

1 TRACKING INFLUENZA A VIRUS INFECTION IN THE LUNG 2 FROM HEMATOLOGICAL DATA WITH MACHINE LEARNING

3 Suneet Singh Jhutti^{1,2,#}, Julia D. Boehme^{3,4,#}, Andreas Jeron^{3,4}, Julia Volckmar^{3,4}, Kristin
4 Schultz^{3,5}, Jens Schreiber⁶, Klaus Schughart^{5, 7,8}, Kai Zhou¹, Jan Steinheimer¹, Horst Stöcker^{1,9,10},
5 Sabine Stegemann-Koniszewski⁶, Dunja Bruder^{4,3*}, Esteban A. Hernandez-Vargas^{1,11*}

- 6 1) Frankfurt Institute for Advanced Studies, 60438, Frankfurt am Main, Germany
7 2) Faculty of Biological Sciences, Goethe University, 60438 Frankfurt am Main, Germany
8 3) Immune Regulation Group, Helmholtz Centre for Infection Research, 38124 Braunschweig,
9 Germany
10 4) Infection Immunology Group, Institute of Medical Microbiology, Infection Control and
11 Prevention, Health Campus Immunology, Infectiology and Inflammation, Otto-von-Guericke-
12 University Magdeburg, 39120 Magdeburg, Germany
13 5) Department of Infection Genetics, Helmholtz Centre for Infection Research, 38124 Braunschweig,
14 Germany
15 6) Department of Pneumology, Health Campus Immunology, Infectiology and Inflammation, Otto-
16 von-Guericke University Magdeburg, 39120 Magdeburg, Germany
17 7) Department of Microbiology, Immunology, and Biochemistry, University of Tennessee Health
18 Science Center, Memphis, TN 38163, United States
19 8) University of Veterinary Medicine Hannover, 30559 Hannover, Germany
20 9) Institut für Theoretische Physik, Goethe Universität Frankfurt, 60438 Frankfurt am Main, Germany
21 10) GSI Helmholtzzentrum für Schwerionenforschung GmbH, 64291 Darmstadt, Germany
22 11) Instituto de Matemáticas, Universidad Nacional Autónoma de México, 76230, Juriquilla, México

23

24

25 Contributions

26

27 # These authors contributed equally: Suneet Singh Jhutti, Julia D. Boehme

28 * These authors jointly supervised this work: Dunja Bruder, Esteban A. Hernandez-Vargas

29 Correspondence to: dunja.bruder@med.ovgu.de, esteban@im.unam.mx

30 **ABSTRACT**

31

32 The tracking of pathogen burden and host responses with minimal-invasive methods during
33 respiratory infections is central for monitoring disease development and guiding treatment
34 decisions. Utilizing a standardized murine model of respiratory Influenza A virus (IAV) infection,
35 we developed and tested different supervised machine learning models to predict viral burden and
36 immune response markers, *i.e.* cytokines and leukocytes in the lung, from hematological data. We
37 performed independently *in vivo* infection experiments to acquire extensive data for training and
38 testing purposes of the models. We show here that lung viral load, neutrophil counts, cytokines
39 like IFN- γ and IL-6, and other lung infection markers can be predicted from hematological data.
40 Furthermore, feature analysis of the models shows that blood granulocytes and platelets play a
41 crucial role in prediction and are highly involved in the immune response against IAV. The
42 proposed *in silico* tools pave the path towards improved tracking and monitoring of influenza
43 infections and possibly other respiratory infections based on minimal-invasively obtained
44 hematological parameters.

45

46

47

48

49

50

51

52

53 INTRODUCTION

54 Respiratory infections by influenza (flu) viruses cause 3 to 5 million cases of severe illness every
55 year¹. Influenza A virus (IAV) is especially severe among high-risk groups like the elderly, infants,
56 pregnant women, and immunocompromised people². Next to its high prevalence in annual
57 epidemics, IAV has led to high mortality during several pandemics including the Spanish flu in
58 1918 and more recently the swine flu in 2009^{3,4}. Generally, the outcome of flu disease highly
59 depends on viral factors as well as host immunity. Accordingly, a fatal course of infection can
60 result from either insufficient control of viral spread, hyperinflammation, and/or a secondary
61 bacterial infection⁵. Thus, tracking of viral burden as well as host responses in the lungs is
62 important for monitoring IAV pathogenesis and tailoring targeted therapies.

63 Methodologically, diagnosis and tracking of acute IAV infection can be performed by assessing
64 viral antigen, nucleic acid, or infectious particles from upper or lower airway lavages, aspirates, or
65 swabs. Likewise, monitoring of lower airway immune responses is accomplished by *e.g.*
66 quantification of inflammatory cytokines or leukocytes in bronchoalveolar lavage fluid (BALF).
67 Next to several obvious disadvantages of these methods (low sensitivity, costly and time-
68 consuming analyses, and/or high technical requirements), the biggest hurdle lies in the invasive
69 sampling procedure that poses a risk to the acutely infected patient⁶. Accordingly, the development
70 of non- or minimally invasive approaches that allow observation of the disease status during IAV
71 infection remains an unmet medical need.

72 Besides inducing acute inflammatory responses in the airways, influenza infection results in
73 peripheral immune activation manifesting in altered blood cell composition⁷, transcriptional
74 signatures, and cytokine and chemokine levels in mice and humans⁸⁻¹⁰. While the severity and
75 longitudinal analyses revealed distinct molecular and/or cellular characteristics in peripheral blood

76 of IAV infected hosts, the suitability of each of these markers (and blood parameters in general)
77 for predicting the disease status is still unknown.

78 Here, we propose for the first time a framework intertwining *in vivo* experiments and machine
79 learning methods to forecast IAV infection parameters in the lung from blood sample data that can
80 be accessed minimal-invasively. To this end, we employed different machine-learning models for
81 blood-lung mapping. Experimentally, we utilized an established mouse model of sublethal
82 respiratory IAV infection^{11,12} and simultaneously assessed the kinetics of pathogen burden, lung
83 inflammation, as well as systemic cellular changes following infection. Ultimately, several
84 independent *in vivo* experiments were used to validate the applicability of the proposed framework.
85 Our primary computational approach was deep learning, which represents a class of machine
86 learning algorithms that uses multiple layers of information processing for feature extraction and
87 pattern analysis^{13,14}. These methods have already been successfully applied in several biological
88 fields including the prediction of transcriptional enhancers¹⁵, protein secondary structure¹⁴, and the
89 pathogenicity of genetic variants¹⁷. The image recognition abilities of machine learning algorithms
90 have already been tested for their diagnostic value and show promising results^{18,19}.

91

92

93

94

95

96

97 **RESULTS**

98 Our methodology consisted of three consecutive stages: *in vivo* experimental data acquisition,
99 model training, and independent experimental model validation (Figure 1). For the acquisition of
100 experimental data, mice were infected with a sublethal dose of PR8 (H1N1) followed by the
101 quantification of lung viral load, pulmonary innate and adaptive leukocyte subsets, pulmonary
102 cytokine levels, and hematological parameters over a total period of 11 days. Blood and lung
103 parameters were measured (Figure 1) in two independent experiments, mice ($n = 4-6$ per time
104 point) were sacrificed on days 1, 2, 3, 4, 5, 7, 9, and 11 post-infection (pi). These experimental
105 data built the basis for model training using different machine learning approaches to identify the
106 relationship between hematological and pulmonary parameters and to train and optimize the model
107 accordingly. To validate the predictive value of our model, we performed two additional
108 independent infection experiments with mice sacrificed on days 2, 4, 6, 9, and 11 pi and used the
109 mathematical models to predict the lung viral burden, leukocyte composition, and cytokine levels,
110 respectively, based on experimental hematological parameters.

111 **Blood and Lung Data Analysis.** In a first step, we conducted a correlation analysis to uncover
112 potential linear relationships between selected hematological and pulmonary parameters (Figure 2
113 and Supplementary Figures S10-12). We found a very strong correlation (Pearson correlation
114 coefficient > 0.9) between blood leukocytes and lymphocytes, which can be attributed to the fact
115 that lymphocytes constitute the largest leukocyte fraction in mice²⁰. Likewise, a strong correlation
116 observed between hematocrit as well as hemoglobin and erythrocyte numbers (Pearson correlation
117 coefficient > 0.8) can be attributed to the fact that hematocrit, as well as hemoglobin, are red blood
118 cell-associated parameters²¹.

119 Strikingly, we did not observe any strong correlation between blood and lung compartment
120 variables. However, there were strong correlations within cells of the lung compartment, such as
121 CD4⁺ T cells and CD8⁺ T cells as well as neutrophils and natural killer (NK) cells. As a linear
122 correlation does not include higher-order nonlinear temporal relations, we next employed machine
123 learning algorithms and compared the obtained results with the linear regression model.

124 **Tracking IAV Infection in the Lungs from Blood-Derived Parameters.** To predict influenza
125 virus levels and immunological markers in the lungs from hematological parameters, we employed
126 feedforward neural networks, gradient boosted regression trees, and a linear regression model.
127 These algorithms considered 14 hematological parameters (see Table 1) as features to predict the
128 respective target lung variable. The main scores used for comparing and evaluating different
129 machine learning algorithms were based on the average squared difference between the estimated
130 values and the actual value (mean squared error). The proportion of the variance in the dependent
131 variable that is predictable from the independent variables was based on the R² score. Overfitting
132 was reduced with regularization techniques presented in methods.

133 Figure 3A illustrates the best model prediction of viral levels in the lung, *i.e.* the feedforward
134 neural network, based on experimental haematological data. We observed that the model had a
135 good qualitative behavior of the lung viral load (measured as copy numbers of NP transcripts) over
136 the experimental time-frame of eleven days. Day 0 represents the control group. The prediction
137 seemed to be most accurate around the peak of viral replication, *i.e.* days 4 and 5 pi. Quantitatively,
138 there was a high variation in the predictions. This was attributed to the variation found in our
139 infection experiments (Figure S1-S9), which is a common observation *in vitro* and *in vivo* viral
140 infection experiments^{22–27}. While for some animals there was a large difference between the actual

141 experimental data and the respective predictions in the testing experiments, overall, the qualitative
142 performance on the testing set was good (Figure 3B).

143 In addition to predicting the lung viral load, we tested several machine learning algorithms to
144 predict target lung leukocytes and cytokines from the hematological parameters (Table 2). The
145 comparison in Figure 4A shows how the respective best model performed for different targets (for
146 all results please refer to the supplemental material, Supplementary Figures S13-15). As a
147 benchmark, we used the mean for each target variable calculated from the training data. In almost
148 all cases, the best model performed better than the benchmark. A positive R^2 score demonstrates
149 the explanatory power of the model. Predictions for lung IFN- γ , viral load, IL-6, and neutrophils
150 were able to outperform the benchmark (Figure 4). Predictions for other lung target parameters are
151 presented in Supplementary Figure S13-S15. Notably, the accuracy of the model predictions was
152 dependent on the stage of the infection. For example, neutrophil numbers in the lungs were better
153 predicted in the later days of infection, while IL-6 and IFN- γ levels were more accurately predicted
154 at the peak of infection (Figure 4B-D). This was also observed for other immune cells such as
155 CD4⁺ and CD8⁺ T cells (Figure S13). Table 2 presents a summary of the best models for predicting
156 the different lung target parameters.

157 To determine the role of each feature from the hematological data for the prediction of lung
158 outcomes, we performed a feature importance analysis. For this, we calculated the permutation
159 importance by swapping out features and evaluating the performance of our testing data. The
160 results of the feature importance analysis are organized from top to bottom in the level of
161 importance in Figure 5. For example, from all hematological parameters, granulocytes showed the
162 greatest impact on the performance of the machine learning models for predicting the viral load,

163 neutrophils, IFN- γ , and IL-6. Interestingly our analyses revealed a pivotal role of blood platelets
164 for predicting both pathogen burden and lung inflammatory milieu (Figure 5). Granulocytes and
165 erythrocytes ranked second and third place, respectively. The feature importance plots of the
166 additional lung target immune cells and cytokines can be found in the supplemental material
167 (Figure S16).

168

169

170

171

172

173

174

175

176

177

178

179 **DISCUSSION**

180 Mathematical modeling of host immune responses has largely contributed to improving our
181 understanding of the overall course of influenza infection²⁸⁻³⁷ as well as the personalization of
182 therapies and vaccines³⁸⁻⁴⁰. Mathematical models consist of systems of ordinary differential
183 equations describing the viral dynamics within the host. However, computational tools for the
184 diagnosis and tracking of respiratory diseases remain a public health challenge.

185 Here, we progress from the state of the art showing for the first time that minimal-invasively
186 acquired haematological parameters can be used to infer lung viral burden, leukocytes, and
187 cytokines following IAV infection in mice. Nevertheless, despite standardized experimental
188 procedures, our analysis showed a large variance in the computational predictions. These can be
189 attributed to the relatively high variances of our experimental data due to biological or
190 experimental variations. For instance, we found differences in some hematological parameters
191 between the training and the testing experiments, which possibly explain the differences in
192 performance between training and testing prediction of the lung viral load for day 2 post-infection.

193 The clinical potential of the framework proposed here consists of a new qualitative vision of the
194 disease processes in the lung compartment. We show that the accumulation and decline of multiple
195 cell types involved in the anti-viral immune response in the lung can successfully be predicted
196 with data derived from peripheral blood analyses. The boosted regression tree, with some
197 modifications, provided the best results for many of the lung immune target cells. On the other
198 hand, some target variables proved to be difficult to predict from hematological data. For instance,
199 alveolar macrophage (AM) numbers could not be predicted with any of the tested algorithms and
200 showed the worst score. This is likely the result of the weak correlations of AMs with the different

201 blood cells analyzed. AMs also demonstrate a quite different behavior than the other cells in terms
202 of their abundance over the course of infection as their number peaked rather late during the
203 infection, *i.e.* when the virus was mostly cleared⁴¹. This phenomenon is most likely the result of
204 inflation of the alveolar macrophage pool by self-renewal⁴² and/or monocyte recruitment
205 processes⁴³. Regarding the less accurate predictions of cytokines within the airways from blood
206 parameters (Supplementary Figure S13), a possible explanation is that a large portion of these
207 cytokines (especially during the early infection stage) originated from lung-resident leukocytes as
208 well as non-leukocytes. Therefore, their temporal quantity and composition are largely determined
209 by local constituents of the lung's immune cell response.

210 Our results show an active reaction chain between peripheral blood parameters and immune cells
211 in the lungs of the mice following IAV infection. Interestingly, we found that peripheral blood
212 platelets play an important role in predicting lung immune cell numbers in IAV infection. In line
213 with our finding of increased numbers of platelets in the blood during the acute phase of IAV
214 infection (see Supplementary Figure S1), platelet accumulation in the pulmonary capillaries is a
215 hallmark of murine IAV H1N1 infection⁴⁴ and contributes to pathogenesis⁴⁵. Importantly, platelet-
216 derived cytokines such as IL-1 β can directly increase endothelial permeability and the expression
217 of important vascular adhesion molecules^{46,47}. In line with this, airway IL-1 β levels were elevated
218 during acute IAV infection (see Figures S8-9). Increased platelet-mediated transendothelial
219 migration of CD4⁺ T cells, CD8⁺ T cells, NK cells, and neutrophils could thus be one conceivable
220 mechanism contributing to the observed strong positive correlation between peripheral blood
221 platelets and the aforementioned lung leukocyte subsets. This is relevant to the anti-viral host
222 response, as the CD4⁺ T cells are central in the activation and maturation of virus-specific CD8⁺
223 T cells⁴⁸, while neutrophils are required for proper NK cell maturation⁴⁹.

224 It should be noted, that although the estimation is called “prediction” in the machine learning
225 domain, and blood-derived data can here be used in a practical way to "predict" the viral load or
226 the number of certain lung immune cells or cytokines in IAV infection, we cannot establish a
227 direction of causality in this case. In other words, we cannot state *e.g.* that platelets are involved
228 in raising the total amount of CD8⁺ T cells or if CD8⁺ T cells drive the increase in platelets. What
229 we can learn from the correlations between variables is that CD8⁺ T cells, CD4⁺ T cells, NK cells,
230 and neutrophils have their strongest positive correlation with platelets between the blood cells
231 analyzed (see Figure 2). The feature importance analysis confirms that platelets play the most
232 important part in the estimation of lung CD4⁺ T and CD8⁺ T cells, followed by erythrocytes.
233 However, the viral load inside the lungs has also a strong correlation with platelets but an even
234 stronger one with granulocytes. The variable importance analysis suggests that platelets and
235 granulocytes do not strongly contribute to the prediction as can be seen in Figure 5.

236 The weak predictive results obtained with the linear regression model signifies that these relations
237 have a high order of complexity. While contributing to the viral clearance, the innate immune
238 system can also exacerbate the lung injury^{50,51}. In this context, tissue injury can be a cause of
239 platelet activation during influenza infection. The role of platelets in human influenza infection
240 has been stressed in recent years^{45,52,53}. Thrombosis, controlled by the innate immune system has
241 been suggested to support immune defense⁵⁴.

242 Hematological parameters such as neutrophil, lymphocyte, and platelet counts, as well as the
243 neutrophil-to-lymphocyte ratio (NLR) have contributed to diagnosing influenza virus infections⁵².
244 Thus, we also addressed if the use of the granulocyte-lymphocyte ratio (GLR) or the platelet-
245 granulocyte ratio (PGR) improves importance in our model predictions (Table S2). We compared
246 the use of only GLR or PGR with the only use of lymphocytes and granulocytes, as well using

247 additional important IAV infection-associated peripheral blood parameters like erythrocytes and
248 hemoglobin. Performances were evaluated with the MSE and R^2 scores on the testing data set. We
249 also calculated the corrected Akaike Information Criterion (AIC_c) during the training to take the
250 complexity of the models into account. Supplementary Table S2 shows that the model performance
251 increased using GLR, while the use of PGR had only a minor effect. Also, adding erythrocytes and
252 hemoglobin did not improve predictions.

253 In summary, blood platelets, granulocytes, and erythrocytes play an important role in
254 understanding the immune response to influenza infection and can be used in conjunction with
255 other blood components for monitoring the lung viral load and lung immune cells in mice.
256 Importantly, our results indicate that a reduced number of variables does not affect
257 model/prediction accuracy. This can help to further reduce the hematology data needed for
258 successful prediction. While recent efforts show evidence for the diagnosis of COVID-19 from
259 blood compartment⁵⁶, further clinical evidence will be needed to show the potential of how our
260 procedure could be generalized to advance medical care.

261

262

263

264

265

266

267 **METHODS**

268 **Experimental Design.** Mice were intranasally infected with a sublethal dose of the mouse-
269 adapted, strictly pneumotropic H1N1 IAV strain PR/8/34 and the lung viral burden, pulmonary
270 innate and adaptive leukocyte subsets, pulmonary cytokine levels, and peripheral blood cell
271 parameters were assessed for 11 days pi (Figure 1).

272 Initial generation, training, and optimization of the computational algorithms were conducted
273 using data from two independent *in vivo* infection experiments. Here, mice were randomly
274 assigned to the respective experimental groups and were either intranasally inoculated with a
275 sublethal dose of IAV or saline (control groups). Mice ($n = 4-6$ /experimental group) were
276 sacrificed on days 1, 2, 3, 4, 5, 7, 9 and 11 post-infection. Experimental readouts for the first
277 experiment were: hematological parameters and lung tissue viral load. In the second experiment,
278 experimental readouts were: hematological parameters, leukocytes in lung tissue, and airway
279 cytokines.

280 For subsequent model validation, two additional, independent *in vivo* infection experiments were
281 performed using the above-mentioned readouts. In these experiments, mice ($n = 3$ /experimental
282 group) were sacrificed at days 2, 4, 6, 9, and 11 post-infection. Hematological parameters used for
283 model generation and evaluation are listed in Table 1.

284 **Mice.** For all experiments, female C57BL/6J01aHsd mice (age 10-12 weeks) from Envigo were
285 used. All mice were housed in the animal facility at the Helmholtz Centre for Infection Research
286 under specific-pathogen-free (SPF) conditions and in accordance with national and institutional
287 guidelines.

288 **Viral preparation and infection.** For viral infections, a mouse-adapted influenza A virus strain
289 (A/Puerto Rico/8/34, H1N1) was utilized. The virus was produced in Madin-Darby Canine Kidney

290 (MDCK) cells⁵⁷ and quantified by calculating the tissue culture infectious dose (TCID₅₀) as
291 previously described⁵⁸. Mice were anesthetized by intraperitoneal injection of ketamine/xylazine
292 and were infected with viral inoculum (0.31 TCID₅₀ in 25µL PBS). Control animals received PBS
293 only.

294 **Quantification of the lung viral load.** Lungs were perfused using PBS. RNA was extracted from
295 whole lung tissue homogenates using the RNeasy Plus Mini Kit (Qiagen). The absence of genomic
296 DNA in RNA samples was initially confirmed by PCR using a Taq DNA-polymerase and primers
297 for the housekeeping gene *Rps9*. Quantitative real-time RT-PCR (qPCR) for detection of viral
298 burden was performed using the SensiFAST™ SYBR® No-ROX One-Step Kit and an influenza
299 nucleoprotein (NP) plasmid standard. The sequences of the used primers were: 5' '
300 CTGGACGAGGGCAAGATGAAGC, 3' TGACGTTGGCGGATGAGCACA (*Rps9*) and 5' '
301 GAGGGGTGAGAATGGACGAAAAAC, 3' CAGGCAGGCAGGCAGGACTT (*Np*).

302 **Hematology analysis.** Blood samples were obtained from the retrobulbar plexus and EDTA was
303 added to prevent coagulation. Samples were analyzed using a VetScan® HM5 machine (Abaxis).

304 **Cytokine detection.** Bronchoalveolar lavage (BAL) was performed with 1mL PBS, samples were
305 spun down (420 x g, 10 min) and BAL fluid (BALF) supernatants were stored at -70°C until further
306 analyses. Cytokine levels in BALF samples were quantified using the LEGENDplex™ Mouse
307 Inflammation Panel (BioLegend) according to the manufacturer's protocol.

308 **Isolation of leukocytes from lung tissue.** Lungs were perfused using PBS, excised, mechanically
309 homogenized, and enzymatic digestion was performed in Iscove's modified Dulbecco's medium
310 (IMDM) supplemented with 0.2mg/mL collagenase D (Roche), 0.01mg/mL DNase I (Sigma-

311 Aldrich), and 5% fetal bovine serum at 37°C for 45min. Digestion was stopped by the addition of
312 EDTA, the cell suspension was filtered through a 100µm cell strainer and spun down (420 x g,
313 10min). Erythrocyte lysis was performed using ammonium-chloride-potassium (ACK) buffer and
314 leukocytes were isolated by gradient centrifugation using Percoll solution (GE Healthcare). Lung
315 tissue leukocytes were filtered again and antibody staining for flow cytometry was performed.

316 **Flow cytometry.** Lung tissue leukocyte samples were subjected to viability staining and blocking
317 of Fc-receptors using LIVE/DEAD™ Fixable Blue Dead cell stain kit (life technologies) and an
318 anti-CD16/32 antibody (clone 93, BioLegend). Cells were then washed and incubated with a
319 staining mix containing antibodies against the following murine antigens: Siglec-F (PE, clone E50-
320 2440, BD), Ly6G (AlexaFluor700, clone 1A8, BD), CD11c (APC, clone N 418, BioLegend),
321 CD11b (BV421, clone M1/70, BD), CD4 (APC-Fire750, clone GK1.5, BioLegend), CD8a
322 (CyChrome, clone 53-6.7, BD), CD3ε (Biotin, clone 145-2C11, BioLegend), NK1.1 (FITC, clone
323 PK136, BioLegend). Secondary staining was performed using streptavidin-BV605 and
324 streptavidin-BV650, respectively (BioLegend). All reagents and antibodies had been titrated
325 before the experiments for optimal staining results. Flow cytometry data were acquired using
326 LSRII and LSR Fortessa instruments (BD). Data were analyzed using FlowJo software (BD).

327 **Data Processing.** The data obtained from the hematological analysis constitutes of 20 different
328 parameters. Some parameters are given in absolute values and percentages. We considered from
329 these parameters only the absolute values, resulting in 14 different parameters (see Table 1). For
330 some mice, it was not possible to extract hematological data and/or according to target infection
331 marker. These mice were removed for the mapping, although their data was used in the correlation
332 analysis. Furthermore, if some values were lower than the measurable threshold, we used the
333 threshold value. Computational algorithms were implemented in python using the *Keras* and

334 *sklearn* libraries. Our study design yielded separate training and testing data sets. The testing data
335 was obtained approximately one year after the training data. All laboratory conditions were kept
336 as similar as possible.

337 Using the data directly for training the algorithms led to poor results and therefore we used data
338 pre-processing techniques. To conserve the nature of hematological parameter distributions, we
339 used the min-max scaling from the *sklearn.preprocessing.MinMaxScaler* class:

$$340 \quad z = \frac{x - x_{min}}{x_{max} - x_{min}}$$

341 For the target infection markers like viral load, any logarithmic function worked well. For
342 simplicity, we used \log_{10} .

343 **Machine Learning Models.** Different machine learning models were tested for the mapping
344 including feedforward neural networks (FNN), gradient boosted regression trees (GBRT), linear
345 regression (LR), support vector machines (SVM), and random forest regression (RFR). The
346 hyperparameters of the models were estimated via grid-search and adjusted via trial and error.
347 FNN and GBRT showed to superior in most cases and RFR was outperformed in every instance
348 with one of the other algorithms.

349 In many cases using PCA before the mapping yielded improved performance. It was found that a
350 dimensionality reduction to six input blood variables was often best. We used the class
351 *sklearn.decomposition.PCA* for implementation. For the feedforward neural network, the *keras*
352 library was used with a TensorFlow backend. We found that one hidden layer was sufficient most
353 of the time and additional layers were not needed. The number of weights varied from 10 to 50.
354 For regularization, the addition of dropout layers with a rate of 0.2 was helpful to prevent

355 overfitting. As an activation function, we used a rectified linear unit (ReLU). We had one output
356 that uses a linear activation function. This was necessary to map the whole range of possible
357 outcome values. We used the Adam optimizer and minimized the mean squared error to find the
358 optimal fit. The weights were initialized according to a He-uniform distribution.⁵⁹ Following
359 common practice in literature we used for the training of the neural network a validation set of
360 10% of the whole training data.

361 The GBRT, LR, SVM, and RFR algorithms are taken from the python library *sklearn*. The
362 hyperparameters of GBRT and RFR models were searched over a grid from 10 to 2000 estimators,
363 a learning rate from 0.001 to 0.09, and a max depth of 2 to 14. The least-square regression was
364 used for optimization. The kernels used for SVM were '*linear*', '*poly*', '*rbf*', '*sigmoid*' and
365 '*precomputed*'.

366 To determine which variables were the most important in our model predictions, we calculated the
367 permutation importance using the `sklearn.inspection.permutation_importance` implementation.
368 For this, we took our best model, respectively, and trained it on the training data set. After the
369 trained model was evaluated on the hold-out testing data set with the mean squared error as metric,
370 a feature column was permuted and the metric was evaluated again. This procedure was repeated
371 100 times and the permutation importance was given by the difference between the baseline metric
372 and the metric from permuted feature columns.

373 **DATA AND CODE AVAILABILITY**

374 The datasets generated and analyzed during the current study are available:

375 https://github.com/Jhutti/Tracking_IAV_from_Blood

376

377 REFERENCES

- 378
379 1. World Health Organization. Influenza (Seasonal) fact sheet.
- 380
381 2. WHO Recommended Surveillance Standards. Second edition.
- 382
383 3. Kilbourne, E. D. Influenza Pandemics of the 20th Century. *Emerging Infectious Diseases* **12**, 9
384 (2006).
- 385
386 4. World Health Organization. Writing Committee of the WHO Consultation on Clinical aspects of
387 pandemic 2009 influenza A (H1N1) virus infection. *New England Journal of Medicine* 1708–1719
388 (2010).
- 389 5. Sharma-Chawla, N. *et al.* In vivo neutralization of pro-inflammatory cytokines during secondary
390 streptococcus pneumoniae infection post influenza a virus infection. *Frontiers in Immunology* **10**,
391 (2019).
- 392 6. Allwinn, R. *et al.* Laboratory diagnosis of influenza – virology or serology? *Medical Microbiology
393 and Immunology 2002 191:3* **191**, 157–160 (2002).
- 394 7. L, D. *et al.* Cellular changes in blood indicate severe respiratory disease during influenza infections
395 in mice. *PloS one* **9**, (2014).
- 396 8. IE, G. *et al.* Untuned antiviral immunity in COVID-19 revealed by temporal type I/III interferon
397 patterns and flu comparison. *Nature immunology* **22**, 32–40 (2021).
- 398 9. Y, Z. *et al.* Pathway mapping of leukocyte transcriptome in influenza patients reveals distinct
399 pathogenic mechanisms associated with progression to severe infection. *BMC medical genomics*
400 **13**, (2020).
- 401 10. BM, C. *et al.* Inflammatory Monocytes Drive Influenza A Virus-Mediated Lung Injury in Juvenile
402 Mice. *Journal of immunology (Baltimore, Md. : 1950)* **200**, 2391–2404 (2018).
- 403 11. N, S.-C. *et al.* Influenza A Virus Infection Predisposes Hosts to Secondary Infection with Different
404 Streptococcus pneumoniae Serotypes with Similar Outcome but Serotype-Specific Manifestation.
405 *Infection and immunity* **84**, 3445–3457 (2016).
- 406 12. Duvigneau, S. *et al.* Hierarchical effects of pro-inflammatory cytokines on the post-influenza
407 susceptibility to pneumococcal coinfection. *Scientific Reports* **6**, 1–11 (2016).
- 408 13. Ian Goodfellow, Yoshua Bengio & Aaron Courville. Deep learning. *MIT press* (2016)
409 doi:10.1007/S10710-017-9314-Z.
- 410 14. DengLi & YuDong. Deep Learning. *Foundations and Trends in Signal Processing* **7**, 197–387 (2014).
- 411 15. Liu, F., Li, H., Ren, C., Bo, X. & Shu, W. PEDLA: predicting enhancers with a deep learning-based
algorithmic framework. *Scientific Reports 2016 6:1* **6**, 1–14 (2016).
16. Wang, S., Peng, J., Ma, J. & Xu, J. Protein Secondary Structure Prediction Using Deep Convolutional
Neural Fields. *Scientific Reports 2016 6:1* **6**, 1–11 (2016).

- 412 17.Quang, D., Chen, Y. & Xie, X. DANN: a deep learning approach for annotating the pathogenicity of
413 genetic variants. *Bioinformatics* **31**, 761–763 (2015).
- 414 18.Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural networks.
415 *Nature* 2017 542:7639 **542**, 115–118 (2017).
- 416 19.Huang, Q., Zhang, F. & Li, X. Machine Learning in Ultrasound Computer-Aided Diagnostic Systems:
417 A Survey. *BioMed Research International* **2018**, (2018).
- 418 20.Kabak, M., Çil, B. & Hocanlı, I. Relationship between leukocyte, neutrophil, lymphocyte, platelet
419 counts, and neutrophil to lymphocyte ratio and polymerase chain reaction positivity. *International*
420 *Immunopharmacology* **93**, 107390 (2021).
- 421 21.Han, Q. *et al.* Role of hematological parameters in the diagnosis of influenza virus infection in
422 patients with respiratory tract infection symptoms. *Journal of Clinical Laboratory Analysis* **34**,
423 e23191 (2020).
- 424 22.Hernandez-Vargas, E. A. *et al.* Effects of Aging on Influenza Virus Infection Dynamics. *Journal of*
425 *Virology* **88**, 4123–4131 (2014).
- 426 23.Hernandez-Vargas, E. A. & Velasco-Hernandez, J. X. In-host Mathematical Modelling of COVID-19
427 in Humans. *Annual Reviews in Control* **50**, 448–456 (2020).
- 428 24.Pawelek, K. A. *et al.* Modeling Within-Host Dynamics of Influenza Virus Infection Including Immune
429 Responses. *PLOS Computational Biology* **8**, e1002588 (2012).
- 430 25.Dobrovoly, H. M., Gieschke, R., Davies, B. E., Jumbe, N. L. & Beauchemin, C. A. A. Neuraminidase
431 inhibitors for treatment of human and avian strain influenza: A comparative modeling study.
432 *Journal of Theoretical Biology* **269**, 234–244 (2011).
- 433 26.Baccam, P., Beauchemin, C., Macken, C. A., Hayden, F. G. & Perelson, A. S. Kinetics of Influenza A
434 Virus Infection in Humans. *Journal of Virology* **80**, 7590–7599 (2006).
- 435 27.Smith, A. M. *et al.* Kinetics of Coinfection with Influenza A Virus and *Streptococcus pneumoniae*.
436 *PLOS Pathogens* **9**, e1003238 (2013).
- 437 28.Miao, H., Xia, X., Perelson, A. S. & Wu, H. On Identifiability of Nonlinear ODE Models and
438 Applications in Viral Dynamics. *SIAM Review* **53**, 3–39 (2011).
- 439 29.Canini, L. & Perelson, A. S. Viral kinetic modeling: state of the art. *Journal of Pharmacokinetics and*
440 *Pharmacodynamics* **41**, 431–443 (2014).
- 441 30.Canini, L. & Carrat, F. Population modeling of influenza A/H1N1 virus kinetics and symptom
442 dynamics. *Journal of Virology* **85**, 2764–2770 (2011).
- 443 31.Handel, A. & Antia, R. A simple mathematical model helps to explain the immunodominance of
444 CD8 T cells in influenza A virus infections. *Journal of virology* **82**, 7768–72 (2008).
- 445 32.Beauchemin, C. & Handel, A. A review of mathematical models of influenza A infections within a
446 host or cell culture: lessons learned and challenges ahead. *BMC public health* **11**, S7 (2011).

- 447 33.Hancioglu, B., Swigon, D. & Clermont, G. A dynamical model of human immune response to
448 influenza A virus infection. *Journal of Theoretical Biology* **246**, 70–86 (2007).
- 449 34.Smith, A. M. Host-pathogen kinetics during influenza infection and coinfection: insights from
450 predictive modeling. *Immunological Reviews* vol. 285 97–112 (2018).
- 451 35.Smith, A. M. & Perelson, A. S. Influenza A virus infection kinetics: quantitative data and models.
452 *Wiley interdisciplinary reviews. Systems biology and medicine* **3**, 429–445 (2011).
- 453 36.Baccam, P., Beauchemin, C., Macken, C. a, Hayden, F. G. & Perelson, A. S. Kinetics of influenza A
454 virus infection in humans. *Journal of virology* **80**, 7590–9 (2006).
- 455 37.Harper, S. A. *et al.* Seasonal Influenza in Adults and Children—Diagnosis, Treatment,
456 Chemoprophylaxis, and Institutional Outbreak Management: Clinical Practice Guidelines of the
457 Infectious Diseases Society of America. *Clinical Infectious Diseases* **48**, 1003–1032 (2009).
- 458 38.Hernandez-Mejia, G. & Hernandez-Vargas, E. A. Uncovering antibody cross-reaction dynamics in
459 influenza A infections. *bioRxiv* 2020.01.06.896274 (2020) doi:10.1101/2020.01.06.896274.
- 460 39.Hernandez-Mejia, G., Alanis, A. Y. & Hernandez-Vargas, E. A. Inverse Optimal Impulsive Control
461 Based Treatment of Influenza Infection. in *IFAC World Congress 2017* vol. 50 12696–12701 (2017).
- 462 40.Parra-Rojas, C., Messling, V. & Hernandez-Vargas, E. A. Adjuvanted influenza vaccine dynamics.
463 *Scientific Reports* **9**, (2019).
- 464 41.Toapanta, F. R. & Ross, T. M. Impaired immune responses in the lungs of aged mice following
465 influenza infection. *Respiratory Research* **10**, 1–19 (2009).
- 466 42.Yao, Y. *et al.* Induction of Autonomous Memory Alveolar Macrophages Requires T Cell Help and Is
467 Critical to Trained Immunity. *Cell* **175**, 1634-1650.e17 (2018).
- 468 43.Aegerter, H. *et al.* Influenza-induced monocyte-derived alveolar macrophages confer prolonged
469 antibacterial protection. *Nature immunology* **21**, 145–157 (2020).
- 470 44.Rommel, M. G. E., Milde, C., Eberle, R., Schulze, H. & Modlich, U. Endothelial–platelet interactions
471 in influenza-induced pneumonia: A potential therapeutic target. *Anatomia, Histologia,*
472 *Embryologia* **49**, 606–619 (2020).
- 473 45.Lê, V. B. *et al.* Platelet Activation and Aggregation Promote Lung Inflammation and Influenza Virus
474 Pathogenesis. <https://doi.org/10.1164/rccm.201406-1031OC> **191**, 804–819 (2015).
- 475 46.Rossaint, J., Margraf, A. & Zarbock, A. Role of Platelets in Leukocyte Recruitment and Resolution
476 of Inflammation. *Frontiers in Immunology* **0**, 2712 (2018).
- 477 47.Hawrylowicz, C. M., Howells, G. L. & Feldmann, M. Platelet-derived interleukin 1 induces human
478 endothelial adhesion molecule expression and cytokine production. *Journal of Experimental*
479 *Medicine* **174**, 785–790 (1991).
- 480 48.Zens, K. D. & Farber, D. L. Memory CD4 T Cells in Influenza. *Current Topics in Microbiology and*
481 *Immunology* **386**, 399–421 (2014).

- 482 49.Jaeger, B. N. *et al.* Neutrophil depletion impairs natural killer cell maturation, function, and
483 homeostasis. *Journal of Experimental Medicine* **209**, 565–580 (2012).
- 484 50.Herold, S., Becker, C., Ridge, K. M. & Budinger, G. R. S. Influenza virus-induced lung injury:
485 pathogenesis and implications for treatment. *European Respiratory Journal* **45**, 1463–1478 (2015).
- 486 51.Kuiken, T., Riteau, B., Fouchier, R. A. M. & Rimmelzwaan, G. F. Pathogenesis of influenza virus
487 infections: the good, the bad and the ugly. *Current Opinion in Virology* **2**, 276–286 (2012).
- 488 52.Koupenova, M. *et al.* The role of platelets in mediating a response to human influenza infection.
489 *Nature Communications* 2019 10:1 **10**, 1–18 (2019).
- 490 53.Assinger, A. Platelets and Infection – An Emerging Role of Platelets in Viral Infection. *Frontiers in*
491 *Immunology* **0**, 649 (2014).
- 492 54.Engelmann, B. & Massberg, S. Thrombosis as an intravascular effector of innate immunity. *Nature*
493 *Reviews Immunology* 2012 13:1 **13**, 34–45 (2012).
- 494 55.Han, Q. *et al.* Role of hematological parameters in the diagnosis of influenza virus infection in
495 patients with respiratory tract infection symptoms. *Journal of Clinical Laboratory Analysis* **34**,
496 e23191 (2020).
- 497 56.Kukar, M. *et al.* COVID-19 diagnosis by routine blood tests using machine learning. *Scientific*
498 *Reports* 2021 11:1 **11**, 1–9 (2021).
- 499 57.Stegemann, S. *et al.* Increased Susceptibility for Superinfection with *Streptococcus pneumoniae*
500 during Influenza Virus Infection Is Not Caused by TLR7-Mediated Lymphopenia. (2009)
501 doi:10.1371/journal.pone.0004840.
- 502 58.REED, L. J. & MUENCH, H. A SIMPLE METHOD OF ESTIMATING FIFTY PER CENT ENDPOINTS.
503 *American Journal of Epidemiology* **27**, 493–497 (1938).
- 504 59.He, K., Zhang, X., Ren, S. & Sun, J. Delving Deep into Rectifiers: Surpassing Human-Level
505 Performance on ImageNet Classification.
- 506
- 507
- 508
- 509
- 510
- 511
- 512
- 513
- 514

515 **ACKNOWLEDGEMENTS**

516 We thank Tatjana Hirsch, Hanna Shkarlet and Karin Lammert for expert technical assistance in
517 infection experiments. This work was supported by the Deutsche Forschungsgemeinschaft with
518 the project HE7707/5-1 and BR2221/6-1; the Universidad Nacional Autonoma de Mexico
519 (UNAM) – PAPIIT with the number IA102521; and the Alfons und Gertrud Kassel-Stiftung.

520

521 **ETHICS DECLARATIONS**

522 All the experiments were approved and conducted in accordance with the guidelines set by the
523 local animal welfare and ethics committee (Niedersächsisches Landesamt für Verbraucherschutz
524 und Lebensmittelsicherheit).

525

526 **COMPETING INTERESTS**

527 The authors declare no competing interests.

528

529

530

531

532

533

534

535

536

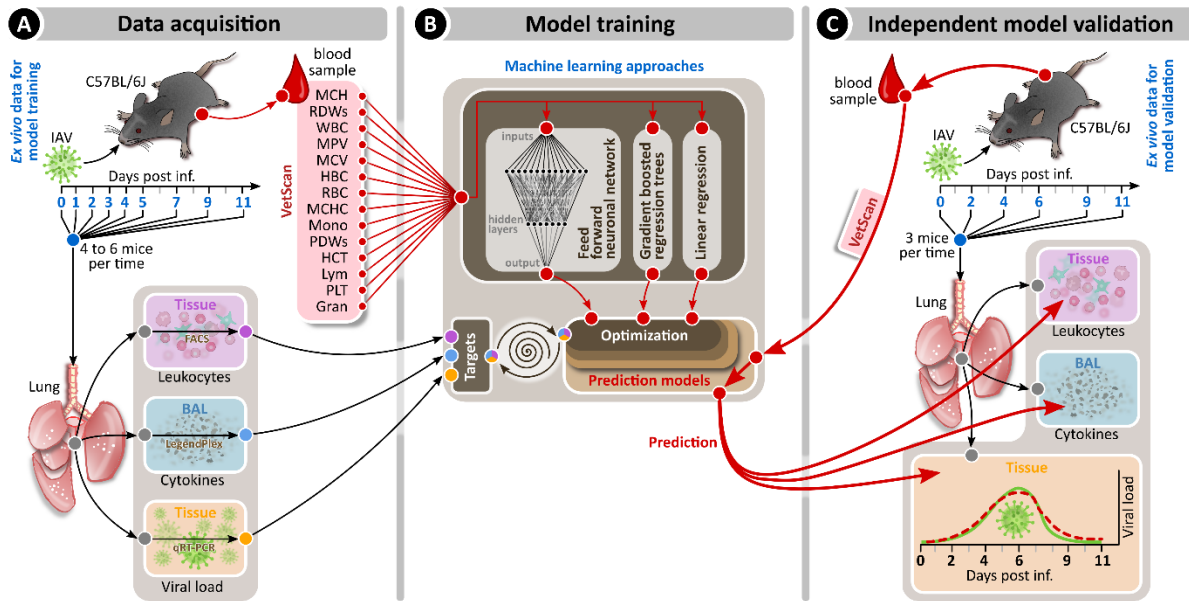
537

538

539

540 **FIGURES**

541



542

543

544 **Figure 1: Experimental scheme for the machine learning approaches of the respiratory IAV infection.**

545 Mice were intranasally infected with a sublethal dose of IAV PR/8/34 on day 0 and sacrificed on the
546 indicated days. Blood was collected for hematology analyses (Supplementary Figures S1-2),
547 bronchoalveolar lavage was performed to analyze lung cytokines and lung tissue samples were used to
548 monitor viral load (experiment 1, Supplementary Figure S5) or pulmonary leukocyte subsets
549 (experiment 2, Supplementary Figures S6, S8) (A). The hematological data from this initial set of
550 experiments were used to build and train different machine learning models (B). Data from a separate
551 experiment (Supplementary Figures S3-5, S7, S9) were used for testing and evaluation of machine learning
552 algorithms (C).

553

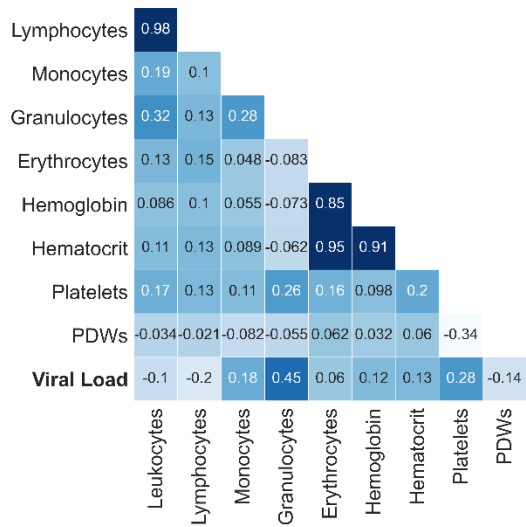
554

555

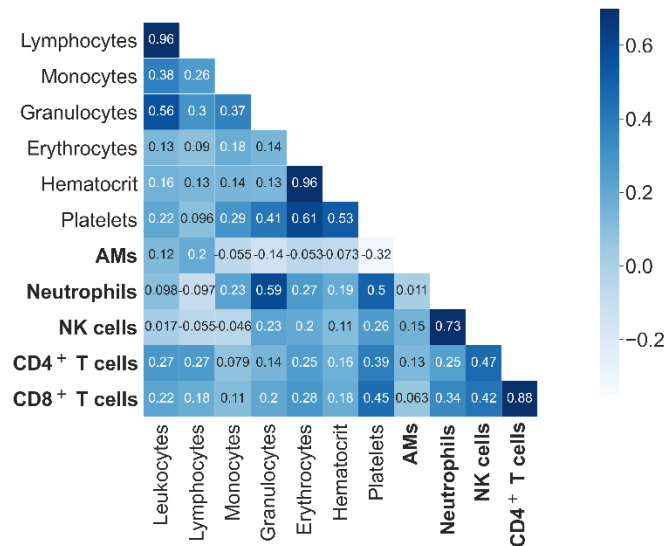
556

557

558 A)



559 B)



560

561 **Figure 2. Selected correlations of blood cells, lung leukocytes, and lung viral load for influenza infection.**

562 (A) shows the correlation of blood cells with lung viral load. (B) shows the correlation of blood cells with
 563 lung leukocytes. The matrices depict the respective Pearson correlation coefficients from the initial
 564 experiments used as training data for the machine learning models. We observed some strongly related
 565 clusters like erythrocytes, haemoglobin, and hematocrit or NK, CD4⁺ T, and CD8⁺ T cells. IAV-associated
 566 lung markers that were later predicted are shown in bold letters. All other parameters were provided to
 567 the algorithm to make the estimation. Here, only a small subset of the data and its correlations are shown.
 568 To view all the data with its correlations, please refer to the supplemental material (Supplementary
 569 Figures S10-12).

570

571

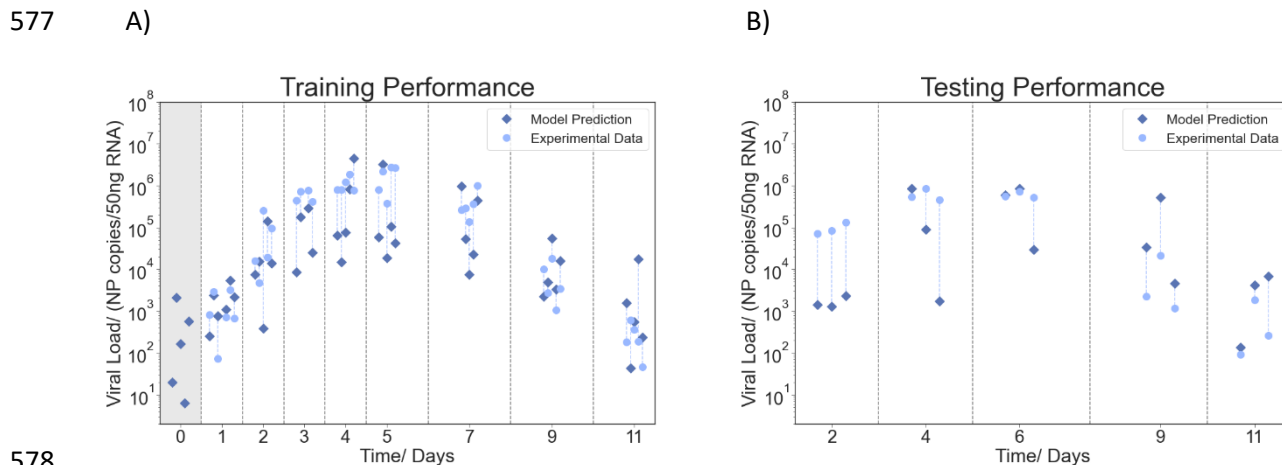
572

573

574

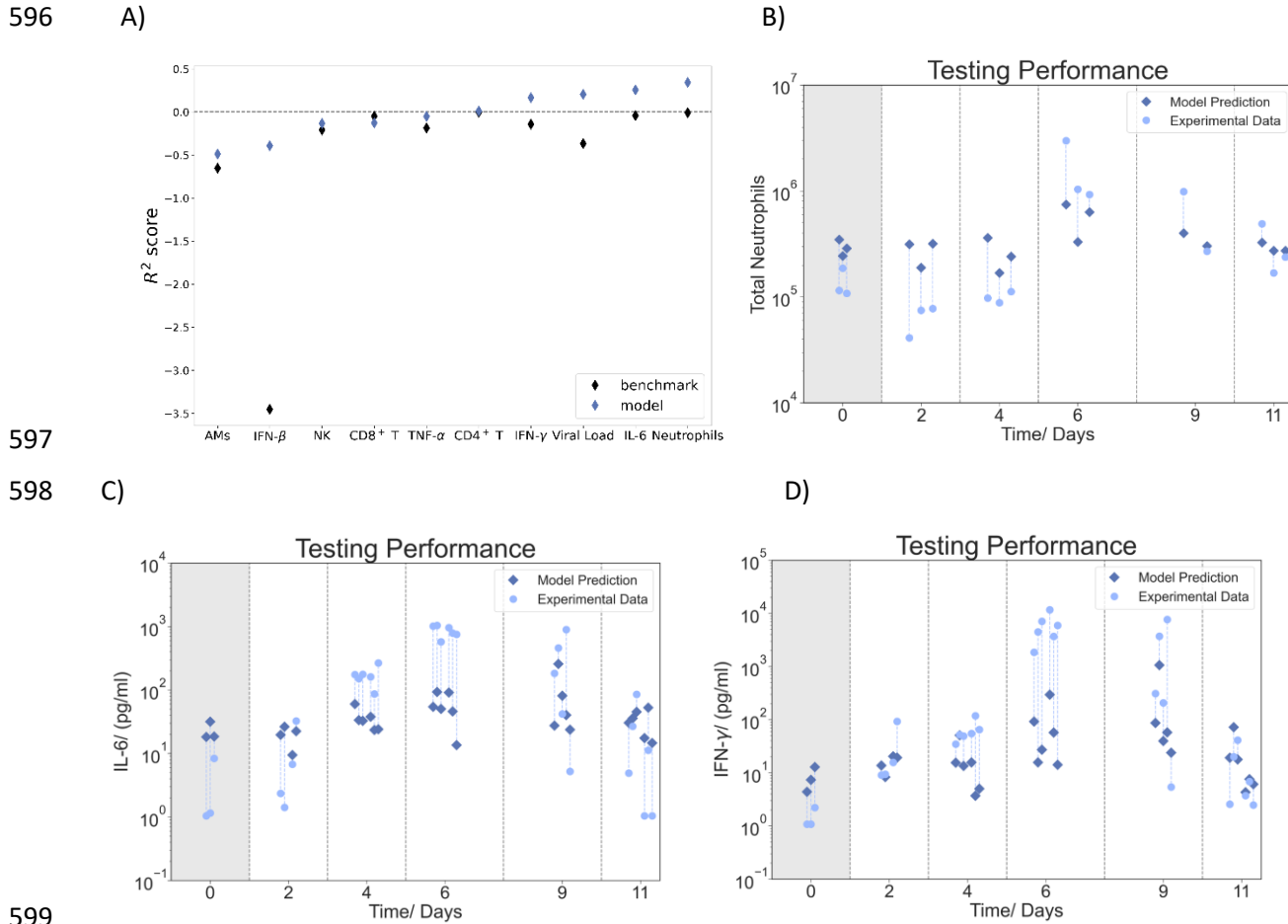
575

576



578
579 **Figure 3. Mapping of the lung viral load from blood data.** The plots show the training and testing
580 performance of the neural network model. (A) shows the performance of the model on the training data.
581 (B) shows the performance of the testing data obtained from a second experiment. Each circle represents
582 one mouse, with its matching individual prediction indicated by a connected (blue dashed line) diamond.
583 The vertical lines divide experimental days. Day 0 marks the control group (viral load below measurable
584 threshold) and is highlighted in grey.

585
586
587
588
589
590
591
592
593
594
595

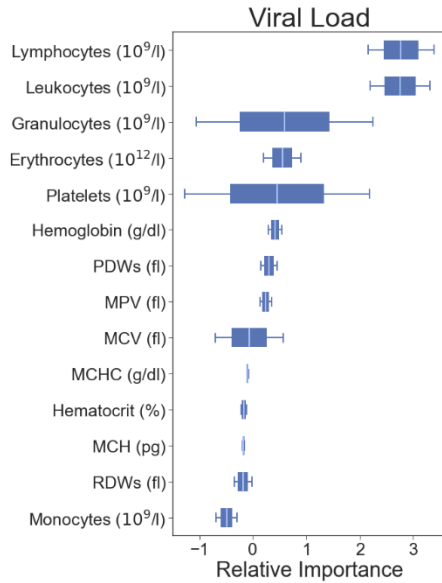


600 **Figure 4. Summary of model predictions for various lung leukocytes and cytokines from hematological**
 601 **data.** (A) R^2 score for different target variables. Blue diamonds show the R^2 score of the best-performing
 602 model. Black diamonds indicate the mean of the target variable obtained from the training data set and
 603 serve as a benchmark. Models that perform better than the benchmark and have a positive R^2 score
 604 indicate the model can make successful predictions. Mapping of (B) neutrophils, (C) IL-6, and (D) IFN- γ
 605 from blood data. Each blue-light circle represents one mouse, with its matching individual prediction
 606 indicated by a connected (blue dashed line) diamond. Gradient boosted regression trees and linear
 607 regression with the aid of PCA worked best for these mappings. We observed that the quality of the
 608 predictions was dependent on the stage of the infection. Neutrophils were estimated more precisely in
 609 the advanced stage of infection, while for IL-6 more accurate estimations were yielded at the peak of the
 610 infection. For complete results, please refer to the supplemental material (Supplementary Figures S13-
 611 15).

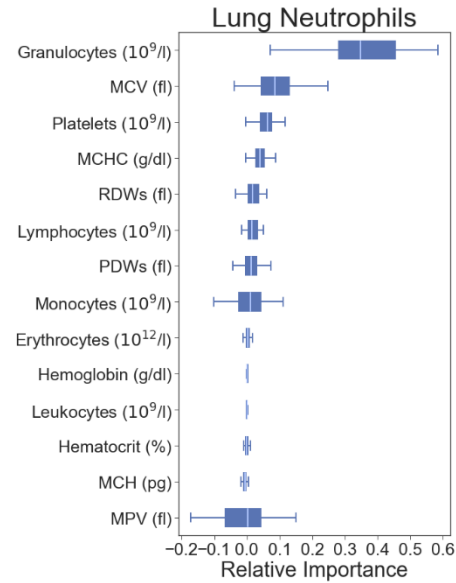
612

613

A)



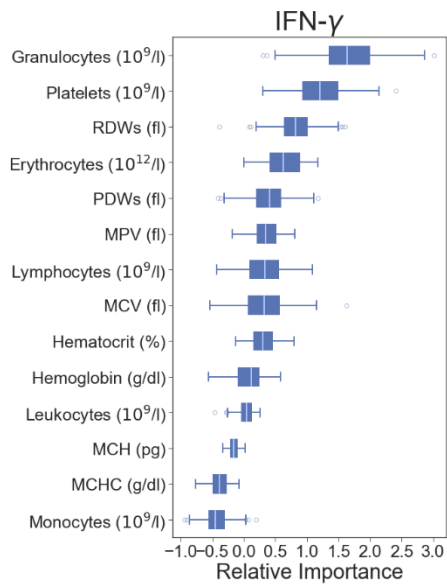
B)



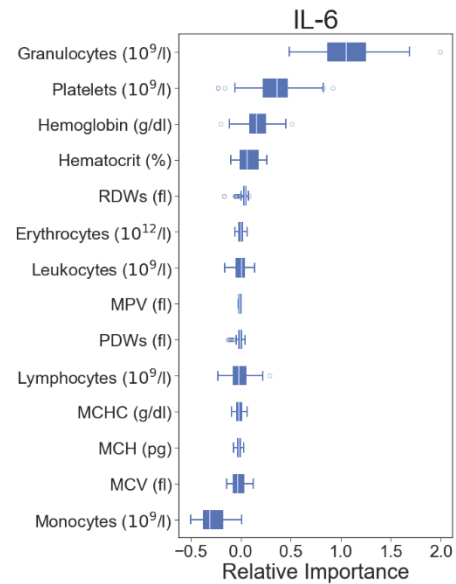
614

615

C)



D)



616

617 **Figure 5. Permutation importance.** The permutation importance for the mapping of (A) the viral load, (B)
 618 neutrophils, (C) IFN- γ , and (D) IL-6 is indicated. We calculated the permutation importance using the
 619 testing data.

620 TABLES

621 **Table 1: Hematological parameters used to estimate the viral load in the lungs of mice.** Blood variables
622 are listed in order of their Pearson correlation with viral load (lowest coefficient: MCH, highest coefficient:
623 granulocytes).

624

Blood variable	Interpretation
MCH (mean corpuscular hemoglobin)	average amount of hemoglobin per red blood cell
RDWs (red cell distribution width)	degree of variation in size and shape of red blood cells
Leukocytes (white blood cells)	protect against infectious diseases and foreign bodies
MPV (mean platelet volume)	average size of platelets in blood
MCV (mean corpuscular volume)	average volume of red cells
Hemoglobin	oxygen carrier in red blood cells
Erythrocytes (red blood cells)	oxygen transportation to the tissue
MCHC (MCH concentration)	concentration of hemoglobin in red blood cells per volume
Monocytes	subtype of white blood cells
PDWs (platelet distribution width)	indicates variation in platelet size
Hematocrit	volume percentage of red blood cells in blood
Lymphocytes	subtype of white blood cells
Platelets	blood component that helps stopping bleeding
Granulocytes	subtype of white blood cells

625

626 **Table 2. Best performing model and respective scores for different targets from the lung milieu.**
627 Hematological data was used in all cases as input to the algorithms. For each target variable to estimate,
628 different machine learning models were tested.

Target Variable	Model	MSE	R ²
Viral Load	Feedforward NN	7.53	0.18
Neutrophils	LR with PCA	0.98	0.25
IFN- γ	Feedforward NN	7.18	0.15
Interleukin-6	GBRT	4.55	0.21

629