

1 **Neighbour Typing Using Long Read Sequencing Provides Rapid Prediction of**
2 **Sequence Type and Antimicrobial Susceptibility of *Klebsiella pneumoniae***

3
4
5 Mabel Budia-Silva^{1,2}, Amanda C. Carroll³, Hiren Ghosh¹, Allison McGeer⁴ Tommaso
6 Giani⁵, Gian Maria Rossolini⁵, Karel Brinda⁶, William P Hanage⁷, Hajo Grundmann¹,
7 Derek R MacFadden³, Sandra Reuter¹

- 8
9 1. Institute for Infection Prevention and Control, Medical Centre – University of
10 Freiburg, Freiburg, Germany
11 2. Faculty of Biology, University of Freiburg, Freiburg, Germany
12 3. The Ottawa Hospital Research Institute, Ottawa, Ontario, Canada
13 4. Sinai Health System, Toronto, Ontario, Canada
14 5. Department of Experimental and Clinical Medicine, University of Florence, and
15 Microbiology and Virology Unit, Florence Careggi University Hospital, Florence,
16 Italy
17 6. Inria, Irisa, Univ. Rennes, Rennes, France
18 7. Center for Communicable Disease Dynamics, Harvard T. H. Chan School of
19 Public Health, Boston, Massachusetts, United States
20

21
22 Corresponding Author:
23 Sandra Reuter
24 Breisacher Straße 115B · 79106 Freiburg
25 sandra.reuter@uniklinik-freiburg.de
26

27 Key Words: RASE, prediction, antibiotic, susceptibility, *Klebsiella pneumoniae*
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 **ABSTRACT**

43 The rapid genome-based diagnostic approach of long read sequencing coupled with
44 neighbor typing offers the potential to improve empiric treatment of infection. However,
45 this approach is still in development, and clinical validation is needed to support its use.
46 In this study, we present an assessment of a neighbour typing method (RASE - resistance
47 associated sequence elements) to predict lineage or sequence type (ST) and
48 antimicrobial susceptibility in real time for *Klebsiella pneumoniae* sensu lato. We analysed
49 the initial reads generated during the early phase of long read sequencing from pure
50 culture (n=99), mock communities (n=20) and metagenomic samples (n=20). RASE
51 accurately identified 69.7% and 70% of STs in pure culture and metagenomes,
52 respectively, and identified the STs of the isolates representing the highest proportion in
53 mock communities. Regarding antimicrobial susceptibility prediction, the probability of
54 susceptibility increased to 72% (95% CI 63%-80%) across all tested antibiotics, when
55 RASE predicted susceptibility, and decreased probability of susceptibility to 8.9% (95%
56 CI 6.4%-9.6%) when was indicative of a resistant phenotype. Our study confirmed that
57 genomic neighbor typing in *K. pneumoniae* sensu lato is capable of providing informative
58 predictions of ST and antibiotic susceptibility in less than ten minutes (after the start of
59 sequencing) with 200-500 reads.

60

61 **IMPORTANCE**

62 The growing burden of antimicrobial resistance is leading to high rates of mortality and
63 morbidity worldwide. This situation has made the selection of empirical antibiotic therapy
64 challenging, due to the risk of treatment failure and the overuse of last-resort antibiotics.
65 The development of new sequencing technologies is helping to reduce the waiting time
66 for a microbiological diagnosis, providing information in the early phase of bacterial
67 infections, which could help improve clinical outcomes in a time of rising antimicrobial
68 resistance. In this context, we assessed the performance of RASE (resistance associated
69 sequence elements) in *Klebsiella pneumoniae*, an opportunistic pathogen frequently
70 associated with nosocomial infections, which can rapidly acquire antibiotic resistance
71 genes. Thus, in our study we provide insights that may aid in the validation of RASE for
72 clinical use.

73 INTRODUCTION

74 The increasing rate of antimicrobial resistance (AMR) is leading to a dearth of therapeutic
75 options for the treatment of bacterial infections (1,2). In 2019, an estimated 1.27 million
76 deaths worldwide were attributed to bacterial AMR, with *Klebsiella pneumoniae* ranked
77 as the third leading pathogen associated with resistance-related mortality (3). Among
78 these cases, *K. pneumoniae* producing extended-spectrum beta-lactamases (ESBL)
79 and/or carbapenemases (CPE) were responsible for an estimated 50,000 and 100,000
80 deaths (3). Many of these multidrug-resistant *K. pneumoniae* are distributed across
81 diverse clonal lineages (4), however, only some specific clones such as sequence type
82 (ST) 258/512, ST11/347/420, ST101, ST307, ST15 and ST147 are disseminated globally,
83 making a major contribution to AMR dissemination (5).

84 To address the expansion of antibiotic-resistant organisms, the World Health
85 Organization has developed a global action plan (6). Many measures have been
86 explored, including the improvement of diagnostic tools to reduce the time required to
87 determine the etiological agent of an infection (6). Given the extended waiting times for
88 detailed pathogen identification and the urgent need for antibiotic therapy in cases of
89 systemic infections (sepsis, meningitis, etc), broad-spectrum antibiotics are used as
90 empirical treatment in order to cover more potential pathogens (7). A rapid microbiological
91 diagnosis could improve antimicrobial stewardship by an earlier selection of accurate
92 treatment, decreasing the incidence of morbidity and mortality associated with bacterial
93 infections (8), and reducing unnecessary antibiotic selective pressure arising from the use
94 of broad spectrum agents (9).

95 Currently, the average time to obtain a microbiologic diagnosis using the classic culture-
96 dependent tools is around ~54 hours (10). New sequencing techniques such as long read
97 sequencing (11), together with genomic neighbour typing, can reduce this waiting time to
98 4 hours from the moment the clinical sample is collected (12). Genomic neighbour typing
99 allows for rapid prediction of lineage and antimicrobial susceptibility of pathogens from
100 clinical specimens by comparison with known and characterized isolates. This method is
101 a two-step algorithm, which first matches sequences of samples generated in real time

102 against a database of reference genomes with a known sequence type, reference
103 phylogeny and associated antimicrobial susceptibility phenotype, and thereafter predicts
104 the probable phenotype based on the closest match and the matching quality. Since
105 closely related isolates typically share similar properties, this provides a reliable heuristic
106 for rapidly inferring the phenotype of the query pathogen. To enable this approach, the
107 software application resistance-associated sequence element (RASE) was created to
108 compare the *k-mer* content of nanopore reads to reference genomes, then calculate
109 similarity weights (13). It predicts the phenotypic profile by identifying the best-matching
110 lineage and comparing it with resistant and susceptible neighbors to determine a
111 susceptibility score (12). RASE was evaluated previously on *Streptococcus pneumoniae*
112 and *Neisseria gonorrhoeae* (12), demonstrating the ability to predict their susceptibility
113 profile within ten minutes of starting sequencing. In this study, we assessed RASE on *K.*
114 *pneumoniae* because of its critical importance to public health. We first developed a
115 RASE database and then evaluated its accuracy, sensitivity and specificity in lineage
116 calling, and explored different approaches to identify strains and resistance genes during
117 sequencing.

118 RESULTS

119 ***Building a K. pneumoniae sensu lato reference Database for RASE***

120 A reference database comprising genome sequences, sequence type (ST) and
121 antimicrobial susceptibility data for 8 antimicrobial agents of 1511 *K. pneumoniae* sensu
122 lato strains from the EuSCAPE study (14) was built using RASE software. The selection
123 was based on isolates that were recovered from several European countries, thereby
124 covering an extensive geographical range and ensuring representation of a large and
125 diverse population. The database comprises strains belonging to 319 different STs within
126 14 clonal lineages, also including *Klebsiella variicola* and *Klebsiella quasipneumoniae*.
127 The predominant STs were high-risk clones ST512, ST11, ST101, ST15, with more than
128 100 isolates aside from ST258 with 72 isolates, and covered strains with variable
129 antimicrobial susceptibility profiles, ranging from susceptibility to all tested antibiotics to
130 complex multidrug resistant profiles (Fig 1, Supplementary table 1).

131 Strains belonging to the high-risk clones (15) were typically resistant to most antibiotics
132 when compared to other clones (Fig.1), with more than 80% resistance to ciprofloxacin
133 and piperacillin-tazobactam. Similarly, a high percentage demonstrated resistance to
134 beta-lactams, with the exception of ertapenem, where susceptibility rates reached 50%
135 and 56% in ST307 and ST405, respectively. For meropenem, resistance rates exceeded
136 94% in ST512, 75% in ST258, and 56% in ST101. Regarding aminoglycosides, more
137 than 80% of ST258 and ST512 were resistant to amikacin, while a similar percentage of
138 the same clones were susceptible to gentamicin. Notably, over 60% of all high-risk clones
139 were susceptible to colistin. Other STs did not demonstrate a consistent susceptibility
140 pattern.

141 ***Comparison of the Performance of Read Extraction Methods and Evaluation of*** 142 ***Preprocessing Tools***

143 As the ultimate goal is to predict susceptibility/resistance from a small number of first
144 reads obtained, we evaluated two strategies for read extraction (nanotimes, first fastq)
145 and one strategy preprocessing (nanotimes+porechop+filtlong [NPF]) for improved
146 quality, as well as their impact on the resulting prediction time. We extracted the first reads
147 generated at the beginning of nanopore sequencing from 88 pure cultures (see Methods).
148 We compared the results obtained by RASE, specifically the phenotypic prediction of
149 eight antibiotics, with the MIC values (considered as true values), to calculate
150 performance metrics, including sensitivity and specificity of susceptibility for each
151 antibiotic. These values (Fig 2a) indicate a comparable overall performance. The
152 sensitivity and specificity values were similar across approaches, although the NPF
153 approach demonstrated the highest median sensitivity (0.75) and specificity (0.65),
154 confirming improved data quality after applying Porechop and Filtlong. No statistically
155 significant differences were found (Friedman test, sensitivity: $\chi^2 = 4.36$, $p = 0.113$;
156 specificity: $\chi^2 = 1.93$, $p = 0.381$). Post-hoc Wilcoxon tests also showed no significant
157 pairwise differences after multiple testing corrections ($p > 0.05$ for all comparisons;
158 Supplementary Table 2).

159 In addition, we analyzed computational performance regarding the time and number of
160 reads required to achieve a stable call. As shown in Fig 2b, nanotimes and first fastq
161 required approximately 5 minutes on average (mean = 4.77 minutes, 95% CI: 4–5.44
162 minutes), while the NPF approach required considerably more time (mean = 8.84
163 minutes; ~10 minutes when considering the upper range). Regarding the number of reads
164 needed (Fig 2c), nanotimes achieved the stable call using the highest number of reads
165 (mean = 433), followed by first fastq (mean = 405), and *NPF* (mean = 223).

166 Taken together, these findings highlight that while all three approaches yield similar
167 sensitivity and specificity for phenotypic prediction, the nanotimes approach
168 demonstrates slightly better sensitivity, albeit requiring more reads. Based on a favorable
169 balance between prediction accuracy and processing time, nanotimes was selected as
170 the default method for subsequent analyses.

171 ***Performance for Lineage Calling for Pure Cultures, Mock Metagenomics and*** 172 ***Metagenomic Samples***

173 To validate RASE lineage calling, we evaluated its performance across three groups of
174 samples: pure cultures (n=99), mock communities (n=20) with known composition,
175 allowing assessment under mixed-population conditions, and metagenomic samples
176 (n=20). For pure culture isolates (Fig. 3a), RASE correctly predicted the ST in 69.7% of
177 cases, with 82.6% of those recognized as the best match. In 14.1% of cases, the wrong
178 ST was predicted, however 35.7% of these were in the same clonal lineage and therefore
179 closely related (Fig 3a). For 16.2% of isolates, the ST was not present in the database,
180 thus RASE was unable to make a correct call, however, classification was possible at the
181 level of the clonal lineage for 37.5% of those isolates. The concordance was calculated
182 using Cohen's Kappa test to assess the level of agreement between the true ST (MLST,
183 short read data) and the ST predicted by RASE, the concordance was 0.823 excluding
184 those isolates belonging to STs not present in the database, and when including the
185 complete collection, the concordance was 0.686. Both had a $p = 0$, suggesting that the
186 observed agreement provides strong validation for the accuracy of RASE in predicting
187 STs represented in the database and related ones.

188 Additionally, we verified the true genomic neighbor. For this analysis, we constructed
189 three phylogenetic trees including isolates from our database and isolates used for
190 validation (pure culture and isolates forming mock communities). Since ST258/512 was
191 one of the most prevalent clonal lineage in our database, we reconstructed a separate
192 tree focused on this lineage and we confirmed that the tested isolates were accurately
193 predicted by RASE as belonging to the same ST (indicated by orange arrows in
194 Supplementary Figure 1a). For those STs not present in the database, RASE correctly
195 identified them as belonging to the same clonal lineage. Only in three cases did RASE
196 predict a different ST. In the phylogenetic tree comprising the STs considered high-risk
197 clones, RASE predicted the correct ST in most cases (Supplementary Figure 2a). The
198 final tree, which includes the remaining STs, RASE also demonstrated accurate
199 predictions. However, for several of these STs not represented in the reference database,
200 RASE was still able to predict the ST of a closely related isolate (Supplementary Figure
201 3a).

202 To assess the effect of mixed samples, we prepared 20 mock communities: 15 with
203 isolates of different STs and 5 with isolates of the same ST, yet distinct, i.e. not closely
204 related, each subset with distinct ratios (Supplementary table 3). Among the subset with
205 the same ST, the STs of 3 mock communities were accurately predicted, while the STs
206 of the remaining mock communities were classified into the same clonal lineage. Within
207 the mock communities with different STs, in the 50:50 ratio subset, one of the two STs
208 present in the mock community was identified, while the second ST was identified as an
209 alternative match. The other collections were mixed with different ratios, in which one of
210 the STs was predominant. In these communities, RASE predicted the STs of the isolates
211 with the highest proportion, except for two communities where the ST of the isolate with
212 the highest proportion was recognized as a second match, and another case where the
213 isolate with the lower proportion was identified as a second match (Table 1, Fig 3b).

214 Among the 20 metagenomic samples, 70% of STs were accurately predicted, regardless
215 of the abundance of *K. pneumoniae*, identified as the best match. 10% were classified in
216 another lineage and 20% were not present in the database (Fig 3c). The concordance of
217 this evaluation was 0.62 and increased to 0.81 excluding those samples belonging to STs

218 not present in our database. Based on the presented results, the best match and second
219 match were often related, frequently belonging to the same clonal lineages; this outcome
220 was consistently influenced by the representativeness of the database.

221 ***Phenotypic Prediction from Pure Cultures and Mock Metagenomics***

222 The RASE performance to predict antibiotic susceptibility phenotype was assessed using
223 the same datasets for the evaluation of lineage prediction, however, samples lacking
224 available MIC data were excluded, resulting in a total of 88 pure culture isolates and 10
225 mock communities that were analyzed. The overall sensitivity and specificity for
226 susceptibility across eight antibiotics in pure culture were 0.71 (95% CI 0.67-0.75) and
227 0.68 (95% CI 0.65-0.73), respectively. Analyzing only the isolates with a lineage score
228 >0.5 (n=25), both parameters increased to 0.91 (95% CI 0.84-0.97) for sensitivity and
229 0.72 (95% CI 0.64-0.80) for specificity. Certain antibiotics, such as amikacin and
230 meropenem, exhibited a sensitivity of 1, while ciprofloxacin demonstrated a specificity of
231 1 (Supplementary Figure 4).

232 After RASE predicted a phenotype in pure cultures as susceptible or resistant, we
233 measured the probability of susceptibility to individual or all antibiotics. This was
234 compared to the empiric treatment thresholds of 80% for mild infections and 90% for
235 moderate to severe infections (16). Across all antibiotics, the baseline probability of
236 susceptibility was 44%, but it increased to 72% (95% CI 63%-80%), when RASE
237 predicted a susceptible phenotype and decreased to 8.9% (95% CI 6.4%-9.6%) when
238 RASE predicted resistant phenotype. With a probable susceptible call from RASE, each
239 tested antibiotic exceeds the pre-test probability value. Nevertheless, ciprofloxacin
240 improved to a 100% along with cefotaxime and colistin, which have post-RASE
241 susceptibility probabilities exceeding 80% (Fig 4), suggesting that RASE predictions for
242 susceptibility are more likely to be accurate.

243 In this step, we evaluated RASE performance to predict antibiotic susceptibility phenotype
244 in 10 mock communities as well, which were prepared with ratios of 50:50 (n=4), 30:70
245 (n=3) and 70:30 (n=3). The selection of isolates was based on their antibiotic susceptibility
246 phenotype, one of the isolates of each mock community was multidrug resistant (R) while

247 the other was multidrug susceptible (S) to eight antibiotics including amikacin (AMK),
248 gentamicin (GEN), piperacillin-tazobactam (TZP), cefotaxime (CTX), ceftazidime (CAZ),
249 meropenem (MER), ciprofloxacin (CIP) and colistin (COL). The mock communities 12 and
250 17 were prepared using the same isolates. Community 12 was composed at a 50:50 ratio,
251 while 17 used a 30:70 ratio, the resistant isolate being present in higher proportion. In
252 both cases, the best match was identified with the same sample. For communities 13 and
253 20, we similarly mixed two of the same isolates in a 50:50 or 70:30 ratio respectively, the
254 best match corresponded to a susceptible isolate for these two collections. For
255 communities 15, 16, and 19, RASE accurately predicted the susceptibility of the
256 susceptible isolate at the 50:50 ratio. When the resistant isolate was present in a higher
257 proportion, RASE identified the best match as a multidrug-resistant isolate. Conversely,
258 when the susceptible isolate was in higher proportion, RASE identified as best match an
259 isolate only resistant to one antibiotic. In most cases of our mock collections, RASE
260 tended to predict the isolates present in the highest proportion (Table 2).

261 After we verified the true neighbors in lineage prediction, we contextualized RASE results
262 for phenotypic prediction with the phenotypic profiles of database isolates, across the
263 three phylogenetic trees. The pure culture isolates and isolates forming mock
264 communities were grouped according to their ST, and the predicted antibiotic phenotype
265 were similar to antibiotic phenotype of the isolates in the same cluster. However, isolates
266 belonging to high-risk clones had more consistent patterns of resistance (Supplementary
267 Figure 1c, 2c, 3b).

268 ***Comparative Analysis of RASE performance with Gene-based predictive Methods***

269 In order to compare RASE performance with gene-based predictive methods, we
270 employed datasets from previous studies that analyzed AMR genes after 15 hours of
271 sequencing. These samples were recovered in the USA (n=40) (17) and Germany (n=2)
272 (18). The reported metrics show high sensitivity with the gene-based methods, however,
273 the reported specificity was low for all or individual antibiotics. RASE demonstrated values
274 over 0.9 for sensitivity and in most cases the specificity was higher for RASE
275 (Supplementary Figure 5a) than for the gene-based method. Additionally, we evaluated

276 the probability of susceptibility for 5 antibiotics, for all of them the values increased when
277 RASE predicted susceptibility (Supplementary Figure 5b).

278 ***Species Identification and Susceptibility Prediction in Metagenomic Samples***

279 As the next step in our analysis, we aimed to assess the applicability of RASE in the
280 metagenome of neonates. For this purpose, we retrieved 20 metagenomic samples
281 collected within the framework of a study “Tracking the Acquisition of Pathogens In Real
282 time” (TAPIR) which recovered nasal and anal swabs from premature neonates in the
283 intensive care unit. We extracted the reads as described above, analyzing only the initial
284 reads generated at the beginning of the sequencing process, and assessed whether this
285 amount was sufficient to identify bacterial species. For this purpose, we employed
286 Kraken2 in conjunction with Bracken to confirm the presence of *K. pneumoniae* and the
287 bacterial proportions within 20 metagenomes. The results revealed consistent proportions
288 of bacteria compared to those identified through the complete 72-hour sequencing
289 analysis (Fig 5, Supplementary table 4).

290 As expected for this sample cohort, the predicted susceptibility by RASE indicates that
291 most of the samples (n=14) are susceptible to all tested antibiotics, while two samples
292 are only susceptible to MER, AMK and Col (Fig 5).

293 **Discussion**

294 The timeframe required to obtain a microbiological diagnosis through culture-dependent
295 techniques takes days (10), impacting the timely and appropriate treatment of severe and
296 common bacterial infections (19). This delay not only hinders the immediate initiation of
297 effective antibiotic therapy but also contributes to the rise of AMR. Whereas appropriate
298 empirical treatment helps to control infections effectively, inappropriate ones exert
299 unnecessary selective pressure on pathogens. We present a promising approach
300 combining long read sequencing technologies with algorithmic neighbor typing (12), to
301 guide empirical antibiotic decisions and significantly reduce diagnostic turnaround time.
302 In this study, we demonstrated that this approach improves the prediction of antibiotic
303 susceptibility in comparison with traditional diagnostic methods (Fig 4). Moreover, the

304 inclusion of lineage or ST prediction provides additional epidemiological insights to inform
305 clinical decision-making, notably in the event of an outbreak. Finally, we observed no
306 statistically significant differences in predictive performance among the read extraction
307 strategies evaluated, supporting the robustness of the proposed method regardless of the
308 input preprocessing pipeline.

309 The *K. pneumoniae* sensu lato used in this study are representative of a broader
310 population at the continental level, as they include susceptible and resistant strains
311 belonging to a wide range of STs, including globally circulating high-risk clones. This
312 diversity enhances the relevance of the findings and supports the applicability of the
313 database beyond regional contexts. In a clinical setting, a timely result would be highly
314 beneficial. Thus, to ensure that the reads evaluated with RASE correspond to the initial
315 reads during a sequencing run, we employed the nanotimes method for read extraction,
316 which selects reads sequenced at the beginning of the sequencing process. Although the
317 predictive performance of reads extracted by nanotimes and first fastq was comparable,
318 nanotimes was preferred due to its alignment with the study's focus on early diagnostic
319 utility. Furthermore, our results demonstrate that the application of preprocessing tools
320 such as Porechop and Filtlong (NPF) resulted in improved read quality but at the cost of
321 increased processing time (Fig 2b), which may delay downstream analysis in time-
322 sensitive contexts.

323 The validation of lineage calling and antibiotic susceptibility prediction using our database
324 was conducted across three sample types: pure culture, mock communities and
325 metagenomes. Concordance rates were high, with 0.82 in pure culture and 0.81 in
326 metagenomic samples, underscoring the robustness of the approach. In mock
327 communities containing more than one isolate, the pipeline consistently predicted the ST
328 of the dominant strain (Table 1). In all tested conditions, RASE correctly identified the ST
329 either as the best or second-best match. When the exact ST was not present in the
330 database, RASE was still able to associate the isolate with its corresponding clonal
331 lineage. In the absence of both ST and clonal lineage, the tool linked the strain to its
332 nearest phylogenetic neighbor (Supplementary Fig 1-3). These results show the
333 importance of maintaining a representative and frequently updated database, particularly

334 at the regional level, to ensure accuracy in lineage and phenotype prediction as novel
335 clones continue to emerge. Neighbor typing can complement existing approaches to
336 surveillance, but not fully replace them.

337 Regarding the prediction of antibiotic susceptibility phenotypes, our analysis yielded an
338 overall sensitivity of 0.91 (95% CI 0.84-0.97) and specificity of 0.72 (95% CI 0.64-0.80)
339 across eight antibiotics and an increase in the predicted probability of susceptibility was
340 observed for all tested antibiotics, supporting their potential applicability in treatment
341 decisions. Moreover, RASE was able to accurately predict the susceptibility profile of the
342 dominant strain in mock communities, however, this was not consistent across all
343 antibiotics tested (Table 2). These results support the utility of RASE, when paired with a
344 suitable database, as a rapid tool for simultaneous lineage and phenotypic prediction
345 using a limited number of reads—typically between 300 and 500—generated in under 10
346 minutes for *K. pneumoniae* (Fig 2b,c).

347 Our findings align with previous studies evaluating RASE for other bacterial pathogens,
348 including *Streptococcus pneumoniae* and *Neisseria gonorrhoeae*, demonstrating
349 prediction times under 10 minutes. In those studies, specificity reached 100% for both
350 species, while sensitivity was reported at 91% for *S. pneumoniae* and 81% for *N.*
351 *gonorrhoeae* (12). Additionally, applications of RASE to *Escherichia coli* and *Klebsiella*
352 *spp.* in Canada have emphasized the necessity of using regionally curated databases to
353 ensure predictive accuracy (20) .

354 The development and integration of novel technologies aim to enable rapid
355 microbiological diagnostics. Among the strategies are methods based on the direct
356 detection of genetic material or proteins, and these approaches have shown high
357 sensitivity, but they still require some hours of sequencing, representing a limitation. In
358 contrast, RASE, as demonstrated through our comparative analysis (Supplementary Fig
359 5-6) in this study, exhibited promising performance with significantly reduced times.

360 A known limitation of the RASE-based approach is its dependency on a species-specific
361 reference database (12,20). This requirement may hinder performance when the
362 causative pathogen is unknown at the time of sequencing, and thus an additional step to

363 identify putative relevant pathogens in a sample is necessary prior to then selecting
364 appropriate RASE databases. Another challenge arising from metagenomic sequencing
365 from swabs is the microbial content of clinical samples. Some swab-based screening
366 samples may contain low biomass, limiting the amount of microbial DNA available for
367 analysis. In our study, pre-enrichment using BHI medium improved microbial yield but
368 resulted in extended processing times. By contrast, biological samples such as blood,
369 urine, and sputum can be sequenced directly without enrichment, as demonstrated in
370 previous studies (21–23).

371 In summary, our study demonstrates that genomic neighbor typing using the *K.*
372 *pneumoniae* RASE database enables rapid prediction of both sequence type or lineage
373 complex (concordance: 0.82) and antibiotic susceptibility phenotype (sensitivity: 0.91,
374 95% CI: 0.84–0.97; specificity: 0.72, 95% CI: 0.64–0.80). This combined approach has
375 the potential to significantly reduce the time to microbiological diagnosis, facilitating more
376 accurate empirical antibiotic selection and contributing to improved treatment outcomes
377 in *K. pneumoniae* infections.

378

379 **METHODS**

380 *Reference Klebsiella pneumoniae sensu lato Database*

381 The RASE reference database was generated using assemblies, genotypic multi-locus
382 sequence typing and antibiotic susceptibility phenotype of *K. pneumoniae* isolates from
383 the EuSCAPE study. This study aimed to describe the epidemiology of carbapenem
384 resistant Enterobacterales, by recovering susceptible and resistant *K. pneumoniae*
385 strains between 2013 to 2014 across Europe (n=1511 strains) (14). We included 1511
386 fasta files of good quality, considering their N50, length and number of contigs, and for
387 which antimicrobial susceptibility testing had been defined (Supplementary table 1). Broth
388 microdilution was performed by reference methodology (24) to determine minimum
389 inhibitory concentration (MIC) (Supplementary table 1). We categorized the MIC values
390 according to the European Committee on Antimicrobial Susceptibility Testing (EUCAST)
391 v15 breakpoints (24), assigning to each strain its respective sequence type (ST) and an

392 antibiotic-specific resistance category (susceptible or resistant) for the following
393 antibiotics: AMK, GEN, TZP, CFP, CTX, CAZ, MER, CIP and COL.

394 *Datasets for Evaluation*

395 Validation was performed with susceptible and resistant *K. pneumoniae* strains from
396 clinical primary specimens and isolates. These included 99 pure cultures selected based
397 on their resistance profiles and STs, as well as the availability of the data we needed to
398 analyze, such as MIC values (isolates from EuSCAPE [n=51] (14), EURECA [n=40] (25)
399 and a project containing isolates from Greece [n=8] (26)) (Supplementary table 5). 20
400 mock communities were prepared from selected pure cultures, following the criteria
401 described below (Supplementary table 3). 20 metagenomic samples of anal and nasal
402 swab recovered from neonates through the TAPIR project were chosen according to the
403 presence and abundance of *K. pneumoniae* (Supplementary table 4). All samples were
404 previously short read sequenced and the antibiotic susceptibility profiles are available for
405 the majority, with the exception of metagenomic samples and 11 pure cultures. Therefore,
406 all phenotypic prediction analyses were performed using only 88 pure cultures. The
407 previously determined STs found using short read data and MICs values were updated
408 based on EUCAST (24) v15, where “susceptible increased exposure” values were
409 reinterpreted as susceptible.

410 *DNA Extraction, Library Preparation and Nanopore Sequencing*

411 Pure culture isolates, including those used for the construction of mock communities,
412 were cultured on blood agar and incubated overnight at 37°C, then species identification
413 was performed using MALDI-TOF mass spectrometry, and a single colony was selected
414 for subsequent DNA extraction. Metagenomic samples were enriched with BHI medium
415 and incubated overnight at 37°C, 210 rpm in a proportion 1:2. To remove the human DNA,
416 enriched samples were centrifuged at 10.000 rpm for 5 minutes, the microbial pellet was
417 resuspended with 1 mL of clean dH₂O at room temperature for 5 minutes and a final
418 resuspension in 200 µL PBS 1X; 5 µL lysozyme and 5 µL lysoplastin.

419 DNA was extracted using the Roche High Pure Template Preparation Kit, and long-read
420 sequencing was performed using the Oxford Nanopore Technology (ONT). Sequencing

421 libraries were prepared using the ligation protocol SQK-LSK114 with native barcode kits
422 (EXP-NBD114 and SQK-LSK 114) for R9 Chemistry, and SQK-NBD114-96 barcoding kit
423 for R10 Chemistry in accordance with the manufacturer's protocol including optional
424 steps. DNA input and other measurements were calculated based on the assumption that
425 1 fmol is equivalent to 5 ng of DNA. Thus, for the final R10 library, 100 ng corresponded
426 to 20 fmol. For the R10 protocol, the initial volume per sample was 11 μ L without water
427 dilution (for clinical samples). Optional NEBNext FFPE DNA Repair and DNA Control
428 Sample were omitted during DNA repair and end-prep. In addition, Bovine Serum Albumin
429 (BSA) was used to improve the sequencing performance as recommended in the
430 protocol.

431 Final libraries were quantified using Qubit 4 Fluorometer, and loaded onto a FLO-MIN114
432 R10.4 flow cell type and sequenced on a GridION X5 Mk1 sequencing platform.
433 Sequence data acquisition, real-time base-calling, and demultiplexing of barcodes were
434 conducted using the graphical user interface MinKNOW (v23.11.7) and the dorado
435 basecaller (v7.2.13).

436 Sequencing data has been deposited in the following ENA projects: Mock Communities
437 PRJEB82665 (supplementary table 3), TAPIR metagenomic clinical samples subsampled
438 to first 60min PRJEB96334 and pure cultures PREJB95992 (supplementary table 4).
439 Long read ONT sequencing used to evaluate RASE and the reference database from
440 projects Daikos Greece PRJEB58216, EuSCAPE and EURECA PRJEB96336
441 (supplementary table 5).

442 *Analysis of Culture Specimens*

443 To evaluate long-read sequencing from pure culture, extraction of reads from isolates
444 were initially performed based on the sequencing time (60 min) using nanotimes
445 (<https://github.com/angelovangel/nanotimes>; first dataset). Along with that, the first fastq
446 sequenced for each isolate was extracted (second dataset). Porechop (v0.2.4;
447 <https://github.com/rrwick/Porechop/tree/master>) was used to trim off adapters and filtlong
448 (v0.2.1; <https://github.com/rrwick/Filtlong>) to filter the long sequences with a minimum
449 read length threshold of 1000, only for reads extracted with nanotimes (third dataset)

450 [NPF]).Therefore, three pure culture sets of sequences of each sample were analyzed
451 with RASE (version 0.1.0.0).

452 *Analysis of Mock Metagenomics*

453 A total of 20 mock communities were created to evaluate the ability of long read
454 sequencing coupled with neighbour typing to predict ST and phenotype in the setting of
455 controlled consortiums of bacterial strains. Isolates to be mixed were chosen based on
456 their ST and antibiotic susceptibility phenotype. Each mock community corresponds to
457 the DNA pool of two (n=16) or three isolates (n=4), with different proportions, resulting in
458 six subsets 50:50 (n=6), 40:60 (n=2), 30:70 (n=5), 20:30:50 (n=2), 10:25:65 (n= 2) and
459 70:30 (n=2) (Supplementary table 3). At least one mock community of each subset was
460 prepared with isolates of the same ST (n=5), except the subset 70:30. Only ten of these
461 communities were used to evaluate the antibiotic susceptibility prediction, including the
462 subsets 50:50 (n=4), 30:70 (n=3) and 70:30 (n=3). Each mock community was prepared
463 with one multidrug resistant (R) and one multidrug susceptible (S) isolate to eight
464 antibiotics AMK, GEN, TZP, CTX, CAZ, MER, CIP and COL. The 20 mock communities
465 were sequenced using nanopore, after the first reads sequenced during the beginning of
466 the sequencing of each mock community were extracted using nanotimes (60min). The
467 number of reads differed between 2551 to 17205 (Supplementary table 3), and then these
468 sequences were assessed with RASE.

469 *Analysis of Metagenomic Samples*

470 Extracted reads employing nanotimes of metagenomic samples were mapped to the
471 human genome (GRCh38) using Minimap2 (v2.24; using the map-ont parameter) (27)
472 and human DNA contaminants were removed from reads using Samtools (v1.14) (28).
473 Species identification was performed using Kraken2 (13). Microbial species abundance
474 was estimated using Bracken (29). Subsequent analysis was done with RASE.

475 *Analysis with RASE*

476 Rase algorithm is based on matching the nanopore reads to a reference database using
477 Prophyle and increasing the weight of the most similar strains. First, the lineage is

478 identified by finding the best match reference genome. The lineage score is calculated by
479 comparing the two best-matching lineages. In the next step, the best match within the
480 lineage is identified, and resistance is predicted from the nearest resistant and susceptible
481 neighbors and the susceptible score is a result of the comparison of their weights. To use
482 the software we follow the steps described in the documentation, then to analyze our
483 results, we first define the stable call based on the fluctuation of the lineage score (LS),
484 when it did not change by more than 0.2 for at least 100 reads, as LS was created as a
485 potential indicator of confidence in lineage prediction. When we performed long read
486 sequencing of isolates that were already contained within EuSCAPE (n=51) (14), we
487 created a new database for each sample by removing the sample being evaluated.

488 *Analysis of external datasets*

489 Additionally, we analyzed two external dataset from two previous studies (17,18) that
490 assessed gene-based methods for identifying antimicrobial resistance (AMR) genes. We
491 downloaded 42 genome assemblies from the ENA database and obtained the
492 corresponding raw sequencing reads generated after 60 minutes of Nanopore
493 sequencing using nanotimes. The reads were processed and analyzed with RASE to
494 predict AMR profiles. We then compared the sensitivity and specificity of our predictions
495 against those reported in the original studies to evaluate performance consistency.

496 *Statistical Methods and Visualization*

497 The concordance between the true ST and the one predicted by RASE was assessed
498 with the Cohen' s Kappa test using the irr package (30) (v2.0.60). Friedman test was
499 performed to compare the three read extraction methods. The performance of RASE
500 predicting susceptibility and non-susceptibility was evaluated by calculating its sensitivity
501 and specificity. Sensitivity was determined by measuring the proportion of true positive
502 results. This was calculated as the number of true positives divided by the sum of true
503 positives and false negatives. Specificity, conversely, was calculated by assessing the
504 proportion of true negative results. It was calculated as the number of true negatives
505 divided by the sum of true negatives and false positives. Plots were visualized using
506 ggplot from Tidyverse (31) (v2.0.0).

507 *Phylogenetic analyses*

508 Phylogenetic trees of the database and tested isolates were estimated using RAxML
509 v8.2.12 (32) based on SNP alignments after mapping to reference genome MGH78578
510 (GenBank accession CP000647) and removal of recombinant regions using Gubbins
511 v2.4.1 (33). Phylogenetic tree visualization was done with iTOL (34).

512 *Ethical approval*

513 This study has been approved by the local ethics committee (22-1039 and 22-1040).

514

515 **ACKNOWLEDGEMENTS**

516 This project was funded by JPIAMR call (K-STaR 01KI1910), and the Federal Ministry of
517 Research, Technology and Space (BMFTR), formerly the Federal Ministry of Education
518 and Research (BMBF; TAPIR 01KI2018) funding to SR. KB was supported by the French
519 National Research Agency (ANR) under Grant ANR-24-CE45-1226 (REALL project). We
520 thank technical support by Leonardo Duarte dos Santos.

521

522 **Conflicts of interest**

523 SR has received travel funds and speaker remuneration from Illumina. All other authors
524 have no conflict of interest to declare

525 **REFERENCES**

- 526 1. Urban-Chmiel R, Marek A, Stępień-Pyśniak D, Wieczorek K, Dec M, Nowaczek A, et al.
527 Antibiotic Resistance in Bacteria-A Review. *Antibiotics (Basel)* [Internet]. 2022 Aug 9;11(8).
528 Available from: <http://dx.doi.org/10.3390/antibiotics11081079>
- 529 2. Mancuso G, De Gaetano S, Midiri A, Zummo S, Biondo C. The Challenge of Overcoming
530 Antibiotic Resistance in Carbapenem-Resistant Gram-Negative Bacteria: “Attack on Titan.”
531 *Microorganisms* [Internet]. 2023 Jul 27;11(8). Available from:
532 <http://dx.doi.org/10.3390/microorganisms11081912>
- 533 3. Antimicrobial Resistance Collaborators. Global burden of bacterial antimicrobial resistance
534 in 2019: a systematic analysis. *Lancet*. 2022 Feb 12;399(10325):629–55.

- 535 4. Arcari G, Carattoli A. Global spread and evolutionary convergence of multidrug-resistant
536 and hypervirulent high-risk clones. *Pathog Glob Health*. 2023 Jun;117(4):328–41.
- 537 5. Wyres KL, Lam MMC, Holt KE. Population genomics of *Klebsiella pneumoniae*. *Nat Rev*
538 *Microbiol*. 2020 Jun;18(6):344–59.
- 539 6. Price R. O’Neill report on antimicrobial resistance: funding for antimicrobial specialists
540 should be improved. *Eur J Hosp Pharm Sci Pract*. 2016 Jul;23(4):245–7.
- 541 7. Hu W, Chen H, Wang H, Peng Q, Wang J, Huang W, et al. Identifying high-risk phenotypes
542 and associated harms of delayed time-to-antibiotics in patients with ICU onset sepsis: A
543 retrospective cohort study. *J Crit Care*. 2023 Apr;74:154221.
- 544 8. Bisarya R, Song X, Salle J, Liu M, Patel A, Simpson SQ. Antibiotic Timing and Progression
545 to Septic Shock Among Patients in the ED With Suspected Infection. *Chest*. 2022
546 Jan;161(1):112–20.
- 547 9. Muteeb G, Rehman MT, Shahwan M, Aatif M. Origin of Antibiotics and Antibiotic
548 Resistance, and Their Impacts on Drug Development: A Narrative Review.
549 *Pharmaceuticals* [Internet]. 2023 Nov 15;16(11). Available from:
550 <http://dx.doi.org/10.3390/ph16111615>
- 551 10. MacFadden DR, Leis JA, Mubareka S, Daneman N. The opening and closing of empiric
552 windows: the impact of rapid microbiologic diagnostics. *Clin Infect Dis*. 2014 Oct
553 15;59(8):1199–200.
- 554 11. Wang Y, Zhao Y, Bollas A, Wang Y, Au KF. Nanopore sequencing technology,
555 bioinformatics and applications. *Nat Biotechnol*. 2021 Nov;39(11):1348–65.
- 556 12. Břinda K, Callendrello A, Ma KC, MacFadden DR, Charalampous T, Lee RS, et al. Rapid
557 inference of antibiotic resistance and susceptibility by genomic neighbour typing. *Nat*
558 *Microbiol*. 2020 Mar;5(3):455–64.
- 559 13. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact
560 alignments. *Genome Biol*. 2014 Mar 3;15(3):R46.
- 561 14. David S, Reuter S, Harris SR, Glasner C, Feltwell T, Argimon S, et al. Epidemic of
562 carbapenem-resistant *Klebsiella pneumoniae* in Europe is driven by nosocomial spread.
563 *Nat Microbiol*. 2019 Nov;4(11):1919–29.
- 564 15. Tryfinopoulou K, Linkevicius M, Pappa O, Alm E, Karadimas K, Svartström O, et al.
565 Emergence and persistent spread of carbapenemase-producing high-risk clones in Greek
566 hospitals, 2013 to 2022. *Euro Surveill* [Internet]. 2023 Nov;28(47). Available from:
567 <http://dx.doi.org/10.2807/1560-7917.ES.2023.28.47.2300571>
- 568 16. Cressman AM, MacFadden DR, Verma AA, Razak F, Daneman N. Empiric Antibiotic
569 Treatment Thresholds for Serious Bacterial Infections: A Scenario-based Survey Study.
570 *Clin Infect Dis*. 2019 Aug 30;69(6):930–7.
- 571 17. Tamma PD, Fan Y, Bergman Y, Perteau G, Kazmi AQ, Lewis S, et al. Applying rapid whole-
572 genome sequencing to predict phenotypic antimicrobial susceptibility testing results among
573 carbapenem-resistant *Klebsiella pneumoniae* clinical isolates. *Antimicrob Agents*

- 574 Chemother [Internet]. 2019 Jan;63(1). Available from: [http://dx.doi.org/10.1128/AAC.01923-](http://dx.doi.org/10.1128/AAC.01923-18)
575 18
- 576 18. Sauerborn E, Corredor NC, Reska T, Perlas A, Vargas da Fonseca Atum S, Goldman N, et
577 al. Detection of hidden antibiotic resistance through real-time genomics. *Nat Commun.*
578 2024 Jun 28;15(1):5494.
- 579 19. Cassini A, Högberg LD, Plachouras D, Quattrocchi A, Hoxha A, Simonsen GS, et al.
580 Attributable deaths and disability-adjusted life-years caused by infections with antibiotic-
581 resistant bacteria in the EU and the European Economic Area in 2015: a population-level
582 modelling analysis. *Lancet Infect Dis.* 2019 Jan;19(1):56–66.
- 583 20. Carroll AC, Mortimer L, Ghosh H, Reuter S, Grundmann H, Brinda K, et al. Rapid inference
584 of antibiotic susceptibility phenotype of uropathogens using metagenomic sequencing with
585 neighbor typing. *Microbiol Spectr.* 2025 Jan 7;13(1):e0136624.
- 586 21. Cheng H, Sun Y, Yang Q, Deng M, Yu Z, Zhu G, et al. A rapid bacterial pathogen and
587 antimicrobial resistance diagnosis workflow using Oxford nanopore adaptive sequencing
588 method. *Brief Bioinform [Internet].* 2022 Nov 19;23(6). Available from:
589 <http://dx.doi.org/10.1093/bib/bbac453>
- 590 22. Serpa PH, Deng X, Abdelghany M, Crawford E, Malcolm K, Caldera S, et al. Metagenomic
591 prediction of antimicrobial resistance in critically ill patients with lower respiratory tract
592 infections. *Genome Med.* 2022 Jul 12;14(1):74.
- 593 23. Ruppé E, d’Humières C, Armand-Lefèvre L. Inferring antibiotic susceptibility from
594 metagenomic data: dream or reality? *Clin Microbiol Infect.* 2022 Sep;28(9):1225–9.
- 595 24. ESCMID-European Society of Clinical Microbiology, Diseases I. eucast: Clinical
596 breakpoints and dosing of antibiotics [Internet]. [cited 2024 Aug 1]. Available from:
597 https://www.eucast.org/clinical_breakpoints
- 598 25. Budia-Silva M, Kostyanev T, Ayala-Montañó S, Bravo-Ferrer Acosta J, Garcia-Castillo M,
599 Cantón R, et al. International and regional spread of carbapenem-resistant *Klebsiella*
600 *pneumoniae* in Europe. *Nat Commun.* 2024 Jun 14;15(1):5092.
- 601 26. Afolayan AO, Rigatou A, Grundmann H, Pantazatou A, Daikos G, Reuter S. Three lineages
602 causing bloodstream infections variably dominated within a Greek hospital over a 15 year
603 period. *Microb Genom [Internet].* 2023 Aug;9(8). Available from:
604 <http://dx.doi.org/10.1099/mgen.0.001082>
- 605 27. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018 Sep
606 15;34(18):3094–100.
- 607 28. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
608 Alignment/Map format and SAMtools. *Bioinformatics.* 2009 Aug 15;25(16):2078–9.
- 609 29. Lu J, Salzberg SL. Ultrafast and accurate 16S rRNA microbial community analysis using
610 Kraken 2. *Microbiome.* 2020 Aug 28;8(1):124.
- 611 30. Gamer M, Lemon J, Fellows I, Singh P. irr: Various Coefficients of Interrater Reliability and
612 Agreement. R package version 0. 2019;84.

- 613 31. Wickham H. Ggplot2. 2nd ed. Basel, Switzerland: Springer International Publishing; 2016.
614 260 p. (Use R!).
- 615 32. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
616 phylogenies. *Bioinformatics*. 2014 May 1;30(9):1312–3.
- 617 33. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid
618 phylogenetic analysis of large samples of recombinant bacterial whole genome sequences
619 using Gubbins. *Nucleic Acids Res*. 2015 Feb 18;43(3):e15.
- 620 34. Letunic I, Bork P. Interactive Tree of Life (iTOL) v6: recent updates to the phylogenetic tree
621 display and annotation tool. *Nucleic Acids Res*. 2024 Jul 5;52(W1):W78–82.

622 **List of Figures, Tables and Supplementary material**

623 **Table 1** Summary of predicted lineage for mock collections.

624 **Table 2** Predicted phenotype for mock collections, displaying the true or actual phenotype
625 and the best match predicted by RASE of 10 mock collections, S represent susceptible
626 and R non-susceptible, and the concordance is coloured green for susceptible and
627 orange for non-susceptible.

628 **Figure 1.** Overview of *Klebsiella pneumoniae* sensu lato RASE database, the dots
629 represent an isolate coloured by their antimicrobial susceptibility for antibiotics in x, and
630 are grouped in y according to the Sequence Type. Barplots show the number of isolates
631 susceptible, and resistant by antibiotic.

632 **Figure 2.** Comparative performance of two read extraction strategies and pre-processing
633 reads in terms of phenotypic prediction accuracy, time to stable call, and number of reads
634 required. a) Sensitivity (left) and specificity (right) of predictions for eight antibiotics were
635 evaluated across two read extraction strategies: nanotimes (blue), first fastq (red), and
636 nanotimes+Porechop+Filtlong (NPF, orange). b) Time (in minutes) and c) number of
637 reads to achieve a stable call.

638 **Figure 3.** Validation of lineage calling of *K. pneumoniae* sensu lato RASE database
639 across three conditions studied: a) pure culture b) mock communities and c)
640 metagenomic samples.

641 **Figure 4.** Probability of susceptibility of eight antibiotics individually and together
642 predicted by RASE in 25 pure culture samples ($L_s > 0.5$), shapes and colors differentiate
643 the pre-test probability in yellow, predicted non-susceptible in orange and green the
644 predicted susceptibility.

645 **Figure 5.** Analysis of metagenomic samples comparing the abundance after 72 hours
646 and 1 hour of sequencing is shown in the first and second bartplots. The abundance of
647 *K. pneumoniae* (Bracken analyses) is coloured in mustard yellow and in green other
648 bacterias, considering that only the three most abundant bacteria per sample are

649 displayed. Antibiotic susceptibility is depicted in the following columns, blue represents
650 susceptible and red resistance isolates.

651 **Supplementary Figure 1** Phylogenetic tree of database isolates, a) pure culture and b)
652 isolates forming mock communities within the ST258/512 clonal lineage, clades are
653 coloured by ST. Arrows link each tested isolate to the match predicted by RASE. Orange
654 arrows indicate predictions within the same ST, light blue for matches in a different ST
655 within the same clonal complex and grey shows isolates of STs absent from the database
656 but matched within the same clonal lineages. c) Phenotypic prediction by RASE of
657 evaluation samples, along with the true antibiotic phenotype of database isolates for eight
658 antibiotics.

659 **Supplementary Figure 2** Phylogenetic tree of database isolates, a) pure culture and b)
660 isolates forming mock communities belonging to the high-risk clones. Each tested isolate
661 is linked to its predicted match according to RASE results. Arrows colours: orange
662 indicates matches within the same ST, light blue matches in a different ST within the
663 same clonal lineage, and purple indicates matches with a different ST. c) Phenotypic
664 prediction by RASE of evaluation samples, along with the true antibiotic phenotype of
665 database isolates for eight antibiotics

666 **Supplementary Figure 3** Phylogenetic tree of remaining isolates belonging to different
667 STs, the most prevalent STs are coloured. Orange arrows indicate matches within the
668 same ST, light blue arrows denote matches in a different ST within the same clonal
669 lineage, and grey arrows correspond to isolates of STs that are not present in the
670 reference database but are matched with the closest isolate. c) Phenotypic prediction by
671 RASE of evaluation samples, along with the true antibiotic phenotype of database isolates
672 for eight antibiotics.

673 **Supplementary Figure 4** Phenotypic validation of the *K. pneumoniae* sensu lato RASE
674 database in 25 pure culture samples, stratified by sensitivity and specificity.

675 **Supplementary Figure 5 a)** Comparative analysis of RASE (blue) and gene-based
676 methods (red) in terms of sensitivity and specificity across six different antibiotics. Each

677 dot represents the predictive performance for a specific antibiotic. **b)** Probability of
678 susceptibility of six antibiotics predicted by RASE in 42 external isolates, shapes and
679 colors differentiate the pre-test probability in yellow, predicted resistant in orange and
680 green the predicted susceptible.

681 **Supplementary Table 1** Reference database information, including quality control of
682 assemblies, multilocus sequence types, and antibiotic susceptibility profiles for the
683 evaluated antibiotics.

684 **Supplementary Table 2** Results of statistical test (Post-hoc Wilcoxon, Shapiro-Wilk and
685 Friedman) used for comparing read extraction methods and pre-processing tools.

686 **Supplementary Table 3** Proportions of isolates forming mock communities and their true
687 STs. Also includes RASE ST predictions and ENA accession numbers.

688 **Supplementary Table 4** Metagenomic samples data including true ST, RASE ST
689 predictions along with bacterial proportions in culture, sequencing after 1hr and 72 hrs.
690 ENA accession numbers are included as well.

691 **Supplementary Table 5** This table presents information of pure culture including true ST
692 and antibiotic susceptibility profiles for evaluated antibiotics. Additionally provides ENA
693 accession number and RASE predictions (stable call) of each sample.

694 **Table 1** Summary of the accurately predicted lineage for mock communities

	Lineage prediction	Ratio					
		10-25-65 (n=2)	20-30-50 (n=2)	30-70 (n=5)	40-60 (n=2)	50-50 (n=6)	70-30 (n=3)
1 ST-high proportion	Best match	1		3	2	5	3
	Best match different ST same CC	1	1			2	
	Second match	1	1	1		1	
2 ST-low proportion	Best match	1	1	1	1		
	Best match different ST same CC		1				
	Second match			2			1
3 ST-lowest proportion	Best match	1					
	Best match different ST same CC	1					
	Second match						

695

696 **Table 2** Predicted phenotype for mock communities

Mock communities	RATIO	MLST match	CC match	AMK		CTX		CAZ		CIP		Col		GEN		MEM	
				Actual	Best match	Actual	Best match	Actual	Best match	Actual	Best match	Actual	Best match	Actual	Best match	Actual	Best match
Mix 12	50%	No	No	S	S	S	R	S	R	S	R	S	S	S	S	S	R
	50%	Yes	Yes	R	S	R	R	R	R	R	R	R	S	R	S	R	R
Mix 17	30%	No	No	S	S	S	R	S	R	S	R	S	S	S	S	S	R
	70%	Yes	Yes	R	S	R	R	R	R	R	R	R	S	R	S	R	R
Mix 13	50%	Yes	Yes	S	S	S	S	S	S	S	S	S	S	S	S	S	S
	50%	No	No	R	S	R	S	R	S	R	S	R	S	R	S	R	S
Mix 20	70%	Yes	Yes	S	S	S	S	S	S	S	S	S	S	S	S	S	S
	30%	No	No	R	S	R	S	R	S	R	S	R	S	R	S	R	S
Mix 15	50%	Yes	Yes	S	S	S	S	S	S	S	S	S	S	S	S	S	S
	50%	No	No	R	S	R	S	R	S	R	S	R	S	R	S	R	S
	30%	No	No	S	S	S	R	S	R	S	R	S	R	S	R	S	R

Mix 16	70%	No	No	R	S	R	R	R	R	R	R	R	R	R	R	R	R
Mix 19	70%	Yes	Yes	S	S	S	S	S	S	S	S	S	S	S	S	S	S
	30%	No	No	R	S	R	S	R	S	R	S	R	S	R	S	R	S
Mix 11	70%	Yes	Yes	S	S	S	S	R	S	S	S	S	S	S	S	S	S
	30%	No	No	R	S	R	S	R	R	R	S	R	S	R	S	R	S
Mix 14	50%	No	No	S	S	S	R	S	R	S	R	S	S	S	S	S	R
	50%	Yes	Yes	R	S	R	R	R	R	R	R	R	R	R	R	R	R
Mix 18	30%	No	No	S	R	S	R	S	R	S	S	S	S	S	R	S	S
	70%	No	No	R	R	R	R	R	R	R	S	R	S	R	R	R	S

697 The table displays the true or actual phenotype and the best match predicted by RASE of 10 mock communities, S
698 represents susceptible and R resistant, and the concordance is coloured , green for susceptible and orange for resistant.

Figure 2

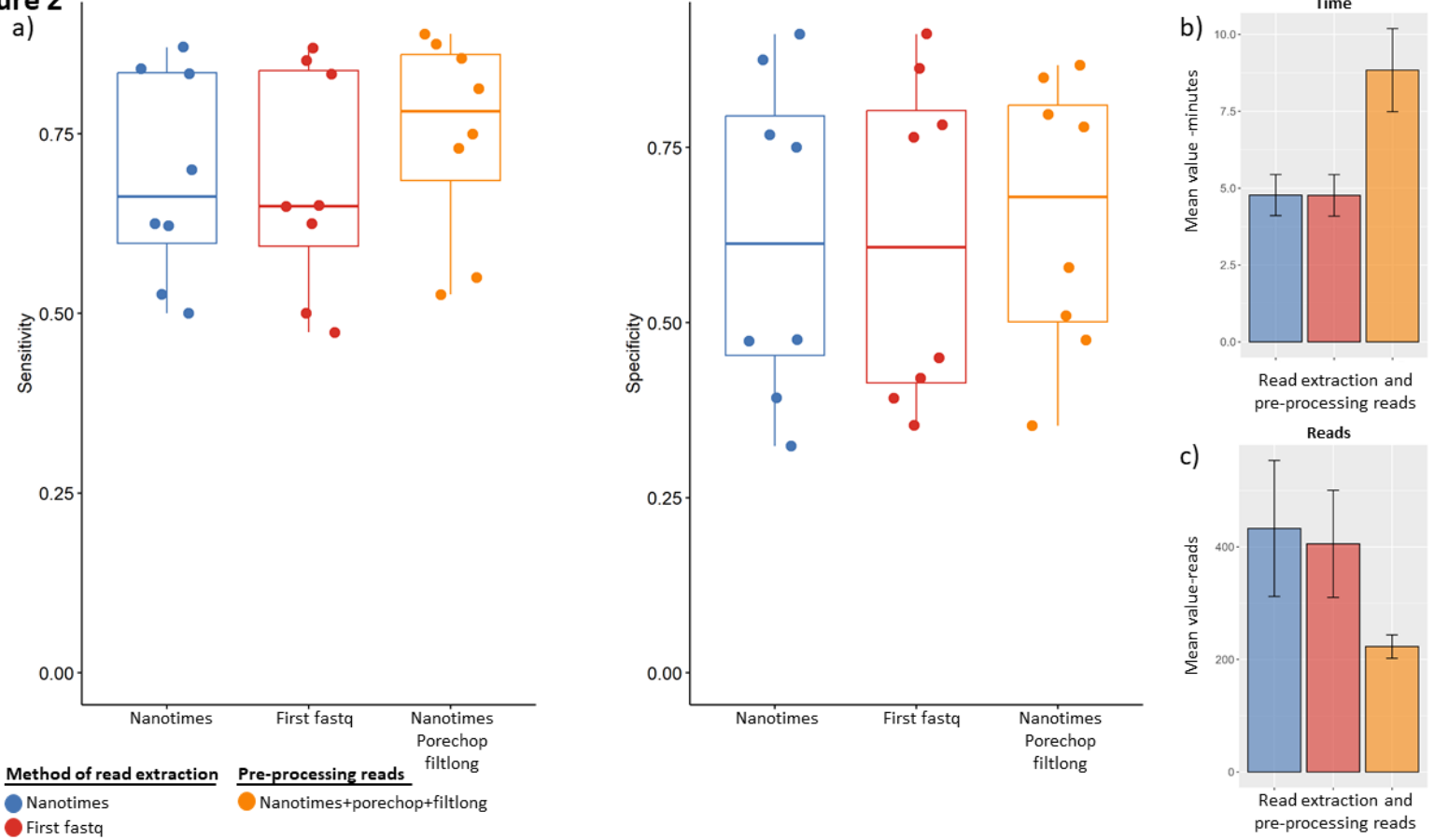


Figure 3

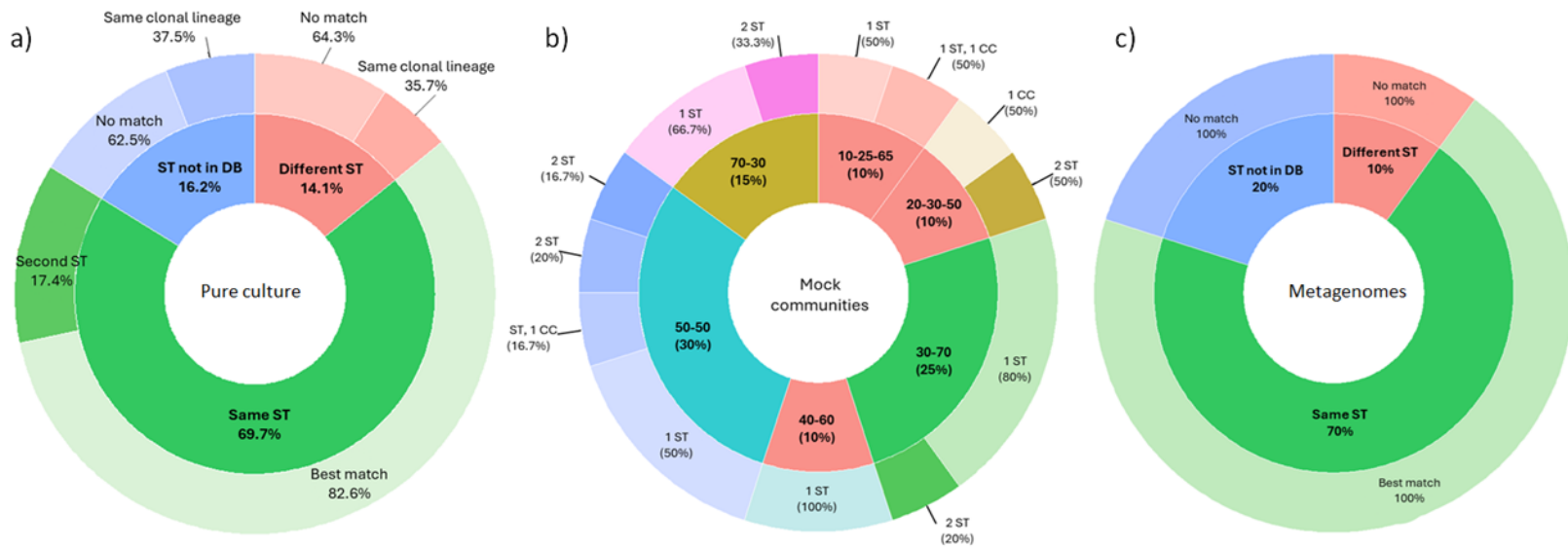


Figure 4

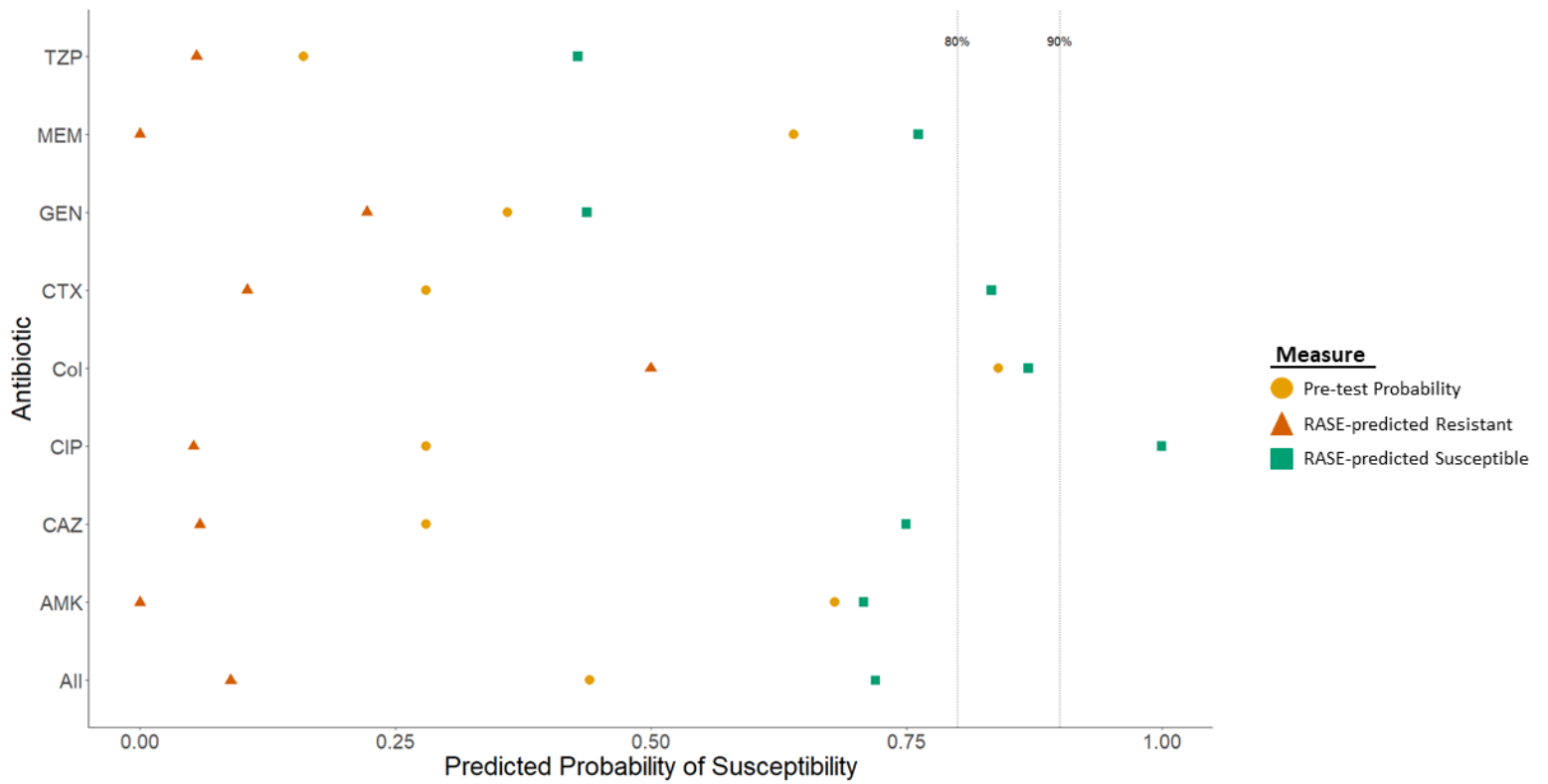
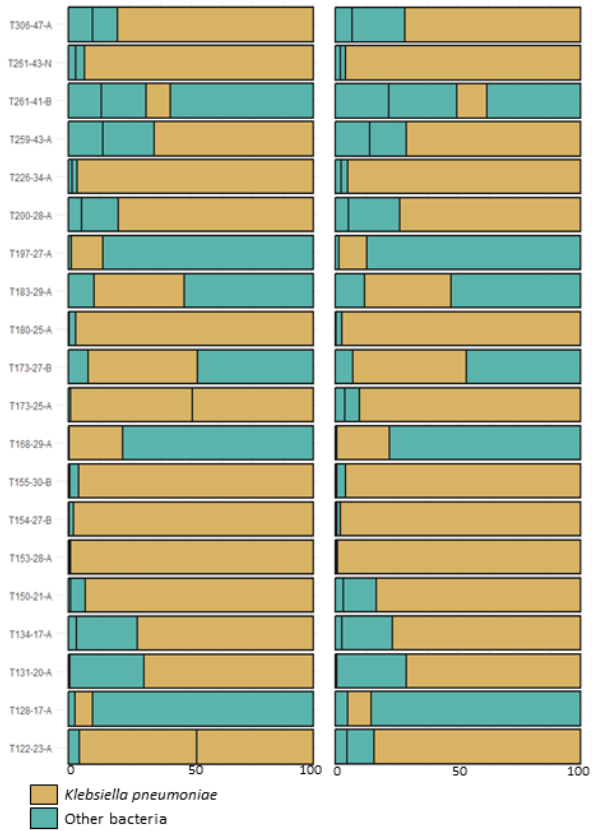


Figure 5

Taxonomy

72 hr

1 hr



Antibiotic Susceptibility RASE prediction

