

1 **Rapid identification of pathogens causing bloodstream infections by Raman**
2 **spectroscopy and Raman tweezers**

3

4 **Katarina Rebrosova^{1*}, Silvie Bernatová², Martin Šiler², Jan Mašek^{3,4}, Ota Samek², Jan**
5 **Ježek², Martin Kizovsky², Veronika Holá¹, Pavel Zemanek², Filip Růžička^{1*}**

6

7 ¹ Department of Microbiology, Faculty of Medicine of Masaryk University and St. Anne's
8 University Hospital, Pekařská 53, Brno 65691, Czech Republic

9 ² Institute of Scientific Instruments of the Czech Academy of Sciences, v.v.i.,

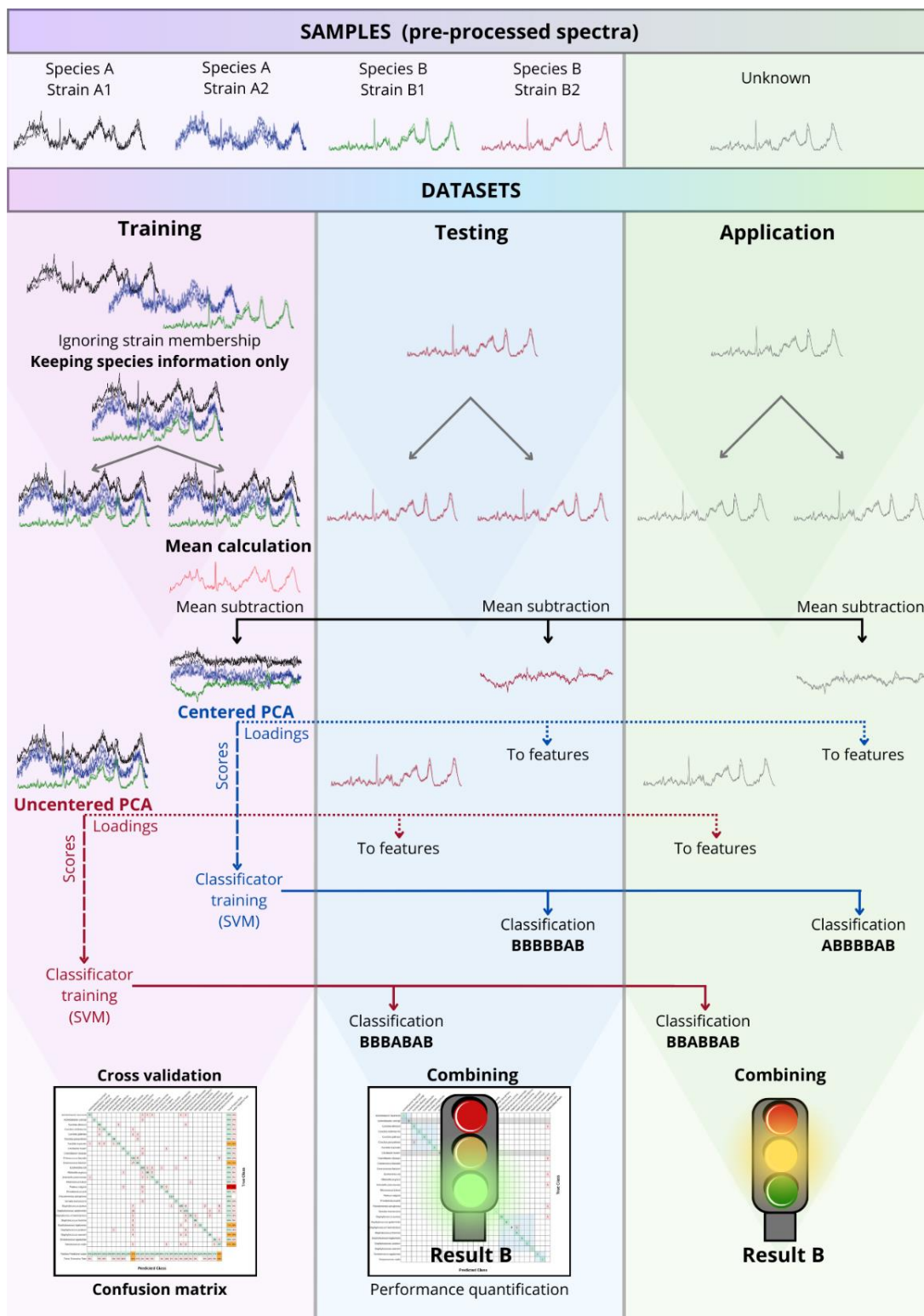
10 Královopolská 147, Brno 61264, Czech Republic

11 ³ National Centre for Biomolecular Research, Faculty of Science, Masaryk University,

12 Kamenice 753/5, Brno 625 00, Czech Republic

13 ⁴ Department of Plant Developmental Genetics, Institute of Biophysics, Academy of Sciences

14 of the Czech Republic, Královopolská 135, Brno 612 65, Czech Republic



15

16 **Figure S1.** Workflow of the species classification process. The processed spectra (Savitzky-
 17 Goly filtered, fluorescence background removed and normalized), which had
 18 been independently identified by MALDI-TOF, are separated into training (~75% of strains)
 19 and testing groups (remaining strains). In the training group (the leftmost column), we employ

*Corresponding authors: k.mlynarikova@gmail.com (K.R.); fruzic@fnusa.cz (F.R.)

20 only the MALDI-TOF confirmed species and do not use strain sub-division. Centered and un-
21 centered Principal Component Analysis (PCA) is applied independently on spectra in the
22 training group giving two sets of PCA scores and loadings. PCA scores of both types
23 represent features that are used to independently train the Support Vector Machines (SVM)
24 classifiers (the quality of training is evaluated by cross-validation confusion matrices, see
25 Fig. 2). The testing data (spectra, the center column), are projected into the same subspace
26 using mean, and loadings of the training data and the obtained features (scores) are used
27 as inputs to the SVM classifiers. The results of single spectra classifications are combined
28 per strain, and the resulting species class is obtained by the semaphore scheme. The overall
29 classification performance is summarized in Fig. 3 and Table 2 of the manuscript. The
30 completely unknown sample consisting of ten spectra is typically processed similarly to
31 testing data (see the rightmost column); for further details, see the main text.

32