

# Automatic discrimination of species within the *Enterobacter cloacae* complex using MALDI-TOF Mass Spectrometry and supervised algorithms

Ana Candela<sup>1\*</sup>, Alejandro Guerrero-López<sup>2\*</sup>, Miriam Mateos<sup>3</sup>, Alicia Gómez-Asenjo<sup>1</sup>, Manuel J. Arroyo<sup>4</sup>, Marta Hernandez-García<sup>3,5</sup>, Rosa del Campo<sup>3,5</sup>, Emilia Cercenado<sup>1,6,7</sup>, Aline Cuénod<sup>8,9</sup>, Gema Méndez<sup>4</sup>, Luis Mancera<sup>4</sup>, Juan de Dios Caballero<sup>3,5</sup>, Laura Martínez-García<sup>3,10</sup>, Desirée Gijón<sup>3,5</sup>, María Isabel Morosini<sup>3,5</sup>, Patricia Ruiz-Garbajosa<sup>3,5</sup>, Adrian Egli<sup>8,9</sup>, Rafael Cantón<sup>3,5</sup>, Patricia Muñoz<sup>1,6,7</sup>, David Rodríguez-Temporal<sup>1†</sup> and Belén Rodríguez-Sánchez<sup>1</sup>.

<sup>1</sup>Clinical Microbiology and Infectious Diseases Department, Hospital General Universitario Gregorio Marañón, Madrid, Spain, and Institute of Health Research Gregorio Marañón (IISGM), Madrid, Spain

<sup>2</sup>Department of Signal Theory and Communication, University Carlos III of Madrid, Madrid, Spain.

<sup>3</sup>Servicio de Microbiología. Hospital Universitario Ramón y Cajal and Instituto Ramón y Cajal de Investigación Sanitaria (IRYCIS), Madrid, Spain.

<sup>4</sup>Clover Bioanalytical Software, Av. del Conocimiento, 41, 18016 Granada, Spain

<sup>5</sup>CIBER en Enfermedades Infecciosas, Madrid, Spain

<sup>6</sup>CIBER de Enfermedades Respiratorias (CIBERES CB06/06/0058), Madrid, Spain.

<sup>7</sup>Medicine Department, Faculty of Medicine, Universidad Complutense de Madrid, Madrid, Spain

<sup>8</sup>Applied Microbiology Research, Department of Biomedicine, University of Basel, Basel, Switzerland.

<sup>9</sup>Division of Clinical Bacteriology and Mycology, University Hospital Basel, Basel, Switzerland.

<sup>10</sup>Centro de Investigación Biomédica en Red de Epidemiología y Salud Pública (CIBERESP), Madrid, Spain.

**Running Title:** *E. cloacae* complex classification with MALDI-TOF MS

**†Corresponding author:**

David Rodríguez-Temporal, PhD.

- 29 Servicio de Microbiología Clínica y Enfermedades Infecciosas. Hospital General Universitario
- 30 Gregorio Marañón.Dr. Esquerdo 46. 28007 Madrid, Spain
- 31 Phone: +34- 91- 426 7163, Fax: +34- 91- 426 9595
- 32 E-mail: [david.rodriquez@iisgm.com](mailto:david.rodriquez@iisgm.com)
- 33 \*These authors contributed equally to the study

## 34 ABSTRACT

35 The *Enterobacter cloacae* complex (ECC) encompasses heterogeneous clusters of  
36 species that have been associated with nosocomial outbreaks. These species may  
37 host different acquired antimicrobial resistance and virulence mechanisms and their  
38 identification are challenging. This study aims to develop predictive models based on  
39 MALDI-TOF MS spectral profiles and machine learning for species-level identification.

40 A total of 198 ECC and 116 *K. aerogenes* clinical isolates from the University Hospital  
41 Ramón y Cajal (Spain) and the University Hospital Basel (Switzerland) were included.  
42 The capability of the proposed method to differentiate the most common ECC species  
43 (*E. asburiae*, *E. kobei*, *E. hormaechei*, *E. roggenkampii*, *E. ludwigii*, *E. bugandensis*)  
44 and *K. aerogenes* was demonstrated by applying unsupervised hierarchical clustering  
45 with PCA pre-processing. We observed a distinctive clustering of *E. hormaechei* and *K.*  
46 *aerogenes* and a clear trend for the rest of the ECC species to be differentiated over  
47 the development dataset. Thus, we developed supervised, non-linear predictive models  
48 (Support Vector Machine with Radial Basis Function and Random Forest). The external  
49 validation of these models with protein spectra from the two participating hospitals  
50 yielded 100% correct species-level assignment for *E. asburiae*, *E. kobei*, and *E.*  
51 *roggenkampii* and between 91.2% and 98.0% for the remaining ECC species. Similar  
52 results were obtained with the MSI database developed recently ([https://msi.happy-](https://msi.happy-dev.fr/)  
53 [dev.fr/](https://msi.happy-dev.fr/)) except in the case of *E. hormaechei*, which was more accurately identified by  
54 Random Forest.

55 In short, MALDI-TOF MS combined with machine learning demonstrated to be a rapid  
56 and accurate method for the differentiation of ECC species.

57  
58 **Keywords:** *Enterobacter* species; mass spectrometry; MALDI-TOF MS; Peak analysis;  
59 Machine Learning

## 60 INTRODUCTION

61 *Enterobacter* is a facultative anaerobic Gram-negative genus that can be found as  
 62 a natural commensal in the gut microbiome of mammals (1). Several species have  
 63 been associated with nosocomial outbreaks causing urinary tract infection, skin and  
 64 soft tissue infection, pneumonia, and bacteremia (2, 3). *Enterobacter cloacae* complex  
 65 (ECC) is of particular clinical interest. This group is composed of 13 heterogenic  
 66 genetic clusters according to *hsp60* gene sequencing: *E. asburiae* (cluster I), *E. kobei*  
 67 (cluster II), *E. hormaechei* subsp. *hoffmannii* (cluster III), *E. roggenkampii* (cluster IV),  
 68 *E. ludwigii* (cluster V), *E. hormaechei* subsp. *oharae*, and subsp. *xiangfangensis*  
 69 (cluster VI), *E. hormaechei* subsp. *hormaechei* (cluster VII), *E. hormaechei* subsp.  
 70 *steigerwaltii* (cluster VIII), *E. bugandensis* (cluster IX), *E. nimipressuralis* (cluster X), *E.*  
 71 *cloacae* subsp. *cloacae* (cluster XI), *E. cloacae* subsp. *dissolvens* (cluster XII), and a  
 72 heterogeneous group of *E. cloacae* sequences are considered as cluster XIII.  
 73 However, the taxonomy of this genus is still under debate (4, 5). In fact, *Enterobacter*  
 74 *aerogenes* has been recently reclassified into the *Klebsiella* genus as *K. aerogenes* (6).  
 75 A more comprehensive study based on whole-genome sequencing (WGS) data from  
 76 ECC isolates yielded a redistribution of the species defined by *hsp60* sequencing (5)  
 77 into 22 clades (7) and allowed the characterization of new ECC species (8).

78 Discrimination of the ECC at the species level is usually performed by sequence-  
 79 based methods. The most commonly targeted gene is *hsp60*, although multi-locus  
 80 sequence typing (MLST) and WGS have also been applied (5, 9, 10). Sequence-based  
 81 diagnostic methods techniques are laborious and require specific equipment.  
 82 Therefore, new emerging techniques such as Matrix Assisted Laser  
 83 Desorption/Ionization Time-of-flight Mass Spectrometry (MALDI-TOF MS) have been  
 84 proposed as an alternative to sequence-based methods. MALDI-TOF MS has shown to  
 85 be an excellent methodology for bacterial identification. It can easily identify *E. cloacae*

86 complex isolates but it showed low discrimination power for the species in this group  
87 when using standard analysis and commercial databases with low resolution (11, 12).

88 This study aimed to develop and validate prediction models for the automatic  
89 species differentiation within the ECC using MALDI-TOF MS and supervised learning  
90 algorithms. This task is important because of the diverse implications of ECC species  
91 in human pathologies and their involvement in nosocomial outbreaks (4). Besides, *E.*  
92 *hormaechei*, the most encountered ECC species in the clinical settings, has been  
93 correlated with the enhanced acquisition of antimicrobial resistance mechanisms and  
94 the expression of virulence factors (13, 14). To achieve this goal, two steps were  
95 conducted in this study. First, we performed an unsupervised clustering to determine  
96 the feasibility of MALDI-TOF MS data for the ECC species identification. Second, we  
97 applied a supervised machine learning algorithm with isolates from University Hospital  
98 Ramón y Cajal, Madrid -Madrid, Spain - and validated our findings with different ECC  
99 isolates from the same hospital and from the University Hospital of Basel, Switzerland.

100

## 101 MATERIALS AND METHODS

### 102 Bacterial isolates

103 Overall, we analyzed 198 clinical isolates belonging to the ECC and nine 116 *K.*  
104 *aerogenes* (formerly *E. aerogenes*). Among them, 164 ECC and 9 *K. aerogenes* were  
105 collected in a surveillance study of antimicrobial resistance in the Hospital Universitario  
106 Ramón y Cajal -UHRC- (Madrid, Spain) between 2005 and 2018 and identified by  
107 partial sequencing of the *hsp60* gene (15). We collected the remaining isolates (34  
108 ECC and 107 *K. aerogenes*) at the University Hospital Basel (UHB; Basel, Switzerland)  
109 between 2016 and 2021 and identified the isolates by whole genome sequencing  
110 (WGS) using KmerFinder 3.2 (16-18). MALDI-TOF MS spectral profiles of these

111 isolates were obtained in Basel and submitted to the Hospital General Universitario  
112 Gregorio Marañón (Madrid, Spain) for further analysis.

113 All isolates from UHRC were incubated overnight at 37°C and metabolically  
114 activated after three subcultures on Columbia Blood Agar (bioMérieux, Marcy l'Etoile,  
115 France) before their analysis with MALDI-TOF MS at the Hospital General Universitario  
116 Gregorio Marañón.

### 117 **Spectra acquisition using MALDI-TOF MS**

118 We identified the isolates using the MBT Smart MALDI Biotyper (Bruker  
119 Daltonics, Bremen, Germany). We spotted all strains from UHRC in duplicate onto the  
120 MALDI target plate and overlaid with 1 µl of 70% formic acid. After drying at room  
121 temperature, we covered and dried the spots with 1 µl HCCA matrix, according to the  
122 manufacturer's indications (Bruker Daltonics). We obtained two spectra in the range of  
123 2,000-20,000 Da on each spot, resulting in 4 spectra per isolate. The isolates from  
124 UHB were analysed in one spot per strain and one spectrum from spot was acquired.

### 125 **Data processing of MALDI-TOF MS protein spectra and development of** 126 **predictive models**

127 For both, feasibility, and supervised studies, we processed all MALDI-TOF MS  
128 spectral profiles with the Clover MS Data Analysis software (Clover Biosoft, Granada,  
129 Spain). We applied pre-processing pipeline to all protein spectra that consisted of: 1)  
130 smoothing -Savitzky-Golay Filter: window length=11; polynomial order=3- and baseline  
131 subtraction -Top-Hat filter method with factor=0.02-; 2) creation of an average  
132 spectrum per isolate; 3) alignment of the average spectra from different isolates -shift:  
133 medium; constant tolerance: 2 Da; linear mass tolerance: 600 ppm-; 4) normalization  
134 by Total Ion Current (TIC).

### 135 Unsupervised feasibility study

136 To study the feasibility of MALDI-TOF MS for differentiation of ECC species, we  
137 proposed an unsupervised study based on Principal Component Analysis (PCA), and t-  
138 distributed Stochastic Neighbour Embedding (t-SNE). For this purpose, an  
139 oversampled balanced dataset of each ECC species was used. We included in total,  
140 126 spectra from the 7 ECC species analysed in this study (sourcing from UHRC and  
141 UHB), as indicated in **Table 1**.

#### 142 Supervised model development

143 Once the feasibility of the study was determined, we proposed the supervised  
144 model development. In this case, three different datasets were created: training  
145 validation set, internal validation set, and external validation set. The details of these  
146 datasets are shown in **Table 1**.

147 Due to the lack of validation samples of *E. ludwigii* and *E. bugandensis*, these  
148 were not included in the development of the supervised model. Therefore, our  
149 supervised model was developed to predict the five ECC species: *E. asburiae* (cluster  
150 I), *E. kobei* (cluster II), and *E. hormaechei* (clusters III, VI, and VIII considered  
151 together), and *E. roggenkampii* (cluster IV). We applied three different supervised  
152 models: Partial Least Square-Discriminant Analysis (PLS-DA), Support Vector Machine  
153 (SVM) with Linear kernel (SVM-L) and with Radial Basis Function (SVM-R) kernel, and  
154 Random Forest (RF). The hyperparameter selection was performed by a 5-fold cross-  
155 validation technique.

156 Finally, we performed two external validations of the predictive models. First,  
157 126 MALDI-TOF MS from UHRC and then 141 MALDI-TOF MS from UHB were blindly  
158 classified by the same predictive models using Clover BioSoft v0.6.1. This software  
159 uses the scikit-learn 0.23.2 python library to implement all statistical methods used in  
160 this study. For reproducibility purposes under FAIR principles, free access to all spectra

161 and to reproduce the analyses in this work can be found at the following url:  
162 <https://platform.clovermsdataanalysis.com/public-repository>.

### 163 MSI Database

164 Recently, an online database has been developed for the rapid differentiation of  
165 ECC species based on their MALDI-TOF MS protein profile (19). This database has  
166 free access (<https://msi.happy-dev.fr/>) and has been built using protein spectra from 42  
167 ECC isolates characterized by sequencing the *hsp60* gene. This identification method  
168 is considered the state-of-the-art method for the identification of ECC isolates at the  
169 species level. Therefore, both external validation datasets were also identified using  
170 the MSI database as a comparison to the methods proposed in this article. As stated  
171 above, MALDI-TOF MS spectra associated to this study have been also made publicly  
172 available.

### 173 **Ethics statement**

174 The Ethics Committee of the Gregorio Marañón Hospital (CEIm) evaluated this  
175 project and considered that all the conditions for waiving informed consent were met  
176 since the study was conducted with microbiological samples and not with human  
177 products. At the University Hospital Basel only anonymized data was used with the  
178 purpose of quality control and assay validation. According to the Swiss Human  
179 Research Act no specific consent is required in this case. Data was either acquired in  
180 routine microbiological diagnostics (excluding cases with a rejected general consent) or  
181 used from a previously published dataset (DRIAMS).

182

## 183 **RESULTS**

### 184 **Feasibility study**



185 To prove the feasibility of MALDI-TOF MS to differentiate ECC species, an  
186 unsupervised hierarchical clustering with PCA and t-SNE pre-processing was  
187 performed (**Figure 1**). The protein spectra of the seven ECC species (*E. asburiae*, *E.*  
188 *kobei*, *E. hormaechei*, *E. roggkampii*, *E. ludwigii*, *E. bugandensis*, and *K.*  
189 *aerogenes*), which are equally represented in the model, were compared. The  
190 dendrogram built with these data showed three main clusters: one containing *E.*  
191 *hormaechei*, a second cluster with *K. aerogenes*, and a third cluster with the rest of the  
192 species. Inside the latter cluster, *E. bugandensis* strains were clustered together and  
193 so did *E. asburiae*, *E. ludwigii*, *E. kobei*, and *E. roggkampii*, although in these four  
194 cases some of the spectra were clustered with the wrong species (**Figure 1A, 1B and**  
195 **1C**).

196 The implementation of PCA to reduce the dimensionality showed that 14  
197 components were needed to explain 95% of the variance (**Figure 1D**). This fact and  
198 the relatively accurate classification of ECC species using an unsupervised algorithm  
199 demonstrated the potentiality of MALDI-TOF MS to differentiate ECC species.

## 200 **Supervised models based on MALDI-TOF MS.**

201 To solve the limitations of unsupervised learning, we added the label knowledge  
202 to the training phase by using supervised algorithms such as PLS-DA, SVM, and RF.  
203 We trained these models using the development dataset shown in **Table 1**, and  
204 selected their hyperparameters by a 5-fold cross-validation technique. This cross-  
205 validation process led to the next hyperparameter selection: for PLS-DA 2 components  
206 were used, for SVM-L the value of C was 10, and for SVM-R the value of C was 10 and  
207 the value of  $\gamma$  was 1000. **Table 2** shows the results obtained for the internal 5-fold  
208 cross-validation, which have been further detailed in **Table S1**.

209 Both *E. hormaechei* and *K. aerogenes* presented the same trend as in the feasibility  
210 study and their differentiation was 100% using non-linear approaches (SVM-R and RF)

211 Only the implementation of a linear approach (SVM-L) yielded lower results for *K.*  
212 *aerogenes* (**Table 2**). For the rest of the analysed ECC species, we also obtained  
213 100% correct classification by the application of non-linear approaches (**Table 2**).

214 In **Figure 2**, the distance between samples calculated by the RF classifier is shown.  
215 We detected a unique cluster for each species. Due to the results presented in **Table**  
216 **2**, only SVM-R and RF were considered for further analysis. **Table 3** shows the results  
217 of SVM-R and RF for the validation dataset collected at UHRC and UHB.

218 Both algorithms, SVM-R and RF, yielded the same results in the external  
219 validation performed on the MALDI-TOF MS spectra from the validation dataset  
220 sourcing from the same hospital (UHRC). In this case, all *K. aerogenes*, *E. asburiae*,  
221 and *E. kobei* isolates were correctly classified meanwhile one *E. hormaechei* strain was  
222 misclassified as *E. kobei* with both algorithms. For *E. roggenkampii*, two isolates were  
223 misclassified as *E. hormaechei* and one as *E. kobei* (**Table S2**). The accuracy of the  
224 model is shown in **Figures 3A and 3B**.

225 Since SVM-R and RF algorithms performed equally, both of them were  
226 considered for external validation with MALDI-TOF MS spectral profiles obtained at the  
227 UHB. In this case, 91.2% of the *E. hormaechei* (n=33), 100% of the *E. roggenkampii*  
228 (n=1), and 98.1% of the *K. aerogenes* isolates were correctly classified by RF, as  
229 shown in **Table 3**. The application of SVM-R yielded lower results for *E. hormaechei*  
230 and *K. aerogenes*. **Figures 3C and 3C** show the accuracy of both SVM-R and RF for  
231 the external validation collection from UHB.

## 232 Identification of the ECC isolates using the MSI Database

233 Finally, the MSI Platform was also used as an identification tool for ECC  
234 species to compare the automatic approach developed in this study versus the current  
235 state-of-the-art method (19). Among the UHRC isolates, the identification rate for *E.*  
236 *asburiae*, *E. hormaechei*, and *E. kobei* was similar to the rates yielded by the predictive

models developed in this study (**Table 3**). The only difference detected between both methods was that the MSI database correctly classified one more *E. roggenkampii* isolate. As for the protein spectra sourcing from the UHB, 84.8% of the *E. hormaechei* and 100% of the *E. roggenkampii* isolates (n=1) were correctly identified with the MSI platform. In this case, the MSI platform provided a lower rate of correct identifications for *E. hormaechei* than the RF algorithm (**Table 3**).

## DISCUSSION

In this study, the implementation of supervised, non-linear algorithms (SVM-R and RF) to MALDI-TOF MS spectra allowed the correct species assignment of 100% isolates belonging to two ECC species (*E. asburiae* and *E. kobei*) and between 91.2% and 98.1% for *E. hormaechei*, *E. roggenkampii*, and *K. aerogenes* (formerly *E. aerogenes*) sourcing from two different hospitals.

Poor discrimination of *E. cloacae* complex species by MALDI-TOF MS has been previously reported either by using commercial (11, 15) or enriched, in-house databases (20). However, a recent study from a research group with broad experience in MALDI-TOF MS and the creation of expanded libraries reported 92.0% correct species-level identification by implementing a specific in-house library enriched with well-characterized ECC strains and correct discrimination of 97.0% *E. hormaechei* isolates (19). This approach can be useful for the discrimination of close-related species, but the construction of a database is cumbersome and requires highly trained personnel. The implementation of the MSI platform allowed 94.9% correct species-level identification of 155 ECC protein spectra in this study. This rate was slightly lower than the obtained with the non-linear algorithms proposed by our approach.

In this study, we demonstrate the feasibility of MALDI-TOF MS to identify species within the ECC. First, hierarchical clustering showed that it is possible to differentiate between species using the information contained in MALDI-TOF MS as

263 reported before (20). Secondly, a supervised study using machine learning algorithms  
264 yielded the correct classification of all ECC species. Therefore, different supervised  
265 classification algorithms were implemented to correctly provide species assignment of  
266 ECC species. The internal validation experiment demonstrated that non-linear  
267 approaches, such as SVM-R or RF, were needed to correctly identify all species. Both  
268 models perfectly classified all samples in internal cross-validation.

269 To further demonstrate that the model can perform in different scenarios with  
270 data different than the spectral profile used for model training, we performed two  
271 validation experiments. First, we carried out a validation with MADI-TOF MS protein  
272 spectra sourcing from UHRC. From a total number of samples of 116, both SVM-R and  
273 RF only misclassified four isolates, i.e., a 96.5% of accuracy in classifying species  
274 within the ECC was yielded. Secondly, we performed an external validation with  
275 MALDI-TOF MS sourcing from UHB to simulated a real-world scenario. These MALDI-  
276 TOF MS protein spectra originated in a different epidemiological scenario and were  
277 processed by operators from the UHB. In this case, SVM-R showed to be overfitted to  
278 the UHRC distribution, which was already pointed out by the value of  $\gamma$  value, scoring  
279 an 83.7% of accuracy. On the other hand, the current state-of-the-art tool -the MSI  
280 database- performed better than SVM-R with 94.9% of accuracy although it was not  
281 able to distinguish the *K. aerogenes* (19). However, RF outperformed both approaches  
282 with over 96.0% of accuracy in identifying the three species. Hence, it is demonstrated  
283 that supervised machine learning algorithms are feasible and, indeed, applicable in  
284 microbiology laboratories

285 One limitation of this study was the fact that all UHRC isolates were  
286 carbapenemase producing isolates, because this was the source of the previously  
287 analysed collection (15). However, the present study provides the first proof of concept  
288 for differentiating ECC species based on machine learning. For a definitive validation,  
289 improvement, and implementation of these predictive models, future studies will involve

290 strains from a more diverse epidemiological and geographical origin and  
291 characteristics. Besides, not all analysed ECC species could be represented in the  
292 external validation dataset due to the lack of isolates from the species *E. ludwigii* and  
293 *E. bugandensis*.

294 The present study provides promising results for differentiating ECC species  
295 based on machine learning and MALDI-TOF MS protein spectra. It also highlights the  
296 facts that MALDI-TOF MS data should be linked to WGS data in order to allow future  
297 work and providing a reference standard. The MALDI-TOF MS and machine learning  
298 approach has been demonstrated to be a rapid and cost-effective method, suitable for  
299 correct species-level assignment of closely-related species, as in the case of ECC. The  
300 use of spectra analysis tools is becoming user-friendly and easy to apply and its use  
301 may provide species-level identification in a fast and inexpensive way. Once the model  
302 is validated with a comprehensive number of ECC species, an open web application  
303 will be deployed to be used by the community freely.

304

## 305 **CONFLICT OF INTERESTS**

306 The authors have no conflict of interests.

307

## 308 **ACKNOWLEDGMENTS**

309 This work was supported by the projects PI15/01073 and PI18/00997 from the Health  
310 Research Fund (FIS. Instituto de Salud Carlos III. Plan Nacional de I+D+I 2013-2016)  
311 of the Carlos III Health Institute (ISCIII, Madrid, Spain) partially financed by the  
312 European Regional Development Fund (FEDER) 'A way of making Europe'. MM was  
313 funded by the Community of Madrid (Programa de Garantía Juvenil, PEJD-2017-  
314 PRE/BMD-5106) and DRT with a postdoc contract from the Intramural Funding

315 Program of the IISGM. BRS (CPII19/00002) is a recipient of a Miguel Servet contract  
316 supported by the FIS. AGL was funded by a predoc contract from the Intramural  
317 Funding Program of the IISGM.

318

## 319 **AUTHOR CONTRIBUTIONS**

320 Ana Candela: experimental part, formal analysis, data collection, validation,  
321 visualization, writing – original draft preparation and review/editing. Alejandro Guerrero-  
322 López: formal analysis, data collection, validation, visualization, writing – original draft  
323 preparation and review/editing. Miriam Mateos and Alicia Gómez-Asenjo: experimental  
324 part, formal analysis, and data collection. Manuel J. Arroyo, Gema Méndez, and Luis  
325 Mancera: data analysis, validation, writing – original draft preparation and  
326 review/editing. Marta Hernández-García, Rosa del Campo: experimental part, formal  
327 analysis, and writing, submission of isolates, original draft preparation, and  
328 review/editing. Aline Cuénod and Adrian Egli: submission of isolates, original draft  
329 preparation, and review/editing. Juan de Dios Caballero, Laura Martínez-García,  
330 Desirée Gijón, María Isabel Morosini, Patricia Ruiz-Garbajosa, Rafael Cantón and  
331 Patricia Muñoz: validation, writing and review/editing. Emilia Cercenado:  
332 conceptualization, formal analysis, validation, writing, and review/editing. David  
333 Rodríguez-Temporal: conceptualization, formal analysis, validation, original draft  
334 preparation and review/editing; Belén Rodríguez-Sánchez: conceptualization, project  
335 administration, formal analysis, supervision, validation, visualization, original draft  
336 preparation and review/ editing.

337

## 338 **FIGURE LEGENDS**

339 **Figure 1. A.** Dendrogram built with 126 MALDI-TOF MS spectra. **B.** PCA of feasibility  
340 study spectra. **C.** t-SNE of study spectra. **D.** PCA Eigenvalues showing the variance of  
341 each component.

342 **Figure 2.** MALDI-TOF MS Euclidean distance between species by Random Forest  
343 classifier.

344 **Figure 3.** ROC Curves and AUC values for both SVM-R (A and C) and RF models (B  
345 and D) were applied to the external validations with MALDI-TOF MS protein spectra  
346 from UHRC (A and B) and UHB (C and D).

347

## 348 REFERENCES

- 349 1. Mezzatesta ML, Gona F, Stefani S. 2012. Enterobacter cloacae complex: clinical impact  
350 and emerging antibiotic resistance. Future Microbiol 7:887-902.
- 351 2. Akbari M, Bakhshi B, Najari Peerayeh S. 2016. Particular Distribution of Enterobacter  
352 cloacae Strains Isolated from Urinary Tract Infection within Clonal Complexes. Iran  
353 Biomed J 20:49-55.
- 354 3. Kremer A, Hoffmann H. 2012. Prevalences of the Enterobacter cloacae complex and its  
355 phylogenetic derivatives in the nosocomial environment. Eur J Clin Microbiol Infect Dis  
356 31:2951-5.
- 357 4. Davin-Regli A, Lavigne JP, Pages JM. 2019. Enterobacter spp.: Update on Taxonomy,  
358 Clinical Aspects, and Emerging Antimicrobial Resistance. Clin Microbiol Rev 32.
- 359 5. Hoffmann H, Roggenkamp A. 2003. Population genetics of the nomenspecies  
360 Enterobacter cloacae. Appl Environ Microbiol 69:5306-18.
- 361 6. Tindall BJ, Sutton G, Garrity GM. 2017. Enterobacter aerogenes Hormaeche and  
362 Edwards 1960 (Approved Lists 1980) and Klebsiella mobilis Bascomb et al. 1971  
363 (Approved Lists 1980) share the same nomenclatural type (ATCC 13048) on the  
364 Approved Lists and are homotypic synonyms, with consequences for the name  
365 Klebsiella mobilis Bascomb et al. 1971 (Approved Lists 1980). Int J Syst Evol Microbiol  
366 67:502-504.
- 367 7. Sutton GG, Brinkac LM, Clarke TH, Fouts DE. 2018. Enterobacter hormaechei subsp.  
368 hoffmannii subsp. nov., Enterobacter hormaechei subsp. xiangfangensis comb. nov.,  
369 Enterobacter roggenkampii sp. nov., and Enterobacter muelleri is a later heterotypic  
370 synonym of Enterobacter asburiae based on computational analysis of sequenced  
371 Enterobacter genomes. F1000Res 7:521.
- 372 8. Wu W, Feng Y, Zong Z. 2019. Characterization of a strain representing a new  
373 Enterobacter species, Enterobacter chengduensis sp. nov. Antonie Van Leeuwenhoek  
374 112:491-500.
- 375 9. Singh NK, Bezdan D, Checinska Sielaff A, Wheeler K, Mason CE, Venkateswaran K.  
376 2018. Multi-drug resistant Enterobacter bugandensis species isolated from the  
377 International Space Station and comparative genomic analyses with human pathogenic  
378 strains. BMC Microbiol 18:175.

- 379 10. Hoffmann H, Stindl S, Ludwig W, Stumpf A, Mehlen A, Heesemann J, Monget D,  
380 Schleifer KH, Roggenkamp A. 2005. Reassignment of enterobacter dissolvens to  
381 Enterobacter cloacae as E. cloacae subspecies dissolvens comb. nov. and emended  
382 description of Enterobacter asburiae and Enterobacter kobei. Syst Appl Microbiol  
383 28:196-205.
- 384 11. Pavlovic M, Konrad R, Iwobi AN, Sing A, Busch U, Huber I. 2012. A dual approach  
385 employing MALDI-TOF MS and real-time PCR for fast species identification within the  
386 Enterobacter cloacae complex. FEMS microbiology letters 328:46-53.
- 387 12. De Florio L, Riva E, Giona A, Dedej E, Fogolari M, Cella E, Spoto S, Lai A, Zehender G,  
388 Ciccozzi M, Angeletti S. 2018. MALDI-TOF MS Identification and Clustering Applied to  
389 Enterobacter Species in Nosocomial Setting. Frontiers in microbiology 9:1885.
- 390 13. Barnes AI, Paraje MG, del CBP, Albesa I. 2001. Molecular properties and metabolic  
391 effect on blood cells produced by a new toxin of Enterobacter cloacae. Cell Biol Toxicol  
392 17:409-18.
- 393 14. Paauw A, Caspers MPM, Leverstein-van Hall MA, Schuren FHJ, Montijn RC, Verhoef J,  
394 Fluit AC. 2009. Identification of resistance and virulence factors in an epidemic  
395 Enterobacter hormaechei outbreak strain. Microbiology (Reading) 155:1478-1488.
- 396 15. Mateos M, Hernandez-Garcia M, Del Campo R, Martinez-Garcia L, Gijon D, Morosini  
397 MI, Ruiz-Garbajosa P, Canton R. 2020. Emergence and Persistence over Time of  
398 Carbapenemase-Producing Enterobacter Isolates in a Spanish University Hospital in  
399 Madrid, Spain (2005-2018). Microbial drug resistance.
- 400 16. Hasman H, Saputra D, Sicheritz-Ponten T, Lund O, Svendsen CA, Frimodt-Moller N,  
401 Aarestrup FM. 2014. Rapid whole-genome sequencing for detection and  
402 characterization of microorganisms directly from clinical samples. J Clin Microbiol  
403 52:139-46.
- 404 17. Larsen MV, Cosentino S, Lukjancenko O, Saputra D, Rasmussen S, Hasman H, Sicheritz-  
405 Ponten T, Aarestrup FM, Ussery DW, Lund O. 2014. Benchmarking of methods for  
406 genomic taxonomy. J Clin Microbiol 52:1529-39.
- 407 18. Clausen P, Aarestrup FM, Lund O. 2018. Rapid and precise alignment of raw reads  
408 against redundant databases with KMA. BMC Bioinformatics 19:307.
- 409 19. Godmer A, Benzerara Y, Normand AC, Veziris N, Gallah S, Eckert C, Morand P, Piarroux  
410 R, Aubry A. 2021. Revisiting Species Identification within the Enterobacter cloacae  
411 Complex by Matrix-Assisted Laser Desorption Ionization-Time of Flight Mass  
412 Spectrometry. Microbiol Spectr 9:e0066121.
- 413 20. Wang YQ, Xiao D, Li J, Zhang HF, Fu BQ, Wang XL, Ai XM, Xiong YW, Zhang JZ, Ye CY.  
414 2018. Rapid Identification and Subtyping of Enterobacter cloacae Clinical Isolates Using  
415 Peptide Mass Fingerprinting. Biomed Environ Sci 31:48-56.

416

417

418

419

420

421



**Table 1.** Number of ECC isolates used for the unsupervised feasibility study and the supervised model development.

ECC species	Unsupervised study		Supervised study	
	Number of balanced samples	Development dataset (UHRC)	Validation dataset (UHRC)	External validation dataset (UHB)
<i>K. aerogenes</i>	18	18	3	107
<i>E. asburiae</i>	18	18	1	0
<i>E. bugandensis</i>	18	0	0	0
<i>E. hormaechei</i>	18	18	51	33
<i>E. kobei</i>	18	18	9	0
<i>E. ludwigii</i>	18	0	0	0
<i>E. roggkampii</i>	18	18	62	1
<b>Total</b>	<b>126</b>	<b>90</b>	<b>126</b>	<b>141</b>

**Table 2.** Accuracy results for internal 5-fold cross-validation over development dataset (90 spectral profiles). PLS-DA: Partial Least Squares-Discriminant Analysis; SVM-L: Support Vector Machine-Linear kernel; SVM-R: Support Vector Machine-Radial Basis Function kernel; RF: Random Forest.

Algorithm	<i>E. asburiae</i>		<i>E. hormaechei</i>		<i>E. kobei</i>		<i>E. roggkampii</i>		<i>K. aerogenes</i>	
<b>PLS-DA</b>	9/18	50%	18/18	100%	5/18	27.8%	4/18	22.2%	18/18	100%
<b>SVM-L</b>	6/18	33.3%	18/18	100%	3/18	16.7%	6/18	33.3%	14/18	77.8%
<b>SVM-R</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>
<b>RF</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>	<b>18/18</b>	<b>100%</b>

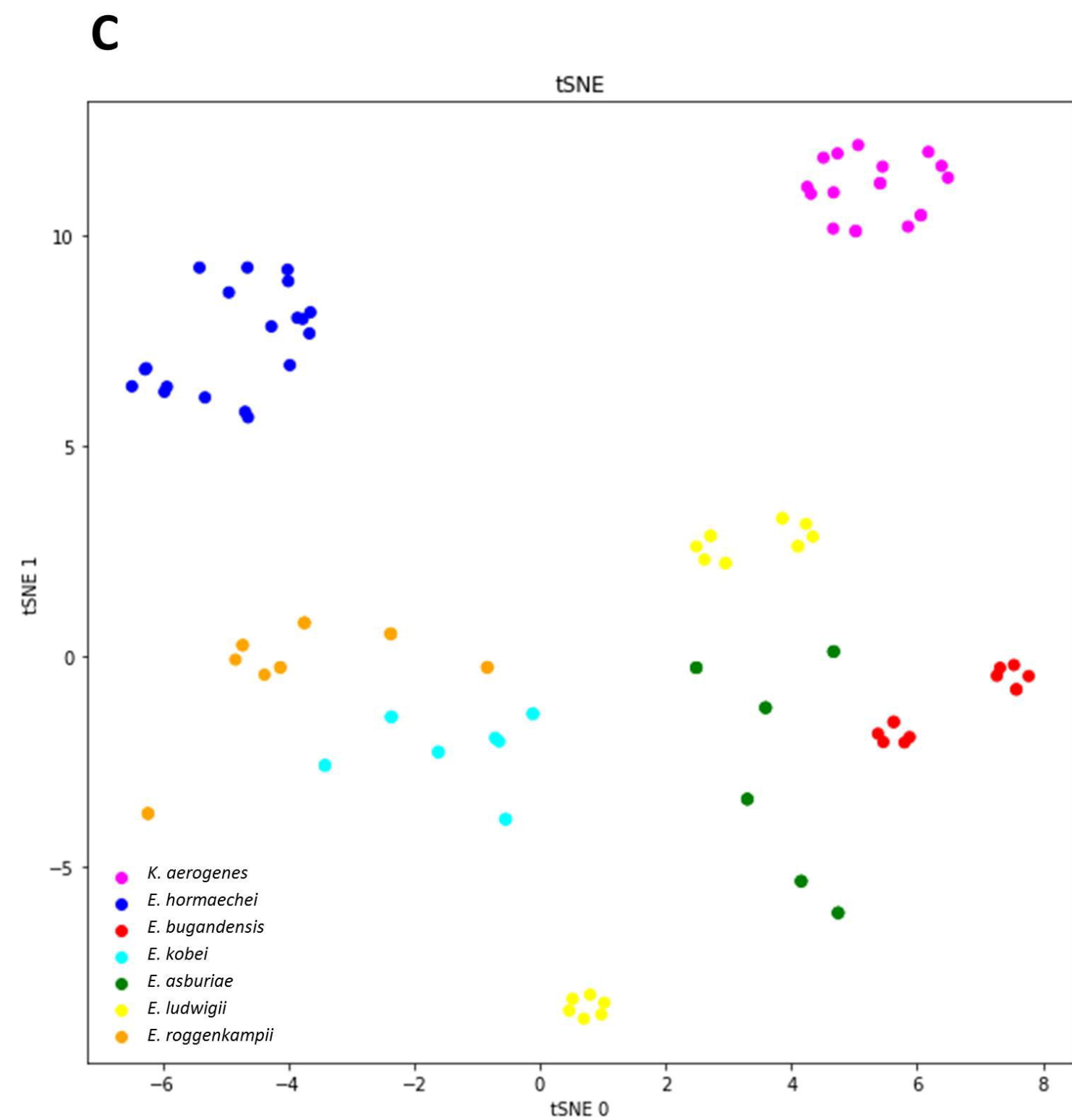
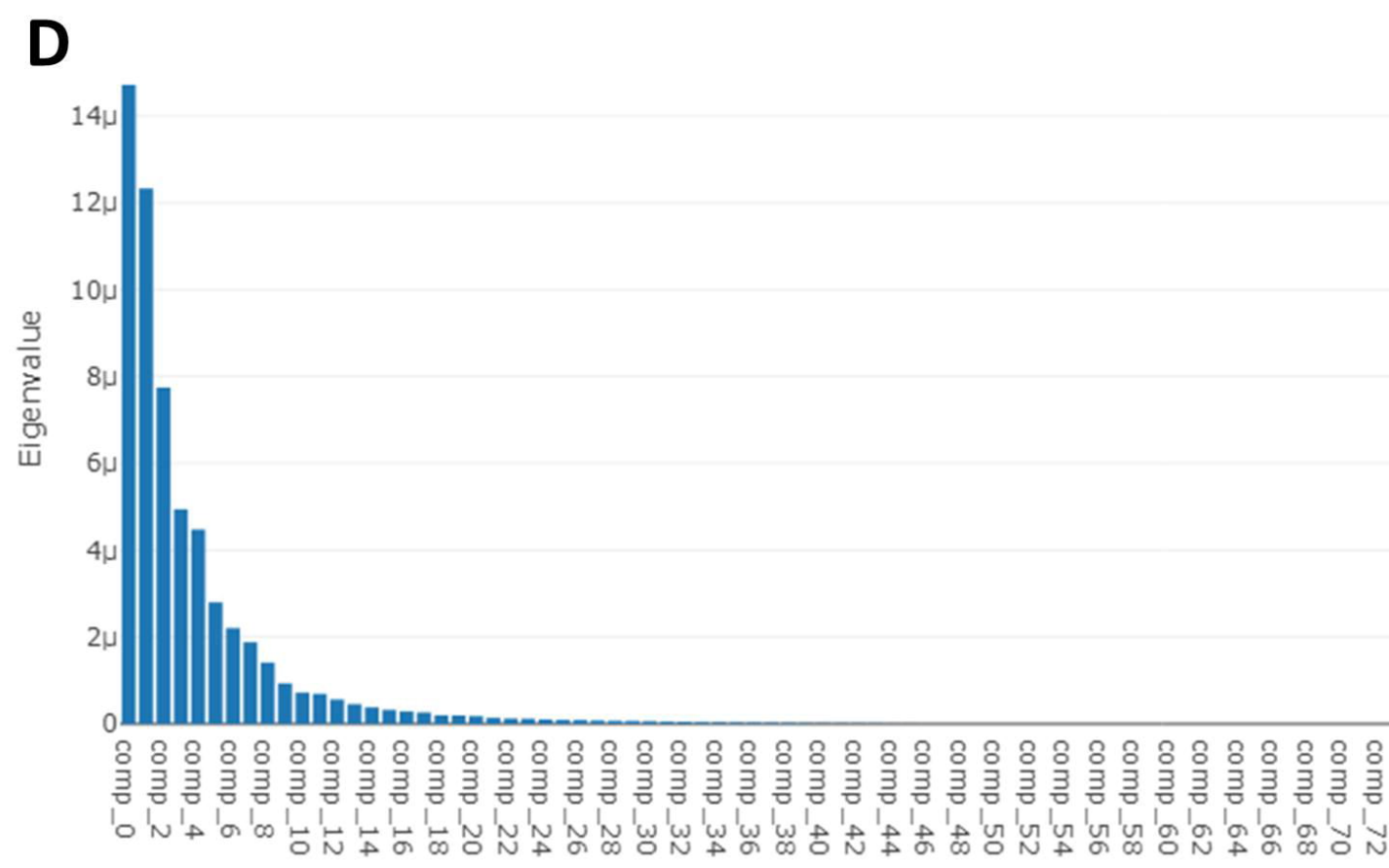
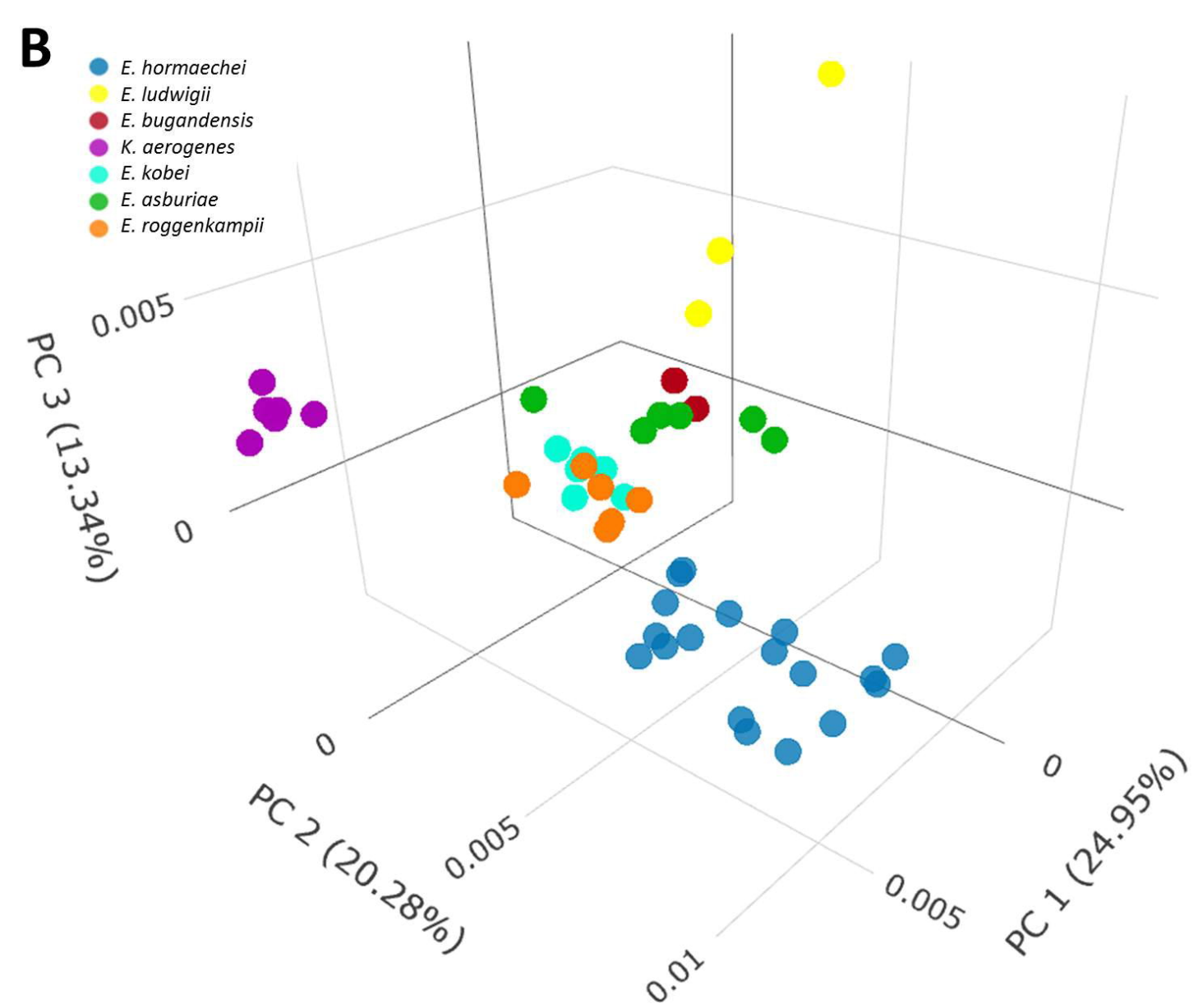
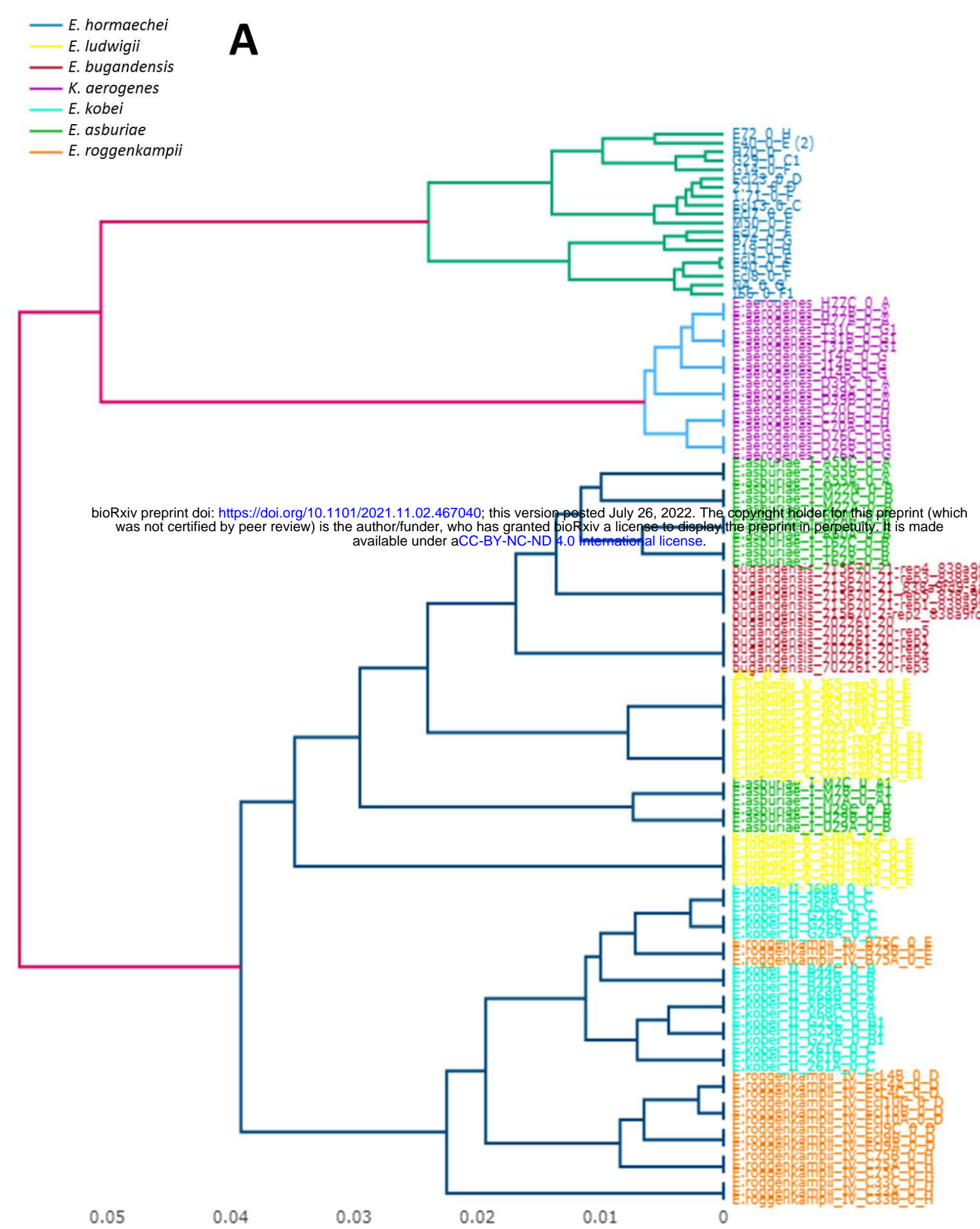
**Table 3.** Accuracy results for the validation dataset from UHRC and UHB, and the identification accuracy obtained by the MSI database. The MSI platform, specific for the identification of ECC species (19), was also applied for the classification of ECC species.

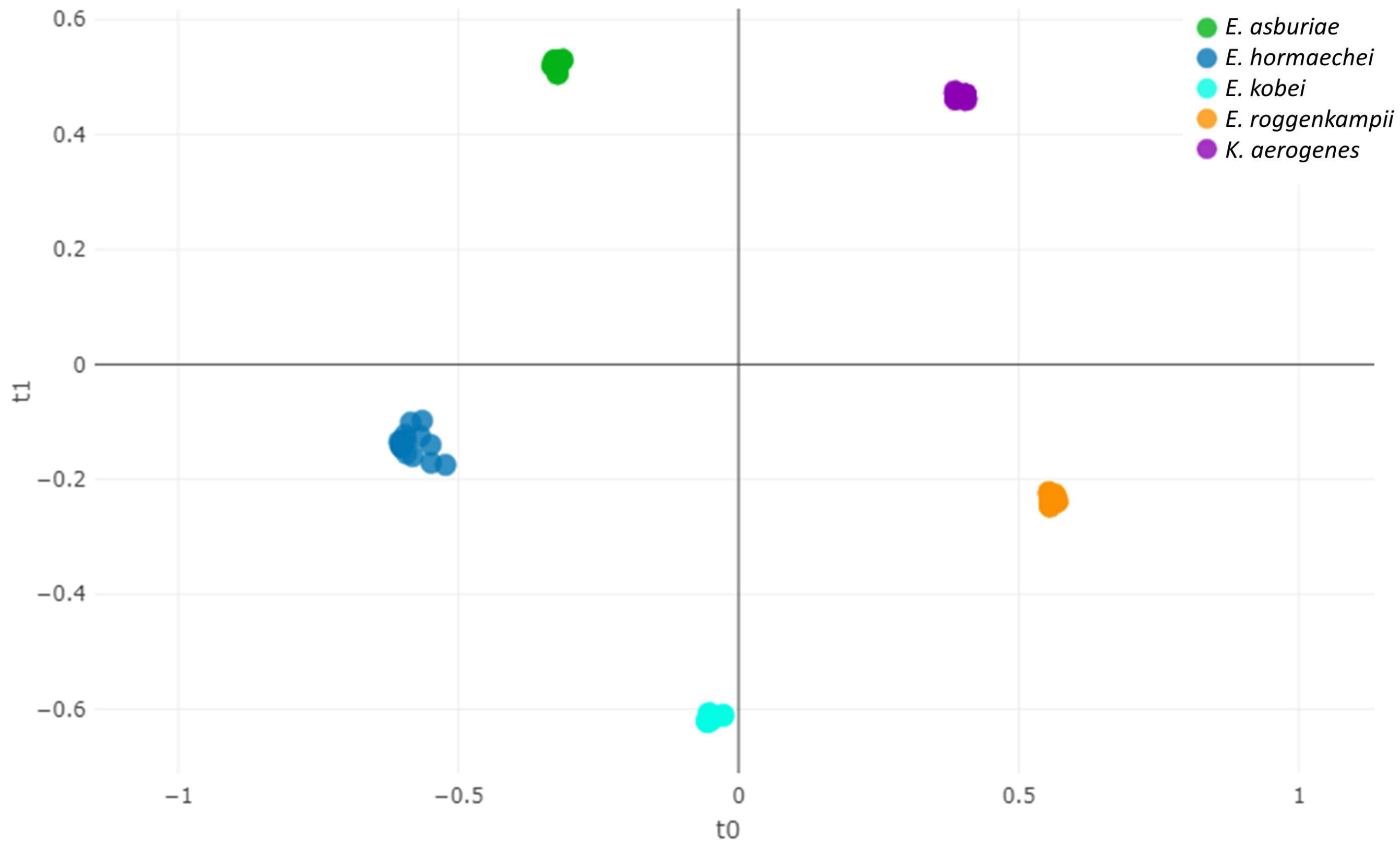
440

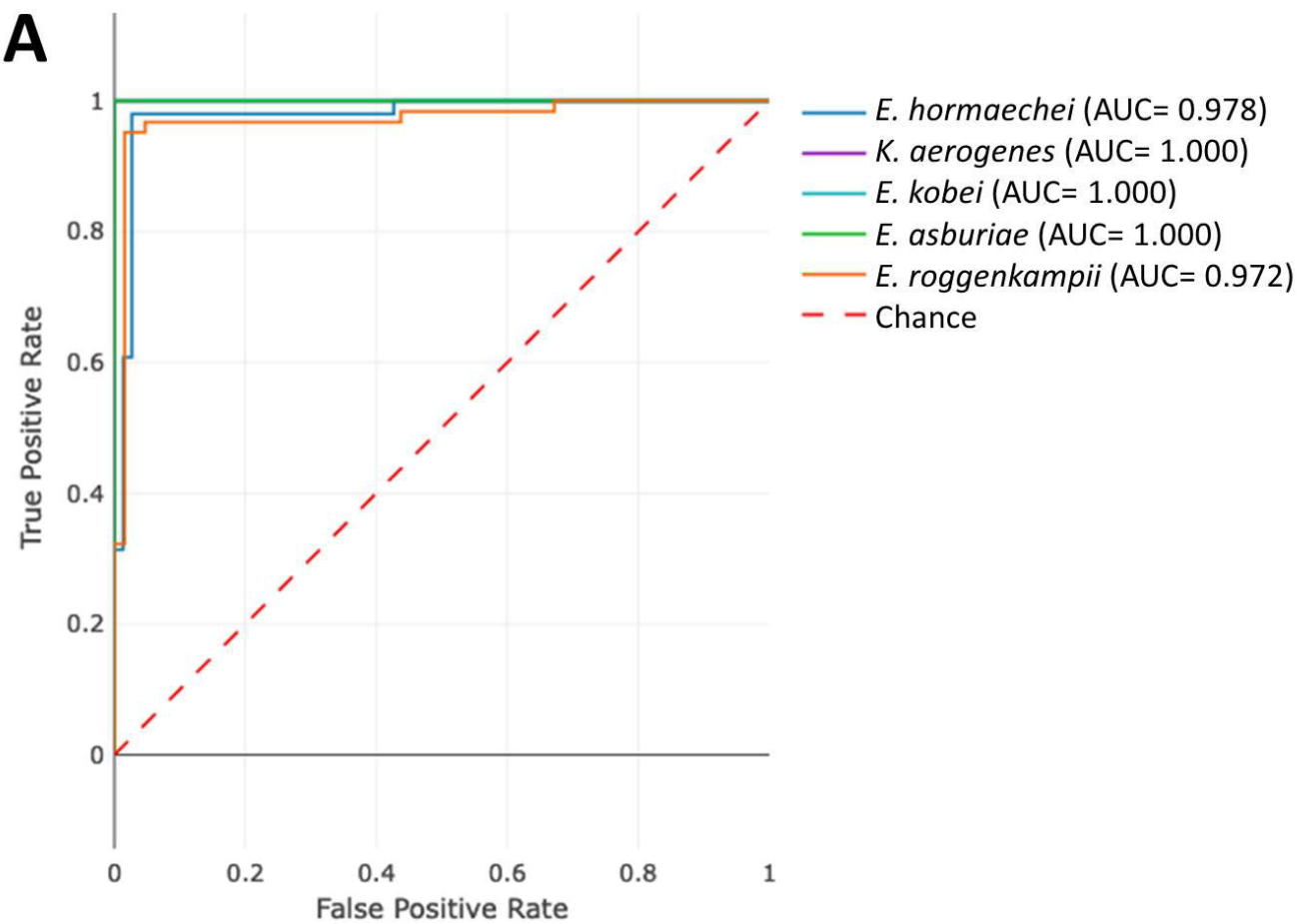
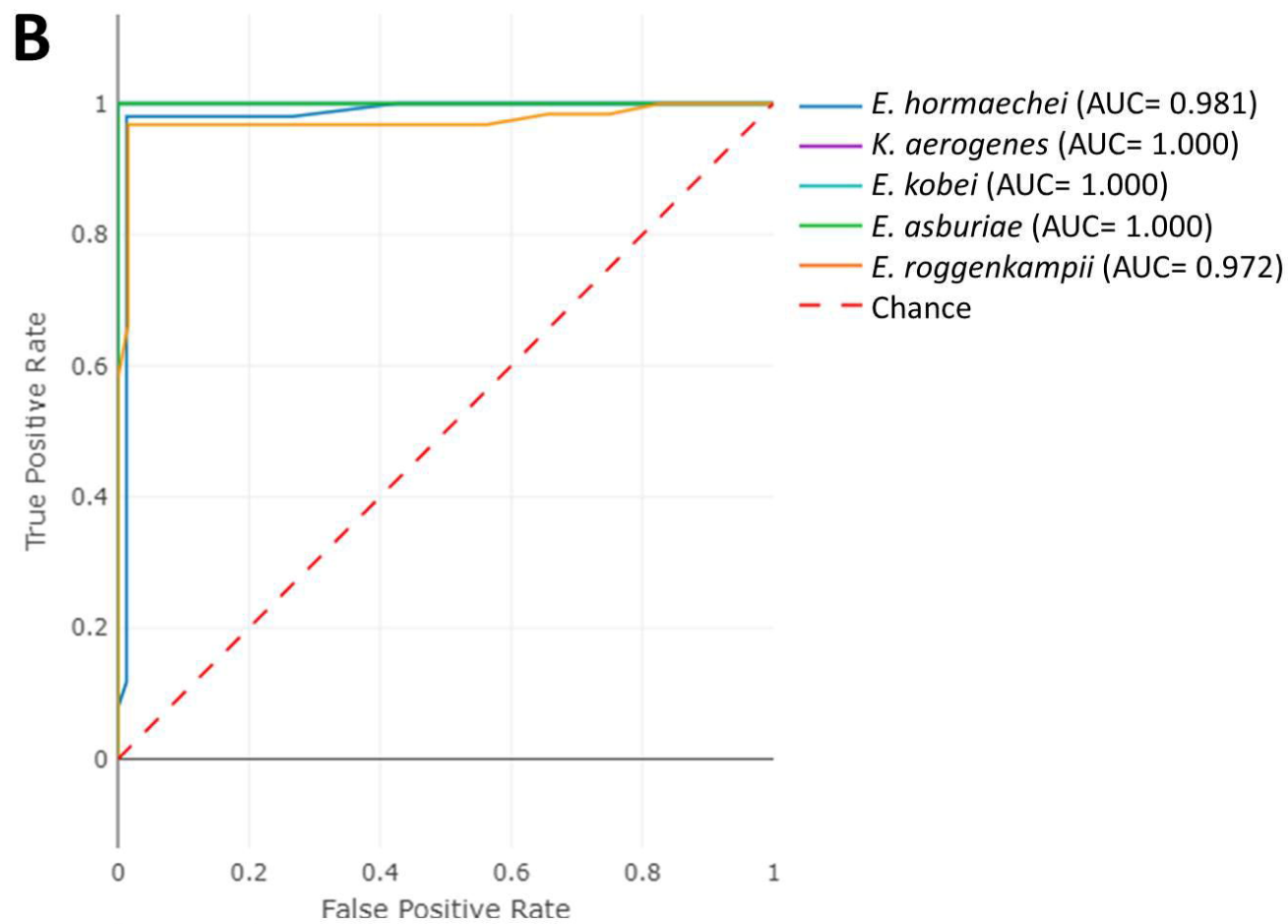
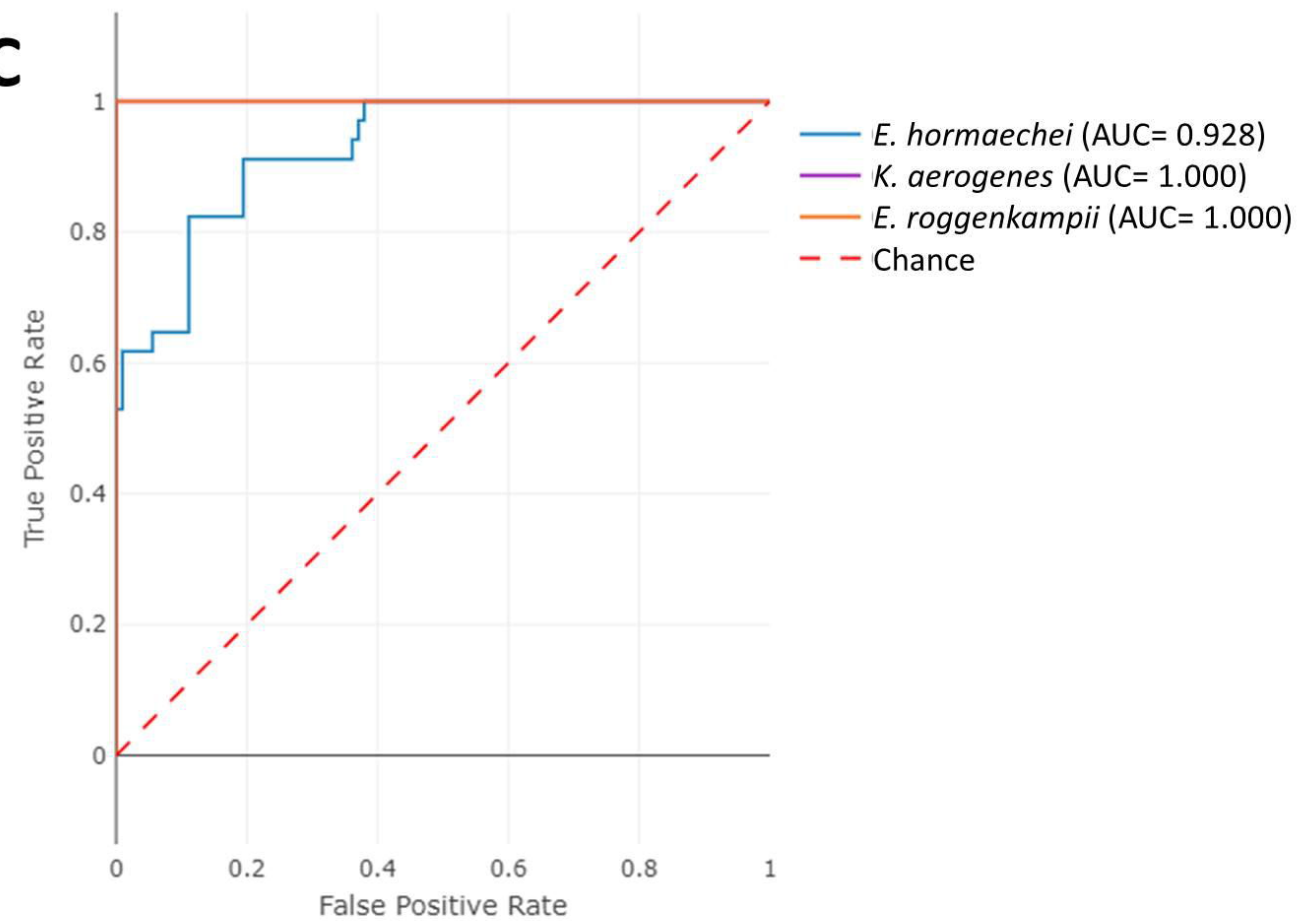
Algorithm	<i>E. asburiae</i>		<i>E. hormaechei</i>		<i>E. kobei</i>		<i>E. roggenkampii</i>		<i>K. aerogenes</i>	
UHRC										
SVM-R	1/1	100%	50/51	98.0%	9/9	100%	59/62	95.2%	3/3	100%
RF	1/1	100%	50/51	98.0%	9/9	100%	59/62	95.2%	3/3	100%
MSI	1/1	100%	50/51	98.0%	9/9	100%	60/62	96.8%	-	-
UHB										
SVM-R	-	-	15/33	44.1%	-	-	1/1	100%	102/107	95.3%
RF	-	-	31/33	91.2%	-	-	1/1	100%	105/107	98.1%
MSI	-	-	28/33	84.8%	-	-	1/1	100%	-	-

441

442





**A****B****C****D**