

1 **Application of MALDI-MS and Machine Learning to Detection of SARS-CoV-2 and non-SARS-CoV-2**
2 **Respiratory Infections.**

3

4 Sergey Yegorov^{1,2*}, Irina Kadyrova^{3*}, Ilya Korshukov³, Aidana Sultanbekova³, Valentina Barkhanskaya³,
5 Tatiana Bashirova⁶, Yerzhan Zhunusov⁷, Yevgeniya Li⁷, Viktoriya Parakhina^{7,8}, Svetlana Kolesnichenko³,
6 Yeldar Baiken^{2,4,5}, Bakhyt Matkarimov⁴, Dmitriy Vazenmiller³, Matthew S. Miller¹, Gonzalo H. Hortelano²,
7 Anar Turmuhambetova³, Antonella E. Chesca⁹, Dmitriy Babenko³.

8

9 ¹ Michael G. DeGroot Institute for Infectious Disease Research; McMaster Immunology Research Centre;
10 Department of Biochemistry and Biomedical Sciences, McMaster University, Hamilton, ON, Canada.

11 ² School of Sciences and Humanities, Nazarbayev University, Astana, Kazakhstan

12 ³ Research Centre, Karaganda Medical University, Karaganda, Kazakhstan.

13 ⁴ National Laboratory Astana, Centre for Life Sciences, Nazarbayev University, Astana, Kazakhstan.

14 ⁵ School of Engineering and Digital Sciences, Nazarbayev University, Astana, Kazakhstan.

15 ⁶ City Centre for Primary Medical and Sanitary Care, Karaganda, Kazakhstan.

16 ⁷ Infectious Disease Centre of the Karaganda Regional Clinical Hospital, Karaganda, Kazakhstan.

17 ⁸ Department of Internal Diseases, Karaganda Medical University, Kazakhstan.

18 ⁹ Faculty of Medicine, Transilvania University, Braşov, Romania

19

20 *** Corresponding authors' contact emails:**

21 yegorovs@mcmaster.ca (SY)

22 ikadyrova@qmu.kz (IK)

23

24

25 **Keywords:** Acute respiratory infection; COVID-19; SARS-CoV-2; MALDI-MS; Machine Learning;
NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

26 Abstract

27 **Background:** Matrix-assisted laser desorption/ionization mass spectrometry (MALDI-MS) could aid the
28 diagnosis of acute respiratory infections (ARI) owing to its affordability and high-throughput capacity.
29 MALDI-MS has been proposed for use on commonly available respiratory samples, without specialized
30 sample preparation, making this technology especially attractive for implementation in low-resource regions.
31 Here, we assessed the utility of MALDI-MS in differentiating SARS-CoV-2 versus non-COVID acute
32 respiratory infections (NCARI) in a clinical lab setting of Kazakhstan.

33 **Methods:** Nasopharyngeal swabs were collected from in- and outpatients with respiratory symptoms and
34 from asymptomatic controls (AC) in 2020-2022. PCR was used to differentiate SARS-CoV-2+ and NCARI
35 cases. MALDI-MS spectra were obtained for a total of 252 samples (115 SARS-CoV-2+, 98 NCARI and 39
36 AC) without specialized sample preparation. In our first sub-analysis, we followed a published protocol for
37 peak preprocessing and Machine Learning (ML), trained on publicly available spectra from South American
38 SARS-CoV-2+ and NCARI samples. In our second sub-analysis, we trained ML models on a peak intensity
39 matrix representative of both South American (SA) and Kazakhstan (Kaz) samples.

40 **Results:** Applying the established MALDI-MS pipeline "as is" resulted in a high detection rate for SARS-
41 CoV-2+ samples (91.0%), but low accuracy for NCARI (48.0%) and AC (67.0%) by the top-performing
42 random forest model. After re-training of the ML algorithms on the SA-Kaz peak intensity matrix, the
43 accuracy of detection by the top-performing Support Vector Machine with radial basis function kernel model
44 was at 88.0, 95.0 and 78% for the Kazakhstan SARS-CoV-2+, NCARI, and AC subjects, respectively with a
45 SARS-CoV-2 vs. rest ROC AUC of 0.983 [0.958, 0.987]; a high differentiation accuracy was maintained for
46 the South American SARS-CoV-2 and NCARI.

47 **Conclusions:** MALDI-MS/ML is a feasible approach for the differentiation of ARI without a specialized
48 sample preparation. The implementation of MALDI-MS/ML in a real clinical lab setting will necessitate
49 continuous optimization to keep up with the rapidly evolving landscape of ARI.

50

51 **Introduction.**

52 The global response to the COVID-19 pandemic has underscored gaps existing in the laboratory-based
53 diagnosis of acute respiratory infection (ARI)(1). In the early stages of the pandemic, a shortage of rapid and
54 inexpensive techniques amenable to modification to adapt to the newly characterized SARS-CoV-2
55 motivated the search for alternative diagnostic tools. Matrix-assisted laser desorption/ionization mass
56 spectrometry (MALDI-MS), a technique traditionally employed in proteomics and metabolomics, has
57 emerged as a promising alternative to molecular and immunochromatography-based assays to detect SARS-
58 CoV-2 (2). Several different MALDI-MS-based approaches involving varied degrees of sample preparation
59 have been described (2).

60 Our clinical laboratory has particularly been interested in the "untargeted" MALDI-MS method, which
61 applies Machine Learning (ML) algorithms to discern SARS-CoV-2 infection using MALDI-MS peak
62 matrices acquired from respiratory samples such as nasopharyngeal swabs (NPS) without a specialized
63 sample preparation (3–5). Therefore, in this study, we explored the feasibility and accuracy of such
64 untargeted MALDI-MS/ML in differentiating SARS-CoV-2 from non-COVID acute respiratory infections
65 (NCARI) in a clinical laboratory setting in Kazakhstan.

66 **Materials and Methods.**

67 **Study setting.**

68 We collected NPS from three participant subgroups: symptomatic SARS-CoV-2+, NCARI and
69 asymptomatic controls (AC). Participants were recruited between May 25, 2020, and December 20, 2022.
70 Written consent was obtained from all adult participants in the presence of a study coordinator; parental
71 consent was obtained for participants under 18 years of age. The ARI diagnosis was made based on the
72 presence of at least one of the following: fever, nasal congestion, cough, sore throat, and/or
73 lymphadenopathy. SARS-CoV-2+ participants were recruited from among patients of the Karaganda
74 regional clinical hospital, hospitalized with a PCR-confirmed SARS-CoV-2 infection. NCARI participants
75 were recruited at the Karaganda regional clinical hospital and the Karaganda City Centre for Primary

76 Healthcare among patients admitted for moderate-severe ARI symptoms. Most (72.4%) NCARI participants
77 were PCR-positive for common respiratory viruses (adenovirus, seasonal coronaviruses, bocavirus,
78 parainfluenza viruses, respiratory syncytial virus, rhinovirus, influenza or metapneumovirus) or bacteria
79 (*Chlamydia pneumoniae* or *Mycoplasma pneumoniae*). Samples were collected around day 3 (median, IQR
80 [2-4]) and day 5 (median, IQR [3-7]) post-symptom onset for the SARS-CoV-2+ and NCARI participants,
81 respectively. The AC sub-group was recruited from amidst the Karaganda University employees. The SARS-
82 CoV-2 infection status was confirmed in the research lab for all samples using SARS-CoV-2 PCR as
83 described earlier [6, 7]. All samples were frozen at -80C until processing.

84 In addition to the MALDI-MS spectra obtained from clinical samples in Kazakhstan, we incorporated into
85 our analysis the publicly available MALDI-MS data from South America (3).

86 **MALDI-MS analysis.**

87 Within feasible limits, we closely followed the published methodology for sample preparation, spectra
88 acquisition and preprocessing (3), with only minor modifications as specified. Spectral acquisition was
89 performed on the MicroFlex LT v. 3.4 instrument (Bruker Daltonics, Bremen, Germany) equipped with a
90 pulsed UV laser (N2 laser with 337 nm wavelength, 150 microJ pulse energy, 3 ns pulse width and 20 Hz
91 repetition rate). After thawing at room temperature, samples were spotted onto the steel target plate at 0.5 µl,
92 covered with 0.5 µl of the HCCA matrix (a solution containing α -cyano-4-hydroxycinnamic acid diluted in
93 acetonitrile, 2.5% trifluoroacetic acid and nuclease-free water) and then air dried. The target plate was then
94 loaded into the instrument. Spectra were generated by summing 500 single spectra (10 * 50 shots) in the
95 range between 3 and 20 kDa, operating in positive-ion linear mode using a 18-20 kV acceleration voltage, by
96 shooting the laser at random positions on the target spot.

97 **Spectral preprocessing.**

98 Raw MALDI-MS files (Bruker) were uploaded and subsequently preprocessed in R (v. 4.3.0) using
99 MALDIquantForeign and MALDIquant (6). To ensure consistency in peak processing with the original
100 untargeted protocol (3), we used the R scripts generously shared by the authors. Briefly, the spectra were

101 trimmed to a 3–15.5 kDa range, square-root transformed, and smoothened via the Savitzky–Golay method.
102 Baseline correction was done using the TopHat algorithm and intensity normalization was done via total ion
103 current calibration as implemented in MALDIquant. Peak detection was performed using a signal-to-noise
104 ratio of 2 and a halfWindowSize of 10, and the peaks were binned with a tolerance of 0.003. Peak binning
105 was performed in two stages to avoid any additional calibration differences. First, each group spectra were
106 binned separately, and peak filtration was performed, keeping only those peaks that were present in 80% of
107 the spectra of each group. Subsequently, all peaks were binned together. The resulting peak intensity matrix
108 was used for the downstream analyses.

109 In Analysis I, to assess the models trained on the South American samples from the source study (3), we
110 made slight modifications to the sample preprocessing protocol as follows. To ensure that we are comparing
111 the same 88 peaks, we employed the "reference" method for peak binning using the median values of the
112 spectra and peaks obtained by Nachtigall *et al.* as a reference and eliminated the filtering procedure for each
113 subgroup. In Analysis II, we constructed a *de novo* peak matrix representative of the combined South
114 America and Kazakhstan dataset using the script provided by Nachtigall *et al.*

115 **Principal component and hierarchical cluster analyses.**

116 PCA was performed using R FactoMineR and factoextra packages. The hierarchical cluster analysis was
117 done by first calculating a distance matrix using the Euclidean method and clustering samples via the
118 unweighted paired group with arithmetic mean (UPGMA) method. Dendrograms were generated using
119 ggtree and ggtreeExtra R packages.

120 **Machine learning and statistical analysis.**

121 We implemented a total of seven ML algorithms, six of which were used in the earlier study [5] (DT
122 (Decision Tree - Quinlan's C5.0 algorithm), KNN (k-Nearest Neighbors), NB (Naive Bayes), RF (Random
123 Forest), SVM-L (Support Vector Machine with linear kernel), SVM-R (Support Vector Machine with radial
124 basis function kernel) plus an additional algorithm XGBoost (eXtreme Gradient Boosting). Analysis 1 was
125 executed by closely following the earlier protocol, with training performed on South American SARS-CoV-

126 2+ and NCARI spectra. Since the training step of analysis II incorporated three sub-groups, i.e. AC samples
127 in addition to the SARS-CoV-2+ and NCARI, the analysis pipeline was modified as outlined below to
128 accommodate this change.

129 Initially, we split the entire sample into two distinct groups: the training dataset, consisting of 80% of
130 samples, and the test group, which accounted for the remaining 20%. In line with Nachtigall *et al* [5], we
131 conducted a training process using a five-fold (outer) nested repeated (five times) ten-fold (inner) cross-
132 validation with a randomized stratified splitting approach. To optimize the performance of each algorithm,
133 we tested 20 hyperparameters in the inner loop of the cross-validation approach, using a random search
134 method. This process was repeated 20 times to ensure robustness and reliability of the model. We selected
135 the best models based on their area under the curve (AUC) score, which is a common metric for evaluating
136 binary-classification model performance, using the Caret R package. In addition, model performance was
137 assessed using several other classification metrics, including F-measure, recall, accuracy, specificity,
138 sensitivity, and positive and negative predictive values in the yardstick R package; differences across the
139 sub-groups were assessed using the Mann-Whitney U non-parametric test in R.

140 **Role of the funding source**

141 The funders had no role in study design, data collection and analysis, decision to publish, or preparation of
142 the manuscript.

143

144 **Results**

145 The primary objective of the study was to assess the capacity of the MALDI-MS approach to detect SARS-
146 CoV-2 infection within a heterogeneous mix of SARS-CoV-2+, NCARI and AC samples (Table 1).

147

148

149

150 **Table 1.** Demographic characteristics of participants.

Characteristic	Overall, N = 252	SARS-COV-2+ 2021, N = 108	SARS-COV-2+ 2022, N = 7	NCARI, N=98	AC, N=39	p-value*
Age, years, median (IQR)	38.0 (18.0,60.0)	61.0 (48.0,69.0)	3 (1.0, 37.0)	8.0 (2.0, 35.0)	34.0 (25.0, 47.0)	<0.001
Male sex, n (%)	114 (45.2%)	49 (45.4%)	6 (85,7 %)	38 (38,8%)	21 (53.8%)	<0.001
Kazakh ethnicity, n (%)	116 (46%)	26 (24.1%)	5 (71,4%)	61 (62.2%)	24 (61.5%)	<0.001
Any comorbidities	104 (41.2%)	72 (66.6%)	1 (14,2%)	15 (15.3 %)	14 (35.9%)	<0.001

151
152 * Differences across the groups were assessed using Kruskal-Wallis or Pearson's Chi-squared tests.
153

154 Therefore, we performed two independent analyses (Figure 1). In the first analysis, we assessed the
155 performance of the Nachtigall *et al.* ML pipeline on the combined pool of samples, both from the original
156 study (data collected from three South American countries in 2020) and Kazakhstan (data collected in 2021
157 and 2022); the ML pipeline in this analysis was trained only on the original South American datasets. In the
158 second analysis, we retrained the ML algorithm, accounting for the spectra contributed by the samples from
159 Kazakhstan and applied this re-trained ML algorithm to the combined pool of samples.

160 **Fig 1. Overall study workflow and description of the analyses.** NCARI: non-COVID acute respiratory
161 infections; AC: asymptomatic controls; ML: machine learning

162

163 **Analysis I.: Applying the "as is" MALDI-MS pipeline to differentiate** 164 **ARI samples collected in Kazakhstan.**

165 To assess how well the original analysis pipeline (3) would differentiate SARS-CoV-2+ samples within the
166 dataset from Kazakhstan, we replicated the steps for i) MALDI-MS peak selection, ii) ML training and iii)
167 ML assessment. Specifically, we focused on the same MALDI-MS peaks that Nachtigall *et al.* (3) used in
168 their analyses (Table S1). These peaks were derived using a six-step spectra processing workflow including

169 spectra transformation and smoothing, baseline removing, spectra calibration, peak detection, and peak
170 processing.

171 We then constructed a peak intensity matrix on the 88 peaks, identical to that used by Nachtigall and
172 colleagues (3), for the downstream analysis of a combined dataset incorporating both the South American
173 (Table S1) and Kazakhstan samples (Figure 2 and Table S2).

174

175 **Fig 2. MALDI-MS peak data generated using nasopharyngeal swabs and processed following the**
176 **MALDI-MS/ML pipeline developed by Nachtigall and colleagues (3).**

177 A-C. representative MALDI-MS spectra from symptomatic SARS-CoV-2+ (A), symptomatic non-SARS-
178 CoV-2 (B) and a healthy control sample from Kazakhstan (C). The central line indicates median value of the
179 spectra, while the shaded region on either side represents the interquartile interval. Insets depict a range from
180 3000 to 5500 m/z encompassing 70% (62/88) of the identified peaks. d. PCA of the combined dataset
181 incorporating MALDI-MS data both from Kazakhstan and South America (2020 SARS-CoV+ and
182 symptomatic SARS-CoV-2-negative)(3).

183

184 We next explored the selected peaks across the comparison groups by reducing the multidimensionality
185 using principal component analyses and dendrograms. Like Nachtigall et al, we did not detect any obvious
186 clustering by sub-group, emphasizing the need for a more sensitive approach to discern subtle differences in
187 the highly multidimensional MALDI-MS peak data (Figure 2D and 2E, Figures S2-S5). Hence, we then
188 applied to our combined Kazakhstan-South America MALDI-MS peak dataset the original Nachtigall *et al.*
189 ML algorithm trained on the original South American samples (3).

190 In keeping with earlier results (3), when tested the South American samples alone, SVM-R provided the
191 highest ROC AUC, although other models had similarly high-performance characteristics (Table S3 and
192 Figure 3A-B) for classifying cases of SARS-CoV-2 and non-SARS-CoV-2.

193

194 **Fig 3. Classification accuracy of the MALDI-ML algorithms assessed on the data from Kazakhstan**
195 **and South America.**

196 A) Accuracy metrics for each of the seven ML models trained on the South American MALDI-MS data
197 (Analysis I in the current study) for the differentiation of study sub-groups. B) ROC curves of the top-
198 performing RF and SVM-L algorithms (Analysis I). C) Accuracy metrics for each of the seven ML models
199 trained on the combined South America-Kazakhstan dataset (Analysis II in the current study) for the
200 differentiation of study sub-groups. D) ROC curves for the top-performing SVM-R and DT algorithms
201 (Analysis II).

202

203 Subsequently, we assessed the performance of the same ML algorithms on samples from Kazakhstan. Here,
204 we observed a broad variation in the ability of the ML models to discern SARS-CoV-2+ samples. RF had the
205 highest percentage of correctly identified 2020-SARS-CoV-2+ samples (91%) (Figure 3A and Table S4,
206 Figure S6). Notably, the accuracy for 2021 SARS-CoV-2 was <60% for all models, similar to the accuracy
207 for identifying NCARI. RF discerned AC with an accuracy of 68%, the highest of all models for this sub-
208 group.

209 **Analysis II: Applying the re-trained MALDI/MS-ML to differentiate**
210 **ARI.**

211 To ensure that we include all relevant MALDI-MS signature peaks representative of all sub-groups, we
212 performed peak selection on the entire pool of samples containing samples from both Kazakhstan and South
213 America (n=615). A total of 120 peaks were identified and a peak intensity matrix was constructed (Table
214 S5). As in Analysis I, PCA and dendrograms did not show any visually apparent clustering of sub-groups
215 (Figures S7-S10). We then proceeded to train ML models on the combined pool consisting of the 120 peaks,
216 of which 53 overlapped with the original 88 peaks.

217 Due to the small sample size of the 2022 subset, the SARS-CoV-2 2021 and 2022 subsets were combined
218 prior to testing the model performance. We then assessed the performance of the trained ML algorithm on

219 the South America-Kazakhstan dataset. All models demonstrated similarly high-performance characteristics
220 in differentiating SARS-CoV-2+ samples. SVM-R and DT slightly outperformed the other five models in
221 discerning SARS-CoV-2 infection from both NCARI and AC with ROC AUC values of 0.983 [0.958, 0.987]
222 and 0.972 [0.966, 0.979], respectively (Figures 3C-D and Table S4). SVM-R, in particular, differentiated the
223 Kazakhstan SARS-CoV-2+, NCARI, and AC subjects with an accuracy of 88.0, 95.0 and 78.0%,
224 respectively (Figure 3C). Both SVM-R and DT were also highly accurate at differentiating NCARI and AC
225 sub-groups (Table S4).

226 **Discussion**

227 Here we aimed to assess the feasibility of deploying MALDI-MS and ML in a clinical lab to differentiate
228 SARS-CoV-2 from other ARI, particularly in the context of minimal specialized sample preparation. Our
229 initial application of the original MALDI-MS/ML pipeline, trained on South American samples (3),
230 demonstrated reduced efficiency in identifying samples from Kazakhstan. Re-training the ML models to
231 incorporate MALDI peak information from a diverse pool of Kazakhstan samples, including SARS-CoV-2+,
232 NCARI subjects, and asymptomatic controls, led to a significant improvement in detection accuracy. Taken
233 as a proof-of-concept, our results support the utility of MALDI-MS/ML, especially in the early phases of
234 respiratory endemics/pandemics, when limited knowledge is available on the infectious pathogen's identity
235 and in low-resource environments, where alternative methods may yet be unavailable.

236 Our replication studies underscore the importance of considering geographical and population-specific
237 variations in the application of MALDI-MS/ML. The observed differences in the performance of the original
238 pipeline trained on South American samples may be attributed to the inherent complexity of NPS, which
239 contains a mixture of host proteins and diverse microbial species (7,8). The sensitivity and specificity of
240 MALDI-MS/ML may also be affected by variability in immune response to different viral loads and the
241 presence of co-infections (9,10). These challenges emphasize the need for careful evaluation and calibration
242 in the application of MALDI-MS/ML.

243 Our study has several limitations. The lack of specialized sample preparation, although advantageous for
244 low-resource settings, may introduce variability and noise into the data, a concern raised by other authors

245 (9,10). Due to a relatively small sample size of the NCARI group, we did not further pursue stratification of
246 this group by the causative agents identified via multiplex PCR. The utility of MALDI-MS/ML in
247 differentiating various NCARI would be important to examine in the context of the changing post-pandemic
248 ARI landscape (11). Temporal variation, spanning samples collected over two years (2020-2022), might
249 have contributed to a high heterogeneity of our results. The differences across groups regarding the basic
250 demographics may also have a confounding effect on the results. Further validation of the method in a
251 broader clinical context would be necessary to fully assess the potential for real-world application.

252 In conclusion, our study provides valuable insights into the potential of MALDI-MS as an accessible
253 laboratory-based diagnostic tool for ARI. While promising, the implementation of MALDI-MS/ML in real
254 clinical lab settings will require further optimization, validation, and continuous adaptation to the evolving
255 epidemiological landscape. Further research is needed to explore the specific components of MALDI-MS
256 spectra that are most informative for differentiating various ARI. Such investigations will contribute to the
257 ongoing refinement of this promising diagnostic tool.

258 **Declarations**

259 **Ethics approval and consent to participate.**

260 All study procedures were approved by the Research Ethics Board of Karaganda Medical University under
261 Protocol 12 (approved 45) from 06.04.2020. Written informed consent was obtained from all participants.

262 **Consent for publication**

263 Authors provide consent for the publication of the manuscript detailed above, including any accompanying
264 images or data contained within the manuscript.

265 **Availability of data and materials**

266 All raw data and R code are available through Github ([https://github.com/dimbage/ML_MALDI-TOF_SARS-](https://github.com/dimbage/ML_MALDI-TOF_SARS-CoV-2)
267 [CoV-2](https://github.com/dimbage/ML_MALDI-TOF_SARS-CoV-2)).

268 Acknowledgements

269 We thank the study participants and clinic staff. We are grateful to Professor Leonardo Santos for sharing the
270 R scripts and associated data from their original study.

271 Supplementary information

272 All supplementary information can be found in the Appendix.
273

274 References.

- 275 1. Yegorov S, Goremykina M, Ivanova R, Good SV, Babenko D, Shevtsov A, et al. Epidemiology, clinical
276 characteristics, and virologic features of COVID-19 patients in Kazakhstan: A nation-wide retrospective
277 cohort study. *The Lancet Regional Health – Europe* [Internet]. 2021 May 1 [cited 2021 Aug 9];4.
278 Available from: [https://www.thelancet.com/journals/lanep/article/PIIS2666-7762\(21\)00073-9/abstract](https://www.thelancet.com/journals/lanep/article/PIIS2666-7762(21)00073-9/abstract)
- 279 2. Spick M, Lewis HM, Wilde MJ, Hopley C, Huggett J, Bailey MJ. Systematic review with meta-analysis
280 of diagnostic test accuracy for COVID-19 by mass spectrometry. *Metabolism - Clinical and Experimental*
281 [Internet]. 2022 Jan 1 [cited 2022 Dec 13];126. Available from:
282 [https://www.metabolismjournal.com/article/S0026-0495\(21\)00222-5/fulltext](https://www.metabolismjournal.com/article/S0026-0495(21)00222-5/fulltext)
- 283 3. Nachtigall FM, Pereira A, Trofymchuk OS, Santos LS. Detection of SARS-CoV-2 in nasal swabs using
284 MALDI-MS. *Nat Biotechnol*. 2020 Oct;38(10):1168–73.
- 285 4. Deulofeu M, García-Cuesta E, Peña-Méndez EM, Conde JE, Jiménez-Romero O, Verdú E, et al.
286 Detection of SARS-CoV-2 Infection in Human Nasopharyngeal Samples by Combining MALDI-TOF
287 MS and Artificial Intelligence. *Front Med (Lausanne)*. 2021;8:661358.
- 288 5. Tran NK, Howard T, Walsh R, Pepper J, Loegering J, Phinney B, et al. Novel application of automated
289 machine learning with MALDI-TOF-MS for rapid high-throughput screening of COVID-19: a proof of
290 concept. *Sci Rep*. 2021 Apr 15;11(1):8219.
- 291 6. Gibb S, Strimmer K. MALDIquant: a versatile R package for the analysis of mass spectrometry data.
292 *Bioinformatics*. 2012 Sep 1;28(17):2270–1.
- 293 7. Rajagopala SV, Bakhom NG, Pakala SB, Shilts MH, Rosas-Salazar C, Mai A, et al. Metatranscriptomics
294 to characterize respiratory virome, microbiome, and host response directly from clinical samples. *Cell*
295 *Rep Methods*. 2021 Oct 25;1(6):100091.
- 296 8. Sandybayev NT, Belousov VY, Stochkov VM, Solomadin MV, Granica J, Yegorov S. The
297 nasopharyngeal virome in adults with acute respiratory infection [Internet]. *bioRxiv*; 2023 [cited 2023
298 Aug 29]. p. 2023.08.21.554191. Available from:
299 <https://www.biorxiv.org/content/10.1101/2023.08.21.554191v1>
- 300 9. Iles RK, Zmuidinaite R, Iles JK, Carnell G, Sampson A, Heeney JL. Development of a Clinical MALDI-
301 ToF Mass Spectrometry Assay for SARS-CoV-2: Rational Design and Multi-Disciplinary Team Work.
302 *Diagnostics (Basel)*. 2020 Sep 24;10(10):746.

Yegorov et al.

MALDI-MS/ML to detect ARI

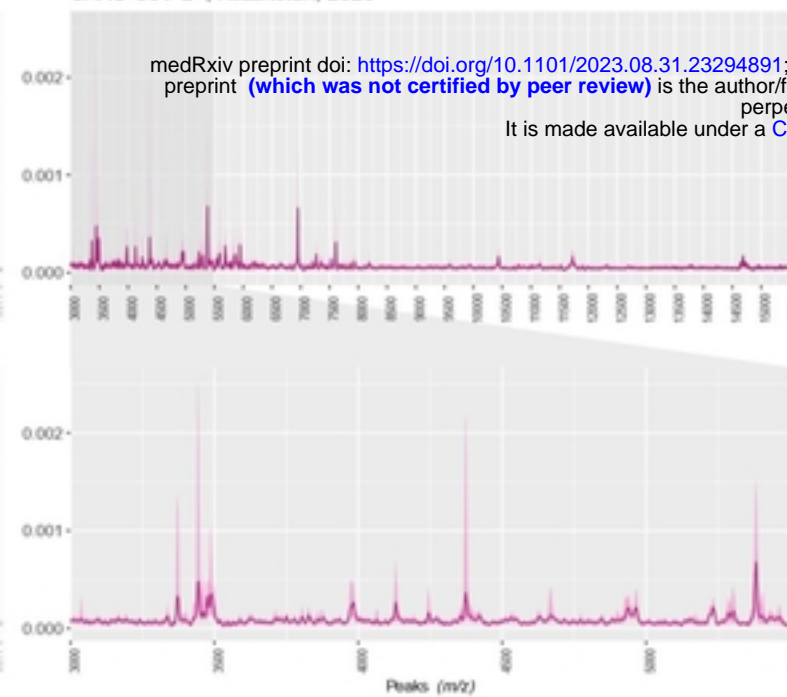
- 303 10. Renuse S, Vanderboom PM, Maus AD, Kemp JV, Gurtner KM, Madugundu AK, et al. A mass
304 spectrometry-based targeted assay for detection of SARS-CoV-2 antigen from clinical specimens.
305 EBioMedicine. 2021 Jul;69:103465.
- 306 11. Sandybayev N, Beloussov V, Strochkov V, Solomadin M, Granica J, Yegorov S. Characterization of
307 viral pathogens associated with symptomatic upper respiratory tract infection in adults during a low
308 COVID-19 transmission period. PeerJ. 2023;11:e15008.

309

310

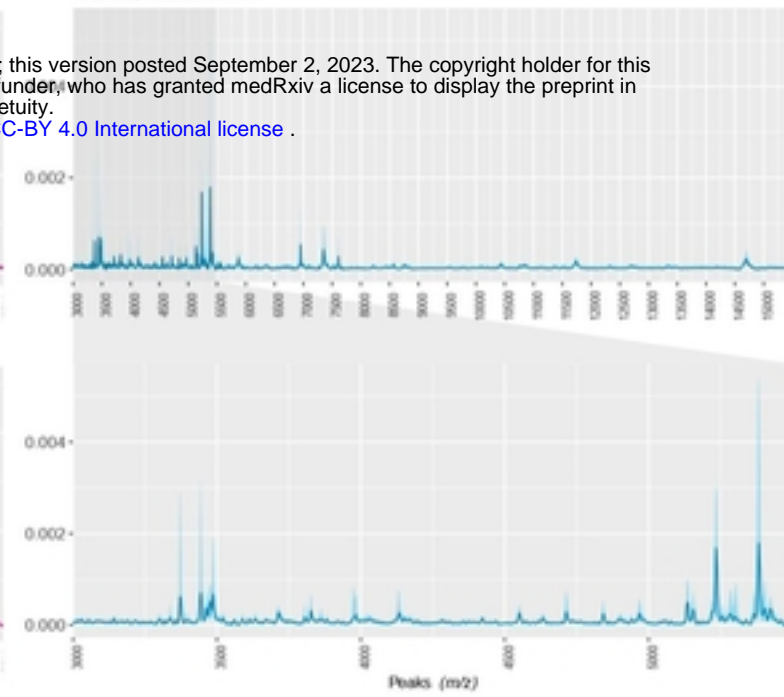
A)

SARS-CoV-2+, Kazakhstan, 2020



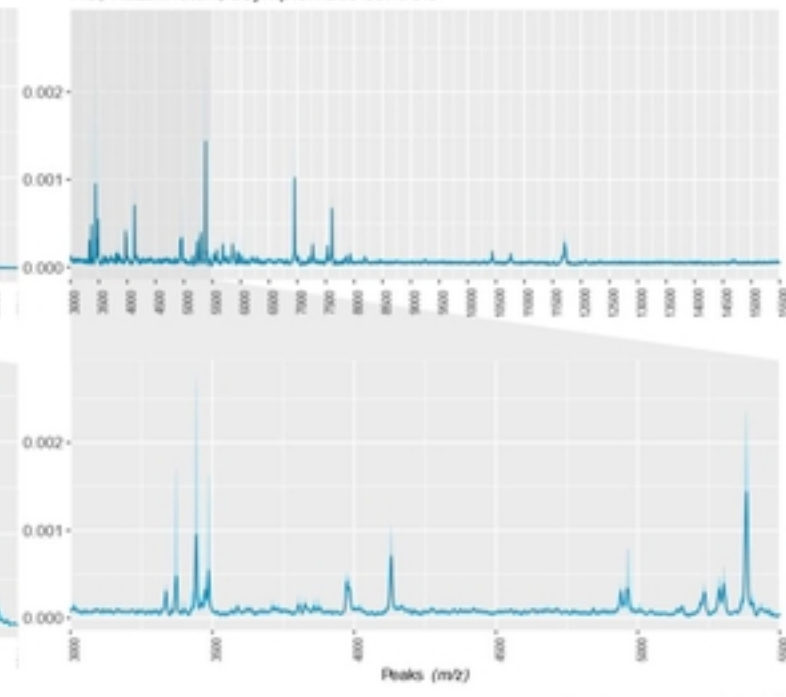
B)

NCARI, Kazakhstan



C)

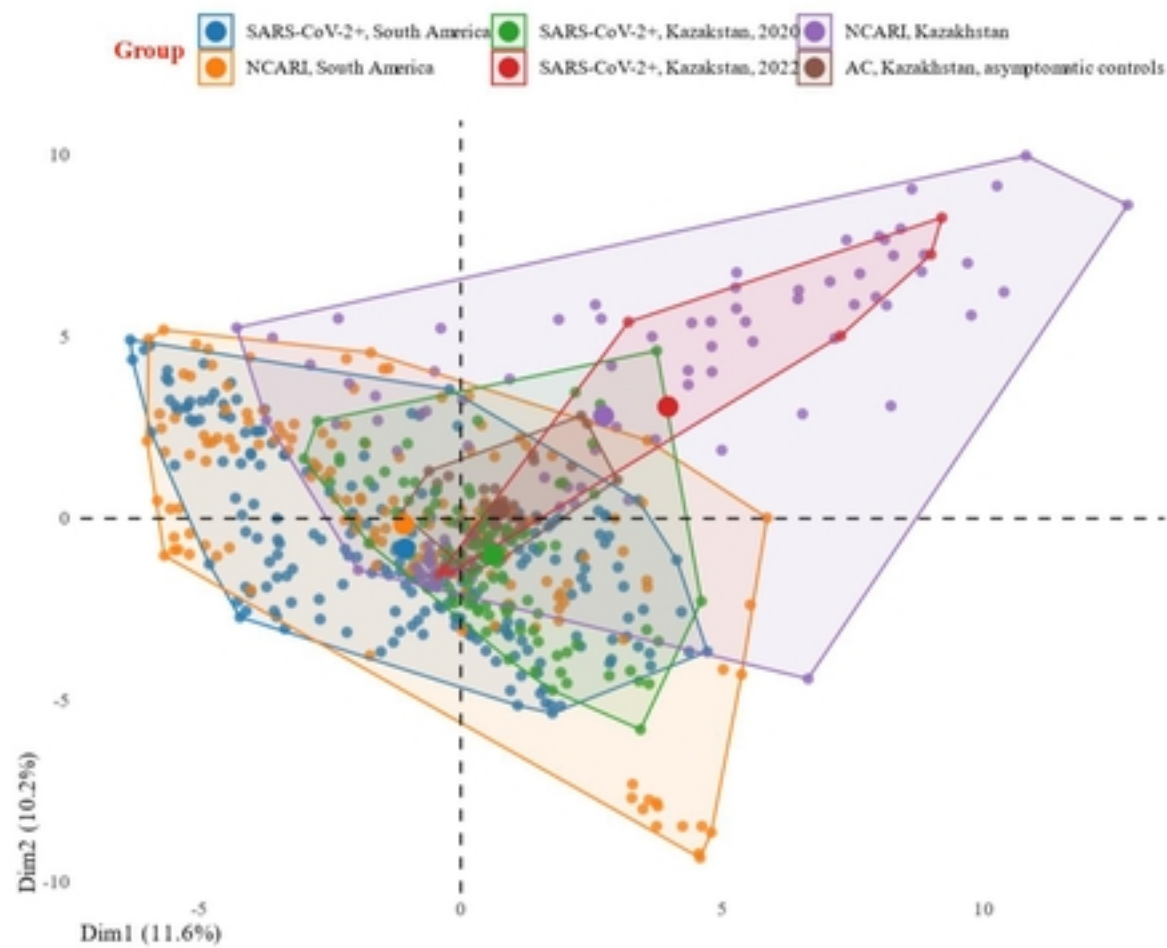
AC, Kazakhstan, asymptomatic controls



medRxiv preprint doi: <https://doi.org/10.1101/2023.08.31.23294891>; this version posted September 2, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

D)



E)

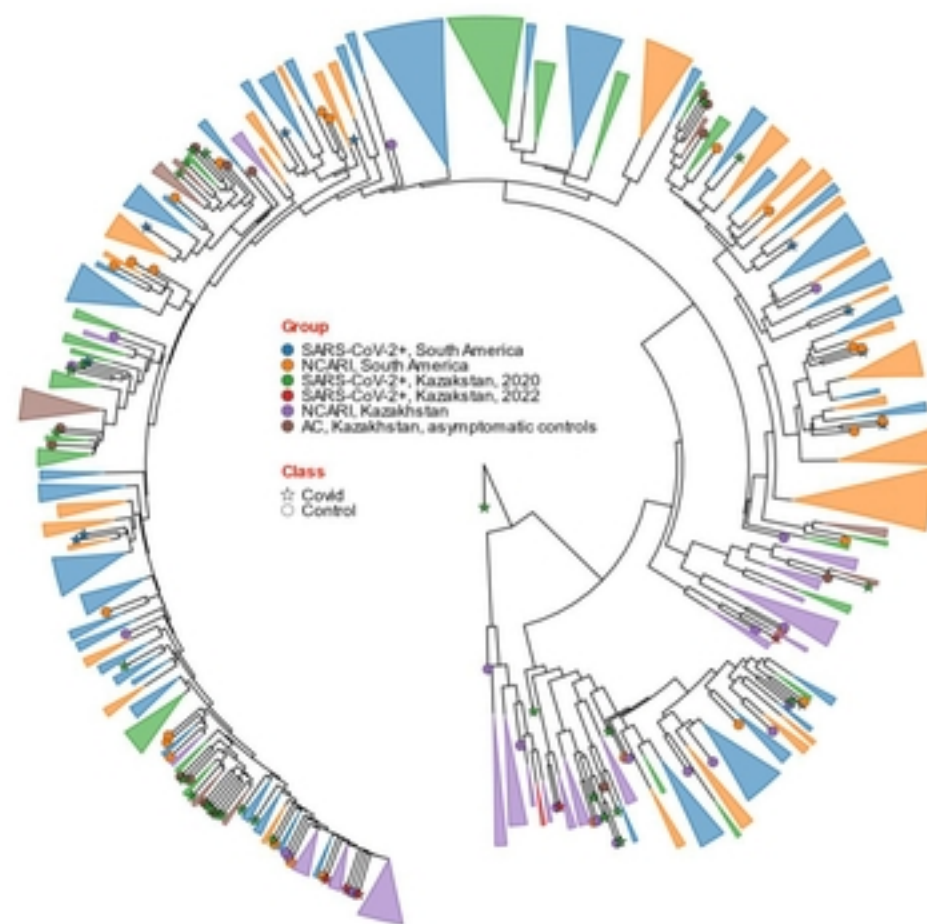
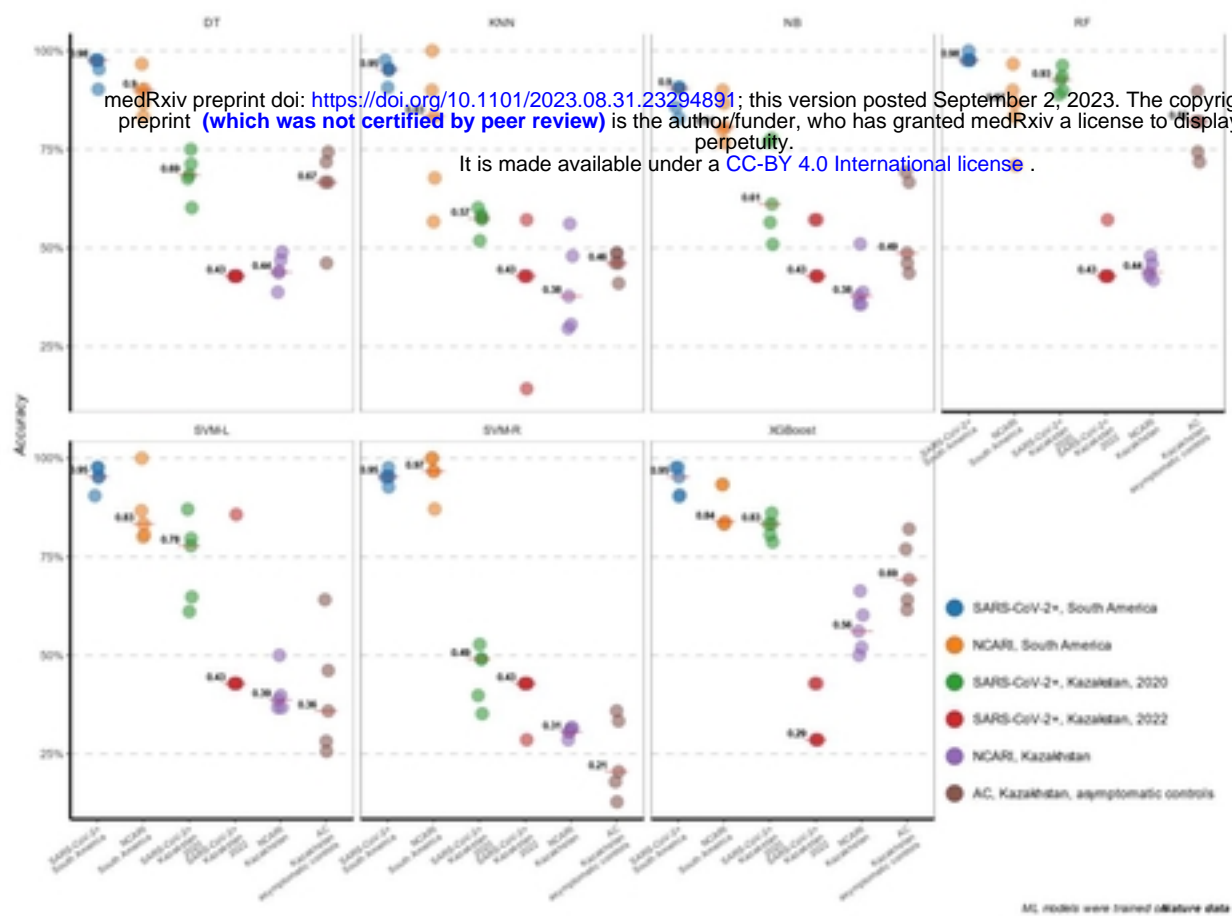
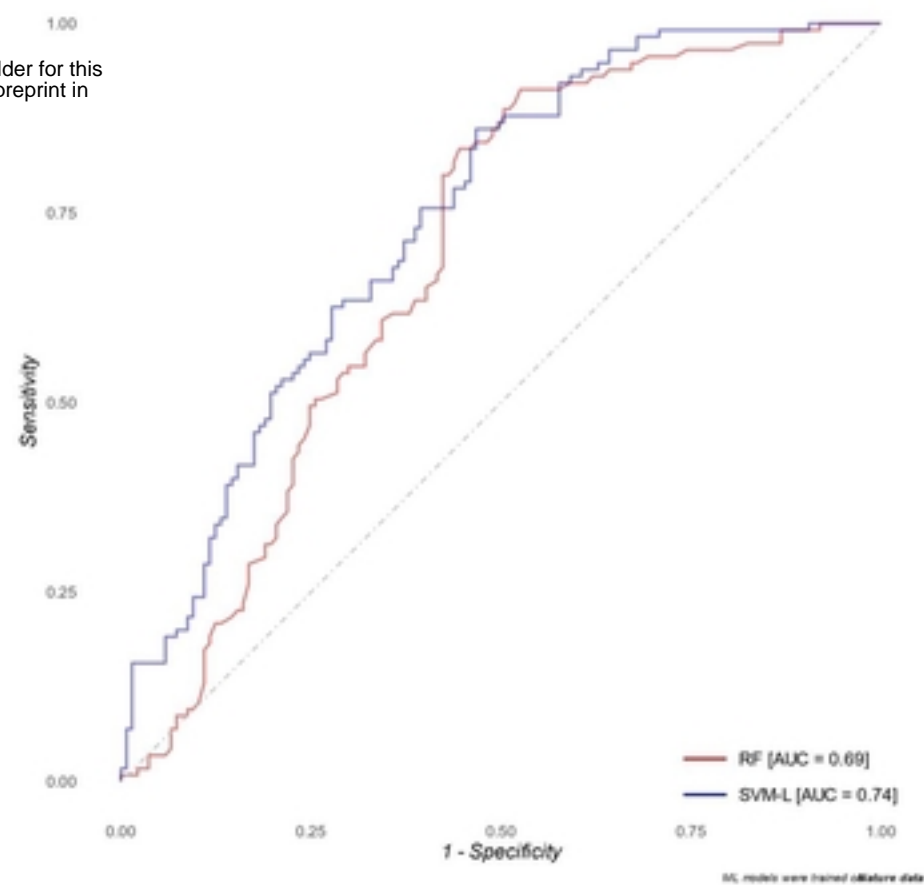


Fig 2
Figure 2

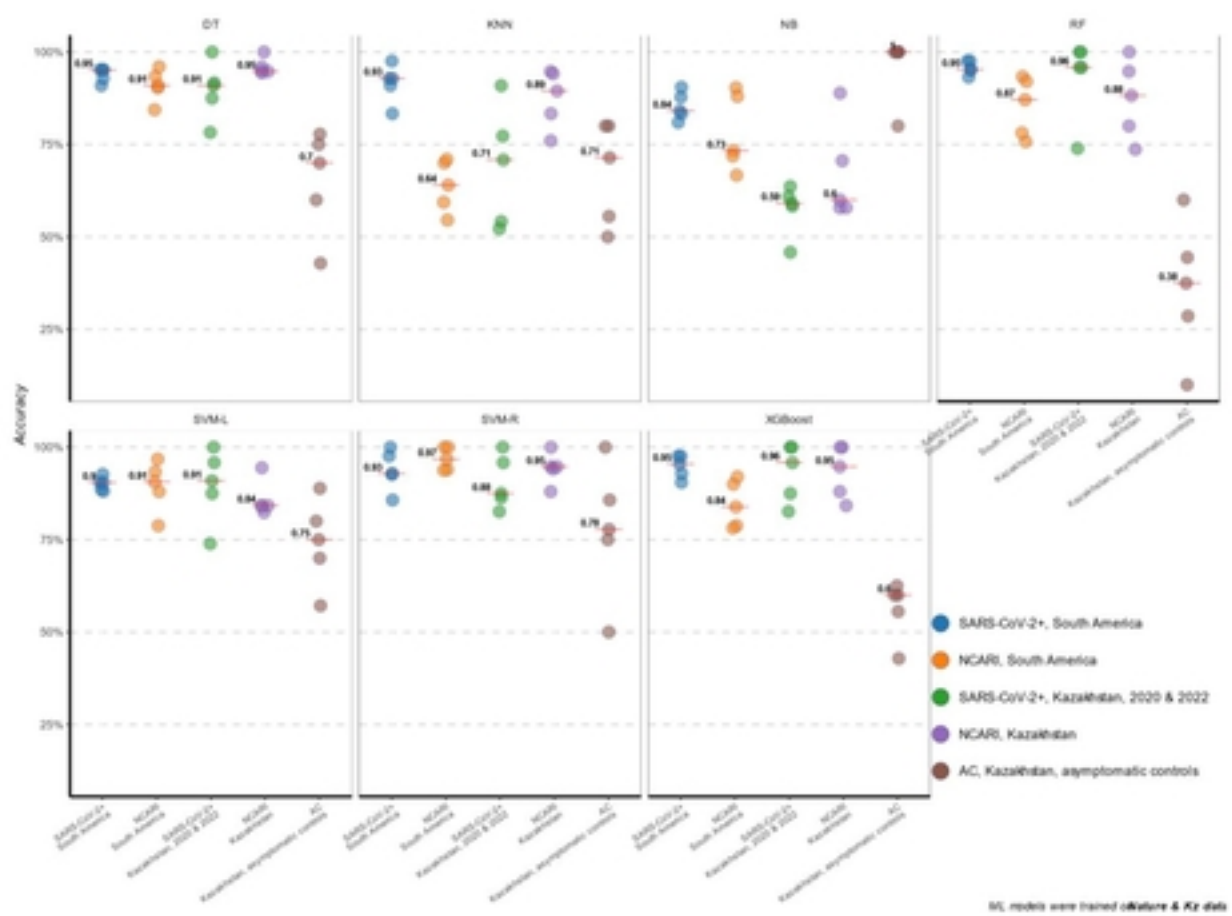
A)



B)



C)



D)

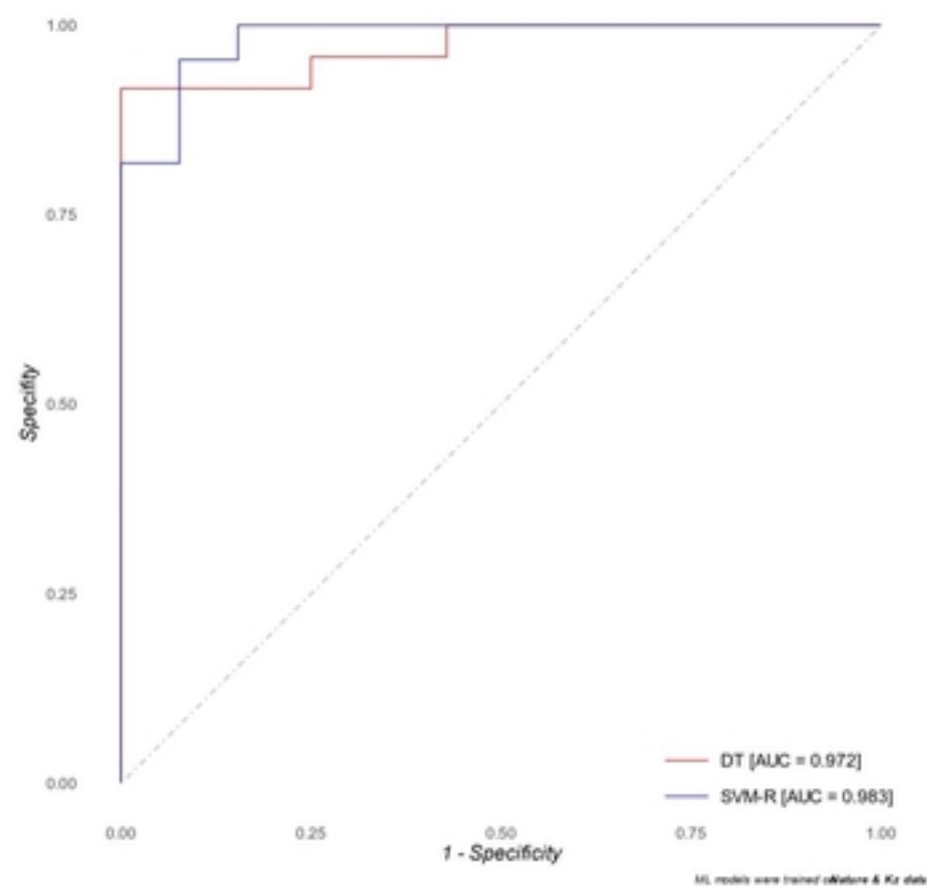


Fig 3
Figure 3

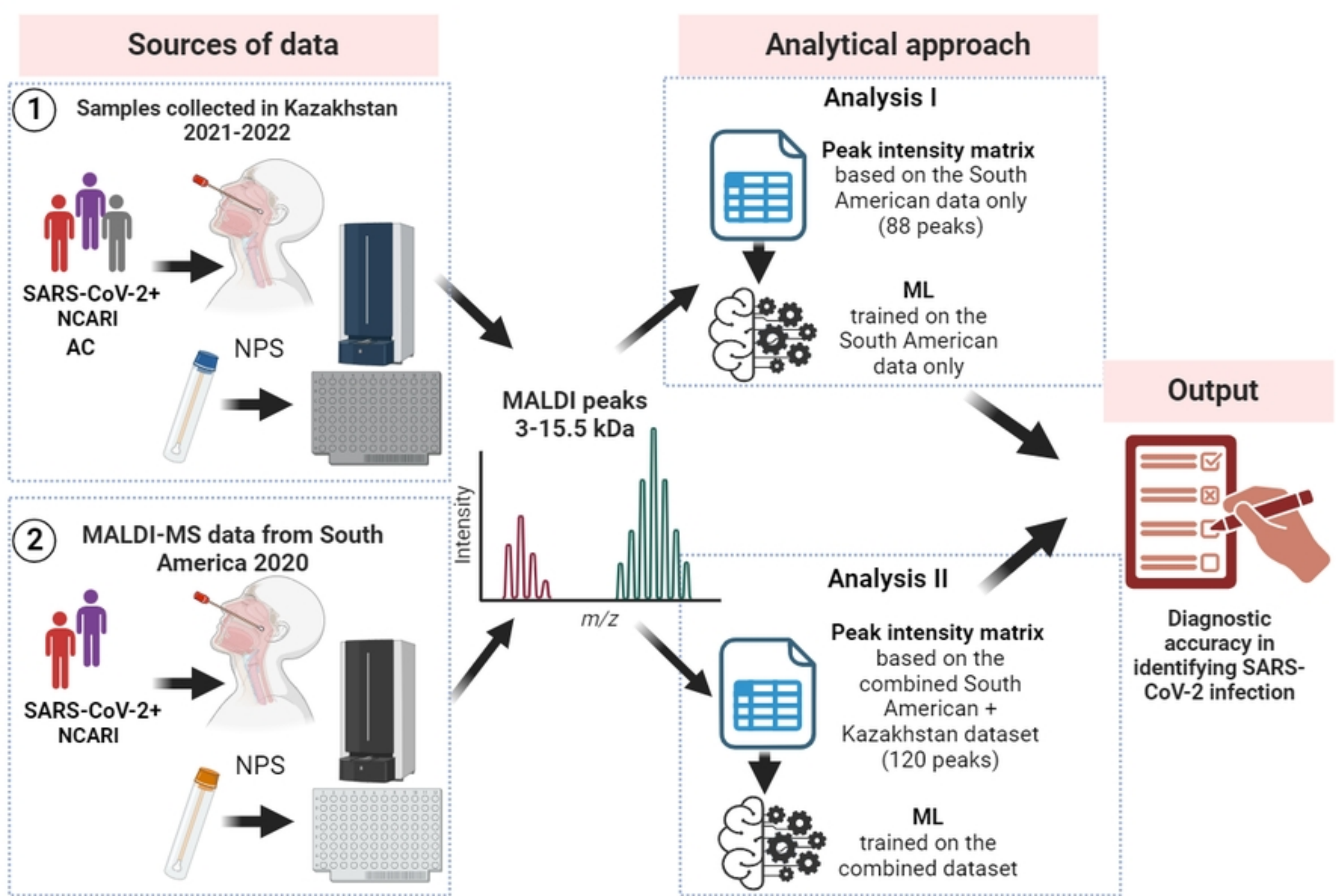


Figure 1