

Stabl: sparse and reliable biomarker discovery in predictive modeling of high-dimensional omic data

Julien Hedou (✉ jhedou@stanford.edu)

Stanford University School of Medicine

Ivana Maric (✉ ivanam@stanford.edu)

Stanford University School of Medicine <https://orcid.org/0000-0002-9441-521X>

Grégoire Bellan (✉ gbellan@surge.care)

Telecom Paris <https://orcid.org/0000-0002-8804-8062>

Jakob Einhaus (✉ jeinhaus@stanford.edu)

Stanford University

Dyani Gaudilliere (✉ dyani.gaudilliere@stanford.edu)

Department of Anesthesiology, Perioperative and Pain Medicine, Stanford University School of Medicine

<https://orcid.org/0000-0003-3017-5939>

Francois-Xavier Ladant (✉ fxladant@gmail.com)

Harvard University

Franck Verdonk (✉ fverdonk@stanford.edu)

Stanford University

Ina Stelzer (✉ istelzer@stanford.edu)

Department of Anesthesiology, Perioperative and Pain Medicine, Stanford University School of Medicine

<https://orcid.org/0000-0002-9974-4661>

Dorien Feyaerts (✉ dfeyaer@stanford.edu)

Stanford University

Amy Tsai (✉ astsai@stanford.edu)

Department of Anesthesiology, Perioperative and Pain Medicine, Stanford University School of Medicine

<https://orcid.org/0000-0002-2380-6424>

Edward Ganio (✉ eganio@stanford.edu)

Stanford University

Maximilian Sabayev (✉ sabayev@stanford.edu)

Stanford University

Joshua Gillard (✉ jgillard@stanford.edu)

Stanford University

Adam Bonham (✉ bonh5363@stanford.edu)

Stanford University

Masaki Sato (✉ satomas@stanford.edu)

Stanford University

Maïgane Diop (✉ mdiop@stanford.edu)

Stanford University

Martin Angst (✉ ang@stanford.edu)

Stanford University School of Medicine <https://orcid.org/0000-0002-1550-8136>

David Stevenson (✉ dks750@stanford.edu)

Stanford University

Nima Aghaeepour (✉ naghaeep@stanford.edu)

Stanford University <https://orcid.org/0000-0002-6117-8764>

Andrea Montanari (✉ montanar@stanford.edu)

Stanford University

Brice Gaudilliere (✉ gbrice@stanford.edu)

Stanford University <https://orcid.org/0000-0002-3475-5706>

Article

Keywords:

DOI: <https://doi.org/>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations: Yes there is potential Competing Interest. Julien Hedou declares advisory board membership in SurgeCare. Brice Gaudilliere declares advisory board membership in SurgeCare and maternica therapeutics.

Stabl: sparse and reliable biomarker discovery in predictive modeling of high-dimensional omic data

Julien Hédou^{1*}, Ivana Marić^{2*}, Grégoire Bellan^{3*}, Jakob Einhaus^{1,4*}, Dyani K. Gaudillière⁵, Francois-Xavier Ladant⁶, Franck Verdonk^{1,7}, Ina A. Stelzer¹, Dorien Feyaerts¹, Amy S. Tsai¹, Edward A. Ganio¹, Maximilian Sabayev¹, Joshua Gillard¹, Thomas A. Bonham¹, Masaki Sato¹, Maïgane Diop¹, Martin S. Angst¹, David Stevenson², Nima Aghaeepour^{1,2,8}, Andrea Montanari^{9,10}, Brice Gaudillière^{1,2,ε}

¹Department of Anesthesiology, Perioperative & Pain Medicine, Stanford University, Stanford, CA

²Department of Pediatrics, Stanford University, Stanford, CA

³Télécom Paris, Paris, France

⁴Department of Pathology and Neuropathology, University Hospital and Comprehensive Cancer Center Tübingen, Tübingen, Germany

⁵Division of Plastic and Reconstructive Surgery, Department of Surgery, Stanford University, Stanford, CA

⁶Department of Economics, Harvard University, Cambridge, MA

⁷Sorbonne University, GRC 29, AP-HP, DMU DREAM, Department of Anesthesiology and Intensive Care, Hôpital Saint-Antoine, Assistance Publique-Hôpitaux de Paris; Paris, France

⁸Department of Biomedical Data Science, Stanford University, Stanford, CA

⁹Department of Statistics, Stanford University, Stanford, CA

¹⁰Department of Electrical Engineering, Stanford University, Stanford, CA

*Julien Hédou, Ivana Marić, Grégoire Bellan, and Jakob Einhaus contributed equally to this manuscript.

εCorrespondence:

Brice Gaudillière, Department of Anesthesiology, Perioperative and Pain Medicine, 300 Pasteur Drive, Room S238, Stanford, CA 94305-5117, USA; email: gbrice@stanford.edu

Abstract

High-content omic technologies coupled with sparsity-promoting regularization methods (SRM) have transformed the biomarker discovery process. However, the translation of computational results into a clinical use-case scenario remains challenging. A rate-limiting step is the rigorous selection of reliable biomarker candidates among a host of biological features included in multivariate models. We propose Stabl, a machine learning framework that unifies the biomarker discovery process with multivariate predictive modeling of clinical outcomes by selecting a sparse and reliable set of biomarkers. Evaluation of Stabl on synthetic datasets and four independent clinical studies demonstrates improved biomarker sparsity and reliability compared to commonly used SRMs at similar predictive performance. Stabl readily extends to double- and triple-omics integration tasks and identifies a sparser and more reliable set of biomarkers than those selected by state-of-the-art early- and late-fusion SRMs, thereby facilitating the biological interpretation and clinical translation of complex multi-omic predictive models.

The complete package for Stabl is available online at <https://github.com/gregbellan/Stabl>.

44 Introduction

45

46 High-content omic technologies, such as transcriptomics, metabolomics, or cytometric immunoassays,
47 are increasingly employed in biomarker discovery studies.^{1,2} The ability to measure thousands of
48 molecular features in each biological specimen provides unprecedented opportunities for development
49 of precision medicine tools across the spectrum of health and disease. Omic technologies have also
50 dictated a shift in statistical analysis of biological data. The traditional univariate statistical framework is
51 maladapted to large omic datasets characterized by a high number of molecular features p relative to the
52 number of samples n . The $p \gg n$ scenario drastically reduces the statistical power of univariate analyses,
53 a problem that cannot be easily overcome by increasing the value of n due to cost or sample availability
54 constraints.^{3,4}

55

56 Statistical analysis in biomarker discovery research comprises three related yet distinct tasks, all of which
57 are necessary for translation into clinical use and impacted by the $p \gg n$ problem: 1) prediction of the
58 clinical endpoint via identification of a multivariate model with high predictive performance (*predictivity*),
59 2) selection of a limited number of features as candidate clinical biomarkers (*sparsity*), and 3) confidence
60 that the selected features are among the true set of features (i.e., truly related to the outcome, *reliability*).

61

62 Several machine-learning methods, including sparsity-promoting regularization methods (SRMs), such
63 as least absolute shrinkage and selection operator (Lasso)⁵ or elastic net (EN),⁶ provide predictive
64 modeling frameworks adapted to $p \gg n$ omic datasets, but the selection of a sparse and reliable set of
65 candidate biomarkers remains an important challenge. Most rely on an L1-regularization to limit the
66 number of features used in the final model. However, the learning phase of the model is often performed
67 on a limited number of samples, such that small perturbations in the training data can yield wide
68 differences in the features selected in the predictive model.⁷⁻⁹ This undermines confidence in the features
69 selected, as current SRMs do not provide objective metrics to determine whether these features are truly
70 related to the outcome. This inherent limitation of SRMs can result in poor sparsity and reliability, thereby
71 hindering the biological interpretation and clinical significance of the predictive model. As such, few omic
72 biomarker discovery studies progress to further clinical development phases.^{1-4,10,11}

73

74 High-dimensional feature selection methods such as stability selection (SS) or Model-X knockoff improve
75 reliability by controlling for false discoveries in the selected set of features.^{12,13} However, in these
76 methods, the threshold for feature selection or the target false discovery rate (FDR) are defined a priori,
77 which uncouples the feature selection from the multivariate modeling process. Without prior knowledge
78 on the data, these methods can lead to suboptimal feature selection, requiring multiple iterations to
79 identify a desired threshold. This limitation also precludes optimal integration of multiple omic datasets
80 into a unique predictive model, as a single fixed selection threshold may not be suited to the specificities
81 of each dataset.

82

83 Here we introduce Stabl, a supervised machine learning framework that bridges the gap between
84 multivariate predictive modeling of high-dimensional omic data and the sparsity and reliability
85 requirements of an effective biomarker discovery process. Stabl combines the injection of knockoff-
86 modeled noise or random permutations into the original data, a data-driven signal-to-noise threshold, and
87 integration of selected features into a predictive model. Systematic benchmarking of Stabl against Lasso,
88 EN, and SS using synthetic datasets, three existing real-world omic datasets, and a newly generated
89 multi-omic clinical dataset demonstrates that Stabl overcomes the shortcomings of state-of-the-art SRMs:
90 Stabl yields highly reliable and sparse predictive models while identifying biologically plausible features
91 amenable to further development into diagnostic or prognostic precision medicine assays.

92

93 The complete package for Stabl is available online at <https://github.com/gregbellan/Stabl>.

94

95 Results

96

97 ***Selection of reliable predictive features using estimated false discovery proportion (FDP)***

98

99 When applied to a single cohort drawn at random from the population, SRMs will select informative
100 features (i.e., truly related to the outcome) with a higher probability, on average, than uninformative
101 features (i.e., unrelated to the outcome).^{5,12} However, as uninformative features typically outnumber
102 informative features in high-dimensional omic datasets,^{1,2,11} the fit of an SRM model on a single cohort
103 can lead to selection of many uninformative features despite a low probability of selection.^{12,14} To address
104 this issue, Stabl implements the following strategy (Fig. 1 and methods):
105

- 106 1. Stabl fits SRM models (e.g., Lasso or EN) on subsamples of the data using a procedure similar
107 to SS.¹² Subsampling mimics the availability of multiple random cohorts and estimates each
108 feature's frequency of selection across all iterations. However, this procedure does not provide
109 an optimal frequency threshold to discriminate between informative and uninformative features
110 objectively.
- 111 2. To define the optimal frequency threshold, Stabl creates artificial features unrelated to the
112 outcome (noise injection) via random permutations¹⁻³ or knockoff sampling,^{13,15,16} which we
113 assume behave similarly to uninformative features in the original dataset¹⁷ (see theoretical
114 guarantees in methods). The artificial features are used to construct a surrogate of the false
115 discovery proportion (FDP₊). We define the "reliability threshold", θ , as the frequency threshold
116 yielding the minimum FDP₊ across all possible thresholds. This method for determining θ is
117 objective, in that it minimizes a proxy for the FDP. It is also data-driven, as it is tailored to individual
118 omic datasets.

119
120 As a result, Stabl provides a unifying procedure that selects features above the reliability threshold while
121 building a multivariate predictive model. Stabl is amenable to classification and regression tasks and
122 extends to integration of multiple datasets of different dimensions and from different omic modalities.
123

124 ***Stabl improves sparsity and reliability while maintaining predictivity: synthetic modeling***

125
126 We benchmarked Stabl against Lasso and EN using synthetically generated training and validation
127 datasets containing known informative and uninformative features (Fig. 2a). Simulations representative
128 of real-world scenarios were performed, including variations in the sample size (n), total features (p), and
129 informative features (S). Models were evaluated using three performance metrics (Fig. 2b):
130

- 131 1. *Sparsity*: the average number of features selected compared to the number of informative
132 features.
- 133 2. *Reliability*: overlap between the features selected by the algorithm and the true set of informative
134 features (Jaccard Index).
- 135 3. *Predictivity*: mean square error (MSE).

136 Before performing benchmark comparisons, we tested whether the FDP₊ defined by Stabl experimentally
137 controls the FDR at the reliability threshold θ , as the true value of the FDR is known for the synthetic
138 dataset. We observed that FDP₊(θ) was indeed greater than the true FDR value (Fig. 2c and S1). These
139 observations experimentally confirmed the validity of Stabl in optimizing the frequency threshold for
140 feature selection. Furthermore, under the assumption that the uninformative features and the artificial
141 features are interchangeable, we bound the probability that FDP exceeds a multiple of the proximity to
142 FDP₊(θ), thus providing a theoretical validation of our experimental observations (see theoretical
143 guarantee in methods).
144

145 Stabl was tested using a random permutation method (Fig. 2 and S2-5) or model-X knockoffs (Fig. S5)
146 for noise generation. In each case, Stabl achieved higher sparsity compared to Lasso or EN (Fig. S6),
147 as the number of features selected by Stabl was lower across all conditions tested and converged
148 towards the true number of informative features (Fig. 2d). The reliability was also higher for Stabl than
149 for Lasso or EN, such that the features selected by Stabl were closer to the true set of informative features
150 (Fig. 2e). Meanwhile, Stabl had similar or better predictivity compared to Lasso or EN (Fig. 2f).
151

152 Further modeling experiments tested the impact of the data-driven computation of θ while building the
153 multivariate model compared to SS (i.e., choosing a fixed frequency threshold a priori). Three

154 representative frequency thresholds were evaluated: 30%, 50%, or 80% (Fig. 2g-i and S7-9). The
155 performance of models built using a fixed frequency threshold varied greatly depending on the simulation
156 conditions. For example, for a small sample size ($n < 75$), the 30% threshold had the best sparsity and
157 reliability. However, for a large sample size ($n > 500$), the 80% threshold resulted in greater
158 performances. In contrast, Stabl models systematically reached optimal sparsity, reliability, and
159 predictivity performances. Further, we show that θ varied greatly with the sample size (Fig. 2j and S10),
160 illustrating how Stabl adapts to datasets of different dimensions to identify an optimal frequency threshold
161 solution.

162
163 In sum, synthetic modeling results show that Stabl achieves better sparsity and reliability compared to
164 Lasso or EN while preserving predictivity and that the set of features chosen by Stabl is closer to the true
165 set of informative features. The results also emphasize the advantage of the data-driven adaptation of
166 the frequency threshold to each dataset's unique characteristics rather than using an arbitrarily fixed
167 threshold.

168 ***Stabl enables effective biomarker discovery in clinical omic studies***

169
170
171 We evaluated Stabl's performance on four independent clinical omic datasets. Three were previously
172 published with standard SRM analyses, while the fourth is a newly generated dataset introduced and
173 analyzed for the first time here. Because clinical omic datasets can vary greatly with respect to
174 dimensionality, signal-to-noise ratio, and technology-specific data preprocessing, we tested Stabl on
175 datasets representing a range of bulk and single-cell omics technologies, including RNA sequencing
176 (RNA-Seq), high-content proteomics (SomaLogic and Olink platforms), untargeted metabolomics, and
177 single-cell mass cytometry.

178
179 For each dataset, Stabl was compared to Lasso and EN on single-omic data or to early fusion and late
180 fusion on multi-omic data over 50 random repetitions using a repeated five-fold cross-validation (CV)
181 strategy. As the true set of informative features is not known for real-world datasets, the performance
182 metrics differed from those used for the synthetic datasets:

- 183
- 184 1. *Sparsity*: determined by the average number of features selected throughout the CV procedure.
- 185 2. *Reliability*: assessed using univariate statistics in the absence of a known true set of features.
- 186 3. *Predictivity*: the area under the receiver operator characteristic curve (AUROC) and the area
187 under the precision-recall curve (AUPRC) for classification tasks or the MSE for regression tasks.
- 188

189 *Identification of sparse, reliable, and predictive candidate biomarkers from single-omic clinical datasets*

190
191 Stabl was first applied to two single-omic clinical datasets featuring a robust biological signal with
192 significant diagnostic potential. The first example is a large-scale plasma cell-free RNA dataset ($p =$
193 $37,184$ cfRNA features) isolated from pregnant patients with the aim of classifying normotensive or
194 preeclamptic (PE) pregnancies (Fig. 3a,b).^{18,19} The second example is a high-plex proteomic dataset (p
195 $= 1,463$ proteomic features, Olink) collected from two independent cohorts (a training and a validation
196 cohort) of SARS-CoV-2-positive patients to classify COVID-19 disease severity (Fig. 3c,d).^{20,21} In these
197 two examples, although both Lasso and EN models achieved very good predictive performance ($AUROC$
198 > 0.80 , Fig. 3, S11-12), the lack of sparsity or reliability hindered the identification of a manageable
199 number of candidate biomarkers, necessitating additional feature selection methods that were decoupled
200 from the predictive modeling process.¹⁸⁻²¹

201
202 Consistent with the results obtained using synthetic data, Stabl achieved comparable predictivity to Lasso
203 (Fig. 3e,f) and EN (Fig. S11a,b) when applied to the single-omic datasets. However, Stabl identified
204 sparser models. For the PE dataset, the average number of features selected by Stabl was reduced over
205 20-fold compared to Lasso (Fig. 3g) or EN (Fig. S11c) respectively. For the classification of patients with
206 mild or severe COVID-19, the number of features selected by Stabl was reduced by a factor of 2.7
207 compared to Lasso (Fig. 3h) and 4.5 compared to EN (Fig. S11d).

208

209 Stabl's reliability performance was also improved compared to Lasso and EN. The univariate p-values
210 (Mann-Whitney test) for the features selected by Stabl were lower than for those selected by Lasso (Fig.
211 3i,j) or EN (Fig. S11e,f). Independent evaluation of the COVID-19 validation dataset confirmed these
212 results (Table S1): 100% of features selected by Stabl passed a 5% FDR threshold (Benjamini-Hochberg
213 correction) on the COVID-19 validation dataset ($mean -\log[p\text{-value}] = 9.0$), compared to 91% for Lasso
214 ($mean -\log[p\text{-value}] = 6.7$, Fig. 3k) and 85% for EN ($mean -\log[p\text{-value}] = 6.2$, Fig. S11g).

215
216 Stabl was also compared to SS using 30%, 50%, and 80% fixed frequency thresholds (Table S2).
217 Consistent with the synthetic modeling analyses, the predictivity and sparsity performances of SS varied
218 greatly with the choice of threshold, while Stabl provided a solution that optimized sparsity while
219 maintaining predictive performance. For example, using SS with a 30% compared to a 50% threshold
220 resulted in a 42% decrease in predictivity for the COVID-19 dataset ($AUROC_{30\%} = 0.85$ vs. $AUROC_{50\%} =$
221 0.49), with a model selecting no features. Conversely, for the PE dataset, fixing the frequency threshold
222 at 30% vs. 50% resulted in a 5.3 fold improvement in sparsity with only a 6% decrease in predictivity
223 ($AUROC_{30\%} = 0.83$ vs. $AUROC_{50\%} = 0.78$).

224
225 Identification of fewer and more reliable features using Stabl facilitated the biomarker discovery process,
226 pinpointing the most informative biological features associated with the clinical outcome. For example,
227 three out of thirteen (23%) cfRNA features (CDK10,²² TRIO,²³ and PLEK2²⁴) selected by the final Stabl
228 PE model encoded proteins with fundamental cellular function, providing biologically-plausible biomarker
229 candidates. Other features were non-coding RNAs or pseudogenes, with yet unknown biological function
230 (Table S3). For the COVID-19 dataset, several features identified by Stabl echoed key pathobiological
231 mechanisms of the host inflammatory response to COVID-19. For example, CCL20 is a known element
232 of the COVID-19 cytokine storm,^{25,26} CRTAC1 is a newly identified marker of lung function,²⁷⁻²⁹ PON3 is
233 a known biomarker decreased during acute COVID-19 infection,³⁰ and MZB1 is a protein associated with
234 high neutralization antibody titers after COVID-19 infection (Fig. 3j).²⁰ The Stabl model also selected
235 MDGA1, a previously unknown biomarker candidate of COVID-19 severity (Table S4).

236
237 Together, the results show that Stabl improves the reliability and sparsity of biomarker discovery in two
238 single-omic datasets of widely different dimensionality while maintaining predictivity performance.

239
240 *Stabl successfully extends to multi-omic data integration*

241
242 We extended the assessment of Stabl to complex clinical datasets combining multiple omic technologies.
243 In this case, the algorithm first selects a reliable set of features at the single-omic level, then integrates
244 the features selected for each omic dataset in a final learner algorithm, such as linear or logistic
245 regression.

246
247 We compared Stabl to early and late fusion Lasso, two commonly employed strategies for multi-omic
248 modeling, on the prediction of a continuous outcome variable from a triple-omic dataset. The analysis
249 leveraged a unique longitudinal biological dataset collected in independent training and validation cohorts
250 of pregnant individuals, together with curated clinical information (Fig. 4a).³¹ The study aimed to predict
251 the difference in days between the time of blood sample collection and spontaneous labor onset (i.e.,
252 time to labor). The study addresses an important clinical need for improved prediction of labor onset in
253 term and preterm pregnancies as standard predictive methods are inaccurate.^{32,33}

254
255 The triple-omic dataset contained a proteomic dataset ($p = 1,317$ features, Somalogic), a metabolomic
256 dataset ($p = 3,529$ untargeted mass spectrometry features), and a single-cell mass cytometry dataset (p
257 $= 1,502$ immune cell features, see methods). When compared to early and late fusion Lasso, Stabl
258 estimated the time to labor with comparable predictivity (Fig. 4b training and validation cohorts), while
259 selecting fewer and more reliable features (Fig. 4c). Importantly, Stabl calculated a different reliability
260 threshold for each omic sublayer ($\theta[\text{Proteomics}] = 36\%$, $\theta[\text{Metabolomics}] = 35\%$, $\theta[\text{mass cytometry}] =$
261 52% , Fig. 4g-i). On the validation dataset, available for the proteomic and mass cytometry data only, 26%
262 of features selected by Stabl passed a 5% FDR threshold (Benjamini-Hochberg correction), compared to
263 4% for early fusion Lasso and 5% for late fusion Lasso, showing that Stabl selected more reliable features
264 (Table S5). These results emphasize the advantage of the data-driven threshold, as fixing a common

265 frequency threshold across all omic layers would have been suboptimal, risking over- or under-selecting
266 features in each omic dataset to be integrated into the final predictive model.

267
268 From a biological standpoint, Stabl streamlined the interpretation of our prior multivariate analyses,³¹
269 honing in on sentinel elements of a systemic biological signature predicting the onset of labor that could
270 be leveraged for development of a blood-based diagnostic test. The Stabl model highlighted dynamic
271 changes in 11 metabolomic, 17 proteomic, and two immune cell features with approaching labor (Fig. 4j-
272 l, Table S6), including a regulated decrease in innate immune cell frequencies (e.g., neutrophils) and
273 their responsiveness to inflammatory stimulation (e.g., pSTAT1 signaling response to IFN α in NK
274 cells^{34,35}), along with a synchronized increase in pregnancy-associated hormones (e.g., 17-
275 Hydroxyprogesterone³⁶), placental-derived (e.g., Siglec-6,³⁷ Angiotensin 2/sTie2³⁸), and immune
276 regulatory plasma proteins (e.g., IL-1R4,³⁹ SLPI⁴⁰).

277
278 *Stabl identifies promising candidate biomarkers from a newly generated multi-omic dataset*

279
280 Application of Stabl to the three existing omic datasets demonstrated the algorithm's performance in the
281 context of biomarker discovery studies with a known biological signal. To complete its systematic
282 evaluation, Stabl was applied to our newly generated multi-omic clinical study performing an unbiased
283 biomarker discovery task. The aim of the study was to develop a model to predict which patients will
284 develop a postoperative surgical site infection (SSI) from analysis of pre-operative blood samples (Fig.
285 5a). A cohort of 274 patients undergoing major abdominal surgery were enrolled and preoperative blood
286 samples were collected. Using a matched, nested case-control design, 93 patients were selected from
287 the larger cohort to minimize the effect of clinical or demographic confounders on identified predictive
288 models (Table S7). These samples were analyzed using a combined single-cell mass cytometry (Fig.
289 S13) and plasma proteomics (Somalogic) approach.

290 Stabl merged all omic datasets into a final model that accurately classified patients with and without SSI
291 ($AUROC_{Stabl} = 0.80 [0.69, 0.89]$). When compared to early and late fusion Lasso, Stabl had comparable
292 predictive performance (Fig. 5b, S14), yet superior sparsity (Fig. 5c) and reliability performance (Fig.
293 5h,i). As a result of the frequency-matching procedure, there were no differences in major demographic
294 and clinical variables between the two patient groups, suggesting that model predictions were primarily
295 driven by pre-operative biological differences in patients' susceptibility to develop an SSI.

296
297 Stabl selected four mass cytometry and 25 plasma proteomic features that were combined into a
298 biologically interpretable immune signature predictive of SSI. Examination of Stabl features revealed cell-
299 type specific immune signaling responses associated with SSI (Fig. 5h) that resonated with circulating
300 inflammatory mediators (Fig. 5i, Table S8). Notably, the STAT3 signaling response to IL-6 in neutrophils
301 was increased before surgery in patients predisposed to SSI. Correspondingly, patients with SSI had
302 elevated plasma levels of IL-1 β and IL-18, two potent inducers of IL-6 production in response to
303 inflammatory stress.^{41,42} Other proteomic features selected by the model included CCL3, which
304 coordinates recruitment and activation of neutrophils, and the canonical stress response protein HSPH1.
305 These findings are consistent with previous studies showing that heightened innate immune cell
306 responses to inflammatory stress, such as surgical trauma,^{43,44} can result in diminished defensive
307 response to bacterial pathogens,³⁹ thus increasing a patient's susceptibility to subsequent infection.

308
309 Altogether, application of Stabl in the setting of a new biomarker discovery study provided a manageable
310 number of candidate biomarkers of SSI, pointing at plausible biological mechanisms that can be targeted
311 for further diagnostic or therapeutic development.

312 313 Discussion

314
315 Stabl is a machine learning method for analysis of high-dimensional omic data designed to unify the
316 biomarker discovery process by identifying sparse and reliable biomarker candidates within a multivariate
317 predictive modeling framework. Application of Stabl to several real-world biomarker discovery tasks
318 demonstrates the versatility of the algorithm across a range of omic technologies, single- and multi-omic
319 datasets, and clinical endpoints. Results from these diverse clinical use cases emphasize the advantage

320 of Stabl's data-driven adaptation to the specificities of each omic dataset, which enables reliable selection
321 of biologically interpretable biomarker candidates conducive to further clinical translation.
322

323 Stabl builds on previous methods, including Bolasso, SS, and Model-X knockoff. These methods improve
324 reliability of sparse learning algorithms by employing a bootstrap procedure, or using artificial
325 features.^{5,12,14,16} However, these methods rely on a fixed or user-defined frequency threshold to
326 discriminate between informative and uninformative features. In practice, in the $p \gg n$ context, objective
327 determination of the optimal frequency threshold is difficult without prior knowledge of the data, as shown
328 by the results from our synthetic modeling. The requirement for prior knowledge impairs the capacity for
329 predictive model building, limiting these previous methods to sole feature selection.
330

331 Stabl improves on these methods by experimentally, and, under certain assumptions, theoretically,
332 generalizing previous false discovery rate control methods devised for model-X knockoffs and random
333 permutation noise.^{13,45,46} Minimization of the FDP surrogate (FDP₊) offers two main benefits. First, it
334 expresses a trade-off between reliability and sparsity, as it is the sum of an increasing and a decreasing
335 function of the threshold. Second, assuming exchangeability between artificial and uninformative features
336 Stabl's procedure guarantees a stochastic upper bound to the FDP using the reliability threshold
337 estimate, which ensures reliability in the optimization procedure. By minimizing this function *ex-ante*,
338 Stabl objectively defines a model fit from the procedure without requiring prior knowledge of the data.
339

340 On a synthetic dataset, we experimentally demonstrate that Stabl selects an optimal reliability threshold
341 by minimizing the FDP₊ and allows for improved reliability and sparsity compared to Lasso or EN at
342 similar predictivity performance. When tested on real-world omic studies, Stabl also performed favorably
343 compared to Lasso and EN. For each case study, the identification of a manageable number of reliable
344 biomarkers facilitated the interpretation of the multivariate predictive model. Prior analyses of similar
345 datasets^{18,20,21,31} required suboptimal analysis frameworks: either post-hoc analyses were performed
346 using user-defined cut-offs for feature selection after an initial model fit, or features associated with the
347 clinical endpoint were selected before modeling, thus risking overfitting. In contrast, Stabl embeds the
348 discovery of reliable candidate biomarkers within the predictive modeling, alleviating the need for
349 separate analyses.
350

351 Stabl extended readily to analysis of multi-omic datasets where a predictive model can utilize features
352 from different biological systems. Here, Stabl offers an alternative that avoids the potential shortcomings
353 of early and late fusion strategies. In the case of early fusion, all omic datasets are first concatenated
354 before applying a statistical learner. This leads to optimization on all omics combined, regardless of the
355 specific properties (e.g., dimensions, correlation structure, underlying noise) of individual omic
356 datasets.⁴⁷⁻⁵⁰ In contrast, the late fusion method trains the learner on each omic data layer independently
357 and merges the predictions into a final dataset.^{19,21,31,51-53} In this case, although a model is adapted to
358 each omic, the resulting model does not weigh features from different omics directly against each other.
359 Stabl analyzes each omic data layer independently and fits specific reliability thresholds before selecting
360 the most reliable features to be merged in a final layer, thus combining the advantages of both methods.
361 Multi-omic data integration with Stabl was particularly useful for analysis of our newly generated dataset
362 in patients undergoing surgery. In this case, the Stabl model comprised several features that were
363 biologically consistent across the plasma and single-cell datasets, revealing a patient-specific immune
364 signature predictive of SSI that appears to be programmed before surgery.
365

366 Our study has several limitations. Although we demonstrate the validity and performances of Stabl
367 experimentally and theoretically under the assumption of exchangeability between artificial and
368 uninformative features, a more general theoretical underpinning of the method will require further
369 guarantee. In addition, our evaluation of Stabl's performance focused on fitting Lasso and EN models as
370 gold standard SRMs. Further development of Stabl will be needed to allow for fitting of any SRM. While
371 Stabl is designed to simultaneously optimize reliability, sparsity, and predictivity performances, other
372 algorithms have been developed to address each of these performance tasks individually, such as double
373 machine learning⁵⁴ for reliability, Boruta⁵⁵ for sparsity, and random forest⁵⁶ or gradient boosting⁵⁷ for
374 predictivity. Additional studies are required to systematically evaluate each method's performance in
375 comparison to, or integrated with, the Stabl statistical framework. Finally, multi-omic data integration is

376 an active area of research. Integrating emerging algorithms such as cooperative multiview learning⁵⁸ may
377 further improve Stabl's performance in multi-omic modeling tasks.

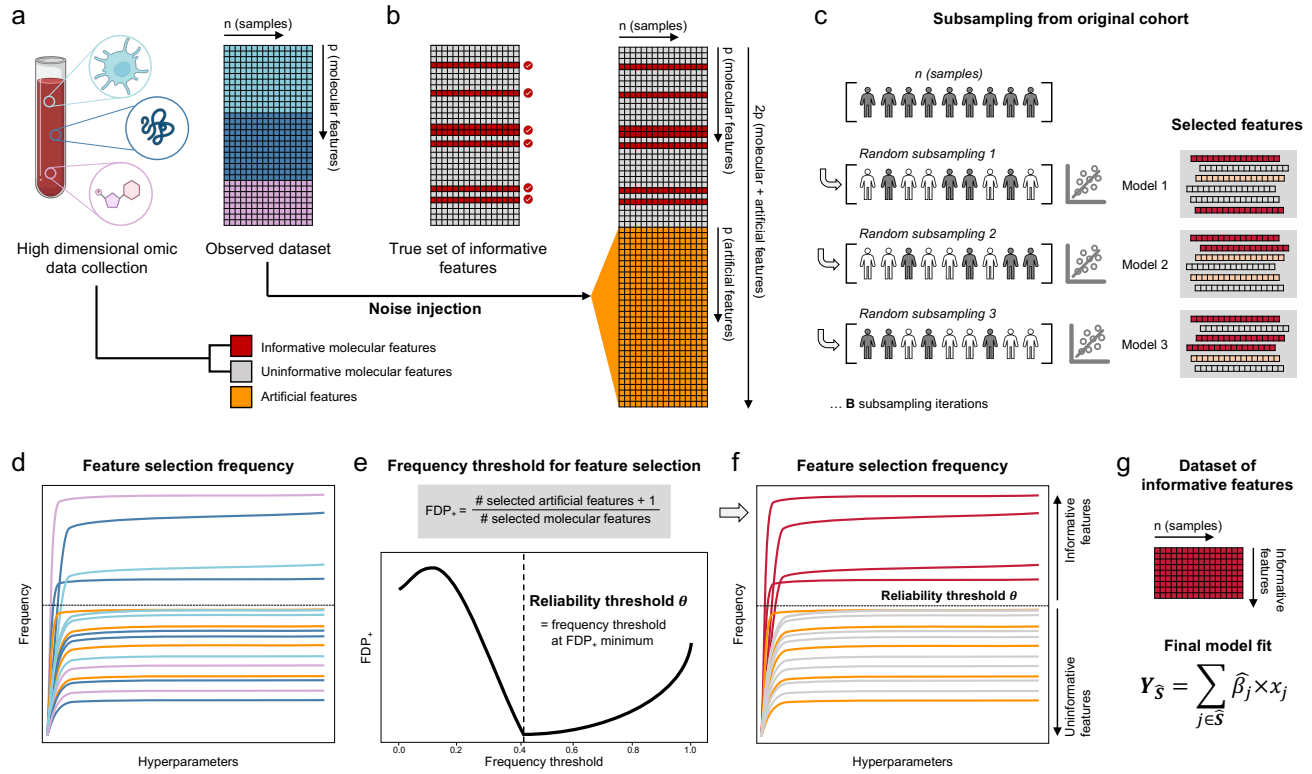
378
379 Analysis of high-dimensional omic data has transformed the biomarker discovery process but
380 necessitates new machine learning methods to facilitate clinical translation. Stabl addresses key
381 requirements of an effective biomarker discovery pipeline offering a unified supervised learning
382 framework that bridges predictive modeling of clinical endpoints with selection of reliable candidate
383 biomarkers. Stabl enabled identification of biologically plausible biomarker candidates across multiple
384 real-world single- and multi-omic datasets, providing a robust machine learning pipeline that we believe
385 can be generalized to all omic data.

386
387 **ACKNOWLEDGEMENT:** We thank Dr. Robert Tibshirani for the thorough and critical reading of the
388 manuscript.

389
390 **FUNDING:** This work was supported by the national institute of health (NIH) R35GM137936 (BG),
391 P01HD106414 (NA, DKS, BG), 1K99HD105016-01 (IAS), the Center for Human Systems Immunology
392 at Stanford (BG); the German Research Foundation (JE); the March of Dimes Prematurity Research
393 Center at Stanford University (#22FY19343); the Bill & Melinda Gates Foundation (OPP1189911); the
394 Stanford Maternal and Child Health Research Institute (DF, DKS, BG, NA, MSA); the Charles and Mary
395 Robertson Foundation (NA, DKS)

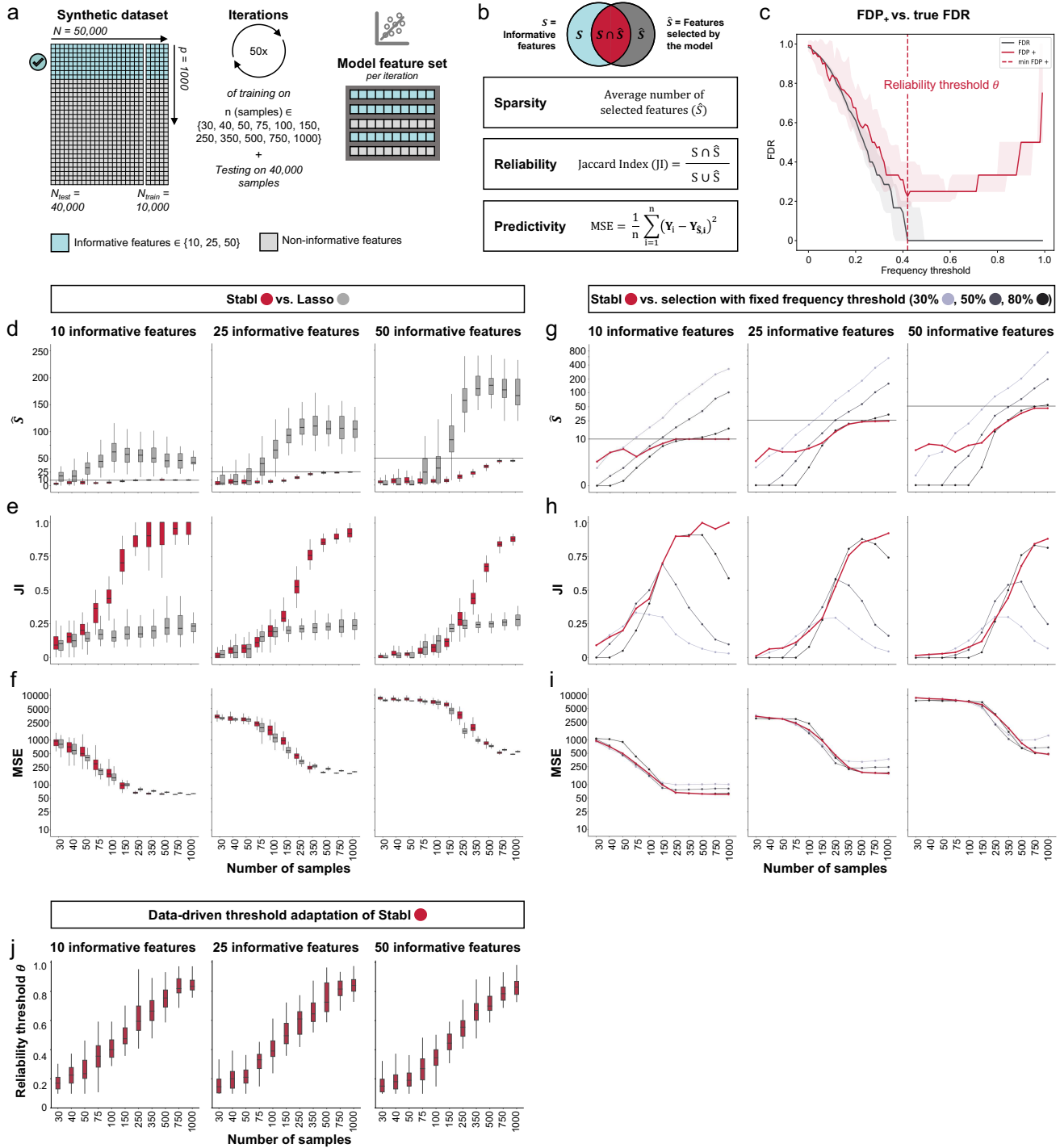
396

Figures



399
400
401
402
403
404
405
406
407
408
409
410

Fig. 1 | Overview of the Stabl algorithm. **a.** An original dataset of size $n \times p$ is obtained from measurement of p molecular features in each one of n samples. **b.** Among the observed features, some are informative (related to the outcome, red), and others are uninformative (unrelated to the outcome, grey). p artificial features (orange), all uninformative by construction, are injected into the original dataset to obtain a new dataset of size $n \times 2p$. **c.** B sub-sample iterations are performed from the original cohort of size n . At each iteration k , Lasso models varying in their regularization parameter λ are fitted on the subsample, which results in a different set of selected features for each iteration. **d.** In total, for a given λ , B sets or selected features are generated. The proportion of sets in which feature i is present defines the feature selection frequency $f_i(\lambda)$. Plotting $f_i(\lambda)$ against $1/\lambda$ yields a stability path graph. Features whose maximum frequency is above a frequency threshold (t) are selected in the final model. **e.** Stabl uses the reliability threshold (θ), obtained by computing the minimum to the false discovery proportion surrogate (FDP_+ , see methods). **f.g.** The set of features with a selection frequency larger than θ (i.e., reliable features) is included in a final predictive model.



411
 412 **Fig. 2 | Synthetic dataset benchmarking.** **a.** A synthetic dataset consisting of $N = 50,000$ samples \times $p = 1,000$ features was generated. Some
 413 features are correlated with the outcome (informative features, light blue), while the others are not (uninformative features, grey). Forty thousand
 414 samples are held out for validation. Out of the remaining 10,000, 50 sets ranging of sample sizes n ranging from 30 to 1,000 are drawn randomly.
 415 **c.** Three metrics are used to evaluate performance: *sparsity* (average number of selected features compared to the number of informative
 416 features), *reliability* (Jaccard Index, JI, comparing the true set of informative features to the selected feature set), and *predictivity* (mean squared
 417 error, MSE). **c.** The surrogate for the false discovery proportion (FDP+, red line) and the experimental false discovery rate (FDR, dotted line) are
 418 shown as a function of the frequency threshold. An example is shown for $n = 150$ samples and 25 informative features (all other conditions are
 419 shown in Fig. S1). The FDP+ estimate approaches the experimental FDR around the reliability threshold, θ . **d-f.** Sparsity (**d**), reliability (JI, **e**),
 420 and predictivity performances (MSE, **f**) of Stabl (red box plots) and least absolute shrinkage and selection operator (Lasso, grey box plots) as
 421 a function of the number of samples (n , x-axis) for 10 (left panels), 25 (middle panels), or 50 (right panels) informative features. **g-i.** Sparsity (**g**),
 422 reliability (**h**), and predictivity (**i**) performances of models built using a data-driven reliability threshold θ (Stabl, red lines) or a fixed frequency
 423 threshold (i.e., SS) of 30% (light grey lines), 50% (Lasso, dark grey lines), or 80% (black lines). The feature set selected by Stabl remains closer
 424 in number (sparsity) and composition (reliability) to the true set of informative features, while achieving a superior or comparable predictive
 425 performance to models built using a fixed threshold. **j.** The reliability threshold chosen by Stabl is shown as a function of the sample size (n , x-
 426 axis) for 10 (left panel), 25 (middle panel), or 50 (right panel) informative features. Benchmarking of Stabl against elastic net (EN) is shown in
 427 Fig. S6.

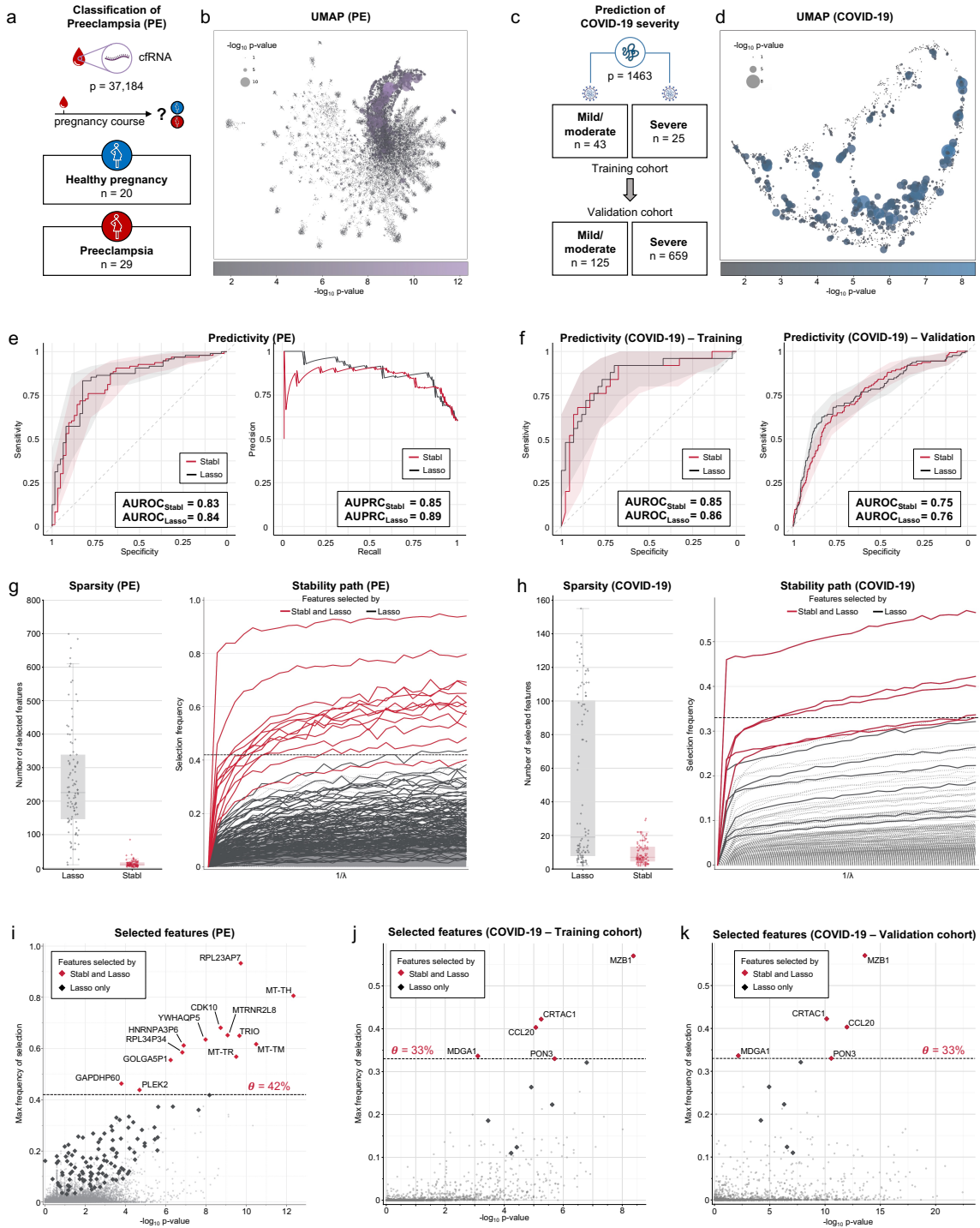
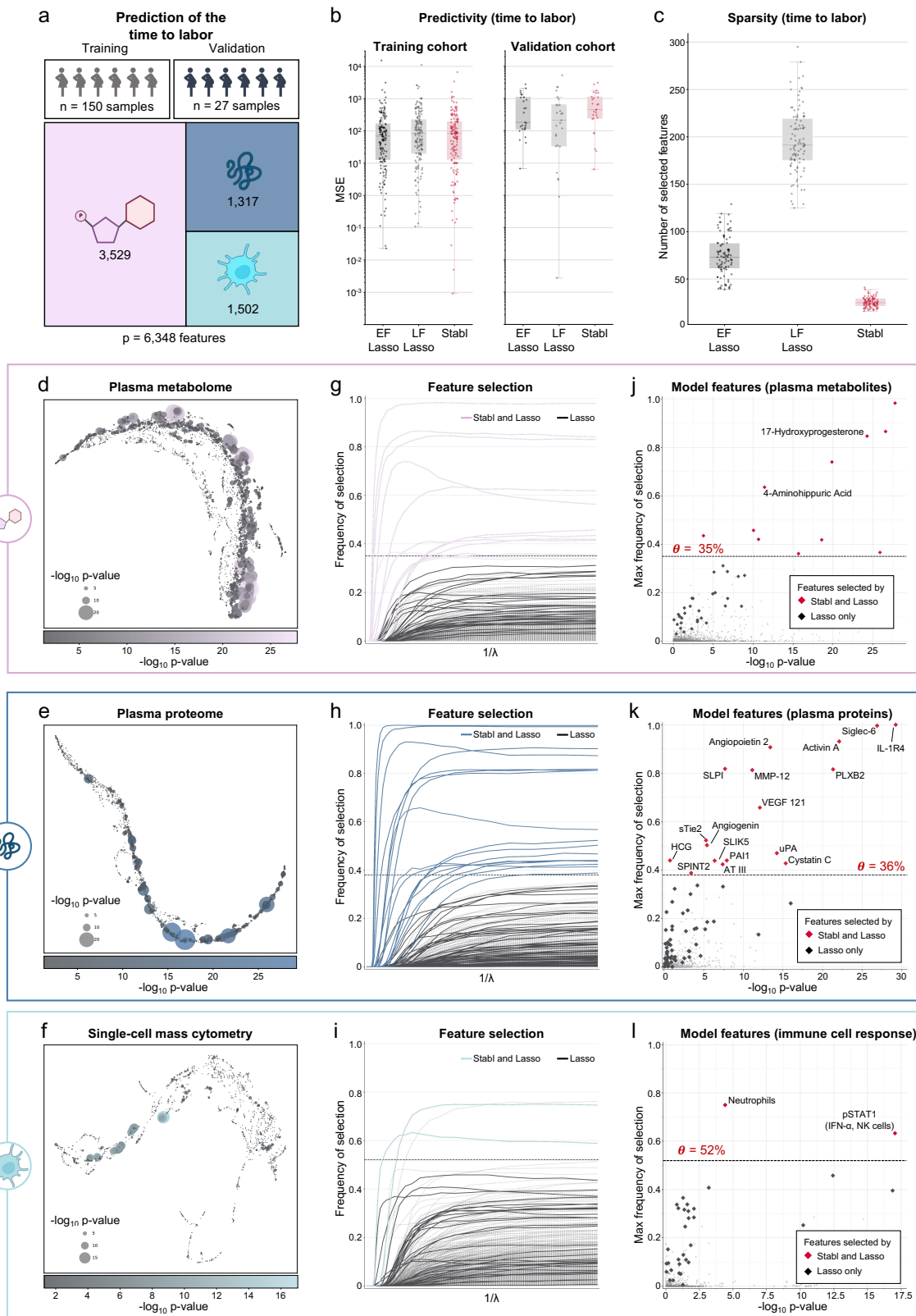


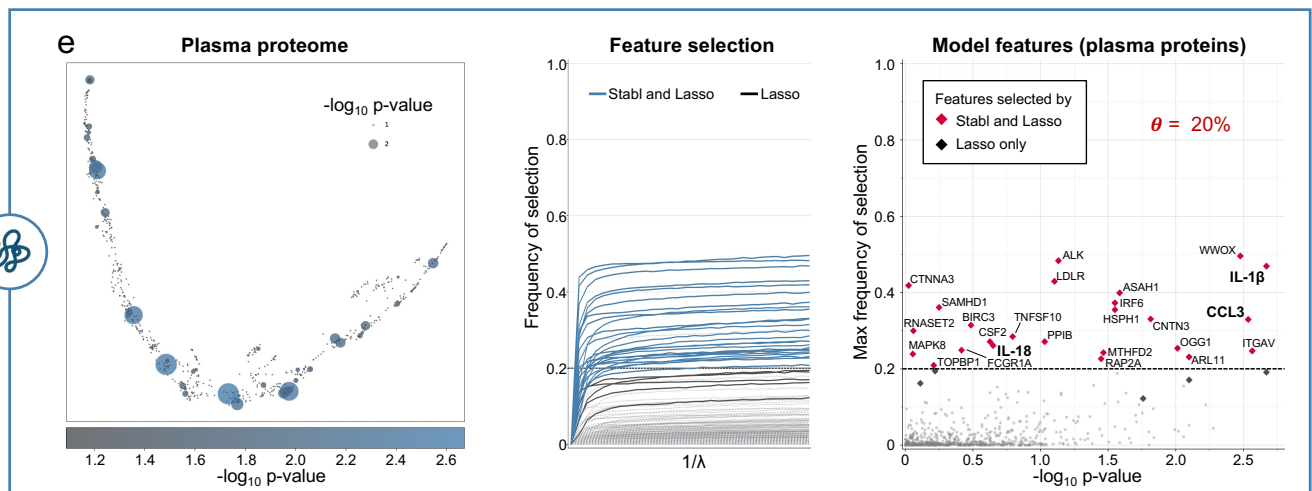
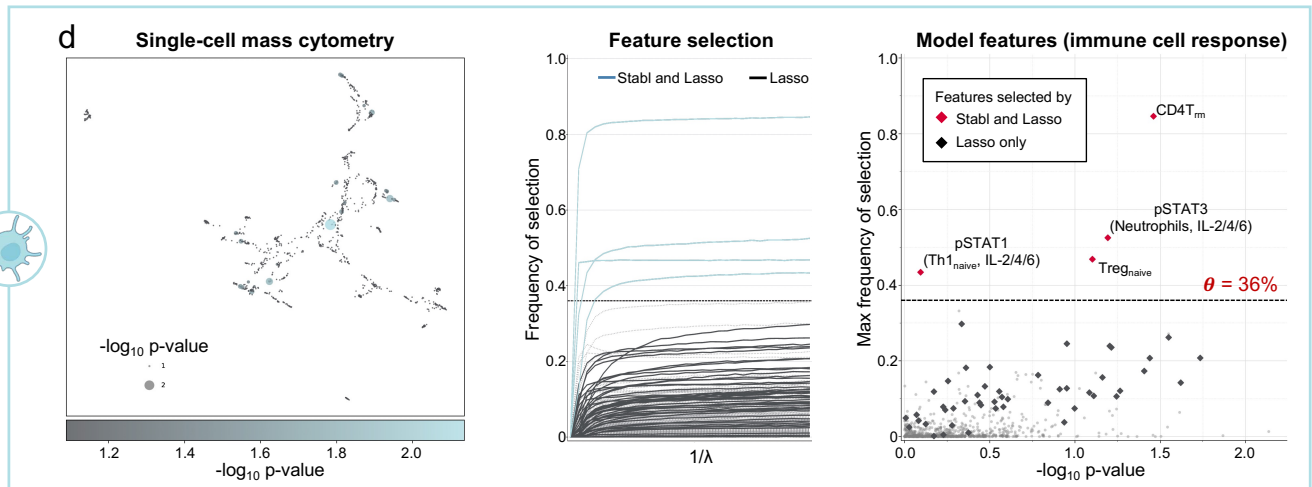
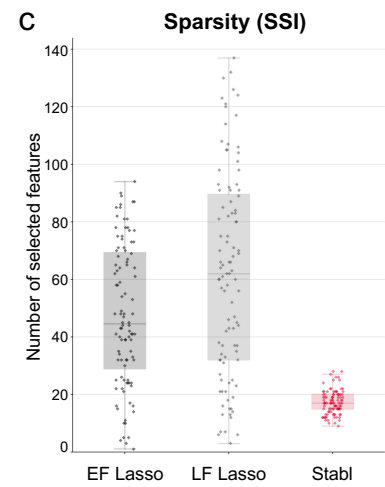
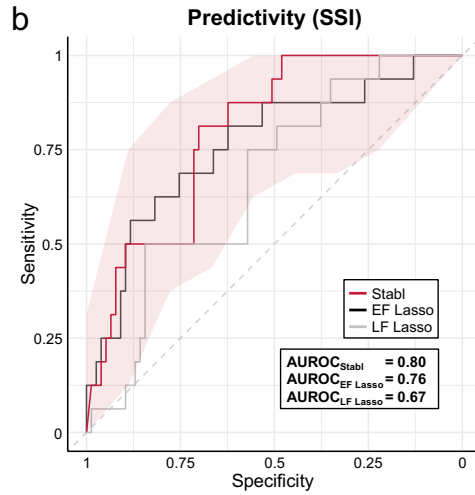
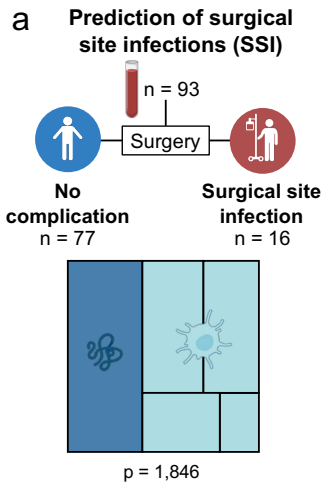
Fig. 3 | Performance of Stabl compared to Lasso on transcriptomic and proteomic data. **a.** Clinical case study 1: Classification of individuals with normotensive pregnancy or preeclampsia (PE) from the analysis of circulating cell-free RNA (cfRNA) sequencing data. Number of samples (n) and features (p) are indicated. **b.** UMAP visualization of the cfRNA transcriptomic features, node size and color are proportional to the strength of the association with the outcome calculated as the p -value in a univariate Mann-Whitney test using a $-\log_{10}$ scale. **c.** Clinical case study 2: Classification of mild vs. severe COVID-19 in two independent patient cohorts from the analysis of plasma proteomic data (Olink). **d.** UMAP visualization of the proteomic data. Node characteristics as in (b). **e.** Predictivity performances of Stabl and Lasso for the PE datasets. $AUROC_{Stabl} = 0.83$ [0.76, 0.90], $AUROC_{Lasso} = 0.84$ [0.78, 0.90] (p -value = 0.28, Bootstrap test); $AUPRC_{Stabl} = 0.85$ [0.77, 0.93], $AUPRC_{Lasso} = 0.89$ [0.83, 0.94] (p -value = 0.18). **f.** AUROC comparing predictive performance of Stabl and Lasso on training (left panel) and validation (right panel) cohorts for the COVID-19 dataset. Training: $AUROC_{Stabl} = 0.85$ [0.74, 0.94], $AUROC_{Lasso} = 0.86$ [0.75, 0.94] (p -value = 0.37). Validation: $AUROC_{Stabl} = 0.75$ [0.71, 0.79], $AUROC_{Lasso} = 0.76$ [0.71, 0.81] (p -value = 0.44). AUPRC are shown in Fig. S12. **g-h.** Left panels. Sparsity performances for the PE (g, number of features selected across cross-validation iterations, median_{Stabl} = 11.0, IQR = [7.8, 16.0], median_{Lasso} = 225.5, IQR = [147.5, 337.5], p -value < $1e-16$) and COVID-19 (h, median_{Stabl} = 7.0, IQR = [4.8, 13.0], median_{Lasso} = 19.0, IQR = [8.0, 100.0], p -value = $4e-10$) datasets. Right panels. Stability path graphs showing the regularization parameter against the selection frequency. The reliability threshold (θ), is indicated (dotted line) **i-k.** Volcano plots depicting the reliability performances of Stabl and Lasso for the PE (i), COVID-19 training (j) and COVID-19 validation (k) datasets. The maximum frequency of selection of each feature is plotted against the $-\log_{10}$ p -value using a univariate Mann-Whitney test. Features selected by Stabl/Lasso only are colored in red/black respectively. Features selected by Stabl are labeled. PE: mean $-\log_{10}(p\text{-value})_{Stabl} = 8.2$; mean $-\log_{10}(p\text{-value})_{Lasso} = 3.3$. COVID-19 training: mean $-\log_{10}(p\text{-value})_{Stabl} = 5.5$; mean $-\log_{10}(p\text{-value})_{Lasso} = 5.2$. COVID-19 validation: mean $-\log_{10}(p\text{-value})_{Stabl} = 9.7$; mean $-\log_{10}(p\text{-value})_{Lasso} = 7.8$. Benchmarking of Stabl against elastic net (EN) is shown in Fig. S11.

428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447



448
449
450
451
452
453
454
455
456
457
458
459
460

Fig. 4 | Stabl's performances on a triple-omic data integration task. **a.** Clinical case study 3. Prediction of the time to labor from the longitudinal assessment of plasma proteomic (Olink), metabolomic (untargeted mass spectrometry), and single-cell mass cytometry datasets in two independent longitudinal cohorts of pregnant individuals. **b.** Predictivity performances (MSE, median, and IQR) for early-fusion (EF), late-fusion (LF) Lasso and Stabl, on the training (left panel) and validation (right panel) cohorts. **c.** Sparsity performances (number of features selected across cross-validation iterations, $median_{Stabl} = 25.0$, $IQR = [22.0, 29.0]$, $median_{EF} = 73.0$, $IQR = [61.8, 87.3]$, p -value $< 1e-16$, $median_{LF} = 191.5$, $IQR = [175.8, 218.8]$, p -value $< 1e-16$). **d-f.** UMAP visualization of the metabolomic (**d**), plasma proteomic (**e**), and single-cell mass cytometry (**f**) datasets. Node size and color are proportional to the strength of the association with the outcome. **g-i.** Stability path graphs depicting the selection of metabolomic (**g**), plasma proteomic (**h**), and single-cell mass cytometry (**i**) features by Stabl. The data-driven reliability threshold θ is computed for individual omic datasets and indicated by a dotted line. **j-l.** Volcano plots depicting the reliability performances of Stabl and Lasso for each independent omic data: the metabolomics (**j**), plasma proteomic (**k**), and single-cell mass cytometry (**l**) datasets. The maximum frequency of selection of each feature is plotted against the $-\log_{10}$ p-value using a univariate Mann-Whitney test. Features selected by Stabl/Lasso only are colored in red/black respectively. Features selected by Stabl are labeled.



461
462
463
464
465
466
467
468
469

Fig. 5 | Candidate biomarker identification using Stabl for analysis of a newly generated multi-omic clinical dataset. a. Clinical case study 4. Prediction of postoperative surgical site infections (SSI) from the combined plasma proteomic and single cell mass cytometry assessment of pre-operative blood samples in patients undergoing abdominal surgery. **b.** Predictivity performances (AUROC) for Stabl, early fusion (EF) and late fusion (LF) Lasso. **c.** Sparsity performances (number of features selected across cross-validation iterations, $median_{Stabl} = 17.0$, $IQR = [15.0, 20.0]$, $median_{EF} = 44.5$, $IQR = [29.0, 69.3]$, $p\text{-value} < 1e-16$, $median_{LF} = 62.0$, $IQR = [32.0, 89.5]$, $p\text{-value} < 1e-16$). **d-e.** UMAP (left panel), stability paths (middle panel), and volcano plots (right panels) visualization of the single-cell mass cytometry (**d**) and plasma proteomics (**e**) datasets. The data-driven reliability threshold θ is computed for individual omic datasets and indicated by a dotted line on the volcano plots.

470	EXTENDED DATA	
471		
472	Extended Data Figure S1	Comparison of FDP ₊ vs. true FDR in synthetic dataset benchmarking.
473	Extended Data Figure S2	Comparison of Stabl and Lasso sparsity performance on synthetic data.
474	Extended Data Figure S3	Comparison of Stabl and Lasso reliability performance on synthetic data.
475	Extended Data Figure S4	Comparison of Stabl and Lasso predictivity performance on synthetic data.
476	Extended Data Figure S5	Comparison of Stabl and Elastic Net (EN) sparsity, reliability and predictivity performances on synthetic data.
477		
478	Extended Data Figure S6	Comparison of Stabl and Lasso sparsity, reliability and predictivity performances on synthetic data using Model-X knockoffs.
479		
480	Extended Data Figure S7	Comparison of Stabl and selection with fixed frequency threshold sparsity performance on synthetic data.
481		
482	Extended Data Figure S8	Comparison of Stabl and selection with fixed frequency threshold reliability performance on synthetic data.
483		
484	Extended Data Figure S9	Comparison of Stabl and selection with fixed frequency threshold predictivity performance on synthetic data.
485		
486	Extended Data Figure S10	Reliability threshold variation with the number of samples.
487	Extended Data Figure S11	Performance of Stabl compared to EN on transcriptomic (Preeclampsia, PE) and proteomic (COVID-19) datasets.
488		
489	Extended Data Figure S12	Predictivity of Stabl and Lasso for the training and validation cohort of the COVID-19 dataset.
490		
491	Extended Data Figure S13	Gating strategy for mass cytometry analyses (SSI dataset).
492	Extended Data Figure S14	Predictive performance of Stabl, Early Fusion and Late Fusion Lasso for the SSI dataset
493		
494	Extended Data Table S1	Univariate p-values for clinical case study 2: COVID-19.
495	Extended Data Table S2	Predictivity and sparsity comparison for Stabl vs. Stability Selection on single omic datasets.
496		
497	Extended Data Table S3	Features selected by Stabl for clinical case study 1: Preeclampsia (PE).
498	Extended Data Table S4	Features selected by Stabl for clinical case study 2: COVID-19.
499	Extended Data Table S5	Univariate p-values for clinical case study 3: Time to labor.
500	Extended Data Table S6	Features selected by Stabl for clinical case study 3: Time to labor.
501	Extended Data Table S7	Clinical information for clinical case study 4: surgical site infections (SSI).
502	Extended Data Table S8	Features selected by Stabl for clinical case study 4: surgical site infections (SSI).
503		

- 506 1. Subramanian, I., Verma, S., Kumar, S., Jere, A. & Anamika, K. Multi-omics Data Integration,
507 Interpretation, and Its Application. *Bioinforma. Biol. Insights* **14**, 1177932219899051 (2020).
- 508 2. Wafi, A. & Mirnezami, R. Translational –omics: Future potential and current challenges in precision
509 medicine. *Methods* **151**, 3–11 (2018).
- 510 3. Dunkler, D., Sánchez-Cabo, F. & Heinze, G. Statistical Analysis Principles for Omics Data. in
511 *Bioinformatics for Omics Data: Methods and Protocols* (ed. Mayer, B.) 113–131 (Humana Press,
512 2011). doi:10.1007/978-1-61779-027-0_5.
- 513 4. Ghosh, D. & Poisson, L. M. “Omics” data and levels of evidence for biomarker discovery. *Genomics*
514 **93**, 13–16 (2009).
- 515 5. Tibshirani, R. Regression Shrinkage and Selection Via the Lasso. *J. R. Stat. Soc. Ser. B Methodol.*
516 **58**, 267–288 (1996).
- 517 6. Zou, H. & Hastie, T. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B*
518 *Stat. Methodol.* **67**, 301–320 (2005).
- 519 7. Xu, H., Caramanis, C. & Mannor, S. Sparse Algorithms are not Stable: A No-free-lunch Theorem. 9.
- 520 8. Roberts, S. & Nowak, G. Stabilizing the lasso against cross-validation variability. *Comput. Stat. Data*
521 *Anal.* **70**, 198–211 (2014).
- 522 9. Homrighausen, D. & McDonald, D. The lasso, persistence, and cross-validation. in *Proceedings of*
523 *the 30th International Conference on Machine Learning* 1031–1039 (PMLR, 2013).
- 524 10. Olivier, M., Asmis, R., Hawkins, G. A., Howard, T. D. & Cox, L. A. The Need for Multi-Omics
525 Biomarker Signatures in Precision Medicine. *Int. J. Mol. Sci.* **20**, 4781 (2019).
- 526 11. Tarazona, S., Arzalluz-Luque, A. & Conesa, A. Undisclosed, unmet and neglected challenges in
527 multi-omics studies. *Nat. Comput. Sci.* **1**, 395–402 (2021).
- 528 12. Meinshausen, N. & Bühlmann, P. Stability Selection. Preprint at
529 <https://doi.org/10.48550/arXiv.0809.2932> (2009).
- 530 13. Candès, E., Fan, Y., Janson, L. & Lv, J. Panning for Gold: Model-X Knockoffs for High-dimensional
531 Controlled Variable Selection. Preprint at <https://doi.org/10.48550/arXiv.1610.02351> (2017).
- 532 14. Bach, F. R. Bolasso: model consistent Lasso estimation through the bootstrap. in *Proceedings of the*
533 *25th international conference on Machine learning - ICML '08* 33–40 (ACM Press, 2008).
534 doi:10.1145/1390156.1390161.
- 535 15. Barber, R. F. & Candès, E. J. Controlling the false discovery rate via knockoffs. *Ann. Stat.* **43**, 2055–
536 2085 (2015).
- 537 16. Ren, Z., Wei, Y. & Candès, E. Derandomizing Knockoffs. Preprint at <http://arxiv.org/abs/2012.02717>
538 (2020).
- 539 17. Weinstein, A., Barber, R. & Candès, E. A Power and Prediction Analysis for Knockoffs with Lasso
540 Statistics. Preprint at <https://doi.org/10.48550/arXiv.1712.06465> (2017).
- 541 18. Moufarrej, M. N. *et al.* Early prediction of preeclampsia in pregnancy with cell-free RNA. *Nature* **602**,
542 689–694 (2022).
- 543 19. Marić, I. *et al.* Early prediction and longitudinal modeling of preeclampsia from multiomics. *Patterns*
544 **3**, 100655 (2022).
- 545 20. Filbin, M. R. *et al.* Longitudinal proteomic analysis of severe COVID-19 reveals survival-associated
546 signatures, tissue-specific cell death, and cell-cell interactions. *Cell Rep. Med.* **2**, (2021).
- 547 21. Feyaerts, D. *et al.* Integrated plasma proteomic and single-cell immune signaling network signatures
548 demarcate mild, moderate, and severe COVID-19. *Cell Rep. Med.* **3**, 100680 (2022).
- 549 22. Kasten, M. & Giordano, A. Cdk10, a Cdc2-related kinase, associates with the Ets2 transcription factor
550 and modulates its transactivation activity. *Oncogene* **20**, 1832–1838 (2001).
- 551 23. Bellanger, J.-M. *et al.* The two guanine nucleotide exchange factor domains of Trio link the Rac1 and
552 the RhoA pathways in vivo. *Oncogene* **16**, 147–152 (1998).
- 553 24. Bach, T. L. *et al.* PI3K regulates pleckstrin-2 in T-cell cytoskeletal reorganization. *Blood* **109**, 1147–
554 1155 (2007).
- 555 25. Markovic, S. S. *et al.* Galectin-1 as the new player in staging and prognosis of COVID-19. *Sci. Rep.*
556 **12**, 1272 (2022).
- 557 26. COvid-19 Multi-omics Blood ATlas (COMBAT) Consortium. Electronic address:
558 julian.knight@well.ox.ac.uk & COvid-19 Multi-omics Blood ATlas (COMBAT) Consortium. A blood

- 559 atlas of COVID-19 defines hallmarks of disease severity and specificity. *Cell* **185**, 916-938.e58
560 (2022).
- 561 27. Mayr, C. H. *et al.* Integrative analysis of cell state changes in lung fibrosis with peripheral protein
562 biomarkers. *EMBO Mol. Med.* **13**, e12871 (2021).
- 563 28. Overmyer, K. A. *et al.* Large-scale Multi-omic Analysis of COVID-19 Severity. *medRxiv*
564 2020.07.17.20156513 (2020) doi:10.1101/2020.07.17.20156513.
- 565 29. Mohammed, Y. *et al.* Longitudinal Plasma Proteomics Analysis Reveals Novel Candidate Biomarkers
566 in Acute COVID-19. *J. Proteome Res.* **acs.jproteome.1c00863** (2022)
567 doi:10.1021/acs.jproteome.1c00863.
- 568 30. Gisby, J. *et al.* Longitudinal proteomic profiling of dialysis patients with COVID-19 reveals markers of
569 severity and predictors of death. *eLife* **10**, e64827 (2021).
- 570 31. Stelzer, I. A. *et al.* Integrated trajectories of the maternal metabolome, proteome, and immunome
571 predict labor onset. *Sci. Transl. Med.* **13**, eabd9898 (2021).
- 572 32. Suff, N., Story, L. & Shennan, A. The prediction of preterm delivery: What is new? *Semin. Fetal.*
573 *Neonatal Med.* **24**, 27–32 (2019).
- 574 33. Marquette, G. P., Hutcheon, J. A. & Lee, L. Predicting the spontaneous onset of labour in post-date
575 pregnancies: a population-based retrospective cohort study. *J. Obstet. Gynaecol. Can. JOGC J.*
576 *Obstet. Gynecol. Can. JOGC* **36**, 391–399 (2014).
- 577 34. Shah, N. M. *et al.* Changes in T Cell and Dendritic Cell Phenotype from Mid to Late Pregnancy Are
578 Indicative of a Shift from Immune Tolerance to Immune Activation. *Front. Immunol.* **8**, 1138 (2017).
- 579 35. Kraus, T. A. *et al.* Characterizing the pregnancy immune phenotype: results of the viral immunity and
580 pregnancy (VIP) study. *J. Clin. Immunol.* **32**, 300–311 (2012).
- 581 36. Shah, N. M., Lai, P. F., Imami, N. & Johnson, M. R. Progesterone-Related Immune Modulation of
582 Pregnancy and Labor. *Front. Endocrinol.* **10**, 198 (2019).
- 583 37. Brinkman-Van der Linden, E. C. M. *et al.* Human-specific expression of Siglec-6 in the placenta.
584 *Glycobiology* **17**, 922–931 (2007).
- 585 38. Kappou, D., Sifakis, S., Konstantinidou, A., Papantoniou, N. & Spandidos, D. A. Role of the
586 angiopoietin/Tie system in pregnancy (Review). *Exp. Ther. Med.* **9**, 1091–1096 (2015).
- 587 39. Huang, B. *et al.* Interleukin-33-induced expression of PIBF1 by decidual B cells protects against
588 preterm labor. *Nat. Med.* **23**, 128–135 (2017).
- 589 40. Li, A., Lee, R. H., Felix, J. C., Minoo, P. & Goodwin, T. M. Alteration of secretory leukocyte protease
590 inhibitor in human myometrium during labor. *Am. J. Obstet. Gynecol.* **200**, 311.e1-311.e10 (2009).
- 591 41. Tosato, G. & Jones, K. D. Interleukin-1 induces interleukin-6 production in peripheral blood
592 monocytes. *Blood* **75**, 1305–1310 (1990).
- 593 42. Lee, J.-K. *et al.* Differences in signaling pathways by IL-1beta and IL-18. *Proc. Natl. Acad. Sci. U. S.*
594 *A.* **101**, 8815–8820 (2004).
- 595 43. Fong, T. G. *et al.* Identification of Plasma Proteome Signatures Associated with Surgery Using
596 SOMAscan. *Ann. Surg.* **273**, 732–742 (2021).
- 597 44. Rumer, K. K. *et al.* Integrated Single-cell and Plasma Proteomic Modeling to Predict Surgical Site
598 Complications: A Prospective Cohort Study. *Ann. Surg.* **275**, 582–590 (2022).
- 599 45. He, K. *et al.* A theoretical foundation of the target-decoy search strategy for false discovery rate
600 control in proteomics. **33**.
- 601 46. He, K., Li, M., Fu, Y., Gong, F. & Sun, X. Null-free False Discovery Rate Control Using Decoy
602 Permutations. *Acta Math. Appl. Sin. Engl. Ser.* **38**, 235–253 (2022).
- 603 47. Yuan, Y. *et al.* Assessing the clinical utility of cancer genomic and proteomic data across tumor types.
604 *Nat. Biotechnol.* **32**, 644–652 (2014).
- 605 48. Gentles, A. J. *et al.* Integrating Tumor and Stromal Gene Expression Signatures With Clinical Indices
606 for Survival Stratification of Early-Stage Non-Small Cell Lung Cancer. *J. Natl. Cancer Inst.* **107**,
607 djv211 (2015).
- 608 49. Perkins, B. A. *et al.* Precision medicine screening using whole-genome sequencing and advanced
609 imaging to identify disease risk in adults. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 3686–3691 (2018).
- 610 50. Chaudhary, K., Poirion, O. B., Lu, L. & Garmire, L. X. Deep Learning-Based Multi-Omics Integration
611 Robustly Predicts Survival in Liver Cancer. *Clin. Cancer Res. Off. J. Am. Assoc. Cancer Res.* **24**,
612 1248–1259 (2018).
- 613 51. Yang, P., Hwa Yang, Y., B. Zhou, B. & Y. Zomaya, A. A Review of Ensemble Methods in
614 Bioinformatics. *Curr. Bioinforma.* **5**, 296–308 (2010).

- 615 52. Zhao, J. *et al.* Learning from Longitudinal Data in Electronic Health Record and Genetic Data to
616 Improve Cardiovascular Event Prediction. *Sci. Rep.* **9**, 717 (2019).
- 617 53. Chabon, J. J. *et al.* Integrating genomic features for non-invasive early lung cancer detection. *Nature*
618 **580**, 245–251 (2020).
- 619 54. Chernozhukov, V. *et al.* Double/debiased machine learning for treatment and structural parameters.
620 *Econom. J.* **21**, C1–C68 (2018).
- 621 55. Kursu, M. B. & Rudnicki, W. R. Feature Selection with the Boruta Package. *J. Stat. Softw.* **36**, 1–13
622 (2010).
- 623 56. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
- 624 57. Friedman, J. H. Stochastic gradient boosting. *Comput. Stat. Data Anal.* **38**, 367–378 (2002).
- 625 58. Ding, D. Y., Li, S., Narasimhan, B. & Tibshirani, R. Cooperative learning for multiview analysis. *Proc.*
626 *Natl. Acad. Sci.* **119**, e2202113119 (2022).
- 627

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [NBTRA59139supplemental.pdf](#)
- [NBTRA59139supplemental.pdf](#)