg hypertension), gestational hypertension (gest hypertension), and fetal abnormalities. We report the total number of women in the delivery cohort at high risk for each clinical risk factor ( $n$ ). Enrichments for additional risk factors are given in Fig. S7.

305 ***Preterm birth prediction algorithm stratifies deliveries into clusters with different preterm birth risk***  
306 ***and distinct comorbidity signatures***

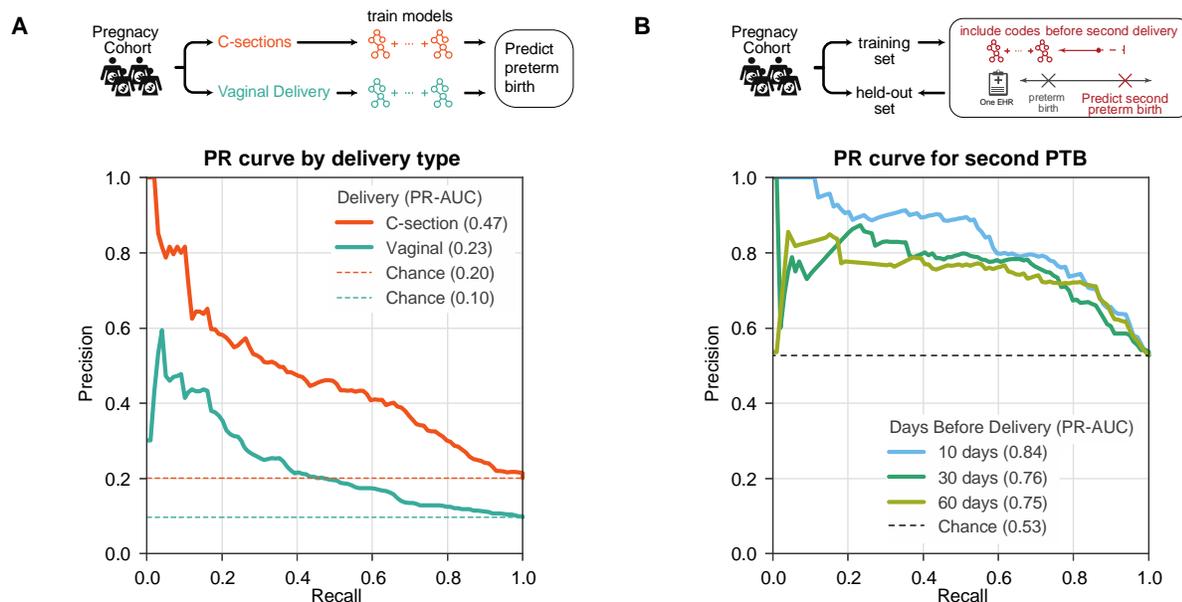
307 Understanding the statistical patterns identified by machine learning models is crucial for their adoption  
308 into clinical practice. Unlike deep learning approaches, decision tree-based models are easier to interpret.  
309 We calculated feature importance as measured by SHapley Additive exPlanation (SHAP) scores[51,52]  
310 for each delivery and feature pair in the held-out cohort for the model using billing codes at 28 weeks of  
311 gestation ('Billing-code-based model', Figure 4A). SHAP scores quantify the marginal additive  
312 contribution of each feature to the model predictions for each individual. Next, we performed a density-  
313 based clustering on the patient by feature importance matrix and visualized clusters using UMAP (Fig.  
314 5A, Methods). This approach focuses the clustering on the features for each individual prioritized by the  
315 algorithm. We identified six clusters with 927 to 102 women. PTB prevalence was the high in the clusters  
316 one to four (blue, pink, green, orange, Fig. 5B) indicating differential risk for preterm birth. Performance  
317 varied across the clusters; the yellow cluster with low PTB prevalence had the highest PPV while clusters  
318 with higher PTB prevalence had higher recall (Fig. 5C, D).

319 To evaluate whether clusters had distinct phenotype profiles, we calculated the enrichment of  
320 demographic and clinical risk factor traits within each cluster using Fisher's exact test (Methods). These  
321 traits were extracted from structured fields in EHRs or ascertained using combinations of billing codes.  
322 Although these billing codes are used to train the model, the combination of codes used to ascertain risk  
323 factor traits are not encoded in the training data. White women are significantly enriched in the cluster 5  
324 (odds ratio, OR = 1.2, p-value = 0.02, Fisher's exact test, Figure 5E). Hispanic women also had  
325 significant positive enrichment in cluster four (OR = 2.5, p-value = 0.0002) and cluster 6 (OR = 1.6, p-  
326 value = 0.008) and were depleted (negative enrichment) in the cluster five (OR= 0.5, p-value = 4.42E-6,  
327 Figure 5D). African American and Asian women also exhibit modest enrichment in different clusters  
328 (Table S1).

329 We also tested for enrichment of clinical risk factors of preterm birth in the clusters. We observed distinct  
 330 patterns of enrichment and depletion for each clinical risk factor (Fig. 5F, Fig. S7). Gestational  
 331 hypertension had strong and enrichment in cluster three (OR = 26.4, p-value = 9.0E-39). Fetal  
 332 abnormalities demonstrated a similar pattern with strong enrichment in cluster one (OR = 8.5, p-value =  
 333 2.07E-10). Extreme age at delivery (>34 or <18 years old) was enriched, though more weakly, (OR = 1.2  
 334 to 2.2) for the all clusters except five and six. Pre-pregnancy BMI, pre-pregnancy hypertension, and  
 335 gestational hypertension had similar patterns with the strongest enrichment in cluster three. The remaining  
 336 clinical risk factors show similar patterns and are provided in Fig. S7 and Table S2.

337 By analyzing the feature importance values through UMAP embeddings, we identify interpretable clusters  
 338 of individuals discovered by the machine learning model that reflect the complex and multi-faceted paths  
 339 to preterm birth. Overall, the learned rules highlight relationships between clinical factors and preterm  
 340 birth prevalence. For example, some risk factors, such as age at delivery, are enriched in all clusters with  
 341 high preterm birth prevalence. Other factors, such as pre-pregnancy BMI and hypertension, are strongly  
 342 enriched only in specific clusters with high preterm birth prevalence. Thus, this approach enables us to  
 343 interpret phenotypic patterns of risk and identify subgroups among cases learned from complex EHR  
 344 features by the prediction model.

345



346  
347 **Figure 6. Preterm birth prediction accuracy is influenced by clinical context.** (A) Preterm birth prediction  
348 models trained and evaluated only on cesarean section (C-section) deliveries perform better (PR-AUC=0.47) than  
349 those trained only on vaginal delivery (PR-AUC=0.23). ROC-AUC patterns were similar (Fig. S8). Billing codes  
350 (ICD-9 and CPT) present before 28 weeks of gestation were used to train a model to distinguish preterm from non-  
351 preterm birth for either C-sections (n=5,475) or vaginal deliveries (n=15,487). (B) Recurrent preterm birth can be  
352 accurately predicted from billing codes. We trained models to predict preterm birth for a second delivery in a cohort  
353 of 1,416 high-risk women with a prior preterm birth documented in their EHR. Three models were trained using  
354 data available at 10 days, 30 days, and 60 days before the date of second delivery. Models accurately predict the  
355 birth type in this cohort of women with a history of preterm birth (PR-AUC $\geq$ 0.75). ROC-AUC varied from 0.82 at  
356 10 days to 0.77 at 60 days before second delivery (Fig. S9). Expected performance by chance is the preterm birth  
357 prevalence in each cohort (dashed lines).  
358

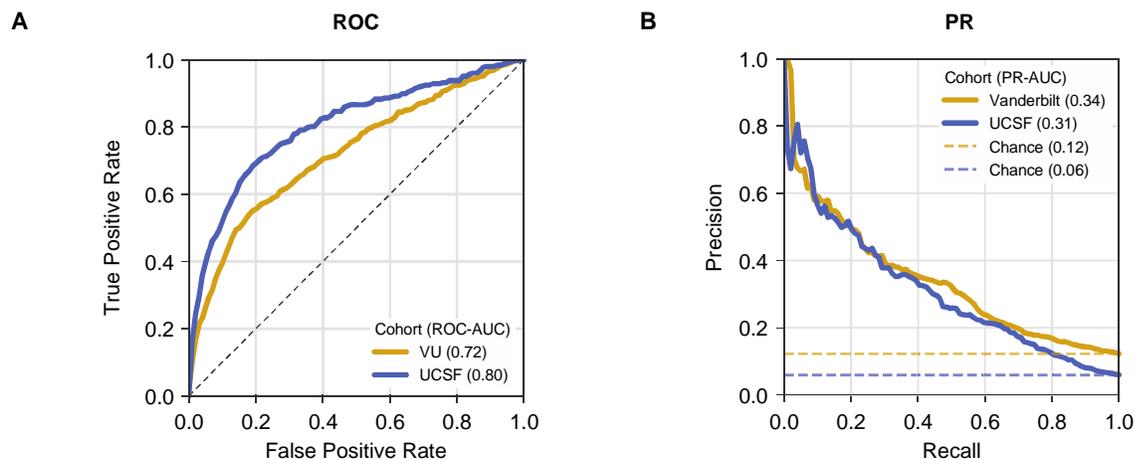
### 359 *Performance varies based on clinical context and delivery history*

360 To further explore the sensitivity of the performance of our approach to clinical context and patient  
361 history, we evaluated how delivery type (vaginal vs. cesarean-section) and a previous preterm birth  
362 influence preterm birth prediction. We trained two classifiers using billing codes (ICD-9 and CPT)  
363 occurring before 28 weeks of gestation: one on a cohort of cesarean-section (n = 5,475) singleton  
364 deliveries and one on vaginal deliveries (n = 15,487). Preterm birth prediction accuracy was higher in the  
365 cesarean-section cohort (PR-AUC = 0.47, chance = 0.20) compared to the vaginal delivery cohort (PR-  
366 AUC = 0.23, chance = 0.10; Fig. 6A). Cesarean-sections also had higher ROC-AUC compared to vaginal  
367 deliveries (0.75 vs. 0.68, Fig. S8). As expected, the preterm birth prevalence was higher in the cesarean-  
368 section cohort.

369 Women with a history of preterm birth are at significantly higher risk for a subsequent preterm birth than  
370 women without a previous history. Therefore, it is particularly important to understand the drivers of risk  
371 in this cohort. We tested if models trained on EHR data of women with a history of preterm birth could  
372 accurately predict the status of their next birth. We assembled 1,416 women with a preterm birth and a  
373 subsequent delivery in the cohort and split them into a training set (80%) and held-out test set (20%) to  
374 evaluate the model performance (Methods). For these women, 53% of the second deliveries were preterm.  
375 Due to limited availability of estimated gestational age data for the recurrent preterm births, which is  
376 necessary to approximate the date of conception, we trained models using billing codes (ICD-9 and CPT)

377 present before each of the following timepoints: 10, 30, and 60 days before the delivery. These models  
378 were all able to discriminate term from preterm deliveries better than chance (Fig. 6B; PR-AUCs $\geq$ 0.75).  
379 The model predicting a second preterm birth as early as 60 days before delivery achieved the high  
380 performance with PR-AUC=0.75 (Fig. 6B, chance=0.53) and ROC-AUC=0.77 (Fig. S9).

381



382  
383

384 **Figure 7. Preterm birth prediction models accurately generalize to an independent cohort.** Performance of  
385 preterm birth prediction models trained at Vanderbilt applied to UCSF cohort. Models were trained using ICD-9  
386 codes present before 28 weeks of gestation at Vanderbilt on 16,857 of women and evaluated on a held-out set at  
387 Vanderbilt (n=4,215, gold) and UCSF cohort (n=5,978, blue). (A) Models accurately predicted preterm birth at  
388 Vanderbilt (ROC-AUC=0.72) and at UCSF (ROC-AUC=0.80). The higher ROC-AUC at UCSF is driven by the  
389 lower prevalence of preterm birth in this cohort. (B) Models performed better than baseline prevalence (chance)  
390 based on the precision-recall curve at Vanderbilt (PR-AUC=0.34) and at UCSF (PR-AUC=0.31). Note that in  
391 contrast to models presented previously this one was trained only on ICD-9 codes, due to the lack of CPT codes in  
392 the UCSF cohort. Feature importance estimates were strongly correlated between the two cohorts (Fig. S10). Cohort  
393 demographics are given in Table S3.

394

### 395 *Models trained at Vanderbilt accurately predict preterm birth in an independent cohort at UCSF*

396 To evaluate whether preterm birth prediction models trained on the Vanderbilt cohort performed well on  
397 EHR data from other databases, we compared their performance on the held-out Vanderbilt cohort  
398 (n=4,215) and an independent cohort from UCSF (n=5,978). The UCSF cohort was ascertained using  
399 similar rules as the Vanderbilt cohort (Methods); age and distribution of race are provided in Table S3.  
400 However, we note that the UCSF cohort has a lower preterm birth prevalence (6%) compared to the  
401 Vanderbilt cohort (13%).

402 To facilitate the comparison, we trained models to predict preterm birth in the Vanderbilt cohort using  
403 only ICD-9 codes present before 28 weeks of gestation. We did not consider CPT codes in this analysis  
404 due to differences in the available billing code data between Vanderbilt and UCSF. As expected from the  
405 previous results, the model accurately predicted preterm birth in the held-out set from Vanderbilt (PR-

406 AUC of 0.34, chance=0.12), but performance was slightly lower than using both ICD and CPT codes  
407 (Fig. 4B). The model trained at Vanderbilt also achieved strong performance in the UCSF cohort. The  
408 classifier had a higher ROC-AUC (0.80) in UCSF cohort compared to the Vanderbilt cohort (0.72; Fig.  
409 7A) and PR-AUC of 0.31 vs 0.34 at Vanderbilt; Fig. 7B). The higher ROC is due to the lower prevalence  
410 of preterm birth in the UCSF cohort and the sensitivity of ROC-AUC to class imbalance[53]. Overall,  
411 these models show striking reproducibility across two independent cohorts.

#### 412 *Similar features are predictive across the independent cohorts*

413 The architecture of boosted decision trees enables straightforward identification of features (ICD-9 codes)  
414 with the largest influence on the model predictions. We used SHAP[54,55] scores to quantify the  
415 marginal additive contribution of each feature to the model predictions for each individual. For each  
416 feature in the ICD-9-based model, we calculated the mean absolute SHAP values across all women in the  
417 held-out set. The mean absolute SHAP value for each feature was highly correlated (spearman  $R=0.93$ ,  $p$ -  
418 value  $< 2.2E-308$ ) between the held-out Vanderbilt set and the UCSF cohort (Fig. S10A). The top 15  
419 features ranked based on the mean absolute SHAP value captured known risk factors (fetal abnormalities,  
420 history of preterm birth, etc.), pregnancy screening and supervision of high-risk pregnancies (Fig. S10B).  
421 Ten of the top 15 features were shared across both cohorts. The full list of SHAP values across all  
422 features are provided in Table S4. This suggests that the model's discovered combination of phenotypes,  
423 including expected risk factors, and the corresponding weights assigned by the machine learning model  
424 are generalizable across cohorts.

425

## 426 **Discussion**

427 Preterm birth is a major health challenge affecting 5-20% of pregnancies[1,2,12] and lead to significant  
428 morbidity and mortality[56,57]. Predicting preterm birth risk could inform clinical management, but no  
429 accurate classification strategies are routinely implemented[24]. Here, we take a step toward addressing  
430 this need by demonstrating the potential for machine learning on dense phenotyping from EHRs to predict  
431 preterm birth in challenging clinical contexts (e.g., spontaneous and recurrent preterm births). However,  
432 we emphasize that more work is needed before these approaches are ready for the clinic. Compared to  
433 other data types in the EHRs, models using billing codes alone had the highest prediction accuracy and  
434 outperformed those using clinical risk factors. Demonstrating the potential broad applicability of our  
435 approach, the model accuracy was similar in an external independent cohort. Combinations of many  
436 known risk factors and patterns of care drove prediction; this suggests that the algorithm builds on  
437 existing knowledge. Thus, we conclude that machine learning based on EHR data has the potential to  
438 predict preterm birth accurately across multiple healthcare systems.

439 Decision tree based models are robust to correlated features, can identify complex non-linear  
440 combinations, and remain transparent for interpretation after training. In addition to these advantages,  
441 decision tree based models have demonstrated strong performance in various clinical prediction tasks[58–  
442 60]. Pregnancy is a clinical context with close monitoring and well defined end-points that may similarly  
443 benefit from machine learning approaches, yet few studies have applied decision tree based machine  
444 learning models to large pregnancy cohorts with rich clinical data[61].

445 Our approach has several distinct advantages compared to published preterm birth prediction models.  
446 First, our models have robust performance. Previous models using risk factors (diabetes, hypertension,  
447 sickle cell disease, history of preterm birth) to predict preterm birth, despite having cohorts up to two  
448 million women[23], have reported ROC-AUCs between 0.69 and 0.74[20–22]. Our models obtain a  
449 ROC-AUC of 0.75 and PR-AUC of 0.40 using data available at 28 weeks of gestation even after

450 excluding multiple gestations. Furthermore, given the unbalanced classification problem (preterm births  
451 are less common than non-preterm), we report high PR-AUCs in addition to high ROC-AUCs. A recent  
452 deep learning model trained using word embeddings from EHRs achieved a high performance (ROC-  
453 AUC = 0.83[61]). This model was evaluated over a stratified high-risk cohort consisting of birth before  
454 28 weeks of gestation. We did not stratify preterm births by severity since more than 85% of preterm  
455 births occur after 32 weeks of gestation[62], however, this is an important topic for future work. Our  
456 models achieve comparable performance with the benefit of easier interpretability, which is an advantage  
457 over deep learning approaches, and we discuss this further below.

458 Second, our models use readily available data throughout pregnancy that do not require invasive  
459 sampling. While some studies have also obtained high ROC-AUCs (e.g., 0.81-0.88), they used serum  
460 biomarkers across small cohorts[17] or acute obstetric changes within days of delivery[16]. The potential  
461 to enable cost-effective and broad application is illustrated by our evaluation of the classifiers on EHR  
462 data from UCSF; however, substantial further work is needed to move from this proof-of-concept analysis  
463 to clinic-ready models. Furthermore, the rich characterization of the phenome provided by EHRs  
464 leveraged by our approach could also complement more invasive biochemical assays.

465 Third, the gradient boosted decision trees we implement are easier to interpret than ‘black-box’ deep  
466 learning models that cannot easily identify features driving predictions. Transparency is an important, if  
467 not necessary, characteristic of machine/artificial learning models deployed in clinical practice[63,64],  
468 and it can facilitate discovery of insights and hypotheses to motivate future work. We reveal the patterns  
469 learned by our model by clustering deliveries using feature importance profiles. The enrichment for  
470 known risk factors in clusters with high preterm birth prevalence establishes confidence in our machine-  
471 learning based prediction models. In addition, we can quantify the strength of enrichment and  
472 combination of risk factors across clusters with distinct comorbid patterns. Since preterm birth is a  
473 heterogenous phenotype[6], and stratifying pregnancies based on clinical features may be critical to

474 uncovering the biological basis of labor[3,65,66], the learned rules from our model offer a possible  
475 method for sub-phenotyping.

476 Finally, our approach generalizes across hospital systems. We demonstrate that billing-code-based models  
477 trained at Vanderbilt achieve similar accuracy in an independent cohort from UCSF. The generalizability  
478 of machine learning models can be constrained by the sampling of the training data. Thus, the accurate  
479 prediction in an independent dataset from an external institution points to several inherent strengths of the  
480 approach. First, successful replication indicates the models' ability to learn predictive signals despite  
481 regional variation in assigning billing codes to an EHR. Second, the large cohorts used to train and  
482 evaluate models at Vanderbilt and UCSF guard against potential weakness of EHRs, such as miscoding or  
483 omission of key data points. These errors are unavoidable in EHRs[67], but the large cohort used to train  
484 our models mitigates these errors and enables the high accuracy in the UCSF dataset, even with its  
485 different demographics. Additionally, idiosyncratic patterns of patient care at the institution used to  
486 develop the algorithm, which would be present in the Vanderbilt training and held-out sets, are unlikely to  
487 be present in the external UCSF cohort and inflate the out-of-sample accuracy. Third, the top features  
488 driving model performance are shared across institutions and reflect combinations of known risk factors  
489 and patterns of care. This aids interpretability of the underlying algorithm and likely reflects underlying  
490 pathophysiology that is innate to preterm birth.

491 We see several avenues for further improving our algorithm. First, some of the top features reflected  
492 routine obstetric care for high-risk pregnancies. Thus, the learning problem could be engineered to force  
493 the algorithm to discover new unappreciated risk factors. Second, we were surprised that the addition of  
494 features beyond billing codes, such as lab values, concepts extracted from clinical notes, and genetic  
495 information did not significantly improve performance. In some cases, any redundant information already  
496 captured by the billing codes would not improve the model's accuracy; this is likely true for clinical  
497 notes. However, other sources, like currently available genetic data and polygenic risk scores, may not  
498 effectively capture underlying etiologies of preterm birth. Thus, these sources may not add more

499 discriminatory power due to limitations in current data. Indeed, the largest published genome-wide study  
500 for preterm birth only explains a very small fraction of the heritability[30], and a polygenic risk score  
501 derived from it was not predictive in our cohort. Further sub-phenotyping of preterm birth will not only  
502 aid in prediction, but also understanding its multifactorial etiology and developing personalized treatment  
503 strategies. More explicit modeling of the temporal dependence between EHR features may further  
504 increase performance. Finally, while we evaluated the ability of our classifiers to discriminate preterm  
505 births, further studies evaluating the calibration of these models are necessary to better risk stratify of  
506 pregnancies.

507 The strong predictive performance of our models suggests that they have the potential to be clinically  
508 useful. Compared to a machine learning model trained using only known risk factors, the billing-code-  
509 based classifier incorporated a broad set of clinical features and predicted preterm birth with higher  
510 accuracy. Furthermore, the superior performance was not driven by the number of risk factors or the total  
511 burden of billing codes. These results indicate the algorithm is not simply identifying less healthy  
512 individuals or those with greater healthcare usage. The models also accurately predicted many preterm  
513 births in challenging and important clinical contexts such as spontaneous and recurrent preterm birth.  
514 Spontaneous preterm births are common[1,12,68], and unlike iatrogenic deliveries, they are more  
515 difficult to predict because they are driven by unknown multifactorial etiologies[12,24]. Similarly, since  
516 a prior history of preterm birth is one of the strongest risk factors[69], distinguishing pregnancies most at  
517 risk for recurrent preterm birth has potential to provide clinical value.

518 However, we emphasize that additional work is needed before this approach is ready for clinical  
519 application. Though it has strong performance, a more comprehensive evaluation of the algorithm against  
520 current clinical practice is needed to determine how early and how much improvement in standard of care  
521 this approach could provide[70]. Furthermore, while our cohorts include diverse individuals and the  
522 algorithm generalizes well, the approach must be evaluated to ensure that it does not introduce or amplify  
523 biases against specific groups or types of preterm birth[71]. In addition, we anticipate further gains in the

524 clinical value of this approach as more modalities of data becomes incorporated in the EHR[72] and  
525 diverse populations become available. Addressing these questions and taking other necessary steps  
526 toward clinical utility will require the close collaboration of diverse experts from basic, clinical, social,  
527 and implementation sciences.

528 Our results provide a proof-of-concept that machine learning algorithms can use the dense phenotype  
529 information collected during pregnancy in EHRs to predict preterm birth. The prediction accuracy across  
530 clinical contexts and compared to existing risk factors suggests such modeling strategies can be clinically  
531 useful. We are optimistic that with the increasingly widespread adoption of, improvement in tools for  
532 extracting meaningful data from them, and integration of complementary molecular data, machine  
533 learning approaches can improve the clinical management of preterm birth.

534

## 535 **Materials and Methods**

### 536 *Ascertaining delivery type and date for Vanderbilt cohort*

537 We identified women with at least one delivery (n=35,282, ‘delivery-cohort’) at Vanderbilt Hospital  
538 based on the presence of delivery-specific billing codes (ICD-9/10 and CPT) or estimated gestational age  
539 (EGA) documented in the EHR. Combining delivery specific ICD-9/10 (‘delivery-ICDs’), CPT  
540 (‘delivery-CPTs’), and EGA data, we developed an algorithm to label each delivery as preterm or not  
541 preterm. Women with multiple gestations (e.g. twins, triplets) were identified using ICD and CPT codes  
542 and exclude for singleton-based analyses. See Supplementary Materials and Methods for exact codes.

543 We demarcate multiple deliveries by grouping delivery-ICDs in intervals of 37 weeks starting with the  
544 most recent delivery-ICD. This step is repeated until all delivery-ICDs in a patient’s EHR are assigned to  
545 a pregnancy. We chose 37-week intervals to maximally discriminate between pregnancies. For each  
546 delivery, we assign a list of labels (preterm, term, or postterm) ascertained using the delivery-ICDs. EGA  
547 values, extracted from structured fields across clinical notes, were mapped to multiple pregnancies using

548 the same procedure. For women with multiple EGA recorded in their EHR, the most recent EGA value  
549 determined the time interval to group preceding EGA values. Based on the most recent EGA value for  
550 each pregnancy, we assigned labels to each delivery (EGA <37 weeks: preterm;  $\geq 37$  and <42 weeks:  
551 term,  $\geq 42$  weeks: postterm). After pooling delivery labels based on delivery-ICDs and EGA, we assigned  
552 a consensus delivery label by selecting the oldest gestational age based classification (i.e. postterm > term  
553 > preterm). By incorporating both billing code and EGA based delivery label and selecting the oldest  
554 gestational classification, we expect this to increase the accuracy of this algorithm, which we evaluate by  
555 chart-review (described in detail below).

556 Since CPT codes do not encode delivery type, we combined the delivery-CPTs with timestamps of  
557 delivery-ICDs and EGAs to approximate the date of delivery. Delivery-CPTs were grouped into multiple  
558 pregnancies as described above. The most recent timestamp from delivery-CPTs, delivery-ICDs, and  
559 EGA values was used as the approximate delivery date for a given pregnancy.

#### 560 *Validating delivery type based on chart review*

561 To validate the delivery type ascertained from billing codes and EGA, we used chart-reviewed labels as  
562 the gold standard. For 104 randomly selected EHRs from the delivery cohort, we extracted the date and  
563 gestational age at delivery from clinical notes. For earliest delivery recorded in the EHR, we assigned a  
564 chart-review based label according to the gestational age at delivery (<37 weeks: preterm; 37 and 42  
565 weeks: term,  $\geq 42$  weeks: postterm). The precision/positive predictive value for the ascertained delivery  
566 type as a binary variable ('preterm' or 'not-preterm') was calculated using the chart reviewed label as the  
567 gold standard. To compare the ascertainment strategy to a simpler phenotyping algorithm, we compared  
568 the concordance of the label derived from delivery-ICDs to one based on the gestational age within three  
569 days of delivery. This simpler phenotyping approach resulted in a lower PPV (85%) and recall (93%; Fig.  
570 S1B) compared to the billing-code-based ascertainment strategy.

#### 571 *Training and evaluating gradient boosted decision trees to predict preterm birth*

572 All models for predicting preterm birth used boosted decision trees as implemented in XGBoost  
573 v0.82[38]. Unless stated otherwise, we trained models to predict the earliest delivery in a woman's EHR  
574 as preterm or not-preterm. The delivery cohort was randomly split into training (80%) and held-out (20%)  
575 sets with equal proportion of preterm cases. For prediction tasks, we used only ICD-9 and excluded ICD-  
576 10 codes to avoid potential confounding effects. The total count of billing codes within a specified time  
577 frame was used as features to train our models; if a woman never had a billing code in her EHR, we  
578 encoded these as '0'. For all models we excluded ICD-9, CPT codes, and EGA used to ascertain delivery  
579 type and date. On the training set, we use tree of Parzen estimators as implemented in hyperopt  
580 v0.1.1[73] to optimize hyperparameters by maximizing the mean average precision. The best set of  
581 hyperparameters was selected after 1,000 trials using 3-fold cross-validation over the training set (80:20  
582 split with equal proportion of preterm cases). We evaluated the performance of all models on the held-out  
583 set using Scikit-learn v0.20.2[74]. All performance metrics reported are on the held-out set. For  
584 precision-recall curves, we define baseline chance for each model as the prevalence of preterm cases. To  
585 ensure no data leaks were present in our training protocol, we trained and evaluated a model using a  
586 randomly generated dataset (n=1,000 samples) with a 22% preterm prevalence. As expected, this model  
587 did not do better than chance (AUC=0.50, PR-AUC=0.22, data not show). All trained models with their  
588 optimized hyperparameters are provided at [https://github.com/abraham-abin13/ptb\\_predict\\_ml](https://github.com/abraham-abin13/ptb_predict_ml).

### 589 *Predicting preterm birth at different weeks of gestation*

590 As a first step, we evaluated whether billing codes could discriminate between delivery types. Models  
591 were trained to predict preterm birth using the total counts of each ICD-9, CPT, or ICD-9 and CPT code  
592 across a woman's EHR. We excluded any codes used to ascertain delivery type or date. All three models  
593 were trained and evaluated on the same cohort of women who had at least one ICD-9 and CPT code (Fig  
594 S2).

595 Next, we evaluated machine learning models at 0, 13, 28, and 35 weeks of gestation by training using  
596 only features present before each timepoint. For the subset of women in our delivery cohort with EGA,

597 we calculated the date of conception by subtracting EGA (recorded within three days of delivery) from  
598 the date of delivery. Next, we trained models using ICD-9 and CPT codes timestamped before different  
599 gestational timepoints with only singleton (Fig. 2B) or including multiple gestations (Fig. S3). The same  
600 cohort of women was used to train and evaluate across models. The sample size varied slightly ( $n =$   
601 11,843 to 10,799) since women who already delivered were excluded at each timepoint.

602 In addition to evaluating models based on the date of conception, we trained models at different  
603 timepoints before the date of delivery (Fig. S4) using the same cohort of women by requiring every  
604 individual in this cohort had to have at least one ICD-9 or CPT code before each timepoint. Evaluating  
605 models before the date of delivery increased the sample size ( $n=15,481$ ) compared to a prospective  
606 conception-based design ( $n=12,410$ ) and yielded similar results.

#### 607 *Evaluating predictive potential of demographic, clinical, and genetic features from EHRs*

608 In addition to billing codes, we extracted structured and unstructured features from the EHRs (Fig. 3A).  
609 We evaluated models using features present before 28 weeks of gestations (Fig 3) and features present  
610 before or after delivery (Fig S6). Structured data included self or third-party reported race (Fig. 1E), age  
611 at delivery, past medical and family history (92 features, see Supplementary Materials and Methods), and  
612 clinical labs. For training models, we only included clinical labs obtained during the first pregnancy and  
613 excluded values greater than four standard deviations from the mean. To capture the trajectory of each  
614 clinical lab's values across pregnancy (307 clinical labs, see Supplementary Materials and Methods), we  
615 trained models using the mean, median, minimum, and maximum lab measurement. For unstructured  
616 clinical text in obstetric and nursing clinical notes, we applied CLAMP[75] to extract UMLS (Unified  
617 Medical Language System) concepts unique identifiers (CUIs and included those with positive assertions  
618 with  $> 0.5\%$  frequency across all EHRs). When training preterm birth prediction models, we one-hot  
619 encoded categorical features. No transformations were applied to the continuous features.

620 A subset of women ( $n=905$ ) was genotyped on the Illumina MEGA<sup>EX</sup> platform. We applied standard  
621 GWAS quality control steps[76] using PLINK v1.90b4s[77]. We calculated a polygenic risk score for

622 each white woman with genotype data based on the largest available preterm birth GWAS [30] using  
623 PRSice-2[78,79]. We assumed an additive model and summed the number of risk alleles at single  
624 nucleotide polymorphisms (SNPs) weighted by their strength of association with preterm birth (effect  
625 size). PRSice determined the optimum number of SNPs by testing the polygenic risk score for association  
626 with preterm birth in our delivery-cohort at different GWAS p-value thresholds. We included date of birth  
627 and five genetic principal components to control for ancestry. Our final polygenic risk score used 356  
628 preterm birth associated SNPs (GWAS p-value < 0.00025).

629 Using the structured and unstructured data derived from the EHR, we evaluated whether adding EHR  
630 features to billing codes could improve preterm birth prediction. Since the number of women varied  
631 across EHR feature, we created subsets of the delivery cohort for each EHR feature. Each subset included  
632 women with at least one recorded value for the EHR feature and billing codes. Then we trained three  
633 models as described above for each subset: 1) using only the EHR feature being evaluated, 2) using ICD-  
634 9 & CPT codes, and 3) using the EHR feature with ICD-9 & CPT codes. Thus, all three models for a  
635 given EHR feature were trained and evaluated on the same cohort of deliveries (Fig. 3A).

### 636 ***Predicting preterm birth using billing codes and clinical risk factors at 28 weeks of gestation***

637 We compared the performance of a model trained using billing codes (ICD-9 and CPT) present before 28  
638 weeks of gestation with a model trained using clinical risk factors to predict preterm delivery (Fig 4).  
639 Both models were trained and evaluated on the same cohort of women (n = 21,099). We selected well-  
640 established obstetric risk factors that included maternal and fetal factors across organ systems, occurred  
641 before and during pregnancy, and had moderate to high risk for preterm birth [3,13,23,43]. For each  
642 individual, risk factors were encoded as high-risk or low-risk binary values. Risk factors such as non-  
643 gestational diabetes status[47], gestational diabetes[47], gestational hypertension, pre-eclampsia or  
644 eclampsia[1,49], fetal abnormalities[13], cervical abnormalities[50], and sickle cell disease[48] status was  
645 defined based on at least one corresponding ICD-9 code occurring before the date of delivery  
646 (Supplementary Materials and Methods). The remaining factors, such as race (Black, Asian, or Hispanic

647 was encoded as higher risk)[20], age at delivery (> 34 or <18 years old)[44–46], pre-pregnancy BMI  $\geq$   
648 35, and pre-pregnancy hypertension (>120/80)[1,49], were extracted from structured fields in EHR. Pre-  
649 pregnancy value was defined as the most recent measurement occurring before nine months of the  
650 delivery date.

### 651 *Density based clustering on feature importance values*

652 To better understand the decision making process of our machine learning models, we calculated feature  
653 importance value for the model predicting preterm birth at 28 weeks of gestation. We used SHapley  
654 Additive exPlanation values (SHAP)[51,52,55] to determine the marginal additive contribution of each  
655 feature for each individual. First, we calculated a matrix of SHAP values of features by individuals from  
656 the held-out cohort. Since the shape of this matrix was too large to perform the density based clustering,  
657 we created an embedding using 30 UMAP components with default parameters as implemented in  
658 UMAPv0.3.8[80]. Next, we performed a density based hierarchical clustering using HDBSCANv0.8.26  
659 [81]. We used default parameters (metric=Euclidean) and tried a range of values for two hyperparameters:  
660 minimum number of individuals in each cluster ('min\_clust\_size) and threshold for determining outlier  
661 individuals who do not belong to a cluster ('min\_samples'). After tuning these two hyperparameters, we  
662 selected the clustering model with the highest density based cluster validity score [81], which measures  
663 the within and between cluster density connectedness. We find a min\_clust\_size = 110 and min\_samples  
664 = 10 had the highest density based cluster validity (DBCV) score with 6 distinct clusters with one cluster  
665 for outliers (Fig. S11). A minority of women (n=16) were not assigned to a cluster ('outliers'). To  
666 visualize the cluster assignments, we performed UMAP on the feature importance matrix with default  
667 settings and two UMAP components and colored each individual by their cluster membership. Finally, we  
668 calculated the preterm birth prevalence and accuracy within each cluster.

### 669 *Comorbidity enrichment within clusters*

670 We tested for enrichment of clinical risk factors within each cluster by using a Fisher Exact test as  
671 implemented in Scipy[82]. For each risk factor, we constructed a contingency table based on a given

672 cluster membership and being high risk for the risk factor. We report enrichment as the odds ratio with  
673 colorbar in log<sub>10</sub> scale of the odds ratio. For sickle cell disease, one cluster did not have any cases of  
674 sickle cell disease.

675 ***Evaluating model performance on spontaneous preterm births, by delivery type, and recurrent preterm***  
676 ***birth***

677 We compared how models trained used billing codes (ICD-9 & CPT) performed in different clinical  
678 contexts. First, we evaluated the accuracy of predicting spontaneous preterm birth using models trained to  
679 predict all types of preterm births. From all preterm cases in the held-out set, we excluded women who  
680 met any of the following criteria to create a cohort of spontaneous preterm births: medically induced  
681 labor, delivery by cesarean section, or preterm premature rupture of membranes. The ICD-9 and CPT  
682 codes used to identify exclusion criteria are provided in Supplementary Materials and Methods. We  
683 calculated recall/sensitivity as the number of predicted spontaneous preterm births out of all spontaneous  
684 preterm births in the held-out set. We used the same approach to quantify performance of models trained  
685 using clinical risk factors (Fig. 4E).

686 We trained models to predict preterm birth among cesarean sections and vaginal deliveries separately  
687 using billing codes (ICD-9 & CPT) as features. Deliveries were labeled as cesarean sections or vaginal  
688 deliveries if they had at least one relevant billing code (ICD-9 or CPT) occurring within ten days of the  
689 date of first delivery in EHR. Billing codes used to determine delivery type are provided in  
690 Supplementary Materials and Methods. Deliveries with billing codes for both cesarean and vaginal  
691 deliveries were excluded. We trained separate models to predict cesarean and vaginal deliveries (Fig. 6A  
692 and Fig. S8).

693 We evaluated how well models using billing codes could predict recurrent preterm birth. From our  
694 delivery cohort, we retained women whose first delivery in the EHR was preterm and a second delivery  
695 for which we ascertained the type (preterm vs. not-preterm) as described above for the first delivery. We  
696 trained models using billing codes (ICD-9 & CPT) at timepoints before the date of delivery because the

697 majority of this cohort did not have reliable EGA at the second delivery. As described earlier, separate  
698 models were trained using billing codes timestamped before timepoint being evaluated (Fig. 6B, Fig. S9).

### 699 *Preterm birth prediction in independent UCSF cohort*

700 We evaluated how well models trained at Vanderbilt using billing codes would replicate in an external  
701 cohort assembled at UCSF. Only the first delivery in the EHR was used for prediction. Women with twins  
702 or multiple gestations, identified using billing codes (Supplementary Materials and Methods), were  
703 excluded. Delivery type (preterm vs. not preterm) was assigned based on the presence of ICD-10 codes.  
704 Term (or not-preterm) deliveries were determined by the presence of an ICD-10 code beginning with the  
705 characters “O80”, specifying an encounter for full-term delivery. Preterm deliveries were determined by  
706 both the absence of ICD-10 codes beginning with “O80” and the presence of codes beginning with  
707 “O60.1”, the family of codes for preterm labor with preterm delivery. We trained models using ICD-9  
708 codes present before 28 weeks of gestation on the Vanderbilt cohort to predict preterm birth. CPT codes  
709 were not used since they were not available from the UCSF EHR system. The 28-week model was  
710 evaluated on the Vanderbilt held-out set and the independent UCSF cohort.

### 711 *Feature interpretation from boosted decision tree models*

712 To determine feature importance, we used SHapley Additive exPlanation values (SHAP)[52,54,55] to  
713 determine the marginal additive contribution of each feature. For the held-out Vanderbilt cohort and the  
714 UCSF cohort, a SHAP value was calculated for each feature per individual. Feature importance was  
715 summarized by taking the mean of the absolute value of SHAP scores across individuals. The top fifteen  
716 features based on the mean absolute SHAP value in either the Vanderbilt or UCSF cohorts values are  
717 reported. To compare how feature importance varies at Vanderbilt and UCSF, we computed the Pearson  
718 correlation of the mean absolute SHAP values.

719

720 **Ethics:** This study exclusively utilized information extracted from medical records in the Vanderbilt  
721 University Medical Center (VUMC) “Synthetic Derivative” database (SD). The SD is a de-identified  
722 copy of the main hospital medical record databases created for research purposes. The de-identification of  
723 SD records was achieved primarily through the application of a commercial electronic program, which  
724 was applied and assessed for acceptable effectiveness in scrubbing identifiers. For instance, if the name  
725 “John Smith” appeared in the original medical record, its corresponding record in the SD does not contain  
726 “John Smith”. Instead, it is permanently replaced with a tag [NAMEAAA, BBB] to maintain the semantic  
727 integrity of the text. Similarly, dates, such as “January 1, 2004” have been replaced with a randomly  
728 generated date, such as “February 3, 2003.”

729 The SD database (which contains over 3 million electronic records, with no defined exclusions) was  
730 accessed through database queries. Searches are logged and audited annually. As no HIPAA identifiers  
731 are available in the SD database, and this work does not plan to re-identify these records using the  
732 identified VUMC database, this study meets criteria for non-human subjects research. Nonetheless, to  
733 ensure confidentiality and appropriate use of the SD, all relevant key personnel for this study entered a  
734 data use agreement, which prohibits any use of the data not described in this application, including the re-  
735 identification of the SD records.

736

737 **Acknowledgments:** We thank the members of the Capra lab and the March of Dimes Prematurity  
738 Research Center Ohio Collaborative for thoughtful discussion on this project.

739 **Funding:** AA was supported by the American Heart Association fellowship 20PRE35080073, National  
740 Institutes of Health (NIH, T32GM007347), the March of Dimes, and the Burroughs Wellcome Fund. MS,  
741 IK, and BL were supported by the March of Dimes and NIH (NLM K01LM012381). JAC was supported  
742 by the NIH (R35GM127087), NIH (1R01HD101669), March of Dimes, and the Burroughs Wellcome  
743 Fund. This work was conducted in part using the resources of the Advanced Computing Center for

744 Research and Education at Vanderbilt University. The dataset(s) used for the analyses described were  
745 obtained from Vanderbilt University Medical Center's BioVU which is supported by numerous sources:  
746 institutional funding, private agencies, and federal grants. These include the NIH funded Shared  
747 Instrumentation Grant S10RR025141; and CTSA grants UL1TR002243, UL1TR000445, and  
748 UL1RR024975. Genomic data are also supported by investigator-led projects that include  
749 U01HG004798, R01NS032830, RC2GM092618, P50GM115305, U01HG006378, U19HL065962,  
750 R01HD074711; and additional funding sources listed at <https://victr.vumc.org/biovu-funding/>. The  
751 content is solely the responsibility of the authors and does not necessarily represent the official views of  
752 the National Institutes of Health, the March of Dimes, or the Burroughs Wellcome Fund.

753  
754 **Author contributions:** *Conceptualization, Methodology:* A.A. and J.A.C. conceived and designed the  
755 study. J.M.N provided clinical interpretation and aided in feature selection. *Data curation:* A.A. and  
756 C.A.B. extracted billing codes, clinical notes, and performed concept extraction on the Vanderbilt cohort.  
757 A.A., P.S. and L.K.D. extracted, cleaned, and provided clinical laboratory data during pregnancy on the  
758 Vanderbilt cohort. *Resources:* D.R.V. provided obstetric and nursing notes on the Vanderbilt cohort. B.L.,  
759 I.K., M.S. extracted the delivery cohort from UCSF. *Formal Analysis, Investigation:* A.A. performed all  
760 analyses on the Vanderbilt cohort under supervision from J.A.C. B.L. and I.K. evaluated models on  
761 UCSF cohorts under supervision from M.S. *Funding acquisition:* J.A.C. *Writing:* A.A. wrote the  
762 manuscript with guidance from J.A.C., J.M.N., M.S., L.M. and A.R.

763

764 **Competing interests:** LJM is a consultant for Mirvie, Inc.

765 **Data and materials availability:** The dataset(s) and code supporting the conclusions of this article  
766 is(are) available in the [https://github.com/abraham-abin13/ptb\\_predict\\_ml](https://github.com/abraham-abin13/ptb_predict_ml) repository.

767

768

769 **References**

770

- 771 1. Goldenberg RL, Culhane JF, Iams JD, Romero R. Epidemiology and causes of preterm birth. *Lancet Lond Engl.*  
772 2008;371: 75–84. doi:10.1016/s0140-6736(08)60074-4
- 773 2. Blencowe H, Cousens S, Oestergaard MZ, Chou D, Moller A-B, Narwal R, et al. National, regional, and  
774 worldwide estimates of preterm birth rates in the year 2010 with time trends since 1990 for selected countries: a  
775 systematic analysis and implications. *Lancet Lond Engl.* 2012;379: 2162–72. doi:10.1016/s0140-6736(12)60820-4
- 776 3. Barros FC, Papageorgiou AT, Victora CG, Noble JA, Pang R, Iams J, et al. The Distribution of Clinical  
777 Phenotypes of Preterm Birth Syndrome. *JAMA Pediatrics.* 2015;169: 220–10.  
778 doi:10.1001/jamapediatrics.2014.3040
- 779 4. Callaghan WM, MacDorman MF, Rasmussen SA, Qin C, Lackritz EM. The Contribution of Preterm Birth to  
780 Infant Mortality Rates in the United States. *Pediatrics.* 2006;118: 1566–1573. doi:10.1542/peds.2006-0860
- 781 5. Liu L, Oza S, Hogan D, Chu Y, Perin J, Zhu J, et al. Global, regional, and national causes of under-5 mortality in  
782 2000–15: an updated systematic analysis with implications for the Sustainable Development Goals. *Lancet.*  
783 2016;388: 3027–3035. doi:10.1016/s0140-6736(16)31593-8
- 784 6. Romero R, Dey SK, Fisher SJ. Preterm labor: One syndrome, many causes. *Science.* 2014;345: 760–765.  
785 doi:10.1126/science.1251816
- 786 7. Iams J, Goldenberg R, Meis P, Mercer B, Moawad A, Das A, et al. The Length of the Cervix and the Risk of  
787 Spontaneous Premature Delivery. *New Engl J Medicine.* 1996;334: 567–573. doi:10.1056/nejm199602293340904
- 788 8. Fuchs F, Monet B, Ducruet T, Chaillet N, Audibert F. Effect of maternal age on the risk of preterm birth: A large  
789 cohort study. *Plos One.* 2018;13: e0191002. doi:10.1371/journal.pone.0191002
- 790 9. Mercer BM, Goldenberg RL, Moawad AH, Meis PJ, Iams JD, Das AF, et al. The Preterm Prediction Study: Effect  
791 of gestational age and cause of preterm birth on subsequent obstetric outcome. *Am J Obstet Gynecol.* 1999;181:  
792 1216–1221. doi:10.1016/s0002-9378(99)70111-0
- 793 10. Mazaki-Tovi S, Romero R, Kusanovic JP, Erez O, Pineles BL, Gotsch F, et al. Recurrent Preterm Birth. *Semin*  
794 *Perinatol.* 2007;31: 142–158. doi:10.1053/j.semperi.2007.04.001
- 795 11. Ananth CV, Kirby RS, Vintzileos AM. Recurrence of preterm birth in twin pregnancies in the presence of a prior  
796 singleton preterm birth. *J Maternal-fetal Neonatal Medicine.* 2008;21: 289–295. doi:10.1080/14767050802010206
- 797 12. Muglia LJ, Katz M. The Enigma of Spontaneous Preterm Birth. *New England Journal of Medicine.* 2010;362:  
798 529–535. doi:10.1056/nejmra0904308
- 799 13. Auger N, Le TUN, Park AL, Luo Z-C. Association between maternal comorbidity and preterm birth by severity  
800 and clinical subtype: retrospective cohort study. *BMC Pregnancy and Childbirth.* 2011;11: 75. doi:10.1186/1471-  
801 2393-11-67
- 802 14. Carter M, Fowler S, Holden A, Xenakis E, Dudley D. The Late Preterm Birth Rate and Its Association with  
803 Comorbidities in a Population-Based Study. *American Journal of Perinatology.* 2011;28: 703–708. doi:10.1055/s-  
804 0031-1280592

- 805 15. Francesca L, Laura M, Giuseppe R, Francesco DA, Ersilia B, Leonardo P, et al. Biomarkers for predicting  
806 spontaneous preterm birth: an umbrella systematic review. *The Journal of Maternal-Fetal & Neonatal Medicine*.  
807 2019;0: 726–734. doi:10.1080/14767058.2017.1297404
- 808 16. Dabi Y, Nedellec S, Bonneau C, Trouchard B, Rouzier R, Benachi A. Clinical validation of a model predicting  
809 the risk of preterm delivery. Terry J, editor. *PloS one*. 2017;12: e0171801. doi:10.1371/journal.pone.0171801
- 810 17. Ngo TTM, Moufarrej MN, Rasmussen M-LH, Camunas-Soler J, Pan W, Okamoto J, et al. Noninvasive blood  
811 tests for fetal development predict gestational age and preterm delivery. *Science*. 2018;360: 1133–1136.  
812 doi:10.1126/science.aar3819
- 813 18. Tarca AL, Pataki BÁ, Romero R, Sirota M, Guan Y, Kutum R, et al. Crowdsourcing assessment of maternal  
814 blood multi-omics for predicting gestational age and preterm birth. *Cell Reports Medicine*. 2021;2: 100323.  
815 doi:10.1016/j.xcrm.2021.100323
- 816 19. Stelzer IA, Ghaemi MS, Han X, Ando K, Hédou JJ, Feyaerts D, et al. Integrated trajectories of the maternal  
817 metabolome, proteome, and immunome predict labor onset. *Sci Transl Med*. 2021;13: eabd9898.  
818 doi:10.1126/scitranslmed.abd9898
- 819 20. Schaaf JM, Ravelli ACJ, Mol BWJ, Abu-Hanna A. Development of a prognostic model for predicting  
820 spontaneous singleton preterm birth. *European journal of obstetrics, gynecology, and reproductive biology*.  
821 2012;164: 150–155. doi:10.1016/j.ejogrb.2012.07.007
- 822 21. Morken NH, Källen K, Jacobsson B. Predicting Risk of Spontaneous Preterm Delivery in Women with a  
823 Singleton Pregnancy. *Paediatric and Perinatal Epidemiology*. 2014;28: 11–22. doi:10.1111/ppe.12087
- 824 22. Weber A, Darmstadt GL, Gruber S, Foeller ME, Carmichael SL, Stevenson DK, et al. Application of machine-  
825 learning to predict early spontaneous preterm birth among nulliparous non-Hispanic black and white women. *Ann*  
826 *Epidemiol*. 2018;28: 783-789.e1. doi:10.1016/j.annepidem.2018.08.008
- 827 23. Baer RJ, McLemore MR, Adler N, Oltman SP, Chambers BD, Kuppermann M, et al. Pre-pregnancy or first-  
828 trimester risk scoring to identify women at high risk of preterm birth. *European Journal of Obstetrics and*  
829 *Gynecology*. 2018;231: 235–240. doi:10.1016/j.ejogrb.2018.11.004
- 830 24. Suff N, Story L, Shennan A. The prediction of preterm delivery: What is new? *Semin Fetal Neonat M*. 2018;24:  
831 27–32. doi:10.1016/j.siny.2018.09.006
- 832 25. Abul-Husn NS, Kenny EE. Personalized Medicine and the Power of Electronic Health Records. *Cell*. 2019;177:  
833 58–69. doi:10.1016/j.cell.2019.02.039
- 834 26. Paquette AG, Hood L, Price ND, Sadovsky Y. Deep phenotyping during pregnancy for predictive and preventive  
835 medicine. *Science Translational Medicine*. 2020;12: eaay1059. doi:10.1126/scitranslmed.aay1059
- 836 27. Artzi NS, Shilo S, Hadar E, Rossman H, Barbash-Hazan S, Ben-Haroush A, et al. Prediction of gestational  
837 diabetes based on nationwide electronic health records. *Nat Med*. 2020;26: 71–76. doi:10.1038/s41591-019-0724-8
- 838 28. Ravizza S, Huschto T, Adamov A, Böhm L, Büsser A, Flöther FF, et al. Predicting the early risk of chronic  
839 kidney disease in patients with diabetes using real-world data. *Nat Med*. 2019;25: 57–59. doi:10.1038/s41591-018-  
840 0239-8
- 841 29. Li R, Chen Y, Ritchie MD, Moore JH. Electronic health records and polygenic risk scores for predicting disease  
842 risk. *Nature Publishing Group*. 2020;31: 1–10. doi:10.1038/s41576-020-0224-1

- 843 30. Zhang G, Feenstra B, Bacelis J, Liu X, Muglia LM, Juodakis J, et al. Genetic Associations with Gestational  
844 Duration and Spontaneous Preterm Birth. *New Engl J Medicine*. 2017;377: 1156–1167.  
845 doi:10.1056/nejmoa1612665
- 846 31. Tomašev N, Glorot X, Rae JW, Zielinski M, Askham H, Saraiva A, et al. A clinically applicable approach to  
847 continuous prediction of future acute kidney injury. *Nature*. 2019;572: 116–119. doi:10.1038/s41586-019-1390-1
- 848 32. Xiao C, Choi E, Sun J. Opportunities and challenges in developing deep learning models using electronic health  
849 records data: a systematic review. *J Am Med Inform Assn*. 2018;25: 1419–1428. doi:10.1093/jamia/ocy068
- 850 33. Zhao J, Feng Q, Wu P, Lupu RA, Wilke RA, Wells QS, et al. Learning from Longitudinal Data in Electronic  
851 Health Record and Genetic Data to Improve Cardiovascular Event Prediction. *Scientific Reports*. 2019;9: 1–10.  
852 doi:10.1038/s41598-018-36745-x
- 853 34. Goldstein BA, Navar AM, Pencina MJ, Ioannidis JPA. Opportunities and challenges in developing risk  
854 prediction models with electronic health records data: a systematic review. *Journal of the American Medical*  
855 *Informatics Association*. 2017;24: 198–208. doi:10.1093/jamia/ocw042
- 856 35. Aung MT, Yu Y, Ferguson KK, Cantonwine DE, Zeng L, McElrath TF, et al. Prediction and associations of  
857 preterm birth and its subtypes with eicosanoid enzymatic pathways and inflammatory markers. *Sci Rep-uk*. 2019;9:  
858 17049. doi:10.1038/s41598-019-53448-z
- 859 36. Rittenhouse KJ, Vwalika B, Keil A, Winston J, Stoner M, Price JT, et al. Improving preterm newborn  
860 identification in low-resource settings with machine learning. *Plos One*. 2019;14: e0198919.  
861 doi:10.1371/journal.pone.0198919
- 862 37. Fergus P, Cheung P, Hussain A, Al-Jumeily D, Dobbins C, Iram S. Prediction of Preterm Deliveries from EHG  
863 Signals Using Machine Learning. *Plos One*. 2013;8: e77154. doi:10.1371/journal.pone.0077154
- 864 38. Krishnapuram B, Shah M, Smola A, Aggarwal C, Shen D, Rastogi R, et al. XGBoost: A Scalable Tree Boosting  
865 System. *Arxiv*. 2016; 785–794. doi:10.1145/2939672.2939785
- 866 39. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning, Data Mining, Inference, and*  
867 *Prediction*. 2009. doi:10.1007/978-0-387-84858-7
- 868 40. Corey KM, Kashyap S, Lorenzi E, Lagoo-Deenadayalan SA, Heller K, Whalen K, et al. Development and  
869 validation of machine learning models to identify high-risk surgical patients using automatically curated electronic  
870 health record data (Pythia): A retrospective, single-site study. *Plos Med*. 2018;15: e1002701.  
871 doi:10.1371/journal.pmed.1002701
- 872 41. Jing L, Cerna AEU, Good CW, Sauers NM, Schneider G, Hartzel DN, et al. A Machine Learning Approach to  
873 Management of Heart Failure Populations. *Jacc Hear Fail*. 2020;8: 578–587. doi:10.1016/j.jchf.2020.01.012
- 874 42. Carter J, Seed PT, Watson HA, David AL, Sandall J, Shennan AH, et al. Development and validation of  
875 predictive models for QUiPP App v.2: tool for predicting preterm birth in women with symptoms of threatened  
876 preterm labor. *Ultrasound Obst Gyn*. 2020;55: 357–367. doi:10.1002/uog.20422
- 877 43. Vogel JP, Chawanpaiboon S, Moller A-B, Watananirun K, Bonet M, Lumbiganon P. The global epidemiology  
878 of preterm birth. *Best Pract Res Cl Ob*. 2018;52: 3–12. doi:10.1016/j.bpobgyn.2018.04.003
- 879 44. Smith GCS, Pell JP. Teenage Pregnancy and Risk of Adverse Perinatal Outcomes Associated With First and  
880 Second Births: Population Based Retrospective Cohort Study. *Obstet Gynecol Surv*. 2002;57: 136–137.  
881 doi:10.1097/00006254-200203000-00002

- 882 45. Waldenström U, Aasheim V, Nilsen ABV, Rasmussen S, Pettersson HJ, Schytt E, et al. Adverse Pregnancy  
883 Outcomes Related to Advanced Maternal Age Compared With Smoking and Being Overweight. *Obstetrics Gynecol.*  
884 2014;123: 104–112. doi:10.1097/aog.0000000000000062
- 885 46. Carolan M. Maternal age  $\geq 45$  years and maternal and perinatal outcomes: A review of the evidence. *Midwifery.*  
886 2013;29: 479–489. doi:10.1016/j.midw.2012.04.001
- 887 47. Ray JG, Vermeulen MJ, Shapiro JL, Kenshole AB. Maternal and neonatal outcomes in pregestational and  
888 gestational diabetes mellitus, and the influence of maternal obesity and weight gain: the DEPOSIT study. *Qjm Int J*  
889 *Medicine.* 2001;94: 347–356. doi:10.1093/qjmed/94.7.347
- 890 48. Whiteman V, Salinas A, Weldeselashe HE, August EM, Mbah AK, Aliyu MH, et al. Impact of sickle cell disease  
891 and thalassemias in infants on birth outcomes. *Eur J Obstet Gyn R B.* 2013;170: 324–328.  
892 doi:10.1016/j.ejogrb.2013.06.020
- 893 49. Umesawa M, Kobashi G. Epidemiology of hypertensive disorders in pregnancy: prevalence, risk factors,  
894 predictors and prognosis. *Hypertens Res.* 2017;40: 213–220. doi:10.1038/hr.2016.126
- 895 50. Koullali B, Oudijk MA, Nijman TAJ, Mol BWJ, Pajkrt E. Risk assessment and management to prevent preterm  
896 birth. *Seminars Fetal Neonatal Medicine.* 2016;21: 80–88. doi:10.1016/j.siny.2016.01.005
- 897 51. Lundberg SM, Lee S-I. A Unified Approach to Interpreting Model Predictions. In: Guyon I, Luxburg U V,  
898 Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. Curran Associates, Inc.; 2017. pp. 4765–4774.
- 899 52. Lundberg SM, Erion G, Chen H, DeGrave A, Prutkin JM, Nair B, et al. From local explanations to global  
900 understanding with explainable AI for trees. *Nat Mach Intell.* 2020;2: 56–67. doi:10.1038/s42256-019-0138-9
- 901 53. Davis J, Goadrich M. The relationship between Precision-Recall and ROC curves. 2006; 233–240.  
902 doi:10.1145/1143844.1143874
- 903 54. Lundberg SM, Lee S-I. A Unified Approach to Interpreting Model Predictions. In: I. Guyon, U. V. Luxburg,  
904 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, et al., editors. *Advances in Neural Information Processing*  
905 *Systems 30.* Curran Associates, Inc.; 2017. pp. 4765–4774. Available: <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
- 907 55. Lundberg SM, Nair B, Vavilala MS, Horibe M, Eisses MJ, Adams T, et al. Explainable machine-learning  
908 predictions for the prevention of hypoxaemia during surgery. *Nat Biomed Eng.* 2018;2: 749–760.  
909 doi:10.1038/s41551-018-0304-0
- 910 56. Creanga AA, Berg CJ, Syverson C, Seed K, Bruce FC, Callaghan WM. Pregnancy-Related Mortality in the  
911 United States, 2006–2010. *Obstetrics Gynecol.* 2015;125: 5–12. doi:10.1097/aog.0000000000000564
- 912 57. Hirshberg A, Srinivas SK. Epidemiology of maternal morbidity and mortality. *Semin Perinatol.* 2017;41: 332–  
913 337. doi:10.1053/j.semperi.2017.07.007
- 914 58. Kopitar L, Kocbek P, Cilar L, Sheikh A, Stiglic G. Early detection of type 2 diabetes mellitus using machine  
915 learning-based prediction models. *Sci Rep-uk.* 2020;10: 11981. doi:10.1038/s41598-020-68771-z
- 916 59. Yan L, Zhang H-T, Goncalves J, Xiao Y, Wang M, Guo Y, et al. An interpretable mortality prediction model for  
917 COVID-19 patients. *Nat Mach Intell.* 2020;2: 283–288. doi:10.1038/s42256-020-0180-7
- 918 60. Couronné R, Probst P, Boulesteix A-L. Random forest versus logistic regression: a large-scale benchmark  
919 experiment. *Bmc Bioinformatics.* 2018;19: 270. doi:10.1186/s12859-018-2264-5

- 920 61. Gao C, Osmundson S, Edwards DRV, Jackson GP, Malin BA, Chen Y. Deep learning predicts extreme preterm  
921 birth from electronic health records. *J Biomed Inform.* 2019;100: 103334. doi:10.1016/j.jbi.2019.103334
- 922 62. Torchin H, Ancel P-Y. [Epidemiology and risk factors of preterm birth]. *J De Gynecol Obstetrique Et Biologie*  
923 *De La Reproduction.* 2016;45: 1213–1230. doi:10.1016/j.jgyn.2016.09.013
- 924 63. Topol EJ. High-performance medicine: the convergence of human and artificial intelligence. *Nat Med.* 2019;25:  
925 44–56. doi:10.1038/s41591-018-0300-7
- 926 64. He J, Baxter SL, Xu J, Xu J, Zhou X, Zhang K. The practical implementation of artificial intelligence  
927 technologies in medicine. *Nat Med.* 2019;25: 30–36. doi:10.1038/s41591-018-0307-0
- 928 65. Esplin MS. The Importance of Clinical Phenotype in Understanding and Preventing Spontaneous Preterm Birth.  
929 *American Journal of Perinatology.* 2016;33: 236–244. doi:10.1055/s-0035-1571146
- 930 66. Manuck TA, Esplin MS, Biggio J, Bukowski R, Parry S, Zhang H, et al. The phenotype of spontaneous preterm  
931 birth: application of a clinical phenotyping tool. *Am J Obstet Gynecol.* 2015;212: 487.e1-487.e11.  
932 doi:10.1016/j.ajog.2015.02.010
- 933 67. Phelan M, Bhavsar NA, Goldstein BA. Illustrating Informed Presence Bias in Electronic Health Records Data:  
934 How Patient Interactions with a Health System Can Impact Inference. *Egems Wash Dc.* 2017;5: 22.  
935 doi:10.5334/egems.243
- 936 68. Moutquin J-M. Classification and heterogeneity of preterm birth. *Bjog Int J Obstetrics Gynaecol.* 2003;110: 30–  
937 33. doi:10.1046/j.1471-0528.2003.00021.x
- 938 69. Phillips C, Velji Z, Hanly C, Metcalfe A. Risk of recurrent spontaneous preterm birth: a systematic review and  
939 meta-analysis. *Bmj Open.* 2017;7: e015402. doi:10.1136/bmjopen-2016-015402
- 940 70. Shah NH, Milstein A, PhD SCB. Making Machine Learning Models Clinically Useful. *Jama.* 2019;322: 1351–  
941 1352. doi:10.1001/jama.2019.10306
- 942 71. Gianfrancesco MA, Tamang S, Yazdany J, Schmajuk G. Potential Biases in Machine Learning Algorithms  
943 Using Electronic Health Record Data. *Jama Intern Med.* 2018;178: 1544. doi:10.1001/jamainternmed.2018.3763
- 944 72. Weng C, Shah N, Hripcsak G. Deep Phenotyping: Embracing Complexity and Temporality—Towards  
945 Scalability, Portability, and Interoperability. *J Biomed Inform.* 2020;105: 103433. doi:10.1016/j.jbi.2020.103433
- 946 73. Bergstra J, Yamins D, Cox D. Making a science of model search: Hyperparameter optimization in hundreds of  
947 dimensions for vision architectures. *International conference on machine learning.* 2013. pp. 115--123.
- 948 74. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in  
949 Python. *Journal of Machine Learning Research.* 2011;12: 2825--2830.
- 950 75. Soysal E, Wang J, Jiang M, Wu Y, Pakhomov S, Liu H, et al. CLAMP – a toolkit for efficiently building  
951 customized clinical natural language processing pipelines. *J Am Med Inform Assn.* 2017;25: 331–336.  
952 doi:10.1093/jamia/ocx132
- 953 76. Marees AT, Kluiver H de, Stringer S, Vorspan F, Curis E, Claire CM, et al. A tutorial on conducting  
954 genome-wide association studies: Quality control and statistical analysis. *International Journal of Methods in*  
955 *Psychiatric Research.* 2018;27: e1608. doi:10.1002/mpr.1608

- 956 77. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising to the  
957 challenge of larger and richer datasets. *GigaScience*. 2015;4: 7. doi:10.1186/s13742-015-0047-8
- 958 78. Euesden J, Lewis CM, O'Reilly PF. PRSice: Polygenic Risk Score software. *Bioinformatics*. 2015;31: 1466–  
959 1468. doi:10.1093/bioinformatics/btu848
- 960 79. Choi SW, O'Reilly PF. PRSice-2: Polygenic Risk Score software for biobank-scale data. *Gigascience*. 2019;8.  
961 doi:10.1093/gigascience/giz082
- 962 80. McInnes L, Healy J, Melville J. UMAP: Uniform Manifold Approximation and Projection for Dimension  
963 Reduction. *Arxiv*. 2018.
- 964 81. Campello RJGB, Moulavi D, Sander J. Advances in Knowledge Discovery and Data Mining, 17th Pacific-Asia  
965 Conference, PAKDD 2013, Gold Coast, Australia, April 14-17, 2013, Proceedings, Part II. 2013; 160–172.  
966 doi:10.1007/978-3-642-37456-2\_14
- 967 82. Virtanen P, Gommers R, Oliphant TE, Haberland M, Reddy T, Cournapeau D, et al. SciPy 1.0: fundamental  
968 algorithms for scientific computing in Python. *Nat Methods*. 2020;17: 261–272. doi:10.1038/s41592-019-0686-2
- 969
- 970
- 971