

1 **Single nucleotide variants in *Pseudomonas aeruginosa* populations from sputum correlate**
2 **with baseline lung function and predict disease progression in individuals with cystic fibrosis**

3
4 Morteza M. Saber^{#1}, Jannik Donner^{#2}, Inès Levade², Nicole Acosta³, Michael D. Parkins^{3,4}, Brian
5 Boyle⁵, Roger Levesque⁵, Dao Nguyen^{1,2,6*} and B. Jesse Shapiro^{1,7*}

6
7 ¹Department of Microbiology and Immunology, McGill University, Montreal, QC, Canada

8 ²Department of Medicine, Research Institute of the McGill University Health Centre, Montreal,
9 Quebec H4A 3J1, Canada

10 ³Department of Microbiology, Immunology and Infectious Disease, University of Calgary,
11 Calgary, AB, Canada.

12 ⁴Department of Medicine, University of Calgary, Calgary, AB, Canada.

13 ⁵Integrative Systems Biology Institute, Université Laval, Québec, Canada

14 ⁶Meakins Christie Laboratories, Research Institute of the McGill University Health Centre,
15 Montreal, QC, Canada

16 ⁷McGill Genome Centre, Montreal, QC, Canada

17
18
19 * Corresponding authors

20 E-mail: dao.nguyen@mcgill.ca, jesse.shapiro@mcgill.ca

21 # These authors contributed equally to this work

22

23 Abstract

24 Complex polymicrobial communities inhabit the lungs of individuals with cystic fibrosis (CF) and
25 contribute to the decline in lung function. However, the severity of lung disease and its progression
26 in CF patients are highly variable and imperfectly predicted by host clinical factors at baseline,
27 CFTR mutations in the host genome, or sputum polymicrobial community variation. The
28 opportunistic pathogen *Pseudomonas aeruginosa* (*Pa*) dominates airway infections in the majority
29 of CF adults. Here we hypothesized that genetic variation within *Pa* populations would be
30 predictive of lung disease severity. To quantify *Pa* genetic variation within whole CF sputum
31 samples, we used deep amplicon sequencing on a newly developed custom Ion AmpliSeq panel of
32 209 *Pa* genes previously associated with the host pathoadaptation and pathogenesis of CF
33 infection. We trained machine learning models using *Pa* single nucleotide variants (SNVs),
34 clinical and microbiome diversity data to classify lung disease severity at the time of sputum
35 sampling, and to predict future lung function decline over five years in a cohort of 54 adult CF
36 patients with chronic *Pa* infection. The models using *Pa* SNVs alone classified baseline lung
37 disease with good sensitivity and specificity, with an area under the receiver operating
38 characteristic curve (AUROC) of 0.87. While the models were less predictive of future lung
39 function decline, they still achieved an AUROC of 0.74. The addition of clinical data to the models,
40 but not microbiome community data, yielded modest improvements (baseline lung function:
41 AUROC=0.92; lung function decline: AUROC=0.79), highlighting the predictive value of the
42 AmpliSeq data. Together, our work provides a proof-of-principle that *Pa* genetic variation in
43 sputum is strongly associated with baseline lung disease, moderately predicts future lung function
44 decline, and provides insight into the pathobiology of *Pa*'s effect on CF.

46 Importance

47 Cystic fibrosis (CF) is among the most common, life-limiting inherited disorder, caused by
48 mutations in the CF transmembrane conductance regulator (CFTR) gene. CF causes progressive
49 damage to the lungs, the major cause of morbidity and mortality in CF patients. However, the rate
50 of lung function decline is highly variable across CF patients, and cannot be fully explained using
51 existing biomarkers in the human genome or patient co-morbidities. *Pseudomonas aeruginosa*
52 (*Pa*) is known to evolve and adapt within chronic CF infections. We hypothesized that within-
53 patient *Pa* diversity could affect lung disease severity. In a CF cohort study, we demonstrate the
54 utility of machine learning tools for predictive modeling of baseline lung function and subsequent
55 decline in CF patients using deep within-patient *Pa* amplicon sequencing. Our findings show the
56 potential of these models to identify high-risk CF patients based on *Pa* diversity within the lung.

57 [Introduction](#)

58 Cystic fibrosis (CF) is an autosomal recessive disorder caused by mutations in the CF
59 transmembrane conductance regulator (CFTR) gene and is the most common lethal Mendelian
60 disease in populations with European ancestry (Welsh, Ramsey et al. 2001). The resulting lung
61 disease is the major cause of morbidity and mortality in CF patients, with lung failure the most
62 common cause of death (Turcios 2020). However, the rate of disease progression and lung function
63 decline is highly variable across CF populations, and cannot be fully explained by variations in
64 CFTR alleles or other modifier genes (Shanthikumar, Neeland et al. 2019).

65
66 While CF airway infections are polymicrobial and microbiome diversity has been associated with
67 lung disease severity in many studies such as (Cox, Allgaier et al. 2010, van der Gast, Flight, Smith
68 et al. 2015, Zhao, Schloss et al. 2012, Coburn, Wang et al. 2015, Cuthbertson, Walker et al. 2020),
69 *Pseudomonas aeruginosa* (*Pa*) is recovered in the majority of adult CF patients and often
70 dominates the CF airway microbiome once established as a chronic infection (Goddard, Staudinger
71 et al. 2012, Zhao, Schloss et al. 2012). Infection with *Pa* in early life is widely recognized to be
72 associated with a greater decline in lung function and mortality (Kosorok, Zeng et al. 2001,
73 Emerson, Rosenfeld et al. 2002, Fothergill, Walshaw et al. 2012). Notably, *Pa* airway infections
74 can persist even with highly effective CFTR-correcting treatment (Hisert, Heltshe et al. 2017,
75 Harris, Wagner et al. 2020).

76
77 Over the course of chronic infection in the CF lung, *Pa* undergoes genetic diversification, selection
78 and adaptive evolution, resulting in a genetically and phenotypically diverse population of
79 clonally-related *Pa* within each patient (Tümmler 2006, Smith, Buckley et al. 2006, Bragonzi,
80 Paroni et al. 2009, Mowat, Paterson et al. 2011, Folkesson, Jelsbak et al. 2012, Marvig, Sommer
81 et al. 2015). How this pathoadaptation affects the clinical course of CF lung disease remains poorly
82 understood. We therefore focused on examining the association between *Pa* genetic variation and
83 the severity and progression of lung disease in CF patients with chronic *Pa* infections. We
84 hypothesized that the within-host genetic variation in *Pa* populations during chronic CF lung
85 infections are associated with baseline lung function and subsequent progression (*i.e.* decline in
86 lung function), as measured by spirometry.

87 Within-host mutations can significantly affect the virulence of *Pa* and host responses to *Pa*,
88 (Marvig, Johansen et al. 2013, Marvig, Sommer et al. 2015, Williams, Evans et al. 2015,
89 Klockgether, Cramer et al. 2018, Dettman and Kassen 2021). Previous studies have examined the
90 genetic variation of *Pa* across cohorts of CF patients by performing whole-genome sequencing
91 (WGS) of one or few *Pa* clones isolated from CF sputum samples – an approach that fails to
92 capture the polyclonal nature of *Pa* in the CF lung and is subject to profound sampling bias. While
93 shotgun metagenomic analysis of CF sputum is increasingly used for microbiome analyses
94 (Nelson, Pope et al. 2019, Whelan, Waddell et al. 2020, Lim et al. 2014.), the overwhelming
95 abundance of host derived DNA in samples continues to hamper the ability to resolve within
96 species genetic variation. To overcome these challenges, here we applied a custom-made amplicon
97 sequencing (AmpliSeq) panel of 209 genes in the *Pa* genome previously known to be involved in
98 the pathoadaptation and pathogenesis of CF infections (Supplementary Data 1). The Ion AmpliSeq
99 platform was selected because it provides a means for quantitative and sensitive measurement of
100 single nucleotide variant (SNV) frequencies within the *Pa* population, directly from CF sputum
101 without the need to culture and sequence hundreds of isolates per individual sample.

102
103 We then used several machine learning (ML) approaches to classify lung disease severity (at the
104 time of sample collection) and to predict future disease progression (over five years) based on the
105 SNV frequency data from a cohort of 54 adult CF patients with chronic *Pa* infection. ML has been
106 successfully applied to predict phenotypes from genotype data in other model systems (Dias and
107 Torkamani 2019). ML models can explicitly include the interactions and correlations between
108 features (in our case, SNVs), which is particularly common in bacterial population structures in
109 which SNVs are often genetically linked on the same clonal genomic background (Lees, Mai et al.
110 2020).

111
112 Our study provides proof-of-principle evidence that the population of *Pa* in CF sputum includes
113 bacterial genetic biomarkers that are associated with lung disease status and could serve to identify
114 individuals at increased risk of future lung function decline. Additionally, this work identified
115 genetic variation in *Pa* genes that merit further investigation for their potential roles in the
116 pathogenesis of CF lung disease.

117 Results

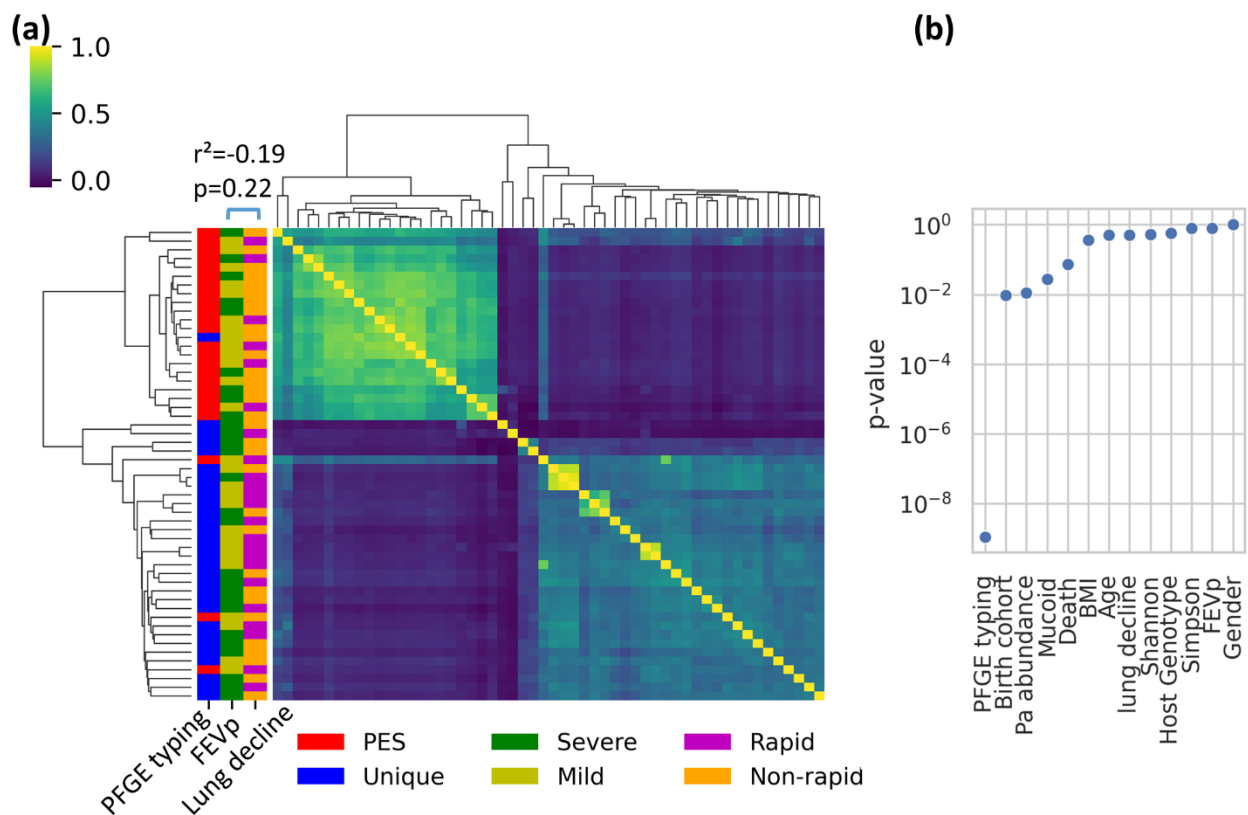
118 We studied a previously described and well-characterized cohort of young adult CF patients aged
119 18 to 22 with chronic *Pa* infection (Acosta, Heirali et al. 2018). Using *Pa* SNV frequencies
120 quantified by AmpliSeq in patient sputum, we sought to predict two measures of lung disease
121 severity: (1) baseline lung function (FEV_p score) at the time of sputum sample collection,
122 classified as severe or mild, and (2) relative lung function decline in the following five years,
123 classified as rapid or non-rapid. After filtering for AmpliSeq sequencing quality, we excluded 10
124 patients with low coverage of *Pa*, leaving 54 patients for further analysis (Methods). The clinical
125 and demographic characteristics of the final cohort are summarized in **Table 1**, and the excluded
126 patients were not apparent outliers in their clinical profiles (data not shown). From the filtered
127 sequence data we identified SNVs within the 209 genes represented in the AmpliSeq panel, and
128 estimated the frequency of each SNV within each patient sputum sample. In total across the 54
129 patient samples, we identified 7,867 synonymous and 4,452 non-synonymous SNVs
130 (Supplementary Data 2). All variants were used for population stratification analysis and only non-
131 synonymous SNVs were used for training ML models.

132

133 Stratification in the *Pa* population

134 We first quantified the extent of *Pa* population stratification, which can be problematic if there are
135 genetic structures (e.g. clonally-related clusters or strains) that can be confounded with the lung
136 disease outcomes of interest. If a particular genetic cluster or lineage is associated with worsened
137 lung disease, it then becomes difficult to pinpoint the most likely SNVs associated with the disease
138 outcome because all mutations (whether related to disease or not) in a cluster are correlated. We
139 know *a priori* based on pulsed-field gel electrophoresis (PFGE) typing that our dataset contains a
140 highly prevalent lineage of *Pa* (called Prairie Epidemic Strain or PES; sequence type (ST)-192;
141 **Table 1**) suspected to be associated with disproportionate lung disease (Somayaji, Lam et al.
142 2017). We confirmed this by hierarchical clustering of the *Pa* AmpliSeq data (n=12,319 SNVs,
143 including both synonymous and non-synonymous variants), which revealed two apparent genetic
144 clusters (**Fig. 1a**), one of which was strongly associated with the PES lineage (Fisher exact test,
145 odds ratio=168.0, $P = 1.1e-09$; **Fig. 1b**). The observed *Pa* genetic clusters are also weakly
146 associated with the birth cohort (Chi square test, $P = 0.0095$) which is likely due to unequal
147 prevalence of PES across birth cohorts (Supplementary table S1). No other clinical factor was

148 significantly associated with either genetic cluster (**Fig. 1b**). Importantly, neither cluster is
149 correlated with either baseline lung function (Fisher exact test, p-value: 0.81) or lung function
150 decline (Fisher exact test, p-value: 0.51) (**Fig. 1b**), indicating that these outcomes are unlikely to
151 be confounded by *Pa* population stratification, and finer-grained predictive modeling is warranted.
152 We also noted that lung disease progression over 5 years (lung function decline) is not significantly
153 correlated with baseline lung function at sample collection (**Fig. 1a**).
154



155
156

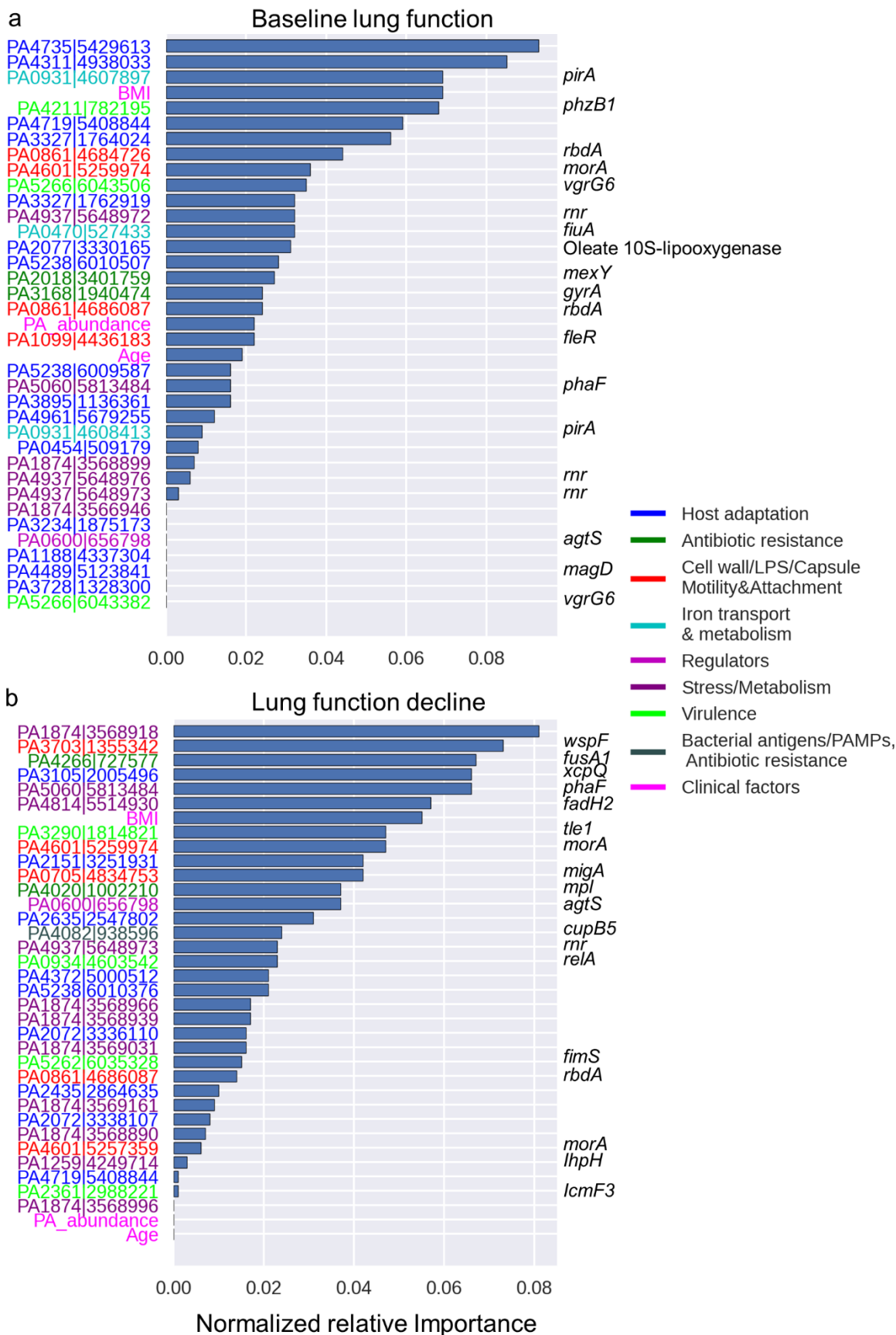
157 **Figure 1. The *Pa* population is stratified into two genetic clusters, neither of which is associated with**
158 **baseline lung function (FEVp) or lung function decline.**

159 (A) Heatmap showing correlations in SNV frequencies between pairs of sputum samples. Strong
160 correlations are yellow; weak correlations in blue. Rows and columns (samples) are ordered by
161 hierarchical clustering. Distribution of baseline lung function measured by FEVp score (27 Severe
162 and 27 Mild individuals), lung function decline (23 Rapid and 31 Non-rapid individuals) and PFGE
163 typing (25 PES and 29 Unique) are presented on the y axis. Baseline lung function and lung function
164 decline over five years are not significantly correlated (Pearson R^2 score=-0.19, p-value=0.22).
165 (B) *P*-values for the association between clinical data and genetic clusters are determined by *t*-test for
166 numerical data and chi-square test for categorical data (Methods). Only the association between
167 PFGE type (PES or non-PES) is significantly associated with the genetic clusters in panel A ($P <$
168 0.0045 after Bonferroni correction for multiple tests).

169 Genetic and clinical features associated with baseline lung function and lung function decline in CF
170 patients

171 A common challenge in predicting outcomes from sequence data is the sparsity of the data, that is,
172 relatively few available samples compared to the large number of genetic markers (called
173 “features” in ML context). To address this problem, feature selection has been used to remove
174 non-informative features (*i.e.*, SNVs and clinical factors) and focus only on the most predictive
175 ones (Mobegi, Cremers et al. 2017, Recker, Laabei et al. 2017, Méric, Mageiros et al. 2018,
176 Macesic, Don’t Walk et al. 2020). We used an ensemble gradient boosting technique for feature
177 selection (Methods). Out of 4,452 non-synonymous SNVs and eleven clinical factors considered,
178 our model selected only 34 SNVs (hereafter called predictor SNVs) and three clinical factors (age,
179 BMI and *Pa* abundance) that account for 99% of the cumulative feature importance (**Fig. 2**). This
180 means that a minimal set of SNVs and clinical factors provides 99% of the information used in
181 predicting baseline lung function at the time of sampling (**Fig. 2a**). An equivalent analysis for lung
182 function decline after five years identified 33 predictor SNVs and the same three clinical factors
183 that contributed to 99% of the cumulative feature importance (**Fig. 2b**). For both baseline lung
184 function and future lung function decline, the phenotype is not simply predicted based on the
185 presence/absence of each SNV, but rather on more subtle information about SNV allele
186 frequencies within patients. In other words, predictive SNVs occur at a range of frequencies, rather
187 than being clustered mainly around 0 or 1 (**Supplementary fig. 1**).

188

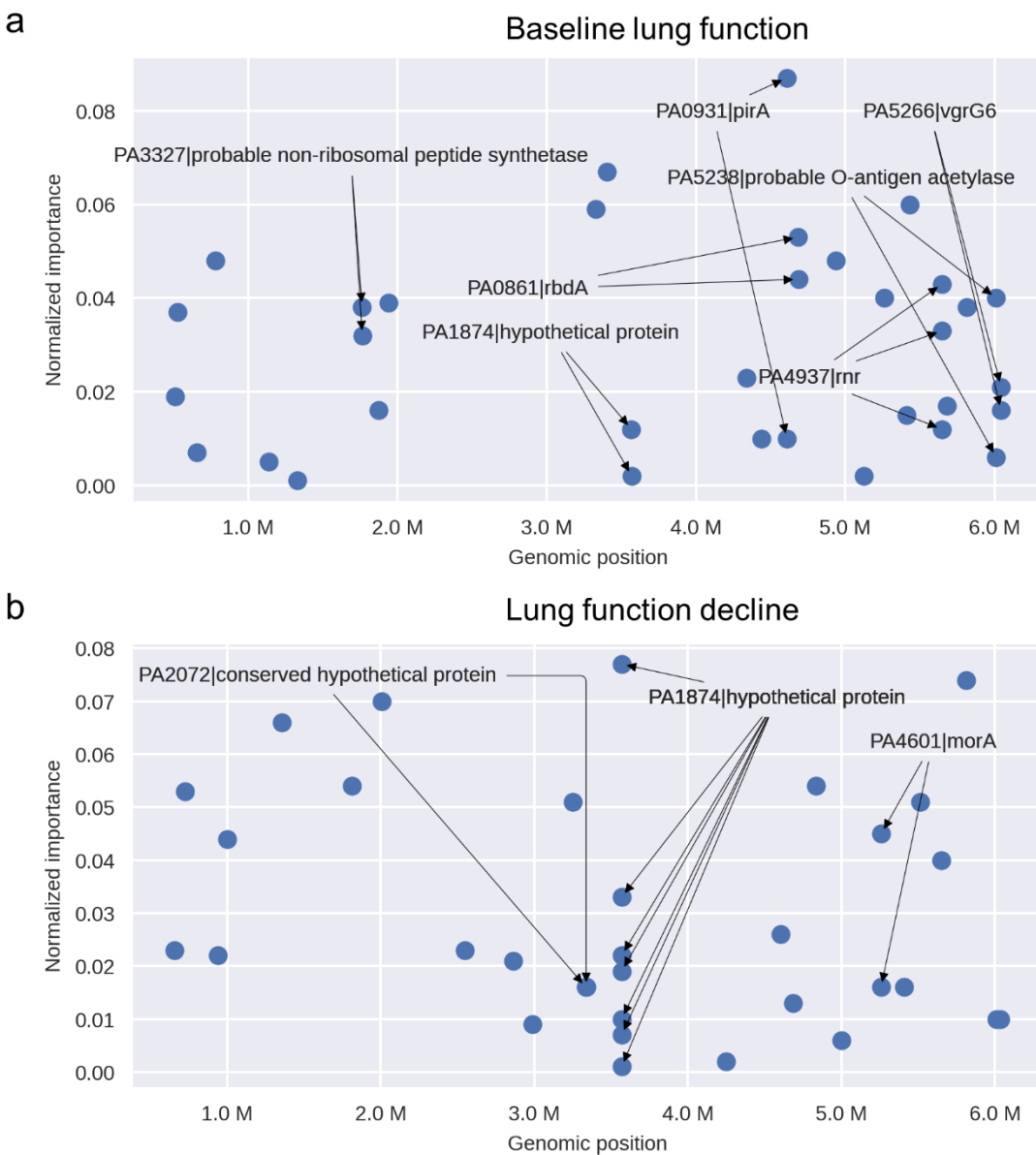


190 **Figure 2. *Pa* genes and clinical factors selected as top predictive features of baseline lung function**
191 **and lung function decline.** Normalized importance of genomic and clinical data that contribute to 99%
192 cumulative relative importance for prediction of (A) baseline lung function at time of sample collection and
193 (B) risk of 5-year progression (lung function decline). On the y-axis, gene identifiers (locus tag|chromosome
194 location based on PES genome) are color-coded based on their functional classification and named genes
195 are shown on the right, when available.

196
197 The three selected clinical factors associated with both baseline lung function and lung function
198 decline are body mass index (BMI), *Pa* relative abundance (from 16S rRNA gene amplicon
199 sequence data from a previous study of the same cohort; (Acosta, Heirali et al. 2018), and age (**Fig.**
200 **2**). Multiple studies have shown an association between poor lung function and low BMI (Snell,
201 Bennetts et al. 1998, Cystic Fibrosis Foundation 2006, Kumru, Emiralioglu et al. 2018), high
202 abundance of *Pa* (Cox, Allgaier et al. 2010) and age (Zhao, Hao et al. 2020, Cox, Allgaier et al.
203 2010). As expected, *Pa* relative abundance also showed a strong negative correlation with Shannon
204 and Simpson microbiome diversity indices (**Supplementary fig. 2**), indicating that *Pa* abundance
205 can be considered as a proxy for lung microbiome diversity in our dataset. However, Shannon and
206 Simpson diversity indices were not selected as predictive features in our model, consistent with a
207 previous work (Acosta et al. 2018; Zhao, Hao et al. 2020). This suggests that, even if low
208 microbiome diversity indices are associated with CF disease progression, the low diversity is likely
209 driven by the dominance of key pathogens such as *Pa*. By identifying previously known clinical
210 determinants of the lung function in CF patients, these results provide validation for the ensemble
211 gradient boosting approach to feature selection.

212
213 To interpret the possible roles of *Pa* SNVs in CF lung disease, we classified the known or predicted
214 function of genes containing predictor SNVs (hereafter called predictor genes) into functional
215 categories manually curated based on existing literature. The predictor SNVs with the highest
216 weighted importance for both baseline lung function and future lung function decline outcomes
217 are located within genes that play a role in seven functional categories (**Table 2**). The distribution
218 of predictor genes is generally similar to the distribution of gene functions included in the
219 AmpliSeq panel (**Supplementary fig. 3**). However, the predictor genes for baseline lung function
220 are enriched in iron transport and metabolism (13.4% in baseline lung function predictor genes vs.
221 1.4% in the AmpliSeq panel, $P = 0.00018$; **Supplementary fig. 3**). The genes encoding the ferric
222 enterobactin receptor (PirA) and the ferrichrome receptor (FiuA) respectively account for 8.9%

223 and 3.7% of the total normalized importance for baseline lung function (**Fig. 2a**), and *pirA* contains
224 multiple predictor SNVs (**Fig. 3a**). In contrast, the predictor genes for lung function decline are
225 enriched in stress/metabolism (33.6 % in lung function decline predictor genes vs. 13.4% in the
226 AmpliSeq panel, $P = 0.002$; **Supplementary fig. 3**). Notably, the hypothetical protein PA1874
227 accounts for 16.9% of the total normalized importance for lung function decline prediction and
228 includes 7 out of 33 predictor SNVs (**Fig. 2b, Fig. 3b**), as well as two predictor SNVs for baseline
229 lung function (**Fig. 3a**). This hypothetical protein has also been shown to play a role in resistance
230 of *Pa* to multiple antibiotics (Zhang and Mah 2008). The PA4937 gene, which encodes an RNase
231 R exoribonuclease involved in stress/metabolism, is also a predictor of both baseline lung function
232 and subsequent lung function decline, with multiple predictor SNVs (**Fig. 2a, Fig. 3a**). The two
233 genes PA0861 (*rbda*) and PA4601 (*morA*), which encode regulators of genes involved in the
234 bacterial cell wall, LPS, capsule, motility and attachment, are also among the important predictor
235 genes for both baseline lung function and lung function decline (**Fig. 2**).
236



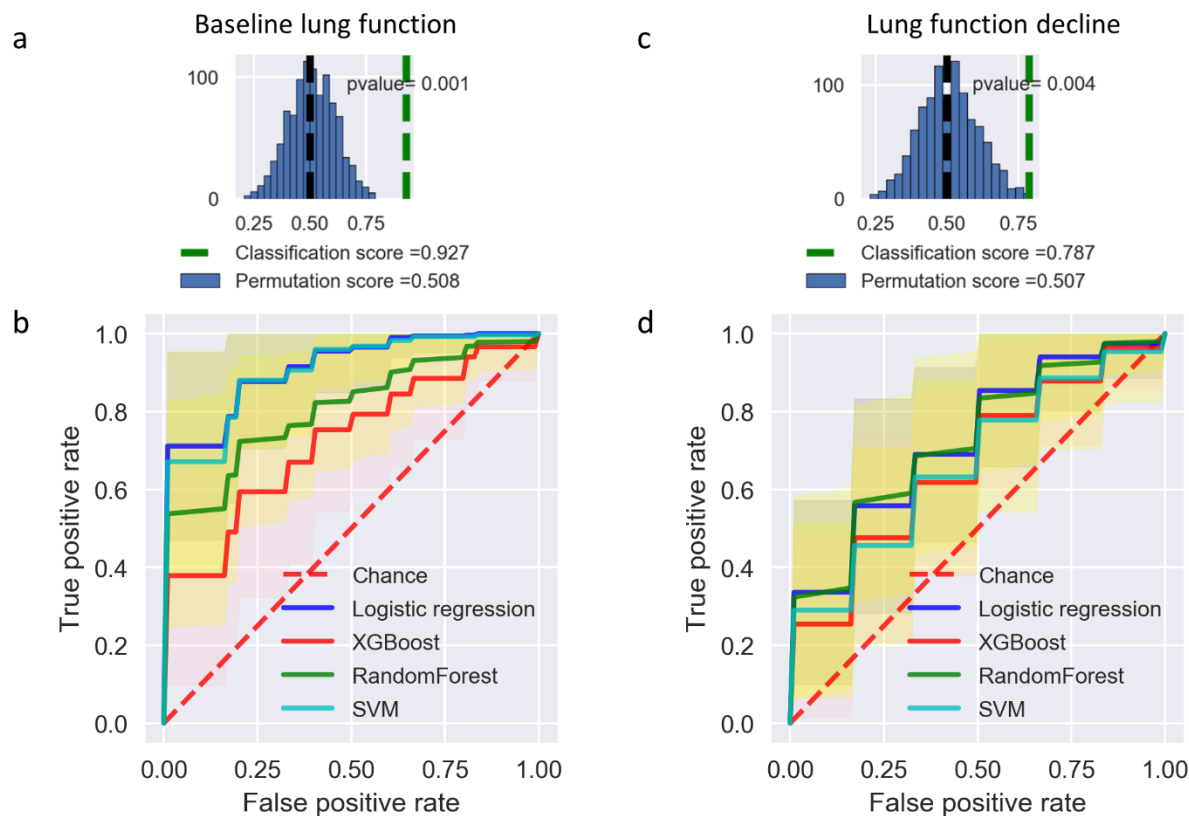
237
238

239 **Figure 3. Genomic locations and importance of genes containing predictor SNVs.**
240 (A) Genes containing SNVs predictive of baseline lung function.
241 (B) Genes containing SNVs predictive of lung function decline. Genes including multiple SNVs are
242 shown with arrows.

243
244

245 Predicting lung disease severity and progression in individuals with CF using genetic and clinical factors
246 Having identified *Pa* SNVs predictive of disease outcomes, we evaluated their predictive
247 performance using the area under receiver operating characteristic (AUROC) curve and other
248 standard metrics (Methods). A genetic programming algorithm identified logistic regression as the
249 best predictive model for both baseline lung function and lung function decline. To confirm this
250 result, we compared the performance of logistic regression with three other common ML
251 algorithms including support vector machines (SVM), random forests, and extreme gradient
252 boosting (XGBoost). It should be noted that while we used an ensemble LGBM model (Methods)
253 for feature selection, we did not use it for predictive modeling to avoid overfitting (*i.e.* to confirm
254 that our feature selection method is not biased toward selecting features that could be used only by
255 one specific ML model). Logistic regression had the best predictive performance for both
256 phenotypes (**Table 3, Fig. 4**). Specifically, cross-validation (Methods) showed that logistic
257 regression was the most accurate and precise model for both baseline lung function, with an
258 average AUROC score of 0.87 (95% CI, 0.84-0.90), and for lung function decline, with a score of
259 0.74 (95% CI: 0.71-0.78; **Table 3, Supplementary table 2**). Linear logistic regression is a simple
260 classification model that makes it reasonably robust against overfitting (Kuhn and Johnson 2013).
261 The second-best method was SVM, also a linear model (**Supplementary table 2**). Across all
262 models, the baseline lung function phenotype was more accurately predicted than lung function
263 decline, consistent with predictions becoming more uncertain further into the future (**Table 3,**
264 **Supplementary table 2**). Importantly, all models could predict both phenotypes significantly
265 better than expected by chance (compared to a permutation test using data with shuffled outcome
266 labels; **Fig. 4a, c, Supplementary fig. 4**).

267



268
269 **Figure 4. Predictive models of baseline lung function and lung function decline perform significantly**
270 **better than expected by chance.**

- 271 (A) Classification score of baseline lung function using logistic regression (green dashed line) is
272 significantly higher than expected based on permuted data (mean shown in black dashed line).
273 (B) Average AUROC scores of different ML models to predict baseline lung function, compared to
274 chance (permuted sample labels). Shading indicates the 95% confidence interval.
275 (C) Classification score of lung function decline using logistic regression compared to permuted data.
276 (D) Average AUROC scores for lung function decline prediction.

277
278 Clinical factors have been traditionally used to predict lung disease progression in CF patients
279 (Alaa and van der Schaar 2018). We therefore assessed if integrating clinical factors could improve
280 upon the predictions based on *Pa* AmpliSeq data alone. The three clinical factors (BMI, *Pa*
281 abundance and age) identified by our feature selection approach are previously recognized as
282 factors affecting lung function in CF patients. Including these three clinical factors in our
283 predictive models led to modest (~5% increase in AUROC) performance increases for both
284 baseline lung function and lung function decline outcomes across the four ML models (**Table 4,**
285 **Supplementary table 2**). We conclude that, while these clinical factors are useful, most of the
286 predictive power comes from the *Pa* genetic data.

287
288 Lack of generalizability is one of the main limiting factors for the translation of prediction models
289 into clinically useful diagnostics. Machine learning models often have low generalizability (i.e.
290 “overfit”) in scenarios where the model performs well on the dataset used to train the model, but
291 fails to achieve similar prediction accuracy on new data. We plotted learning curves to assess how
292 logistic regression predictions improved by training on more data (Raschka 2018). We found that
293 the performance difference between training and testing data decreases as sample size increases
294 (**Supplementary Figure 5**). There were no major differences in prediction accuracy of training
295 and testing datasets (**Supplementary Figure 5**) which suggests the model does not suffer from
296 significant overfitting. We also noted that cross validation scores for both baseline lung function
297 and lung function decline models continued to increase for the testing dataset as more data was
298 used for model training (**Supplementary Figure 5**) which suggests the model could be further
299 improved with more data.

300

301 Discussion

302 Considering the critical role of *Pa* in CF-related morbidity and mortality, here we established a
303 link between *Pa* genetic diversity and lung disease severity in a cohort of CF young adults with
304 chronic *Pa* infections. Despite a modest sample size, our study provides a proof of principle
305 demonstrating the utility of ML models for predictive modeling of lung function severity and
306 decline in CF patients using bacterial genetic and clinical data. Although our models do not appear
307 to be severely overfit, fully validating their predictive performance will require independent
308 cohorts. We also identified potential genetic biomarkers associated with lung disease severity.
309 Overall, our findings provide evidence that ML models can identify CF individuals at high-risk
310 for poor *Pa* infection outcomes using *Pa* genetic data.

311
312 Our work is based on a subset of samples from a previously described cohort study that identified
313 dominance of *Pa* in the sputum microbiome (and the resulting reduction of community diversity)
314 as a predictor of lung function decline in a cohort of young CF adults (Acosta, Heirali et al. 2018).
315 Here we focused on a subset of patients with a lung microbiome dominated by *Pa*. While these
316 patients are already at increased risk of lung disease, we found that the severity of disease at the

317 time of sampling and five years into the future could be predicted based on genetic variation within
318 the infecting *Pa* population. Even within this patient cohort in which *Pa* was present, we replicated
319 the finding that *Pa* relative abundance is associated with disease severity and progression –
320 although it is a less important biomarker than many SNVs within the *Pa* genome. This suggests
321 that genetic variation in dominant pathogens can significantly complement and improve upon
322 predictions of disease status based on the microbiome. Along these lines, another recent study
323 showed that the *Pa* genomic data can predict pathogenicity in mouse models (Pincus, Ozer et al.
324 2020).

325
326 We recognize that, in addition to variation in the host genome, the polymicrobial community
327 inhabiting the CF lung has been identified as an important modifier of disease progression.
328 Numerous studies of the lung microbiome have shown an association between decreasing
329 microbial community diversity and worsening lung function (Cox, Allgaier et al. 2010, van der
330 Gast, Flight, Smith et al. 2015, Zhao, Schloss et al. 2012, Coburn, Wang et al. 2015, Cuthbertson,
331 Walker et al. 2020), as well as progression to end-stage lung disease (Acosta, Heirali et al. 2018).
332 However, microbiome diversity may have limited predictive value as there is high interpersonal
333 variability in lung microbiomes (Cuthbertson, Walker et al. 2020), and a large number of adult CF
334 individuals have microbiomes dominated by pathogens such as *Pa*. Zhao, Hao et al. (2020) also
335 recently showed that microbiome composition data does not improve machine learning (ML)
336 prediction performance compared with using only clinical factors. Our results suggest that genetic
337 diversity within key pathogens like *Pa* could complement or even supersede microbiome
338 community diversity for predicting clinical outcomes in specific patient subsets.

339
340 Limitations of our study include a relatively small sample size of patients ($N=54$) from a single
341 cohort, and a relatively small number of targeted genes ($N=209$) included in the AmpliSeq panel.
342 As such, we consider our work a proof of concept that could be improved upon in larger cohorts
343 and by including more loci in the *Pa* genome. Indeed, learning curves showed that predictive
344 accuracy is likely to improve with more samples. Although we performed cross-validation by
345 subsampling our 54 patients for model training and testing, the model should ideally be tested on
346 a completely independent cohort to assess its real-world predictive value. Despite these limitations,
347 the models made significantly better predictions than expected by chance. As expected, predicting

348 lung function decline five years into the future proved more challenging than doing so at the time
349 of sampling. Still, our results provide a key first step toward clinical diagnostics of patients most
350 at risk of lung function decline.

351
352 As with any genotype-phenotype association method, our approach does not fully guarantee causal
353 relationships, and rather provides candidate genes for further experimental testing. Our study is
354 further complicated by the fact that the phenotypes of interest (i.e., baseline lung function and
355 future lung function decline) are complex host phenotypes, while the genotype data comes from
356 only *Pa*. It is therefore unclear to what extent *Pa* SNVs play a causal role in lung function decline,
357 or merely serve as useful biomarkers. Regardless, we were able to pinpoint SNVs in several genes
358 of interest. This was feasible because the strong population stratification of *Pa* into PES and non-
359 PES lineages was fortunately not associated with the disease outcomes of interest. This allowed
360 us to identify SNVs in several genes that provided independent biomarkers of disease.

361
362 Several candidate genes containing SNVs predictive of disease status and progression were
363 identified as targets for further investigation. We note that genes in the AmpliSeq panel were
364 selected *a priori* for their known involvement in virulence, disease progression, or within-patient
365 evolution. However, it was not known *a priori* which, if any, of these genes would contain SNVs
366 predictive of lung function or decline. For example, we found that baseline lung function predictor
367 SNVs are enriched in genes involved in iron transport and metabolism. The AmpliSeq panel only
368 included three iron-related genes, of which two (*pirA* and *fiuA*) contained SNVs associated with
369 baseline lung function. Updated AmpliSeq panels or whole-genome sequencing, along with
370 targeted experimental studies, could be used to test the hypothesis that variation in these genes
371 plays a role in disease progression. Multiple studies have shown competition for iron to be key for
372 the survival and virulence of many of the pathogens that reside in the CF lung, including *Pa*
373 (Bouvier 2016, Firoz, Haris et al. 2021). We also found that SNVs predictive of lung function
374 decline are enriched in genes involved in stress/metabolism. Notably, the gene PA1874 includes
375 seven predictor SNVs comprising 16.9% of total feature importance for lung function decline, and
376 two predictor SNVs for baseline lung function prediction, suggesting its general importance in
377 disease severity and progression in CF patients. PA1874 encodes a multidrug efflux pump
378 involved in biofilm-dependent resistance to antibiotics including tobramycin, gentamicin, and

379 ciprofloxacin (Zhang and Mah 2008; Poudyal and Sauer 2018), and could be a potentially
380 promising biomarker of CF disease severity, which merits further investigation.

381
382 Among the set of clinical factors studied, BMI, *Pa* abundance, and age were identified as important
383 predictors of both baseline lung function and lung function decline. These are all known risk
384 factors for CF disease severity and progression (Acosta et al. 2018, Snell, Bennetts et al. 1998,
385 Kumru, Emiralioglu et al. 2018, Cox, Allgaier et al. 2010, Zhao, Hao et al. 2020). By including
386 these features in our prediction models, we noted a moderate increase across all the measured
387 metrics relative to using only AmpliSeq data. These results are in line with previous studies
388 showing the improvement of ML-based phenotype prediction by adding relevant clinical data
389 (MacFadden, Melano et al. 2019, Pincus, Ozer et al. 2020). We note that clinical factors only
390 modestly improved the performance of the models (~5%), highlighting the rich information and
391 predictive value of the *Pa* AmpliSeq data alone.

392
393 In summary, our study demonstrates that SNVs in the *Pa* genome, assayed with an AmpliSeq panel
394 and identified by ML models, can be powerful predictors of lung disease severity and progression
395 in CF patients with chronic *Pa* infections. Even though this disease outcome is affected by multiple
396 microbial, host genetic and environmental factors, *Pa* SNVs add complementary predictive value.
397 With additional genetic and clinical data, our ML model could be further fine-tuned and eventually
398 used as a biomarker to preemptively identify individuals with CF at high-risk for more aggressive
399 observation and treatment.

400

401 [Materials and methods](#)

402 [Patient selection, sample and clinical data collection](#)

403 The Calgary biobank includes frozen whole sputum samples prospectively collected from
404 individuals with CF followed at the Calgary Adult CF clinic from 1998 to 2017, as described
405 previously (Acosta, Heirali et al. 2018, Acosta, Whelan et al. 2017). A cohort of 104 individuals
406 between the ages of 18 and 22 with sputum available from the Calgary biobank was previously
407 characterized (Acosta, Heirali et al. 2018). For this study, we selected from this cohort all
408 individuals with sputum cultures positive for *Pa* (64 out of 104 patients). Out of these 64 samples,
409 54 yielded AmpliSeq data of sufficient depth (>10X average depth of coverage of the targeted

410 genes) and were retained for further analysis. Clinical data collected for each patient is outlined in
411 **Tables 1** and **S1** and includes age, gender, body mass index, CFTR genotype, birth cohorts, and
412 microbiology (mucoid phenotype, *Pa* relative abundance, microbiome diversity indices). The
413 study was carried out with the approval from the Research Ethics Boards from the University of
414 Calgary (15-0854) and McGill University Health Centre (15-623).

415
416 As a measure of lung disease severity at the time of sputum collection, we used the spirometric
417 measure of forced expiratory volume in one second, percent predicted (hereafter referred to as
418 ‘baseline lung function’ and noted FEV_p) a standard measure of lung function normalized for age,
419 height, and self-identified gender and ethnicity. Baseline lung function was categorized as severe
420 for FEV_p < 60%, and mild for FEV_p ≥ 60%. Long-term lung function decline (hereafter noted as
421 ‘lung function decline’) was measured using the relative rate of FEV_p decline (determined by
422 subject-specific constructed linear regressions over the 5 years following sputum collection as
423 described in Acosta, Heirali et al. (2018). Lung function decline was categorized as ‘rapid’ when
424 the 5-year FEV_p decline was >5%, and ‘non-rapid’ when less than or equal to 5%.

425
426 [Sputum DNA extraction and microbiome analyses](#)

427 Genomic DNA was extracted from a single biobanked sputum sample per patient as previously
428 described (Acosta, Heirali et al. 2018), and used as template for 16S rRNA gene amplicon and Ion
429 AmpliSeq sequencing. The Prairie Epidemic Strain (PES) genotype, a highly prevalent strain in
430 our study population, was identified by pulse field gel electrophoresis (PFGE) and/or multi-locus
431 sequence typing (MLST) (Parkins, Glezerson et al. 2014). For microbiome analysis, bacterial
432 communities in CF sputum and reagent blanks were characterized by amplification and sequencing
433 the V3-V4 region of the 16S rRNA gene, as previously described (Acosta, Heirali et al. 2018). The
434 sequencing reads were then processed to identify operational taxonomic units (OTUs) (Acosta,
435 Whelan et al. 2017). Relative *Pa* abundance was determined as the proportion of *Pseudomonas*
436 reads relative to the total 16S reads.

437
438 [Ion AmpliSeq panel design and sequencing](#)

439 The AmpliSeq panel targeted 209 *Pa* genes previously implicated in pathogenicity, antimicrobial
440 resistance and within-host pathoadaptation during chronic infection (Supplementary Data 1). The

441 AmpliSeq primer panel (generated by Life Technologies, Carlsbad, CA, U.S.A.) was designed by
442 the AmpliSeq Custom Services (White Glove, Thermo Fisher Scientific) to provide high
443 sequencing coverage of the target genes based on the *Pa* PAO1 genome (NCBI accession number:
444 GCA_000006765.1), with 100% breadth of coverage for 205 genes and >96% in 4 genes, based
445 on the tiling of amplicons. Four additional genome assemblies of *Pa* clinical isolates
446 (GCF_004375495.1, GCF_004374685.1, GCF_004374275.1 and the PES genome (NCBI
447 BioProject: PRJNA750451) were also evaluated along with PAO1 for the optimization of primer
448 design, tiling and pooling to achieve maximal target coverage by the primer panel with minimal
449 misalignments and homology with the human genome.

450 AmpliSeq libraries were constructed using the Ion AmpliSeq™ Library kit 2.0 and IonCode™
451 barcode set with the following modifications. SparQ magnetic beads (Quantabio) were used for
452 purification, and individual libraries were quantified using the Quant-iT™ PicoGreen™ dsDNA
453 Assay Kit (ThermoFisher). Samples were mixed in equimolar proportions and the pooled library
454 (200 pM) was loaded on an Ion Chef for template preparation using HiQ reagents. The P1 v3 chips
455 were sequenced using an Ion Proton sequencer (500 flows) with P1 HiQ sequencing reagents
456 following manufacturer's instructions.

457 [AmpliSeq variant calling](#)

458 The quality of AmpliSeq sequencing was confirmed using TorrentSuite™ software (v.5.2, Thermo
459 Fisher Scientific Inc.). Raw sequencing reads were trimmed based on per-base phred quality score
460 cutoff ('q' flag) of 18, window size of 1 base pair and minimum remaining sequence length ('l'
461 flag) of 19 using fastq-mcf (v.1.04.636) (Aronesty 2013). Reads were aligned to the PES genome
462 (CP080405) using BWA MEM (Li, 2013), and the alignments were sorted and indexed using
463 SAMtools (v.1.9) (Li, Handsaker et al. 2009). Samples with average sequencing depth $\leq 10X$
464 across the target genes were discarded, leaving 54 samples for further analysis. Single-nucleotide
465 variants (SNVs) with minimum mapping quality of 20, minimum base quality of 18 and minimum
466 coverage of 10x were then identified using VarScan 2 (Koboldt, Zhang et al. 2012) and functional
467 consequences of each SNV were inferred using snpEFF (v.2.4.2) (Cingolani, Platts et al. 2012).
468 The SNV allele frequencies (ranging from 0 to 1) at each polymorphic site covered by the
469 AmpliSeq panel were used to generate a SNV frequency matrix, with samples as rows and
470 nucleotide positions as columns. For baseline lung function (measured based on FEVp score) and

471 lung function decline (disease progression) prediction analysis, all synonymous variants were
472 filtered out and only non-synonymous variants (including nonsense and missense mutations,
473 frameshift deletions and insertions) were used (Supplementary Data 2). All SNVs (including
474 synonymous sites) were included for population stratification analyses.

475

476 [Bacterial population stratification](#)

477 Population stratification in *Pa* was evaluated by calculating pairwise Pearson correlation
478 coefficients between sputum samples based on the SNV frequency matrix followed by
479 determination of distinct genome subgroups using hierarchical agglomerative clustering
480 implemented in SciPy (Virtanen, Gommers et al. 2020) and visualized using the python seaborn
481 package (Waskom, Botvinnik et al. 2017). This identified two major subclusters of *Pa*, one of
482 which was significantly enriched in PES strains. To determine if any clinical factors were
483 associated with these subclusters, we used t-tests for continuous variables including age, body
484 mass index (BMI), Shannon and Simpson diversity indices and *Pa* abundance in the sputum
485 sample. For binary variables including PFGE typing (PES or not), gender, host CFTR genotype,
486 death, mucoid presence/absence status, baseline lung function and lung disease progression (lung
487 function decline), we used a Fisher Exact Test. A Chi-square test was used for the multi-categorical
488 birth cohort factor.

489

490 [Feature selection](#)

491 In a machine learning context, a feature is defined as an individual measurable characteristic of an
492 observed phenomenon. In this study, the features considered are *Pa* genetic variants (SNV
493 frequencies) identified by the AmpliSeq panel and the clinical factors linked to the study patients
494 (Supplementary table 1). In order to reduce the high-dimensionality of the dataset (i.e. high ratio
495 of features to sample size), a feature selection approach was applied using feature-selector v1.0.0
496 (Koehrsen 2019). Briefly, 50 rounds of a gradient boosting ensemble method implemented in
497 LightGBM (Ke, Meng et al. 2017) were conducted on the training dataset sampled by a bootstrap
498 approach (43 samples in training set and 11 in test set in each bootstrap). The feature importance
499 (i.e. scores assigned to each input feature indicating the relative importance of the feature when
500 making a prediction) were averaged over the 50 bootstraps. The set of features required to obtain
501 99% cumulative relative importance were kept to perform prediction analysis and the remaining

502 features were discarded. The enrichment of predictor SNVs across functional gene categories
503 relative to the total genes in the AmpliSeq panel were measured using a Fisher exact test with a
504 family-wise error rate of 0.05 adjusted for multiple testing using Bonferroni method.

505

506 [Training predictive model of lung function](#)

507 The selected SNVs and clinical features were then used to train independent prediction models for
508 baseline lung function and disease progression (lung function decline) in the CF patient cohort. To
509 identify the best prediction model, a genetic programming algorithm implemented in TPOT (Le,
510 Fu et al. 2019) was used. TPOT attempts to identify the machine learning prediction pipeline with
511 the best performance (i.e. cross-validation using the area under the receiver operating characteristic
512 curve (AUROC) as the performance metric). We began with 10,000 random prediction pipelines
513 which were evolved over 5 generations with an offspring size of 100 in each generation, using the
514 recommended mutation rate of 0.9 and recombination rate of 0.1. The performance of each pipeline
515 was evaluated using 20-fold stratified shuffled cross validation. This analysis revealed logistic
516 regression with L2 regularization to be the most generalizable prediction model based on the cross-
517 validation scores. To confirm the results of genetic programming, three additional methods
518 including extreme gradient boosting implemented in XGBoost (Chen and Guestrin 2016),
519 ensemble decision trees implemented in random forest (Svetnik, Liaw et al. 2003) and linear
520 support vector machine (SVM) with linear kernel were also tested, using 20-fold stratified shuffled
521 cross validation implemented in scikit-learn (Pedregosa, Varoquaux et al. 2011). The performance
522 of machine learning models were evaluated using six metrics including AUROC, accuracy
523 (number of correct predictions/total number of predictions), precision (True Positives / (True
524 Positives + False Positives)), recall (True Positives / (True Positives + False Negatives)), F1 score
525 ($2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$), and balanced accuracy (bACC), the average of recall
526 obtained on each class (i.e. severe/mild for baseline lung function and rapid/non-rapid for lung
527 function decline).

528

529 To evaluate the statistical significance of the prediction performances (AUROC scores) obtained
530 by ML models in comparison with random expectations, a non-parametric permutation test
531 (Pedregosa, Varoquaux et al. 2011) was performed using 20-fold stratified shuffled cross-
532 validation across hundred rounds of label switching and model training followed by empirical p -

533 value estimation (i.e. the chance that the observed AUROC scores obtained using the data could
534 be obtained by chance alone).

535
536 [Data availability](#)

537 All amplicon sequencing data generating in this project are deposited in NCBI GenBank under
538 BioProject PRJNA763719.
539

540 [Acknowledgements](#)

541 The project was supported by funding from CIHR (PJT-148827 to DN) and a Vertex Research
542 Innovation Award (DN), and salary support from the Cystic Fibrosis Canada Research
543 Fellowship (Award ID 558850 to JD), the Leopoldina Foundation (German National Academy
544 of Sciences Leopoldina, Award ID LPDS 2017-17), the Réseau en Santé respiratoire (IL), and
545 the Fonds de Recherche en Santé Québec (IL, DN). MMS and BJS were supported by a Genome
546 Canada and Genome Québec Bioinformatics and Computational Biology grant.
547 We would like to acknowledge Michael Surette for providing the PES genome sequence and
548 Pradeep K. Singh for input in the Ampliseq design.

549 **Tables**

550

551 **Table 1. Patient clinical data.**

Patient data	Baseline lung function				Lung function decline over 5 years			
	Severe (n=27)	Mild (n=27)	Test statistic	<i>P</i> value	Non-rapid (n=31)	Rapid (n=23)	Test statistic	<i>P</i> value
<i>Pa</i> relative abundance	0.59 (0.32)	0.35 (0.32)	2.48	0.01	0.46 (0.33)	0.49 (0.32)	0.3	0.75
Age (year)	19.05 (1.13)	19.37 (1.16)	1.40	0.16	19.21 (1.1)	19.22 (1.25)	0.24	0.80
Body Mass Index (kg/m ²)	19.01 (2.3)	21.45 (2.3)	3.45	0.0005	20.62 (2.9)	19.70 (2.1)	1.18	0.23
Shannon index	1.12 (0.64)	1.21 (0.68)	0.42	0.66	1.31 (0.66)	1.06 (0.65)	1.44	0.24
Simpson index	0.48 (0.26)	0.5 (0.28)	0.23	0.81	0.54 (0.26)	0.45 (0.27)	1.61	0.20
PES (PFGE typing)	15	9	2.5	0.17	16	8	2.0	0.27
Homozygous ΔF508	15	13	1.34	0.78	14	14	0.52	0.28
Not-deceased (Death)	22	25	0.35	0.42	27	20	1.01	1
Male (Gender)	11	9	0.72	0.77	8	12	3.13	0.86
Mucoid	24	23	0.71	1	27	20	0.98	1
Birth cohort 1978- 1984	11	8	0.47	0.49	12	7	1.31	0.25
Birth cohort 1985- 1990	9	11	0.2	0.65	7	13	1.8	0.17

552

553 Values show the mean or absolute count, with standard deviation in parenthesis where applicable.

554 Baseline lung function is defined as severe when FEV_p < 60 and mild otherwise. Lung function

555 decline is defined as non-rapid when five-year FEV_p decline < 5% and rapid otherwise. The

556 relative abundance of *Pa*, as well as Shannon and Simpson diversity indices, were computed based

557 on 16S rRNA gene sequencing of the lung microbiome community. Homozygous $\Delta F508$ indicates
558 the counts of individuals with a $\Delta F508/\Delta F508$ genotypes; others include heterozygotes or other
559 genotypes. Test statistics are Wilcoxon rank-sum statistic for numerical data (Pa, Age, BMI,
560 Shannon and Simpson indices) and odds ratio for categorical data.

561 **Table 2. Functional classification of predictor genes used for prediction of baseline lung**
 562 **function and lung function decline.**

	Baseline lung function predictor genes	Lung function decline predictor genes	Shared genes
Host adaptation	PA0454 <i>conserved hypothetical protein</i> PA1188 <i>hypothetical protein</i> PA2077 <i>oleate 10S-lipoxygenase</i> PA3234 <i>probable sodium:solute symporter</i> PA3327 <i>probable non-ribosomal peptide synthetase</i> PA3728 <i>hypothetical protein</i> PA3895 <i>probable transcriptional regulator</i> PA4311 <i>conserved hypothetical protein</i> PA4489 <i>magD</i> PA4735 <i>hypothetical protein</i> PA4961 <i>hypothetical protein</i>	PA2072 <i>conserved hypothetical protein</i> PA2151 <i>conserved hypothetical protein</i> PA2435 <i>probable cation-transporting P-type ATPase</i> PA2635 <i>hypothetical protein</i> PA3105 <i>xcpQ</i> PA4372 <i>hypothetical protein</i>	PA4719 <i>probable transporter</i> PA5238 <i>probable O-antigen acetylase</i>
Antibiotic resistance	PA2018 <i>mexY</i> PA3168 <i>gyrA</i>	PA4020 <i>mpl</i> PA4082 <i>cupB5</i> PA4266 <i>fusA1</i>	
Cell wall, LPS, capsule, motility & attachment	PA1099 <i>fleR</i>	PA0705 <i>migA</i> PA3703 <i>wspF</i> PA4082 <i>cupB5</i>	PA0861 <i>rbdA</i> PA4601 <i>morA</i>
Iron transport and metabolism	PA0470 <i>fiuA</i> PA0931 <i>pirA</i>		
Regulators			PA0600 <i>agtS</i>
Stress/ metabolism		PA1259 <i>lhpH</i> PA4814 <i>fadH2</i>	PA1874 <i>hypothetical protein</i> PA4937 <i>rnr</i> PA5060 <i>phaF</i>
Virulence	PA4211 <i>phzB1</i> PA5266 <i>vgrG6</i>	PA0934 <i>relA</i> PA2361 <i>icmF3</i> PA3290 <i>tle1</i> PA5262 <i>fimS</i>	

563

564 **Table 3. Performance of logistic regression in predicting baseline lung function and lung**
 565 **function decline using genomic data only, or a combination of genomic and clinical data.**
 566 See Methods for descriptions of the performance metrics.

567

		Genomic data (95% CI)	Genomic and clinical data (95% CI)
Baseline lung function	AUROC	0.87 (0.84, 0.9)	0.92 (0.84, 1.00)
	bACC	0.81 (0.78, 0.84)	0.83 (0.72, 0.94)
	Accuracy	0.81 (0.78, 0.84)	0.83 (0.72, 0.94)
	F1	0.81 (0.78, 0.83)	0.83 (0.72, 0.94)
	Precision	0.83 (0.81, 0.86)	0.84 (0.73, 0.94)
	Recall	0.81 (0.78, 0.84)	0.83 (0.72, 0.94)
Lung function decline	AUROC	0.74 (0.71, 0.78)	0.79 (0.70, 0.88)
	bACC	0.63 (0.59, 0.66)	0.66 (0.59, 0.74)
	Accuracy	0.64 (0.60, 0.67)	0.67 (0.60, 0.75)
	F1	0.62 (0.58, 0.65)	0.66 (0.58, 0.74)
	Precision	0.65 (0.61, 0.69)	0.69 (0.60, 0.78)
	Recall	0.64 (0.60, 0.67)	0.67 (0.60, 0.75)

568

569 Bibliography

- 570 Acosta, N., A. Heirali, R. Somayaji, M. G. Surette, M. L. Workentine, C. D. Sibley, H. R. Rabin and M. D.
571 Parkins (2018). "Sputum microbiota is predictive of long-term clinical outcomes in young adults with cystic
572 fibrosis." Thorax 73(11): 1016-1025.
- 573 Acosta, N., F. J. Whelan, R. Somayaji, A. Poonja, M. G. Surette, H. R. Rabin and M. D. Parkins (2017).
574 "The evolving cystic fibrosis microbiome: a comparative cohort study spanning 16 years." Annals of the
575 American Thoracic Society, 14(8): 1288-1297.
- 576
577 Alaa, A. M., & van der Schaar, M. (2018). Prognostication and risk factors for cystic fibrosis via automated
578 machine learning. Scientific reports, 8(1), 1-19.
- 579
580 Aronesty, E. (2013). Fastq-mcf: sequence quality filter, clipping and processor.
581 Bouvier, N. M. (2016). Cystic fibrosis and the war for iron at the host-pathogen battlefield. Proceedings of
582 the National Academy of Sciences, 113(6), 1480-1482.
- 583
584 Bragonzi, A., M. Paroni, A. Nonis, N. Cramer, S. Montanari, J. Rejman, C. Di Serio, G. Döring and B.
585 Tümmler (2009). "Pseudomonas aeruginosa microevolution during cystic fibrosis lung infection establishes
586 clones with adapted virulence." American journal of respiratory and critical care medicine 180(2): 138-145.
- 587
588 Chen, T., & Guestrin, C. (2016, August). Xgboost: A scalable tree boosting system. In Proceedings of the
589 22nd acm sigkdd international conference on knowledge discovery and data mining (pp. 785-794).
- 590
591 Cingolani, P., A. Platts, L. L. Wang, M. Coon, T. Nguyen, L. Wang, S. J. Land, X. Lu and D. M. Ruden
592 (2012). "A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff:
593 SNPs in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3." Fly 6(2): 80-92.
- 594
595 Coburn, B., P. W. Wang, J. D. Caballero, S. T. Clark, V. Brahma, S. Donaldson, Y. Zhang, A. Surendra,
596 Y. Gong and D. E. Tullis (2015). "Lung microbiota across age and disease stage in cystic fibrosis."
597 Scientific reports 5(1): 1-12.
- 598
599 Cox, M. J., M. Allgaier, B. Taylor, M. S. Baek, Y. J. Huang, R. A. Daly, U. Karaoz, G. L. Andersen, R.
600 Brown and K. E. Fujimura (2010). "Airway microbiota and pathogen abundance in age-stratified cystic
601 fibrosis patients." PloS one 5(6): e11044.
- 602
603 Cuthbertson, L., A. W. Walker, A. E. Oliver, G. B. Rogers, D. W. Rivett, T. H. Hampton, A. Ashare, J. S.
604 Elborn, A. De Soyza and M. P. Carroll (2020). "Lung function and microbiota diversity in cystic fibrosis."
605 Microbiome 8: 1-13.
- 606
607 Dettman, J. R. and R. Kassen (2021). "Evolutionary genomics of niche-specific adaptation to the cystic
608 fibrosis lung in Pseudomonas aeruginosa." Molecular biology and evolution 38(2): 663-675.
- 609
610 Dias, R., & Torkamani, A. (2019). Artificial intelligence in clinical and genomic diagnostics. Genome
611 medicine, 11(1), 1-12.
- 612
613 Emerson, J., Rosenfeld, M., McNamara, S., Ramsey, B., & Gibson, R. L. (2002). Pseudomonas aeruginosa
614 and other predictors of mortality and morbidity in young children with cystic fibrosis. Pediatric
615 pulmonology, 34(2), 91-100.
- 616

- 617 Eyre, D. W., D. De Silva, K. Cole, J. Peters, M. J. Cole, Y. H. Grad, W. Demczuk, I. Martin, M. R. Mulvey
618 and D. W. Crook (2017). "WGS to predict antibiotic MICs for *Neisseria gonorrhoeae*." Journal of
619 Antimicrobial Chemotherapy 72(7): 1937-1947.
- 620
621 Firoz, A., M. Haris, K. Hussain, M. Raza, D. Verma, M. Bouchama, K. S. Namiq and S. Khan (2021). "Can
622 Targeting Iron Help in Combating Chronic *Pseudomonas* Infection? A Systematic Review." Cureus 13(3).
623
- 624 Flight, W. G., A. Smith, C. Paisey, J. R. Marchesi, M. J. Bull, P. J. Norville, K. J. Mutton, A. K. Webb, R.
625 J. Bright-Thomas and A. M. Jones (2015). "Rapid detection of emerging pathogens and loss of microbial
626 diversity associated with severe lung disease in cystic fibrosis." Journal of clinical microbiology 53(7):
627 2022-2029.
- 628
629 Folkesson, A., Jelsbak, L., Yang, L., Johansen, H. K., Ciofu, O., Høiby, N., & Molin, S. (2012). Adaptation
630 of *Pseudomonas aeruginosa* to the cystic fibrosis airway: an evolutionary perspective. Nature Reviews
631 Microbiology, 10(12), 841-851.
- 632
633 Fothergill, J. L., Walshaw, M. J., & Winstanley, C. (2012). Transmissible strains of *Pseudomonas*
634 *aeruginosa* in cystic fibrosis lung infections. European Respiratory Journal, 40(1), 227-238.
- 635
636 Goddard, A. F., B. J. Staudinger, S. E. Dowd, A. Joshi-Datar, R. D. Wolcott, M. L. Aitken, C. L. Fligner
637 and P. K. Singh (2012). "Direct sampling of cystic fibrosis lungs indicates that DNA-based analyses of
638 upper-airway specimens can misrepresent lung microbiota." Proceedings of the National Academy of
Sciences 109(34): 13769-13774.
- 639
640 Harris, J. K., B. D. Wagner, E. T. Zemanick, C. E. Robertson, M. J. Stevens, S. L. Heltshe, S. M. Rowe and
641 S. D. Sagel (2020). "Changes in airway microbiome and inflammation with ivacaftor treatment in patients
with cystic fibrosis and the G551D mutation." Annals of the American Thoracic Society 17(2): 212-220.
- 642
643 Hisert, K. B., S. L. Heltshe, C. Pope, P. Jorth, X. Wu, R. M. Edwards, M. Radey, F. J. Accurso, D. J. Wolter
644 and G. Cooke (2017). "Restoring cystic fibrosis transmembrane conductance regulator function reduces
645 airway bacteria and inflammation in people with cystic fibrosis and chronic lung infections." American
journal of respiratory and critical care medicine 195(12): 1617-1628.
- 646
647 Javan, A. O., Shokouhi, S., & Sahraei, Z. (2015). A review on colistin nephrotoxicity. European journal of
648 clinical pharmacology, 71(7), 801-810.
- 649
650 Jolly, A. L., D. Takawira, O. O. Oke, S. A. Whiteside, S. W. Chang, E. R. Wen, K. Quach, D. J. Evans and
651 S. M. J. Fleiszig (2015). "*Pseudomonas aeruginosa*-induced bleb-niche formation in epithelial cells is
652 independent of actinomyosin contraction and enhanced by loss of cystic fibrosis transmembrane-
653 conductance regulator osmoregulatory function." MBio 6(2).
- 654
655 Ke, Guolin, et al. "Lightgbm: A highly efficient gradient boosting decision tree." Advances in neural
information processing systems 30 (2017): 3146-3154.
- 656
657 Klockgether, J., N. Cramer, S. Fischer, L. Wiehlmann and B. Tümmler (2018). "Long-term microevolution
658 of *Pseudomonas aeruginosa* differs between mildly and severely affected cystic fibrosis lungs." American
journal of respiratory cell and molecular biology 59(2): 246-256.
- 659
660 Koboldt, D. C., Q. Zhang, D. E. Larson, D. Shen, M. D. McLellan, L. Lin, C. A. Miller, E. R. Mardis, L.
661 Ding and R. K. Wilson (2012). "VarScan 2: somatic mutation and copy number alteration discovery in
cancer by exome sequencing." Genome Research 22(3): 568-576.

662
663 Koehrsen, W. (2019). Feature Selector: Feature Selection in Python.
664 <https://github.com/WillKoehrsen/feature-selector>
665
666 Kosorok, M. R., L. Zeng, S. E. West, M. J. Rock, M. L. Splaingard, A. Laxova, C. G. Green, J. Collins and
667 P. M. J. P. p. Farrell (2001). "Acceleration of lung disease in children with cystic fibrosis after *Pseudomonas*
668 *aeruginosa* acquisition." *Pediatric Pulmonology* **32**(4): 277-287.
669
670 Kresse, A. U., S. D. Dinesh, K. Larbig and U. J. M. m. Römling (2003). "Impact of large chromosomal
671 inversions on the adaptation and evolution of *Pseudomonas aeruginosa* chronically colonizing cystic
672 fibrosis lungs." *Molecular Microbiology* **47**(1): 145-158.
673
674 Kuhn, M. and K. Johnson (2013). Over-Fitting and Model Tuning. *Applied Predictive Modeling*. New
675 York, NY, Springer New York: 61-92.
676
677 Kumru, B., Emiralioglu, N., & Ozel, H. G. (2018). Does body mass index affect lung function in patients
678 with cystic fibrosis?. *Clinical Nutrition*, 37, S91.
679
680 Le, T. T., W. Fu and J. H. Moore (2019). "Scaling tree-based automated machine learning to biomedical
681 big data with a feature set selector." *Bioinformatics* **36**(1): 250-256.
682
683 Lees, J. A., Mai, T. T., Galardini, M., Wheeler, N. E., Horsfield, S. T., Parkhill, J., & Corander, J. (2020).
684 Improved prediction of bacterial genotype-phenotype associations using interpretable pangenome-spanning
685 regressions. *Mbio*, 11(4)
686
687 Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis and R. Durbin
688 (2009). "The Sequence Alignment/Map format and SAMtools." *Bioinformatics* **25**(16): 2078-2079.
689
690 Lim YW, Evangelista III JS, Schmieder R, Bailey B, Haynes M, *et al*. Clinical insights from
691 metagenomic analysis of sputum samples from patients with cystic fibrosis. *Journal of clinical*
692 *microbiology* (2014). 52:425–437.
693
694 Macesic, N., Don't Walk, O. J. B., Pe'er, I., Tatonetti, N. P., Peleg, A. Y., & Uhlemann, A. C. (2020).
695 Predicting phenotypic polymyxin resistance in *Klebsiella pneumoniae* through machine learning analysis
696 of genomic data. *Msystems*, 5(3).
697
698 MacFadden, D. R., Melano, R. G., Coburn, B., Tijet, N., Hanage, W. P., & Daneman, N. (2019). Comparing
699 patient risk factor-, sequence type-, and resistance locus identification-based approaches for predicting
700 antibiotic resistance in *Escherichia coli* bloodstream infections. *Journal of clinical microbiology*, 57(6).
701
702 Mahé, P., & Tournoud, M. (2018). Predicting bacterial resistance from whole-genome sequences using k-
703 mers and stability selection. *BMC bioinformatics*, 19(1), 1-11.
704
705 Marvig, R. L., H. K. Johansen, S. Molin and L. Jelsbak (2013). "Genome analysis of a transmissible
706 lineage of *Pseudomonas aeruginosa* reveals pathoadaptive mutations and distinct evolutionary paths of
707 hypermutators." *PloS genetics* **9**(9): e1003741.
708
709 Marvig, R. L., Sommer, L. M., Molin, S., & Johansen, H. K. (2015). Convergent evolution and adaptation
710 of *Pseudomonas aeruginosa* within patients with cystic fibrosis. *Nature Genetics*, 47(1), 57.
711

712 Méric, G., L. Mageiros, J. Pensar, M. Laabei, K. Yahara, B. Pascoe, N. Kittiwon, P. Tadee, V. Post, S.
713 Lambie, R. Bowden, J. E. Bray, M. Morgenstern, K. A. Jolley, M. C. J. Maiden, E. J. Feil, X. Didelot, M.
714 Miragaia, H. de Lencastre, T. F. Moriarty, H. Rohde, R. Massey, D. Mack, J. Corander and S. K. Sheppard
715 (2018). "Disease-associated genotypes of the commensal skin bacterium *Staphylococcus epidermidis*."
716 Nature Communications 9(1): 5034.
717
718 Mobegi, F. M., Cremers, A. J., De Jonge, M. I., Bentley, S. D., Van Hijum, S. A., & Zomer, A. (2017).
719 Deciphering the distance to antibiotic resistance for the pneumococcus using genome sequencing data.
720 Scientific reports, 7(1), 1-13.
721
722 Mowat, E., S. Paterson, J. L. Fothergill, E. A. Wright, M. J. Ledson, M. J. Walshaw, M. A. Brockhurst and
723 C. Winstanley (2011). "Pseudomonas aeruginosa population diversity and turnover in cystic fibrosis
724 chronic infections." American journal of respiratory and critical care medicine 183(12): 1674-1679.
725
726 Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,
727 R. Weiss and V. Dubourg (2011). "Scikit-learn: Machine learning in Python." The Journal of machine
728 Learning research 12: 2825-2830.
729
730 Pincus, N. B., E. A. Ozer, J. P. Allen, M. Nguyen, J. J. Davis, D. R. Winter, C.-H. Chuang, C.-H. Chiu, L.
731 Zamorano and A. Oliver (2020). "A genome-based model to predict the virulence of *Pseudomonas*
732 *aeruginosa* isolates." Mbio 11(4).
733
734 Poudyal, B., & Sauer, K. (2018). The ABC of biofilm drug tolerance: the MerR-like regulator BrlR is an
735 activator of ABC transport systems, with PA1874-77 contributing to the tolerance of *Pseudomonas*
736 *aeruginosa* biofilms to tobramycin. Antimicrobial agents and chemotherapy, 62(2).
737
738 Raschka, S. (2018). MLxtend: Providing machine learning and data science utilities and extensions to
739 Python's scientific computing stack. Journal of open source software, 3(24), 638.
740
741 Recker, M., M. Laabei, M. S. Toleman, S. Reuter, R. B. Saunderson, B. Blane, M. E. Török, K. Ouadi, E.
742 Stevens and M. Yokoyama (2017). "Clonal differences in *Staphylococcus aureus* bacteraemia-associated
743 mortality." Nature microbiology 2(10): 1381-1388.
744
745 Saber, M. M., & Shapiro, B. J. (2020). Benchmarking bacterial genome-wide association study methods
746 using simulated genomes and phenotypes. Microbial genomics, 6(3).
747
748 Shanthikumar, S., M. N. Neeland, R. Saffery and S. Ranganathan (2019). "Gene modifiers of cystic
749 fibrosis lung disease: a systematic review." Pediatric pulmonology 54(9): 1356-1366.
750
751 Smith, E. E., D. G. Buckley, Z. Wu, C. Saenphimmachak, L. R. Hoffman, D. A. D'Argenio, S. I. Miller, B.
752 W. Ramsey, D. P. Speert and S. M. Moskowitz (2006). "Genetic adaptation by *Pseudomonas aeruginosa* to
753 the airways of cystic fibrosis patients." Proceedings of the National Academy of Sciences 103(22): 8487-
754 8492.
755
756 Snell, G. I., Bennetts, K., Bartolo, J., Levvey, B., Griffiths, A., Williams, T., & Rabinov, M. (1998). Body
757 mass index as a predictor of survival in adults with cystic fibrosis referred for lung transplantation. The
758 Journal of heart and lung transplantation: the official publication of the International Society for Heart
759 Transplantation, 17(11), 1097-1103.
760

- 761 Somayaji, R., J. C. Lam, M. G. Surette, B. Waddell, H. R. Rabin, C. D. Sibley, S. Purighalla and M. D.
762 Parkins (2017). "Long-term clinical outcomes of 'Prairie Epidemic Strain' *Pseudomonas aeruginosa*
763 infection in adults with cystic fibrosis." Thorax 72(4): 333-339.
764
- 765 Stressmann, F. A., G. B. Rogers, C. J. Van Der Gast, P. Marsh, L. S. Vermeer, M. P. Carroll, L. Hoffman,
766 T. W. V. Daniels, N. Patel and B. Forbes (2012). "Long-term cultivation-independent microbial diversity
767 analysis demonstrates that bacterial communities infecting the adult cystic fibrosis lung show stability and
768 resilience." Thorax 67(10): 867-873.
769
- 770 Svetnik, V., Liaw, A., Tong, C., Culberson, J. C., Sheridan, R. P., & Feuston, B. P. (2003). Random forest:
771 a classification and regression tool for compound classification and QSAR modeling. Journal of chemical
772 information and computer sciences, 43(6), 1947-1958.
- 773 Turcios, N. L. (2020). Cystic fibrosis lung disease: An overview. Respiratory care, 65(2), 233-251.
- 774 Tümmler, B. (2006). Clonal variations in *Pseudomonas aeruginosa*. In *Pseudomonas* (pp. 35-68). Springer,
775 Boston, MA.
776
- 777 Van Der Gast, C. J., A. W. Walker, F. A. Stressmann, G. B. Rogers, P. Scott, T. W. Daniels, M. P. Carroll,
778 J. Parkhill and K. D. Bruce (2011). "Partitioning core and satellite taxa from within cystic fibrosis lung
779 bacterial communities." The ISME journal 5(5): 780-791.
780
- 781 Virtanen, P., R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P.
782 Peterson, W. Weckesser and J. Bright (2020). "SciPy 1.0: fundamental algorithms for scientific computing
783 in Python." Nature methods 17(3): 261-272.
784
- 785 Waskom, M. L. (2021). Seaborn: statistical data visualization. Journal of Open Source Software, 6(60),
786 3021.
787
- 788 Welsh MJ, Ramsey BW, Accurso F, Cutting GR. Cystic fibrosis. In: Scriver CR, Beaudet AL, Sly WS,
789 Valle D, eds. *Metabolic & molecular bases of inherited disease*. 8th ed. Vol. 3. New York: McGraw-Hill,
790 2001:5121-5188.
791
- 792 Williams, D., B. Evans, S. Haldenby, M. J. Walshaw, M. A. Brockhurst, C. Winstanley and S. Paterson
793 (2015). "Divergent, coexisting *Pseudomonas aeruginosa* lineages in chronic cystic fibrosis lung infections."
794 American journal of respiratory and critical care medicine 191(7): 775-785.
795
- 796 Zhang, L., & Mah, T. F. (2008). Involvement of a novel efflux system in biofilm-specific resistance to
797 antibiotics. Journal of bacteriology, 190(13), 4447-4452.
- 798 Zhao, C. Y., Y. Hao, Y. Wang, J. J. Varga, A. A. Stecenko, J. B. Goldberg and S. P. Brown (2020).
799 "Microbiome data enhances predictive models of lung function in people with cystic fibrosis." The Journal
800 of Infectious Diseases. 2020; jiaa655
- 801 Zhao, J., P. D. Schloss, L. M. Kalikin, L. A. Carmody, B. K. Foster, J. F. Petrosino, J. D. Cavalcoli, D. R.
802 VanDevanter, S. Murray and J. Z. Li (2012). "Decade-long bacterial community dynamics in cystic fibrosis
803 airways." Proceedings of the National Academy of Sciences 109(15): 5809-5814.

804