

1 Bactabolize: A tool for high- 2 throughput generation of bacterial 3 strain-specific metabolic models 4

5 Authors

6 Ben Vezina^{1*}, Stephen C. Watts^{1*}, Jane Hawkey¹, Helena B. Cooper¹, Louise M. Judd¹, Adam
7 Jenney², Jonathan M. Monk³, Kathryn E. Holt^{1,4}, Kelly L. Wyres^{1*}

8

9 ¹ Department of Infectious Diseases, Central Clinical School, Monash University, Melbourne,
10 Victoria, Australia

11 ² Microbiology Unit, Alfred Health, Melbourne, Victoria, Australia

12 ³ Department of Bioengineering, University of California, San Diego, CA, United States of America

13 ⁴ Department of Infection Biology, London School of Hygiene and Tropical Medicine, London, UK

14 ⁺ These authors contributed equally

15 ^{*} Corresponding authors (benjamin.vezina@monash.edu, kelly.wyres@monash.edu)

16

17 Abstract

18 Metabolic capacity can vary substantially within a bacterial species, leading to ecological niche
19 separation, as well as differences in virulence and antimicrobial susceptibility. Genome-scale
20 metabolic models are useful tools for studying the metabolic potential of individuals, and with the
21 rapid expansion of genomic sequencing there is a wealth of data that can be leveraged for
22 comparative analysis. However, there exist few tools to construct strain-specific metabolic models
23 at scale.

24 Here we describe Bactabolize (github.com/kelwyres/Bactabolize), a reference-based tool which
25 rapidly produces strain-specific metabolic models and growth phenotype predictions. We describe
26 a pan reference model for the priority antimicrobial-resistant pathogen, *Klebsiella pneumoniae*
27 (github.com/kelwyres/KpSC-pan-metabolic-model), and a quality control framework for using draft
28 genome assemblies as input for Bactabolize.

29 The Bactabolize-derived model for *K. pneumoniae* reference strain KPPR1 outperformed the
30 CarveMe-derived model across ≥ 201 substrate and ≥ 1220 knockout mutant growth predictions.
31 Novel draft genomes passing our systematically-defined quality control criteria resulted in models
32 with a high degree of completeness ($\geq 99\%$ genes and reactions captured) and high accuracy
33 (mean 0.97, $n=10$).

34 We anticipate the tools and framework described herein will facilitate large-scale metabolic
35 modelling analyses that broaden our understanding of diversity within bacterial species and inform
36 novel control strategies for priority pathogens.

37

38 Introduction

39 Bacteria exhibit metabolic diversity and can utilise a broad range of substrates for growth. It has
40 become clear amongst pathogens that there is an intertwined relationship between metabolism
41 and nutrient usage with virulence and antimicrobial resistance (1-7). Comparative analyses of
42 metabolic profiles (e.g. substrate usage) are key to fully understanding these relationships.
43 Traditionally, these profiles have been assessed via phenotypic growth on a limited number of
44 substrates, such as those used to delineate between species (8-10) which form the basis of a
45 number of commercial products for species identification. However, these methods are not
46 sufficiently discriminatory for in-depth comparisons within species, and alternative approaches
47 such as the Omnilog Phenotype MicroArray system (Biolog) are too expensive and/or labour
48 intensive for application to large numbers of isolates. Similarly, probing of essential metabolism-
49 associated genes via transposon mutant libraries (e.g. to identify novel virulence factors and
50 therapeutic targets) (4, 11, 12) cannot be easily scaled across diverse bacterial populations.

51 Genome-scale metabolic models or metabolic reconstructions are a computational approach to
52 analysing the metabolic potential of an organism, within which the entire biochemical network is
53 represented as a stoichiometric matrix (13). Metabolic models are constructed programmatically,
54 but typically informed and at least partially validated using phenotypic growth data (14-16). Once
55 constructed, they can be run through simulations and analysed under various contexts, such as *in*
56 *silico* growth experiments (Flux Balance Analysis [FBA]) to predict substrate usage profiles (17),
57 evaluate the impact of single gene knockouts on growth (14, 18), and identify metabolic
58 chokepoints for drug targets (19), among others. Traditionally, metabolic models are strain-specific
59 (i.e. each model represents a unique individual <http://bigg.ucsd.edu/models>) and may not be
60 applicable to other isolates due to unrepresented genetic diversity.

61 We recently described 37 curated strain-specific models for the *Klebsiella pneumoniae* Species
62 Complex (*KpSC*) (14) comprised of *K. pneumoniae* and its close relatives (20). These organisms
63 are a common cause of healthcare-associated infections world-wide, and among the World Health
64 Organization's priority antimicrobial resistant pathogens (21). *KpSC* are highly diverse and gene
65 content can differ substantially between strains (22, 23). Accordingly, our models varied in terms of
66 gene and reaction content, resulting in variable growth substrate usage profiles and metabolic
67 redundancy (14). Similar variation has also been described in other key bacterial pathogens e.g.
68 *Escherichia coli* (24), *Salmonella enterica* (25), *Staphylococcus aureus* (26) and *Pseudomonas*
69 *aeruginosa* (27). This is highly relevant to the use of metabolic models for the exploration of
70 virulence and antimicrobial resistance, and for the identification of novel drug targets. Therefore,
71 such works should seek to include multiple strain-specific models, and in some cases 100s-1000s
72 of models may be required to accurately represent population diversity (22, 28, 29).

73 There are several open source tools currently available that can rapidly produce strain-specific
74 metabolic models, including CarveMe (30), ModelSEED (31) and KBase (32) (see the recent
75 review by Mendoza and colleagues for comparative descriptions (33)), as well as a recently
76 published modelling and analysis pipeline, ChiMera, which leverages CarveMe for model
77 construction (34). In their systematic analysis Mendoza *et al.* indicated CarveMe and ModelSEED
78 to be of particular interest for large-scale studies due to their speed and model quality (33). Like
79 KBase, ModelSEED is a web interface application, limiting its utility for high-throughput analysis of
80 100s – 1000s of bacterial genomes. CarveMe is a command line application; it is open source but
81 is dependent on commercial solvers such as CPLEX (free for academic use). However, its use of a
82 universal reference model may limit specificity of strain-specific models (35), and result in
83 overestimation of model genes. These limitations can be overcome by manual curation of the
84 output models, but such curation is highly labour intensive and not suitable for high-throughput
85 analyses. Furthermore, the CarveMe database (BiGG universal_model) appears to be no longer

86 actively maintained, meaning that there is no opportunity to integrate novel structural and/or
87 biochemical data as these become available in the literature (as discussed in COBRA community
88 forums).

89 Here, we present Bactabolize (available at <https://github.com/kelwyres/Bactabolize>), an easy-to-
90 use tool which allows scalable production of strain-specific draft metabolic models and prediction
91 of growth phenotypes. Bactabolize builds upon the reference-based model reconstruction
92 approach described by Norsigian *et al.* (35), leveraging the COBRAPy framework (36) and BiGG
93 nomenclature (37). We present a pan-metabolic reference model for the *KpSC* (derived from our
94 37 curated strain-specific models (14)), and describe an exemplar quality control framework for the
95 application of Bactabolize to *KpSC* draft genome assemblies. We show that Bactabolize can
96 rapidly produce strain-specific models from draft genomes with a high degree of completeness (as
97 compared to models generated from completed genome assemblies), resulting in highly accurate
98 growth predictions that match or exceed the accuracy of models from CarveMe and manual
99 curation efforts.

100 Results

101 Description of Bactabolize

102 Bactabolize is written in Python 3 and utilises the metabolic modelling library COBRAPy (36).

103 Bactabolize has four main commands:

- 104 i) Draft model generation (*draft_model* command), which generates a strain-specific draft
105 metabolic reconstruction ('model') using the approach outlined previously (35), and
106 uses gap-filling to identify any missing reactions required to simulate growth in the user-
107 specified conditions
- 108 ii) Patching incomplete models (*patch_model* command) by the addition of missing
109 reactions e.g. those identified by the automated gap-filling process
- 110 iii) Substrate usage analysis via Flux Balance Analysis (FBA) (*fba* command) to predict
111 growth outcomes for a specified range of substrates supported by the model(s)
- 112 iv) **Fig. 1**.

113 Additional processing scripts are provided alongside Bactabolize to improve model metadata
114 annotation (*improve_model_annotations.py*), convert models generated using KBase and
115 ModelSEED to Bactabolize/BiGG-compatible format (*SEED_to_BiGG_model_convert.sh*),
116 generate network graph files from models (*model_to_network_graph.py*) and merging output FBA
117 profiles (*merge_fba_profiles_longtable.sh*).

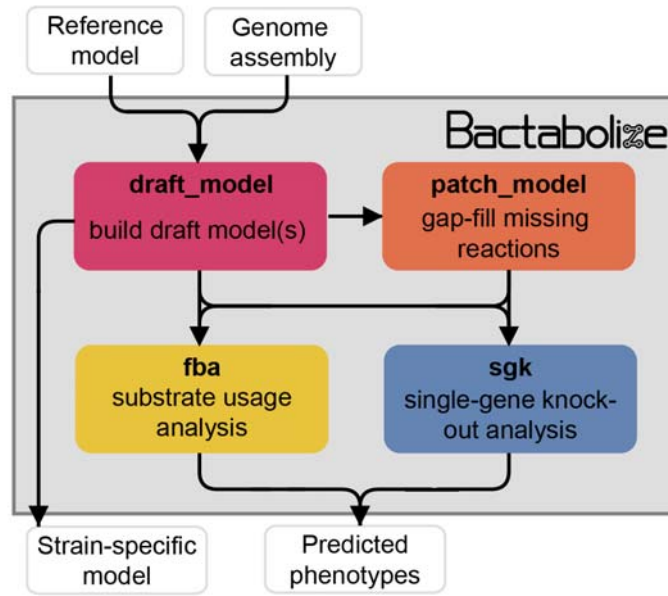
118 For draft model construction, Bactabolize requires users to provide an input assembly (annotated
119 or unannotated FASTA or Genbank format respectively), a reference model (JSON format) and the
120 corresponding reference sequence data (gene and protein sequences in two separate multi-fasta
121 files or a single Genbank annotation in a .gbk file) (**Figure S1**). If the input assembly is
122 unannotated, Bactabolize will identify coding sequences using Prodigal (38) but will otherwise
123 honour the existing coding sequence (CDS) notations and optionally use Prodigal to search for
124 additional CDS. Draft genome-scale metabolic models are output in both SMBL v3.1 (39) and
125 JSON formats (one pair of files for each independent strain-specific model), along with an optional
126 MEMOTE quality report (40). Bactabolize will identify orthologs in the input genome(s) compared
127 to the reference sequence data using Bi-directional BLAST (41) Best Hits (BBH) (42) using
128 BLAST+ (35). Users can parameterise the ortholog finding settings (coverage and identity
129 thresholds) for BBH. Alternatively, there is the option of using protein similarity to identify orthologs
130 instead of identity.

I31 Once a draft model has been constructed, it is validated via a simulated growth experiment on
I32 user-input choice of media and atmosphere (aerobic or anerobic). Predefined media include BG11
I33 (Gibco), M9 + glucose (35), nutrient media (43), Luria-Bertani (LB) (43), Tryptic Soy (TSA) (43),
I34 TSA + sheep blood (43), LB as specified by the CarveMe developers (30), Chemically Defined
I35 Medium (CDM)-like (33), Plantarum Minimal Medium (PMM) PMM5-like (33) and PMM7-like (33).
I36 Users can also define custom media as Bactabolize supports several complex media ingredients,
I37 including peptone (peptic digest of bovine and porcine tissue) (44-46), tryptone (pancreatic digest
I38 of casein) (44, 46, 47), soy peptone/soytone (digest of soymeal) (44, 46, 48, 49), yeast extract (50-
I39 55) and beef extract (44, 46). If the model fails to simulate growth, gap-filling is performed to
I40 indicate missing reactions. Users can add these reactions to a patch JSON file and optionally use
I41 the *patch_model* command to correct the model (**Figure S2**). Bactabolize uses a conservative
I42 gap-filling approach that only adds the minimum number of reactions to enable growth under the
I43 chosen conditions. We recommend testing the models in minimal media and atmosphere expected
I44 to support growth for all isolates of the species of interest, unless the user has access to matched
I45 phenotypic data demonstrating growth for individual isolates in specific conditions. Aggressive
I46 gap-filling will effectively homogenise the models and should be avoided if the goal is to
I47 understand the underlying strain diversity.

I48 Substrate usage analysis (the *fba* command) is performed iteratively for each possible carbon,
I49 nitrogen, sulfur and phosphor substrate supported by the model(s) (**Figure S3**), by replacing the
I50 default substrate in the user specified growth medium (specified in the *fba_spec* JSON file). For
I51 example, in M9 media the default substrates are glucose (carbon), ammonia (nitrogen), sulphate
I52 (sulfur) and phosphate (phosphor). Each substrate can be tested in aerobic and/or anaerobic
I53 conditions. Growth prediction output is recorded in a tab delimited file (one per strain). The
I54 *merge_fba_profiles_longtable.sh* helper script will combine the outputs for multiple strains into a
I55 single file for downstream analysis.

I56 The growth impacts of single-gene knockout mutations can be simulated via the *sgk* command
I57 (**Figure S4**). Bactabolize will iterate through every gene in the model, temporarily removing it and
I58 its associated reactions (unless they are also associated with another gene) and running FBA to
I59 simulate growth in the user-specified conditions. The output is comparable to single-gene knockout
I60 studies such as transposon mutagenesis and can be used to probe gene essentiality.

I61 We recorded the time required for Bactabolize to build draft models and performed 1692
I62 independent growth predictions for each of 35 *KpSC* genomes (tested in triplicate) on a high-
I63 performance computing cluster (Intel(R) Xeon(R) Platinum 8260 CPU @ 2.40GHz and 340 GB of
I64 requested memory on a CentOS Linux release 7.9.2009 environment). The mean CPU time
I65 required for model construction was 98.41 seconds (range 83.72 - 112.55 seconds), while the
I66 mean CPU time for growth predictions was 88.79 seconds (range 85.15 - 103.37 seconds). On a
I67 standard consumer laptop (Intel(R) Core(TM) i5-8365U CPU @ 1.60GHz and 15 GB of memory on
I68 Windows Subsystem for Linux (WSL1) environment), the mean CPU time for model construction
I69 was 87.27 seconds (range 76.133 – 102.694 seconds), while growth predictions took 82.19
I70 seconds (range 72.24 – 103.37 seconds).



171

172 **Figure 1:** Simplified overview of Bactabolize's main commands. In pink is the *draft_model* command, which builds a draft
173 strain-specific metabolic model using an input reference model and an input target assembly (approach adapted from
174 (35)). If the model fails to simulate growth, Bactabolize will attempt automated gap-filling and produce a model patch file.
175 The *patch_model* command (orange) allows the addition of missing reactions to produce a valid draft model that can
176 simulate growth in a user-specified growth environment. A functioning model can be passed to the *fba* command
177 (yellow), which performs Flux Balance Analysis to simulate growth in the user specified conditions, across all carbon,
178 nitrogen, phosphorus and sulfur metabolite sources supported by the model under aerobic and anaerobic conditions. The
179 *sgk* command (blue) shows the Single Gene Knockout analysis, which outputs a predicted phenotype. User inputs and
180 outputs are shown in white boxes while Bactabolize commands are shown inside the grey box.

181 **KpSC pan-metabolic reference model**

182 We constructed a species complex-specific pan-metabolic reference model by combining a
183 collection of 37 manually curated models for which we have previously demonstrated high
184 accuracy (range 88.3%–96.8% for prediction of 94 distinct growth phenotypes (14)). These models
185 represent a diverse collection of *KpSC* (14) (including at least one each of the seven major taxa in
186 the complex; *K. pneumoniae*, *Klebsiella variicola subsp variicola*, *Klebsiella variicola subsp tropica*,
187 *Klebsiella quasipneumoniae subsp quasipneumoniae*, *Klebsiella quasipneumoniae subsp*
188 *similipneumoniae*, *Klebsiella quisvariicola*, *Klebsiella africana*). The combined pan-model, known
189 as *KpSC*-pan v1, comprises a total of 1265 distinct genes, 2319 reactions and 1696 metabolites,
190 and is available at github.com/kelwyres/KpSC-pan-metabolic-model.

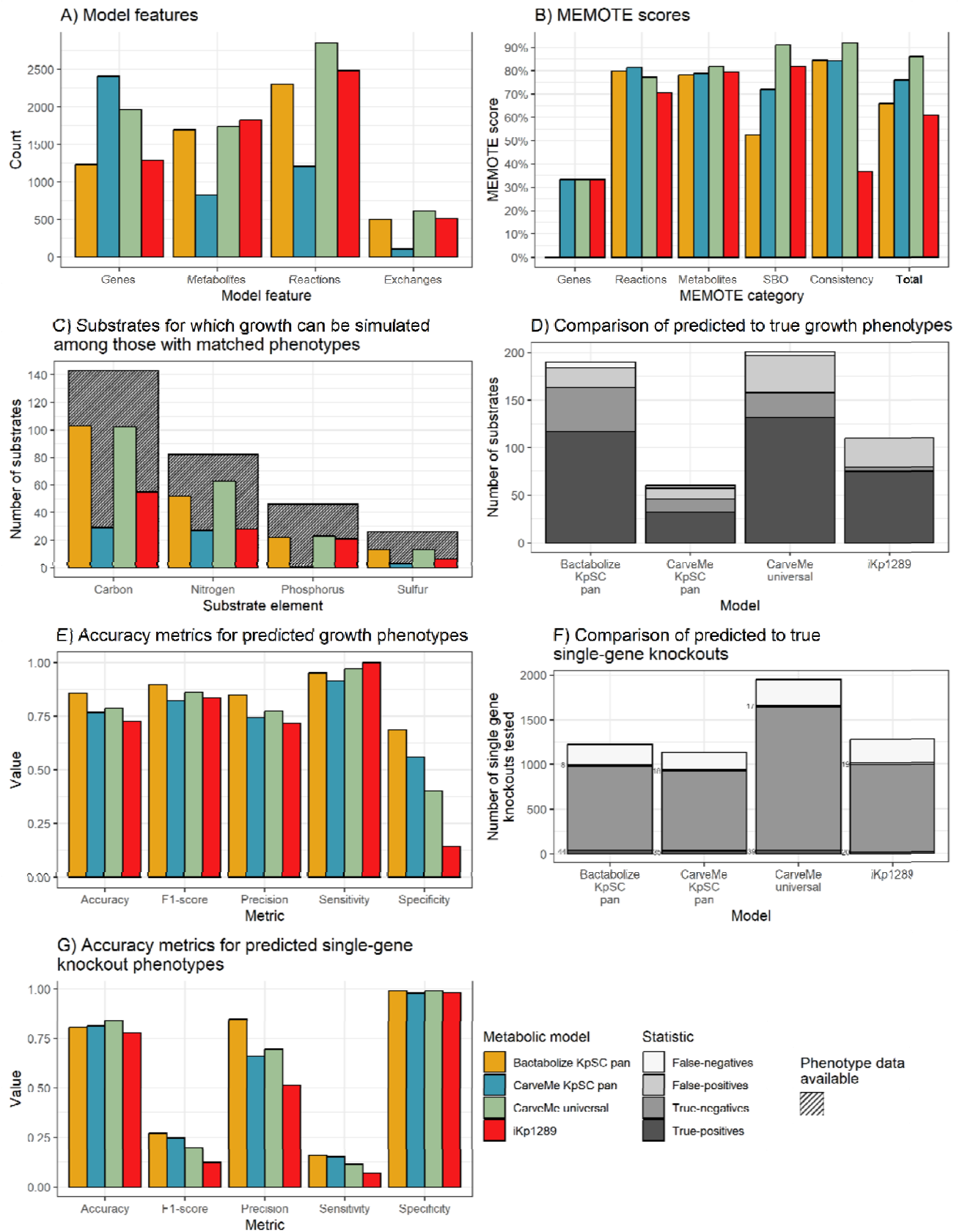
191 **Performance comparison**

192 We compared the output and performance of Bactabolize to CarveMe (30) and a manually curated
193 metabolic reconstruction of *K. pneumoniae* strain KPPR1 (also known as VK055 and ATCC
194 43816, metabolic model named iKp1289) (15). This isolate was chosen as there is a completed
195 genome sequence (Genbank accession: CP009208), single-source growth phenotype (15) and
196 single-gene knockout growth essentiality data available (56). Draft models were built using; i)
197 Bactabolize with the *KpSC* pan v1 reference; ii) CarveMe, with its universal reference model
198 (CarveMe universal); and iii) CarveMe, with *KpSC*-pan v1 reference (CarveMe *KpSC* pan).
199 Importantly, neither *K. pneumoniae* KPPR1 nor its genetic lineage (7 gene multi-locus sequence
200 type, ST493), are represented in the *KpSC* pan-reference model, meaning these benchmarking
201 comparisons were on equal footing.

202 The Bactabolize draft model captured a comparable number of genes and reactions (n = 1233 and 2307, respectively) to
203 the manually curated model (n = 1289 and 2484, respectively) but fewer than the CarveMe universal model (n = 1960
204 and 2857)

205 **Fig. 2A).** In contrast, the number of metabolites represented in the Bactabolize and CarveMe
206 universal models were similar (1696 vs 1737) and both were lower than the number represented in
207 iKp1289 (n = 1827). The CarveMe *KpSC* pan model method captured considerably more genes
208 than any of the other models (n = 2407), but these were associated with many fewer unique
209 reactions and metabolites (1206 and 825, respectively). Upon further investigation we determined
210 that this method resulted in the over prescription of gene reaction rules (GPRs) to multiple
211 reactions (mean 2.2 GPRs per reaction when compared to Bactabolize using the same pan
212 reference model: 1.94 GPRs per reaction; and CarveMe Universal: 2.12 GPRs per reaction).
213 MEMOTE scores, (produced by the MEMOTE report (40)) indicate the quality of the model
214 metadata annotations, with the scores ranging between 0 – 100%. These provide a measure of
215 model portability and the level of connected databases available to support the metabolite, reaction
216 and genetic information represented in the model, but bear no reflection on model accuracy.
217 Bactabolize performs on the lower end, with CarveMe universal performing the best (**Fig. 2B**).
218 However, Bactabolize using the *KpSC*-pan model outperforms the model propagation mode of
219 CarveMe using the same reference model (**Fig. 2B**). Work is ongoing to improve the annotations
220 in the *KpSC*-pan reference model, to improve large-scale model propagation.

221 We assessed the performance of each model for *in silico* prediction of growth phenotypes
222 compared to the previously published experimental data (15). Accuracy, sensitivity, specificity,
223 precision and F1 scores were calculated (57). Note that the specific set of growth substrates and
224 gene knockouts that can be simulated is determined by the sets of genes and metabolites
225 captured by each model and is therefore model-dependent (**Data S1 and S2**). Among those with
226 matched experimental phenotype data, the Bactabolize and CarveMe universal models were able
227 to predict growth for a greater number of carbon, nitrogen, phosphorous and sulfur substrates than
228 both the iKp1289 model and the CarveMe *KpSC* pan models (**Fig. 2C, Data S1**). While the
229 CarveMe universal model had the highest number of true-positive growth predictions overall, it
230 also had a comparably high number of false-positive predictions (**Fig. 2D**). In contrast, the
231 Bactabolize model had fewer false-positive predictions, resulting in the highest overall accuracy
232 metrics (**Fig. 2E, Data S1**). Similarly, while the CarveMe universal model resulted in the highest
233 absolute number of true-positive gene essentiality predictions, driving a high accuracy, the
234 Bactabolize model was associated with the greatest overall precision, sensitivity, and specificity
235 (**Figs. 2F & 2G**).



238 **Figure 2 (previous page):** *K. pneumoniae* KPPR1 metabolic model benchmarking comparisons. **A)** Counts of model
239 features; genes, metabolites and reactions captured by each model. Exchanges refers to number of exchange reactions,
240 a subset of reactions involved in substrate uptake, which determine the number of distinct growth substrates for which
241 phenotypes can be predicted with the model. **B)** MEMOTE scores indicating the richness of annotations and metadata
242 for metabolic model features according to database outlinks. SBO refers to score of Systems Biology Ontology (SBO), a
243 controlled vocabulary for systems biology. Consistency refers to the score of stoichiometric consistency and chemical
244 formulae annotation. Total refers to total MEMOTE score, as a combination of all previous scores, and is shown in bold.
245 **C)** Counts of carbon, nitrogen, phosphorus and sulfur growth substrates that can be simulated by models and for which
246 matched phenotypes were available for comparison (15). Hatched columns indicate the total number of substrates for
247 which phenotypic data for *K. pneumoniae* KPPR1 were described (15). **D)** and **E)** Accuracy metrics for predicted to true
248 phenotypes for the growth substrates shown in D and E, respectively False-negatives, true-negatives, false-positives
249 and true-positives are coloured as shown in legend. **F)** and **G)** Accuracy metrics for the KPPR1 single-gene knockout
250 mutant library described in (56) shown in F and G, respectively. Numbers of true positives and false positives are shown
251 to the left of the respective columns.

252 Quality control framework for input genome assemblies

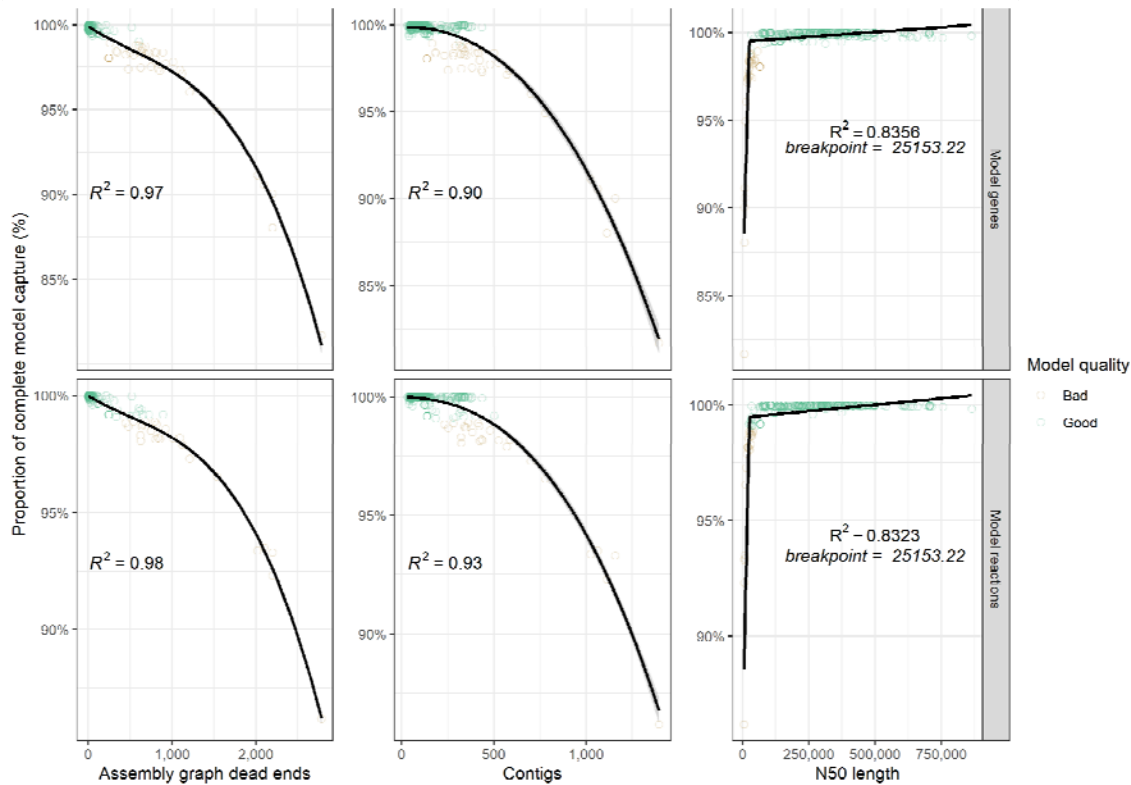
253 There are now thousands of bacterial genomes available in public databases, the majority of which
254 are in draft form. If we are to use these data for high-throughput metabolic modelling studies, it is
255 essential to evaluate the expected model accuracies and understand the minimum input genome
256 quality requirements. Here we performed a systematic analysis leveraging our published curated
257 *KpSC* models (n=37, (14)), which were generated using completed genome sequences and were
258 therefore considered to represent 'complete' models. We randomly subsampled the corresponding
259 Illumina read sets to various depths (10 – 100x, increments of 10) in triplicate and generated draft
260 assemblies that were passed to Bactabolize for generation of draft metabolic models (**Data S3**).
261 Due to low read depth ($\leq 30x$), two isolates were removed from this analysis. Additionally, ten 10x
262 depth read samples failed to produce assemblies, leaving 1040 draft genomes for analysis. The
263 resulting draft metabolic models were compared to the complete models to; i) determine the
264 proportions of complete model genes and reactions captured in the draft models; and ii) compare
265 846 *in silico* aerobic growth predictions in M9 minimal media, where growth on 266 carbon, 153
266 nitrogen, 59 phosphorus and 25 sulfur sources were examined. Substrates containing multiple
267 elements were tested as sole sources of each element independently and in combination, e.g. 1,5-
268 Diaminopentane was tested as a sole carbon, sole nitrogen and sole carbon plus nitrogen source.

269 As expected, assembly quality generally increased with increasing sequencing depth i.e.
270 assemblies generated from higher depth read sets were associated with higher N50 values, fewer
271 contigs and fewer assembly graph dead-ends, although the rate of improvement drastically
272 declined beyond 40-50x depth (**Figure S5, Data S3**). We noted that it was rare for draft models to
273 capture 100% of model genes and reactions (just 420 of all 1040 draft assemblies were associated
274 with models that captured 100% model genes) (**Data S3, Figure S6**), even when using the highest
275 quality draft genomes. However, $\geq 99\%$ of genes and reactions were commonly captured, which
276 plateaued from 40x depth onwards (**Figure S6**). Therefore, we sought to evaluate whether $\geq 99\%$
277 model capture would produce functionally accurate models.

278 We used FBA to simulate substrate growth profiles for the 40x depth assemblies, representing a
279 sequencing depth that can be routinely achieved with standard Illumina library preparations. All but
280 one assembly triplicate set (isolate SB4767 98% gene capture, 99% reaction capture) captured
281 $\geq 99\%$ but $\leq 100\%$ model genes and/or reactions. The substrate growth profiles were then
282 compared to those of the complete models. The vast majority of draft models produced accurate
283 growth predictions; 102 of 108 models resulted in predictions with 100% concordance to those
284 from the corresponding complete models. Three models for *K. quasipneumoniae similibpneumoniae*
285 isolate SB164 resulted in predictions with a mean of 99.8% concordance. The remaining three
286 models were for isolate SB4767 and resulted in mean of 80.4% concordance. Notably, these
287 models were those representing $< 99\%$ gene capture. Together, these data suggest that draft

288 models capturing $\geq 99\%$ of the complete model genes/reactions generate highly accurate growth
289 predictions and that these capture rates can be readily achieved from draft genome assemblies.

290



291

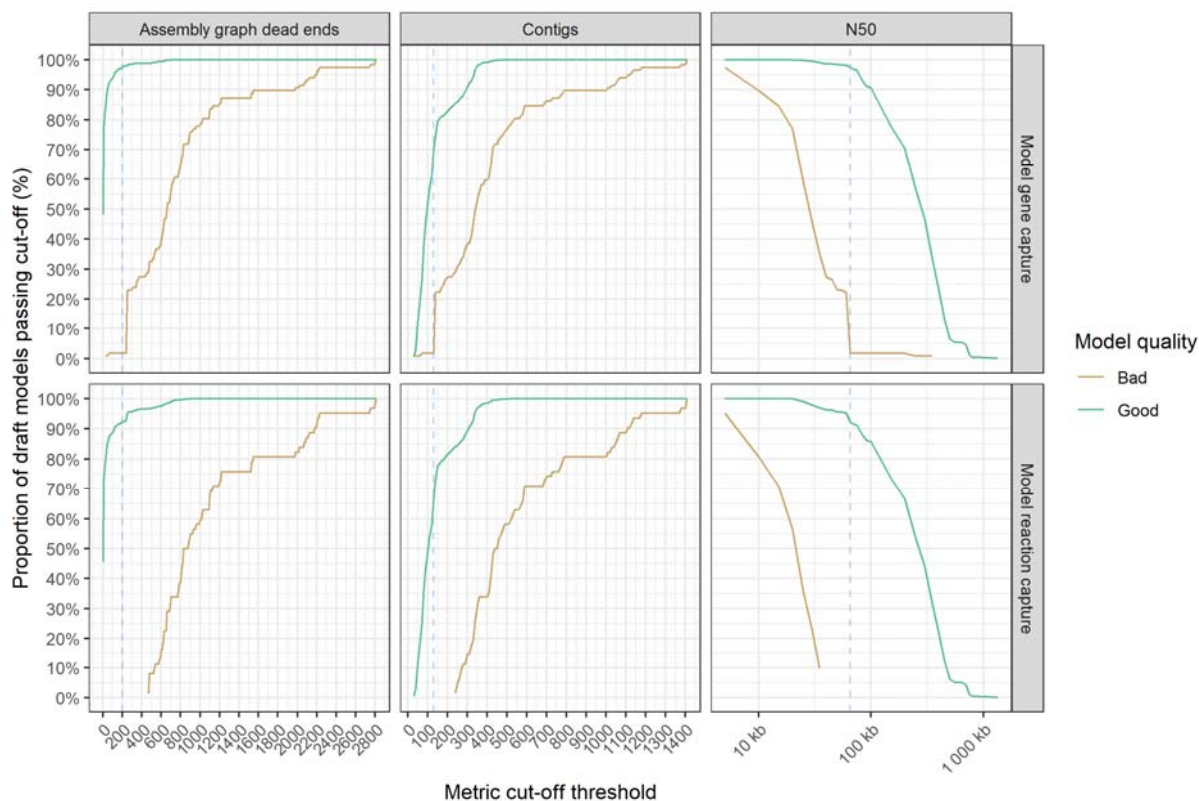
292 **Figure 3:** Scatterplots showing distribution of best performing assembly metrics 'assembly graph dead ends', 'contigs'
293 and 'N50' against model feature capture (genes and reactions). Each point represents the mean values from a single
294 genome (technical triplicate) and is coloured by model quality. 'Good' models capture $\geq 99\%$ of the model metric as
295 compared to the corresponding complete model (shown at each facet), 'Bad' models capture $< 99\%$. Cubic polynomial
296 line plotted for assembly 'graph dead ends', 'contigs', while a segmented linear model was plotted for 'N50'. R^2 is shown
297 on each panel.

298 We investigated the relationships between assembly quality metrics and model gene/reaction
299 capture in more detail. Variation in assembly graph dead-ends accounted for the greatest amount
300 of variation in model capture, closely followed by raw contig counts (cubic polynomial fit, R^2 of
301 ≥ 0.98 for graph dead-ends, R^2 of ≥ 0.9 for contig count). A segmented linear model was fitted to
302 N50 length ($R^2 \geq 0.83$), producing a breakpoint at 25153 bp (**Fig. 3**).

303 To further explore the optimum thresholds for assembly metrics, we tallied the number of draft
304 assemblies resulting in $\geq 99\%$ and $< 99\%$ gene and reaction capture at increasing graph dead-end
305 and contig count count-offs, and decreasing N50 cut-offs. Draft models that captured $\geq 99\%$ of the
306 complete model genes/reactions were considered 'good' models, whereas draft models that
307 captured $< 99\%$ of complete model genes/reactions were considered 'bad' models. The optimum
308 threshold for assembly graph dead end was determined to be ≤ 200 . At this value, 94.44% of 'good'
309 models were captured, and 0% 'bad' models. The optimum threshold for contig counts was
310 determined as ≤ 130 contigs at which 67.92% of 'good' and 0% 'bad' models were captured (**Fig.**
311 **4**). The optimum threshold for N50 was determined to be ≥ 65000 , at which 94.97% of 'good' and
312 1.71% of 'bad' models were captured. The assembly graph dead-end threshold results in
313 comparatively higher sensitivity (i.e. a higher proportion of 'good' models pass the threshold) than
314 contig count and comparatively better specificity (i.e. lower proportion of 'bad' models pass the

315 threshold) than N50, but the underlying metric information is not universally available because
316 many isolate genomes are deposited in public databases only as assemblies without the
317 associated assembly graph. We therefore recommend a three-tier approach, whereby the
318 assembly graph dead-end criterion is preferred if available, followed by N50 and then contig
319 count.

320



321

322 **Figure 4:** Line graphs showing the impact of assembly metric cut-off thresholds on model feature capture ($n = 1040$).
323 'Good' models which captured $\geq 99\%$ of model features are shown in green, while 'bad' models captured $< 99\%$ model
324 features are shown in gold. The blue dotted line shows the metric cut-off thresholds, to minimize the number of models
325 that capture $< 99\%$ model features and maximise models that capture $\geq 99\%$. Metric cut-off statistics are calculated in
326 intervals of 10 for assembly graph dead ends and contigs, and every 5000 for N50.

327

328 Impact of gap-filling models

329 Of the 901 draft genome assemblies which passed our QC criteria (≤ 200 assembly graph dead
330 ends), 23 of the resulting draft models failed to simulate growth in M9 minimal media with glucose
331 (despite capturing $\geq 99\%$ of the genes and reactions in the corresponding complete models). It is
332 expected that all *KpSC* models should be able to simulate growth on M9 media with glucose as a
333 sole carbon source, as this central metabolism is universal amongst *KpSC*. To replace missing,
334 critical reactions required for growth on M9 with glucose, we investigated model gap-filling using
335 the *patch_model* command of Bactabolize. We then assessed the accuracy of the gap-filled
336 models for prediction of growth on the full range of substrates, as compared to the predictions from
337 the corresponding complete models.

338 Gap filling added 1 – 3 missing reactions to each model, with a median of one, fully restoring
339 biomass production in M9 media with glucose in all but two of the 23 failed models. The missing

340 reactions appeared to be random genes across these 23 genomes, likely due to missing
341 information in these assemblies.

342 Substrate usage predictions from the 21 successfully gap-filled models were highly accurate, with
343 18/21 having a prediction concordance of $\geq 99\%$ across all 846 growth conditions (12/21 had 100%
344 concordance) (**Figure S7**). We therefore conclude that models generated for genome assemblies
345 passing our QC criteria, which have been gap-filled to successfully simulate growth on minimal
346 media plus glucose, are suitable for the prediction of growth across a range of substrates.

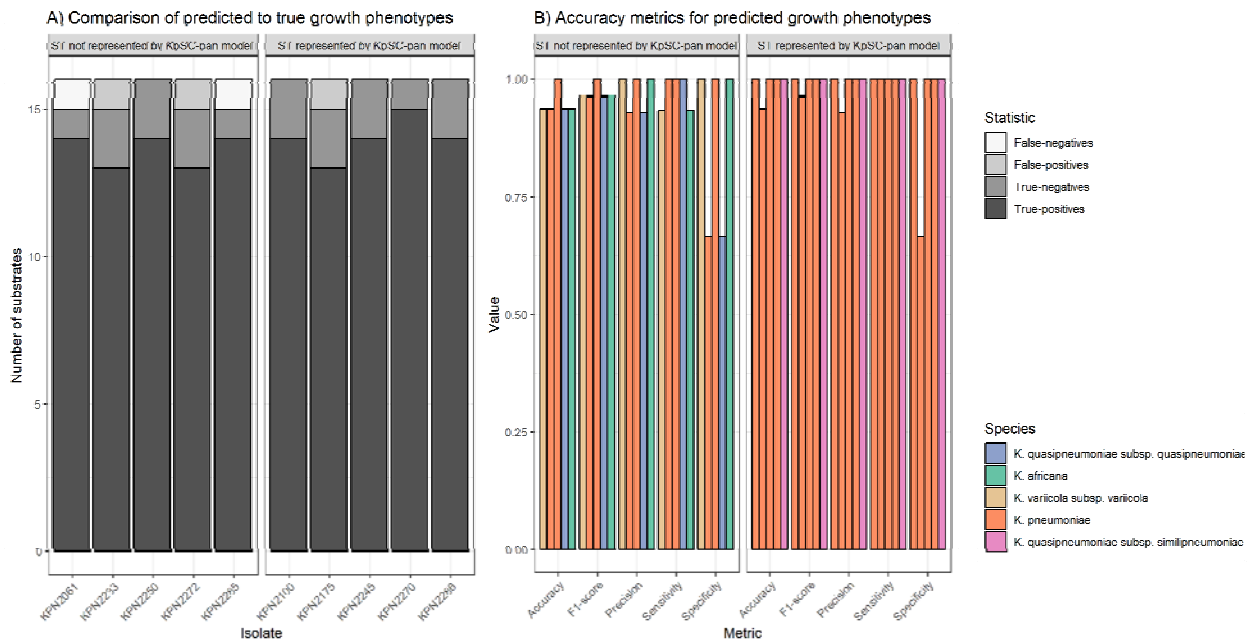
347

348 Predictive accuracy of draft models

349 We assessed the accuracy of Bactabolize for the construction of draft models for 10 novel *KpSC*
350 clinical isolates, representing five of the major taxa in the complex. We included five isolates for
351 which the associated STs were represented in the *KpSC*-pan v1 model and five isolates with STs
352 that were not represented. Whole genome sequence data were generated on the Illumina platform
353 and draft assemblies generated *de novo*. The resultant assemblies had 0-4 graph dead-ends,
354 N50s of 151958-388486 bp and 83-187 contigs (**Data S4**), within the tiered threshold values.

355 FBA was performed, and the predicted growth profiles compared to matched phenotypic growth
356 data for 16 carbon sources derived from Vitek GN ID cards. Though the number of tested carbon
357 sources was limited, all were associated with high accuracy metrics (**Fig. 5, Data S4**). As
358 expected, models for isolates with STs represented in the *KpSC*-pan v1 reference performed
359 slightly better (mean accuracy = 0.98) than those for non-represented STs (mean accuracy =
360 0.95).

361



362

363 **Figure 5: A)** Comparisons of predicted to true phenotypes for 16 carbon source substrates. False-negatives, true-
364 negatives, false-positives and true-positives are coloured as shown in legend. Each column represents a different
365 isolate, separated by ST representation in the *KpSC*-pan model. **B)** Accuracy metrics for predicted vs phenotypic growth
366 comparisons shown in A. Each column represents a different isolate, coloured by taxa and separated by ST
367 representation in the *KpSC*-pan v1 model.

368 Discussion

369 In this work we described Bactabolize, a pipeline for rapid and scalable production of accurate
370 bacterial strain-specific metabolic models and growth phenotype predictions. We describe a pan-
371 reference model for the *KpSC* and demonstrate that a draft strain-specific model generated *de*
372 *novo* via Bactabolize using the *KpSC*-pan v1 reference was highly accurate for growth phenotype
373 prediction (85.79% accuracy for substrate usage across 190 substrates, and 80.57% for gene
374 essentiality across 1220 genes). Importantly, we also described a quality control framework for the
375 use of draft genome assemblies as input for metabolic reconstructions. We used a systematic
376 analysis to; i) evaluate the proportion of gene and reaction capture compared to the corresponding
377 'completed' models; ii) define quality control thresholds for input assemblies (three tier approach
378 for *KpSC*; ≤ 200 assembly graph dead ends, followed by ≥ 65000 N50, followed by ≤ 130 contigs);
379 and iii) estimate the accuracy of the resultant growth predictions. While the quality control
380 thresholds and accuracy estimates are specific to *KpSC*, the conceptual framework can be applied
381 to any organism and is essential to support the confident application of metabolic modelling for
382 large-scale genome datasets. We appreciate that assembly graphs may not be available for dead
383 end count, e.g. for draft genome assemblies accessed via public repositories, however we
384 encourage users to include this information in their quality control procedures wherever possible
385 (e.g. using the recently published counter tool available at [https://github.com/rswick/GFA-dead-](https://github.com/rswick/GFA-dead-end-counter)
386 [end-counter](https://github.com/rswick/GFA-dead-end-counter)), because these counts represent a direct reflection of the completeness of the
387 genome assembly. In contrast, contig counts and N50 are influenced by biological features such
388 as repeat copy numbers as well as the underlying sequence data quality e.g. a bacterial genome
389 harbouring many insertion sequence insertions will result in a draft assembly with a high number of
390 contigs regardless of the sequence data quality and completeness.

391 Bactabolize's reference-based reconstruction approach is reductive, meaning the resultant draft
392 models will comprise only the genes, reactions and metabolites present in the reference, or a
393 subset thereof, and will not include novel reactions unless they are manually identified and curated
394 by the user. This is an important caveat that should be considered carefully for application of
395 Bactabolize to large genome data sets, particularly for genetically diverse organisms such as those
396 in the *KpSC*. The use of a pan-reference derived from multiple curated strain-specific models
397 results in greater representation of the population diversity and partially alleviates the
398 shortcomings of the reference-based approach. However, draft models constructed for strains with
399 a corresponding lineage represented in the reference are likely more accurate. Our analysis
400 indicated that a draft *KpSC* model generated by Bactabolize with the *KpSC* pan v1 reference was
401 equally or more accurate than the current gold standard automated approach, CarveMe with
402 universal model (30), and outperformed a manually curated model (15). The latter was constructed
403 using the KBase pipeline (32), which uses RAST to annotate the sequences with Enzyme
404 Commission numbers. It has been demonstrated several times that the Enzyme Commission
405 scheme has systematic errors (58, 59), leading to a loss in accuracy when compared to the
406 ortholog identification methods used in CarveMe and Bactabolize.

407 The CarveMe universal reference model captures greater diversity than the *KpSC* pan v1
408 reference, which resulted in a comparatively greater number of genes, reactions and metabolites
409 in the corresponding CarveMe draft model, and ability to simulate growth outcomes for a greater
410 number of distinct substrates (**Figure 2**). Overall, CarveMe (with the universal model) performed
411 extremely well, with high numbers of true-positive growth predictions. However, these were also
412 accompanied by comparatively higher numbers of false-positive predictions, which resulted in a
413 lower overall accuracy score for substrate usage analysis compared to Bactabolize with the *KpSC*-
414 pan v1 reference (**Figure 2**), and comparatively lower sensitivity and specificity for the gene
415 essentiality analysis. False-positive predictions may indicate that the relevant metabolic machinery

116 are present in the cell but were not active during the growth experiments (e.g. due to lack of gene
117 expression). In this regard, false-positives are not always a sign of model inaccuracy. However,
118 false-positive predictions can also occur from incorrect gene annotations e.g. due to reduced
119 specificity of ortholog assignment resulting from the use of the universal model without manual
120 curation. Given a key objective here is to facilitate high-throughput analysis for large numbers of
121 genomes, it is not feasible to expect that all models will be manually curated, and therefore we
122 believe that identifying fewer genes with lower overall error rates provides greater confidence in
123 the resulting draft models. We also note that the CarveMe universal reference model is no longer
124 being actively maintained, but in contrast, user defined species- (or genera-) specific references
125 can be iteratively curated and updated to incorporate new knowledge and data as they become
126 available. Accordingly, the accuracy of models derived from such references is expected to
127 continually improve.

128 Bactabolize and the *KpSC* pan v1 model are freely available under open source licenses and
129 satisfy the four features of the FAIR research principles (findability, accessibility, interoperability
130 and reusability) (60). In addition to the *KpSC* pan-reference described here, a pan-reference model
131 has been described previously for *Salmonella enterica* (representing 410 strains (25)). We are
132 actively working to expand and improve the *KpSC* pan reference model and welcome similar
133 efforts to generate high quality references for other organisms. Together these resources will
134 facilitate population wide metabolic analyses for global priority pathogens, which can be used to
135 understand how they transmit, cause disease and evolve drug-resistance, and to identify novel
136 therapeutic targets.

137 **Methods**

138 **Bactabolize pipeline**

139 Bactabolize utilises the existing metabolic modelling library COBRAPy (36) and Python 3 (61). All
140 code is freely available and open source at GitHub (www.github.com/kelwyres/Bactabolize) under
141 a GNU General Public License v3.0. Users should additionally cite COBRAPy (36) if Bactabolize is
142 used.

143 ***Klebsiella pneumoniae* Species Complex-pan metabolic model**

144 The 37 metabolic models from a previous study (14) were combined with the iY1228 model using
145 the `create_master_model.py` script (available at 10.6084/m9.figshare.21728717). Briefly, all GPRs
146 from the iYL1228 model and the associated sequences were included, as well as new GPRs
147 identified from the 36 additional strains by manual curation following comparison to the matched
148 phenotype data (as described in (14)). Additionally, orthologous sequence variants with <75%
149 nucleotide identity to gene sequences associated with these gene reaction rules (GPRs) were
150 added if there was phenotype data supporting the reaction. The biomass reaction was updated,
151 removing the metabolites `udpgalr_c` and `udpgalc_c` as their production was strain-specific.

152 Metadata annotations were improved using the `improve_model_annotations.py` script (also
153 available in the Bactabolize code repository) resulting in the *KpSC*_pan v1 used in this study,
154 available at www.github.com/kelwyres/KpSC-pan-metabolic-model.

155 **Draft model generation**

156 Bactabolize draft models were generated using the `draft_model` command in Bactabolize v1 with
157 the *KpSC*-pan v1 model as a reference, and the following options:

158 `--min_coverage 25 --min_pident 80 --media M9 --atmosphere aerobic`

159 CarveMe draft models were generated firstly using the universal reference with the following
160 commands: '-g M9 -i M9'. The --universe-file mode was also used, so the *KpSC-pan* model could
161 be used as a reference, with the previously described command.

162 **Speed calculations**

163 Bactabolize *draft_model* and *fba* commands were timed via a script using the date +%s.%N
164 command run before and after command on the MASSIVE computing cluster (Intel(R) Xeon(R)
165 Platinum 8260 CPU @ 2.40GHz and 340 GB of memory, CentOS Linux release 7.9.2009
166 environment). Speed tests were also performed on a standard consumer laptop with the following
167 hardware: Intel(R) Core(TM) i5-8365U CPU @ 1.60GHz and 15 GB of memory on Windows
168 Subsystem for Linux (WSL1) environment.

169 **Performance comparison**

170 The genome of *K. pneumoniae* KPPR1 was obtained from Genbank under the accession:
171 CP009208, and draft metabolic models were generated using Bactabolize and CarveMe as
172 described above. The previously described, manually curated model for KPPR1 (iKp1289) was
173 also included for comparison (15). The following KPPR1 phenotype data were retrieved from
174 published studies: BIOLOG Phenotypic Microarray data (15) and single gene knockout data
175 inferred from the outputs of a TraDIS transposon mutagenesis library (56).

176 A list of BIOLOG growth substrates for plates PM1, PM2A, PM3B, and PM4A (62) were converted
177 where possible to BiGG and SEED IDs by manual search of the BiGG (bigg.ucsd.edu) and SEED
178 websites (<https://modelseed.org/biochem/compounds>). An updated BiGG to SEED dictionary can
179 be found in **Data S1**. A total of 143 of 190 carbon, 82 of 95 nitrogen, 46 of 59 phosphor and 26 of
180 35 sulfur substrates were successfully matched to BiGG and SEED IDs (**Data S1**). These growth
181 data were compared to *in silico* predictions generated via FBA using the *fba* command from
182 Bactabolize to optimise the biomass objective function with the following options:

183 --fba_spec_name m9 --fba_open_value -20

184 Gene essentiality was inferred from single gene knockout growth predictions using the *sgk*
185 command from Bactabolize with the following options to mirror the growth conditions of the TraDIS
186 library (LB media grown aerobically):

187 --media_type lb --atmosphere aerobic

188 In all cases, an objective value cut-off of $\geq 10^{-4}$ was used to indicate binarised growth as per
189 previous studies (14, 63).

190 *In silico* predictions were compared to matched phenotype data and the following accuracy metrics
191 were calculated:

$$Precision = \frac{TP}{TP + FP}$$

$$Sensitivity/recall = \frac{TP}{TP + FN}$$

$$Specificity = \frac{TN}{TN + FP}$$

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$F1 - score = 2 \times \frac{Precision \times Sensitivity}{Precision + Sensitivity}$$

192

193 **Quality control framework**

194 Illumina read sets (250 bp paired end) and completed genome sequences for 37 *KpSC* isolates
195 were described previously (14). Here we randomly subsampled the Illumina reads at various
196 depths (10 – 100, by increments of 10) using *rasusa* version 0.3.0 (64) in technical triplicate.
197 Reads were then trimmed using *TrimGalore* version 0.5.0 (65) and assembled *de novo* with
198 *Unicycler* version 0.4.7 (66), default parameters. Assembly statistics and assembly graph dead
199 ends were calculated using the *GFA-dead-end-counter* version 1.0.0
200 (<https://github.com/rwick/GFA-dead-end-counter>) (67). Draft metabolic models were generated
201 with *Bactabolize* using the *KpSC*-pan v1 reference, and growth substrate profiles were predicted
202 as described above. We compared the outputs from models generated for draft genome
203 assemblies to those generated for the corresponding completed genomes. Where necessary
204 models were gap-filled via the *patch_model* command.

205 **Predictive accuracy of draft models**

206 Novel growth phenotype data were generated for 10 *KpSC* clinical isolates from our in house
207 collection using the VITEK 2 GN ID card system as described previously (14). Briefly, isolates
208 were grown on Tryptic Soy (OXOID) agar plates overnight at 37°C, then analysed using VITEK 2
209 GN ID cards (bioMérieux) and read on the VITEK 2 Compact (bioMérieux) as per manufacturer's
210 instructions using software version 8.0. DNA was extracted for whole-genome sequencing via
211 Genfind v3 extraction kit, library preparation performed using Nextera Flex (Illumina) using ¼
212 reagents. Paired-end read data (300 bp) were generated on an Illumina NovaSeq6000 SP v1.0
213 and have been deposited in the European Nucleotide Archive under Bioproject PRJNA777643
214 (individual read accession numbers are given in **Data S4**). Draft genome assemblies were
215 generated with *Unicycler*, and draft metabolic models and growth predictions were generated with
216 *Bactabolize* as described above.

217 **Statistics and visualisation**

218 Statistical analysis and graphical visualisation were performed using R version 4.0.3 [23], RStudio
219 version 1.3.1093 (68), with the following software packages: *tidyverse* version 1.3.1 (69), *viridis*
220 version 0.5.1 (70), *RColorBrewer* version 1.1-2 (71), *ggpubr* version 0.4.0 (72) *ggpmisc* version
221 0.4.4 (73), *aplot* version 0.1.6 (74), *colorspace* version 2.0-2 (75), *ggpattern* version 0.4.3-3 (76),
222 *ggtext* version 0.1.1 (77) and *glue* version 1.4.2 (78).

223 Linear regression analysis was performed in R using the *lm* function in R and a third degree
224 polynomial model was fitted to plots with the following equation: $y \sim \text{poly}(x, 3, \text{raw} = \text{TRUE})$. The
225 segmented linear model was fitted using *segmented* version 1.6-2 (79).

226 All code used to generate results can be found as supplemental material,
227 <https://github.com/kelwyres/Bactabolize> and on Figshare (10.6084/m9.figshare.21728717).

228 **Logo**

229 The *Bactabolize* logo was constructed in *Inkscape* version 1.0.1 (80). The font used is Proportional
230 TFB (81) and Element (82).

231

532 **Author contributions**

533 Conceptualization: BV, JH, JMM, KEH, KLW

534 Methodology: SCW, BV, LMJ, JH, JMM, KLW

535 Software: SCW, BV

536 Validation: BV, HBC, KLW

537 Formal analysis: BV, KLW

538 Investigation: BV, KLW

539 Resources: AJ

540 Writing - Original Draft: BV, KLW

541 Writing - Review & Editing: All authors

542 Visualization: BV, KLW

543 Supervision: JMM, KEH, KLW

544 Project administration: KLW

545 Funding acquisition: JMM, KEH, KLW

546 All authors contributed to and approve of the manuscript in its current form.

547

548 **Acknowledgements**

549 We thank Sylvain Brisse for the isolates used in this study and review of the manuscript.

550

551 **References**

- 552 1. Su Q, Guan T, He Y, Lv H. Siderophore Biosynthesis Governs the Virulence of
553 Uropathogenic *Escherichia coli* by Coordinately Modulating the Differential Metabolism. *Journal of*
554 *Proteome Research*. 2016;15(4):1323-32.
- 555 2. Wu Y, Meng Y, Qian L, Ding B, Han H, Chen H, et al. The Vancomycin Resistance-
556 Associated Regulatory System *VraSR* Modulates Biofilm Formation of *Staphylococcus epidermidis*
557 in an *ica*-Dependent Manner. *mSphere*.6(5):e00641-21.
- 558 3. Mir M, Prusic S, Kang C-M, Lun S, Guo H, Murry JP, et al. Mycobacterial gene *cuvA* is
559 required for optimal nutrient utilization and virulence. *Infection and immunity*. 2014;82(10):4104-17.
- 560 4. Vornhagen J, Sun Y, Breen P, Forsyth V, Zhao L, Mobley HLT, et al. The *Klebsiella*
561 *pneumoniae* citrate synthase gene, *gltA*, influences site specific fitness during infection. *PLOS*
562 *Pathogens*. 2019;15(8):e1008010.
- 563 5. Eberl C, Weiss AS, Jochum LM, Durai Raj AC, Ring D, Hussain S, et al. *E. coli* enhance
564 colonization resistance against *Salmonella* Typhimurium by competing for galactitol, a context-
565 dependent limiting carbon source. *Cell Host & Microbe*. 2021.
- 566 6. Jenior ML, Dickenson ME, Papin JA. Genome-scale metabolic modeling reveals increased
567 reliance on valine catabolism in clinical isolates of *Klebsiella pneumoniae*. *bioRxiv*.
568 2021:2021.09.08.459555.

- 569 7. Hudson AW, Barnes AJ, Bray AS, Zafar MA. *Klebsiella pneumoniae* L-Fucose metabolism
570 promotes gastrointestinal colonization and modulates its virulence determinants. bioRxiv.
571 2022:2022.05.18.492588.
- 572 8. Rodrigues C, Passet V, Rakotondrasoa A, Diallo TA, Criscuolo A, Brisse S. Description of
573 *Klebsiella africanensis* sp. nov., *Klebsiella variicola* subsp. *tropicalensis* subsp. nov. and *Klebsiella*
574 *variicola* subsp. *variicola* subsp. nov. Res Microbiol. 2019;170(3):165-70.
- 575 9. Blin C, Passet V, Touchon M, Rocha EPC, Brisse S. Metabolic diversity of the emerging
576 pathogenic lineages of *Klebsiella pneumoniae*. Environmental Microbiology. 2017;19(5):1881-98.
- 577 10. Brisse S, Fevre C, Passet V, Issenhuth-Jeanjean S, Tournebize R, Diancourt L, et al.
578 Virulent Clones of *Klebsiella pneumoniae*: Identification and Evolutionary Scenario Based on
579 Genomic and Phenotypic Characterization. PLOS ONE. 2009;4(3):e4982.
- 580 11. Mobegi FM, van Hijum SAFT, Burghout P, Bootsma HJ, de Vries SPW, van der Gaast-de
581 Jongh CE, et al. From microbial gene essentiality to novel antimicrobial drug targets. BMC
582 Genomics. 2014;15(1):958.
- 583 12. Hogan AM, Scoffone VC, Makarov V, Gislason AS, Tesfu H, Stietz MS, et al. Competitive
584 Fitness of Essential Gene Knockdowns Reveals a Broad-Spectrum Antibacterial Inhibitor of the
585 Cell Division Protein FtsZ. Antimicrobial Agents and Chemotherapy. 2018;62(12):e01231-18.
- 586 13. Edwards JS, Palsson BO. Systems Properties of the Haemophilus influenzae Rd Metabolic
587 Genotype. Journal of Biological Chemistry. 1999;274(25):17410-6.
- 588 14. Hawkey J, Vezina B, Monk JM, Judd LM, Harshegyi T, López-Fernández S, et al. A
589 curated collection of *Klebsiella* metabolic models reveals variable substrate usage and gene
590 essentiality. Genome Research. 2022.
- 591 15. Henry CS, Rotman E, Lathem WW, Tyo KEJ, Hauser AR, Mandel MJ. Generation and
592 Validation of the iKp1289 Metabolic Model for *Klebsiella pneumoniae* KPPR1. The Journal of
593 Infectious Diseases. 2017;215(suppl_1):S37-S43.
- 594 16. Liao Y-C, Huang T-W, Chen F-C, Charusanti P, Hong JSJ, Chang H-Y, et al. An
595 Experimentally Validated Genome-Scale Metabolic Reconstruction of *Klebsiella pneumoniae* MGH
596 78578, γ L1228. Journal of Bacteriology. 2011;193(7):1710-7.
- 597 17. Orth JD, Thiele I, Palsson BØ. What is flux balance analysis? Nature Biotechnology.
598 2010;28(3):245-8.
- 599 18. Stanway RR, Bushell E, Chiappino-Pepe A, Roques M, Sanderson T, Franke-Fayard B, et
600 al. Genome-Scale Identification of Essential Metabolic Processes for Targeting the Plasmodium
601 Liver Stage. Cell. 2019;179(5):1112-28.e26.
- 602 19. Ramos PIP, Fernández Do Porto D, Lanzarotti E, Sosa EJ, Burguener G, Pardo AM, et al.
603 An integrative, multi-omics approach towards the prioritization of *Klebsiella pneumoniae* drug
604 targets. Scientific Reports. 2018;8(1):10755.
- 605 20. Wyres KL, Lam MMC, Holt KE. Population genomics of *Klebsiella pneumoniae*. Nature
606 Reviews Microbiology. 2020;18(6):344-59.
- 607 21. WHO. WHO publishes list of bacteria for which new antibiotics are urgently needed 2017
608 [Available from: [https://www.who.int/news/item/27-02-2017-who-publishes-list-of-bacteria-for-](https://www.who.int/news/item/27-02-2017-who-publishes-list-of-bacteria-for-which-new-antibiotics-are-urgently-needed)
609 [which-new-antibiotics-are-urgently-needed](https://www.who.int/news/item/27-02-2017-who-publishes-list-of-bacteria-for-which-new-antibiotics-are-urgently-needed)].
- 610 22. Gorrie CL, Mirčeta M, Wick RR, Judd LM, Lam MMC, Gomi R, et al. Genomic dissection of
611 *Klebsiella pneumoniae* infections in hospital patients reveals insights into an opportunistic
612 pathogen. Nature Communications. 2022;13(1):3017.
- 613 23. Holt KE, Wertheim H, Zadoks RN, Baker S, Whitehouse CA, Dance D, et al. Genomic
614 analysis of diversity, population structure, virulence, and antimicrobial resistance *Klebsiella*
615 *pneumoniae*, an urgent threat to public health. Proceedings of the National Academy of Sciences.
616 2015;112(27):E3574.
- 617 24. Monk JM, Charusanti P, Aziz RK, Lerman JA, Premyodhin N, Orth JD, et al. Genome-scale
618 metabolic reconstructions of multiple *Escherichia coli* strains highlight strain-specific adaptations to
619 nutritional environments. Proceedings of the National Academy of Sciences. 2013;110(50):20338.
- 620 25. Seif Y, Kavvas E, Lachance J-C, Yurkovich JT, Nuccio S-P, Fang X, et al. Genome-scale
621 metabolic reconstructions of multiple *Salmonella* strains reveal serovar-specific metabolic traits.
622 Nature Communications. 2018;9(1):3771.

- 323 26. Bosi E, Monk JM, Aziz RK, Fondi M, Nizet V, Palsson BØ. Comparative genome-scale
324 modelling of *Staphylococcus aureus* strains identifies strain-specific metabolic capabilities linked to
325 pathogenicity. *Proceedings of the National Academy of Sciences*. 2016;113(26):E3801-E9.
- 326 27. Bartell JA, Blazier AS, Yen P, Thøgersen JC, Jelsbak L, Goldberg JB, et al. Reconstruction
327 of the metabolic network of *Pseudomonas aeruginosa* to interrogate virulence factor synthesis.
328 *Nature Communications*. 2017;8(1):14631.
- 329 28. Croucher NJ, Coupland PG, Stevenson AE, Callendrello A, Bentley SD, Hanage WP.
330 Diversification of bacterial genome content through distinct mechanisms over different timescales.
331 *Nature Communications*. 2014;5(1):5471.
- 332 29. Cummins EA, Hall RJ, Connor C, McInerney JO, McNally A. Distinct evolutionary
333 trajectories in the *Escherichia coli* pangenome occur within sequence types. *Microbial Genomics*.
334 2022;8(11).
- 335 30. Machado D, Andrejev S, Tramontano M, Patil KR. Fast automated reconstruction of
336 genome-scale metabolic models for microbial species and communities. *Nucleic Acids Res*.
337 2018;46(15):7542-53.
- 338 31. Seaver SMD, Liu F, Zhang Q, Jeffryes J, Faria JP, Edirisinghe JN, et al. The ModelSEED
339 Biochemistry Database for the integration of metabolic annotations and the reconstruction,
340 comparison and analysis of metabolic models for plants, fungi and microbes. *Nucleic Acids Res*.
341 2021;49(D1):D575-D88.
- 342 32. Arkin AP, Cottingham RW, Henry CS, Harris NL, Stevens RL, Maslov S, et al. KBase: The
343 United States Department of Energy Systems Biology Knowledgebase. *Nature Biotechnology*.
344 2018;36(7):566-9.
- 345 33. Mendoza SN, Olivier BG, Molenaar D, Teusink B. A systematic assessment of current
346 genome-scale metabolic reconstruction tools. *Genome Biology*. 2019;20(1):158.
- 347 34. Tamasco G, Kumar M, Zengler K, Silva-Rocha R, da Silva RR. ChiMera: an easy to use
348 pipeline for bacterial genome based metabolic network reconstruction, evaluation and
349 visualization. *BMC Bioinformatics*. 2022;23(1):512.
- 350 35. Norsigian CJ, Fang X, Seif Y, Monk JM, Palsson BO. A workflow for generating multi-strain
351 genome-scale metabolic models of prokaryotes. *Nature Protocols*. 2020;15(1):1-14.
- 352 36. Ebrahim A, Lerman JA, Palsson BO, Hyduke DR. COBRApy: COntstraints-Based
353 Reconstruction and Analysis for Python. *BMC Systems Biology*. 2013;7(1):74.
- 354 37. Schellenberger J, Park JO, Conrad TM, Palsson BØ. BiGG: a Biochemical Genetic and
355 Genomic knowledgebase of large scale metabolic reconstructions. *BMC Bioinformatics*.
356 2010;11(1):213.
- 357 38. Hyatt D, Chen G-L, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal: prokaryotic
358 gene recognition and translation initiation site identification. *BMC Bioinformatics*. 2010;11(1):119.
- 359 39. Keating SM, Waltemath D, König M, Zhang F, Dräger A, Chaouiya C, et al. SBML Level 3:
360 an extensible format for the exchange and reuse of biological models. *Molecular Systems Biology*.
361 2020;16(8):e9110.
- 362 40. Lieven C, Beber ME, Olivier BG, Bergmann FT, Ataman M, Babaei P, et al. MEMOTE for
363 standardized genome-scale metabolic model testing. *Nature Biotechnology*. 2020;38(3):272-6.
- 364 41. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+:
365 architecture and applications. *BMC Bioinformatics*. 2009;10(1):421.
- 366 42. Hernández-Salmerón JE, Moreno-Hagelsieb G. Progress in quickly finding orthologs as
367 reciprocal best hits: comparing blast, last, diamond and MMseqs2. *BMC Genomics*.
368 2020;21(1):741.
- 369 43. BD. Difco™ & BBL™ Manual, 2nd Edition. In: BD, editor. 2009.
- 370 44. Biosciences BD. BD Bionutrients Technical Manual: BD Biosciences – Advanced
371 Bioprocessing. BD Biosciences; 2015.
- 372 45. Loginova LI, Manuilova VP, Tolstikov VP. Content of free amino acids in peptone and the
373 dynamics of their consumption in the microbiological synthesis of dextran. *Pharmaceutical
374 Chemistry Journal*. 1974;8(4):249-51.
- 375 46. ThermoFisherScientific. Technical guide to peptones, supplements, and feeds: Enhancing
376 performance of mammalian and microbial bioprocesses. ThermoFisherScientific; 2019.

- 377 47. Clausen E, Gildberg A, Raa J. Preparation and testing of an autolysate of fish viscera as
378 growth substrate for bacteria. *Applied and Environmental Microbiology*. 1985;50(6):1556-7.
- 379 48. Hagely KB, Palmquist D, Bilyeu KD. Classification of Distinct Seed Carbohydrate Profiles in
380 Soybean. *Journal of Agricultural and Food Chemistry*. 2013;61(5):1105-11.
- 381 49. Choct M, Dersjant-Li Y, McLeish J, Peisker M. Soy Oligosaccharides and Soluble Non-
382 starch Polysaccharides: A Review of Digestion, Nutritive and Anti-nutritive Effects in Pigs and
383 Poultry. *Asian-Australasian Journal of Animal Sciences*. 2010;23.
- 384 50. Tomé D. Yeast Extracts: Nutritional and Flavoring Food Ingredients. *ACS Food Science &*
385 *Technology*. 2021;1(4):487-94.
- 386 51. Plata MR, Koch C, Wechselberger P, Herwig C, Lendl B. Determination of carbohydrates
387 present in *Saccharomyces cerevisiae* using mid-infrared spectroscopy and partial least squares
388 regression. *Anal Bioanal Chem*. 2013;405(25):8241-50.
- 389 52. Liu Y, Huang G, Lv M. Extraction, characterization and antioxidant activities of mannan
390 from yeast cell wall. *International Journal of Biological Macromolecules*. 2018;118:952-6.
- 391 53. Blagović B. Lipid Composition of Brewer's Yeast. *Food Technology and Biotechnology*.
392 2001;39:175-81.
- 393 54. Blagović B, Mesarić M, Marić V, Rupčić J. Characterization of lipid components in the
394 whole cells and plasma membranes of baker's Yeast. *Croatica Chemica Acta*. 2005;78:479-84.
- 395 55. Avramia I, Amariei S. Spent Brewer's Yeast as a Source of Insoluble β -Glucans.
396 *International Journal of Molecular Sciences*. 2021;22(2).
- 397 56. Short Francesca L, Di Sario G, Reichmann Nathalie T, Kleanthous C, Parkhill J, Taylor
398 Peter W, et al. Genomic Profiling Reveals Distinct Routes To Complement Resistance in *Klebsiella*
399 *pneumoniae*. *Infection and Immunity*. 2020;88(8):e00043-20.
- 700 57. Powers DMW. Evaluation: from precision, recall and F-measure to ROC, informedness,
701 markedness and correlation. *arXiv*. 2020;arXiv:2010.16061.
- 702 58. Green ML, Karp PD. Genome annotation errors in pathway databases due to semantic
703 ambiguity in partial EC numbers. *Nucleic Acids Res*. 2005;33(13):4035-9.
- 704 59. Rembeza E, Engqvist MKM. Experimental and computational investigation of enzyme
705 functional annotations uncovers misannotation in the EC 1.1.3.15 enzyme class. *PLOS*
706 *Computational Biology*. 2021;17(9):e1009446.
- 707 60. Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, et al. The FAIR
708 Guiding Principles for scientific data management and stewardship. *Scientific Data*.
709 2016;3(1):160018.
- 710 61. Van Rossum G, & Drake, F. L. Python 3 Reference Manual. Scotts Valley, CA2009.
- 711 62. BIOLOG. Phenotype MicroArrays. In: BIOLOG, editor. 2020.
- 712 63. Norsigian CJ, Attia H, Szubin R, Yassin AS, Palsson BØ, Aziz RK, et al. Comparative
713 Genome-Scale Metabolic Modeling of Metallo-Beta-Lactamase-Producing Multidrug-Resistant
714 *Klebsiella pneumoniae* Clinical Isolates. *Frontiers in Cellular and Infection Microbiology*.
715 2019;9(161).
- 716 64. Hall MB. Rasusa: Randomly subsample sequencing reads to a specified coverage 2019
717 [Available from: <https://github.com/mbhall88/rasusa>].
- 718 65. Krueger F. Trim Galore GitHub2012 [Available from:
719 <https://github.com/FelixKrueger/TrimGalore>].
- 720 66. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies
721 from short and long sequencing reads. *PLOS Computational Biology*. 2017;13(6):e1005595.
- 722 67. Wick RR. Dead-end count for QC of short-read assemblies. 2023.
- 723 68. RStudio-Team. RStudio: Integrated Development for R. Boston, MA RStudio; 2020.
- 724 69. Wickham H, Averick M, Bryan J, Chang W, McGowan LDA, François R, et al. Welcome to
725 the Tidyverse. *Journal of Open Source Software*. 2019;4(43):1686.
- 726 70. Garnier S. viridis: Default Color Maps from 'matplotlib'. 2018.
- 727 71. Neuwirth E. RColorBrewer: ColorBrewer Palettes. 1.1-2 ed2014.
- 728 72. Kassambara A. ggpubr: 'ggplot2' Based Publication Ready Plots. 0.5.0 ed2022.
- 729 73. Aphalo PJ, Slowikowski K, Mouksassi S. ggpmisc: Miscellaneous Extensions to 'ggplot2'.
730 0.4.1 ed2021.
- 731 74. Yu G. aplot: Decorate a 'ggplot' with Associated Information. 0.0.6 ed2020.

- 732 75. Zeileis A, Fisher JC, Hornik K, Ihaka R, McWhite CD, Murrell P, et al. colorspace: A
733 Toolbox for Manipulating and Assessing Colors and Palettes. *Journal of Statistical Software*.
734 2020;96(1):1 - 49.
- 735 76. F.C. M, Davis TL. Package 'ggpattern'. 2022.
- 736 77. Wilke CO. ggtext: Improved Text Rendering Support for 'ggplot2'. 0.1.1 ed2020.
- 737 78. Hester J, Bryan J, RStudio. glue: Interpreted String Literals. 1.6.2 ed2022.
- 738 79. Muggeo VMR. segmented: Regression Models with Break-Points / Change-Points (with
739 Possibly Random Effects) Estimation. 1.6-2 ed2022.
- 740 80. The-Inkscape-Team. Inkscape. 1.0.1 (3bc2e813f5, 2020-09-07) ed2020.
- 741 81. zanatlija. Proportional TFB. 2012.
- 742 82. weknow. Element. 2015.

743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769

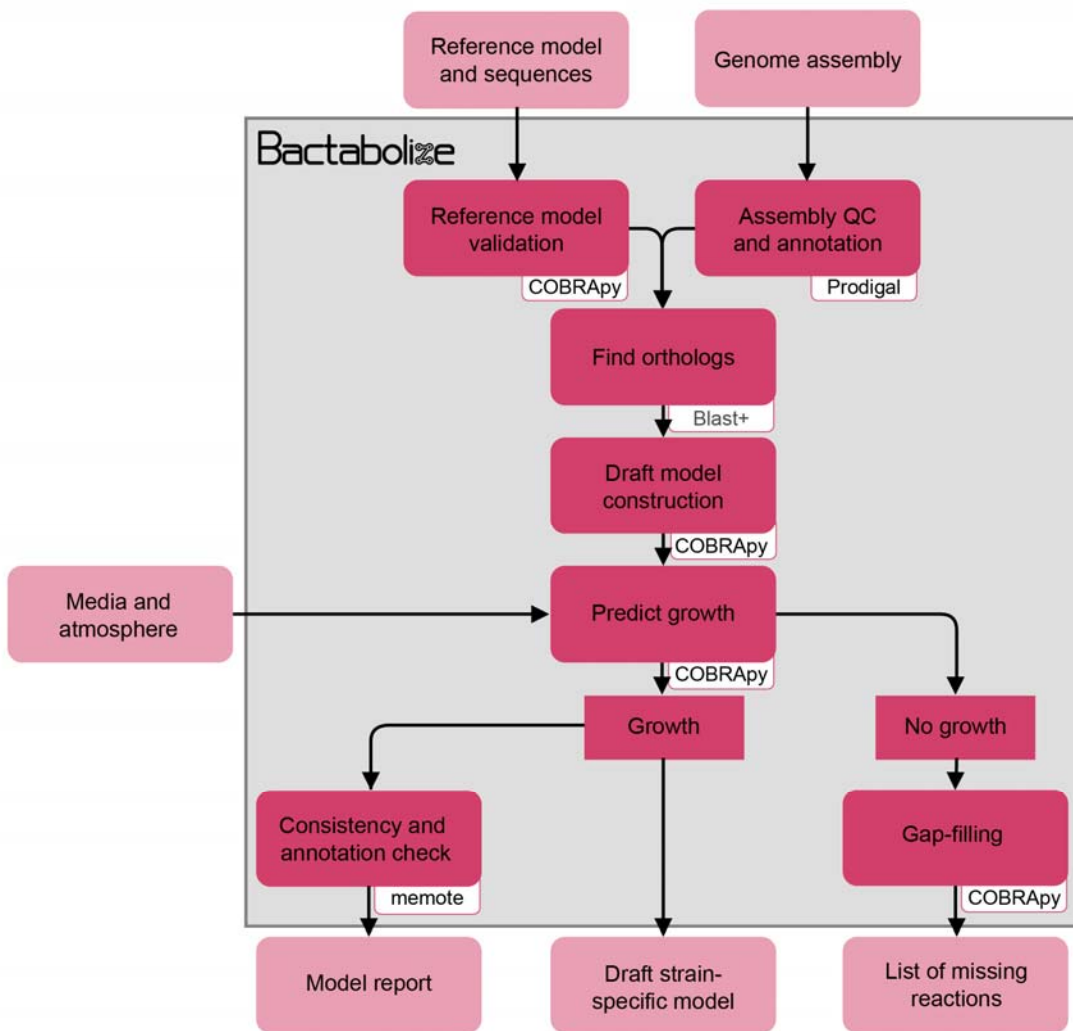
770 Supplementary figures

771

772 Figure S1

773 Flow diagram showing the overview of the draft_model module from Bactabolize, which produces
774 draft metabolic models. Input and output files are shown in light pink while Bactabolize processes
775 are shown in dark pink. Third-party dependencies are indicated within the white boxes.

776



777

778

779

780

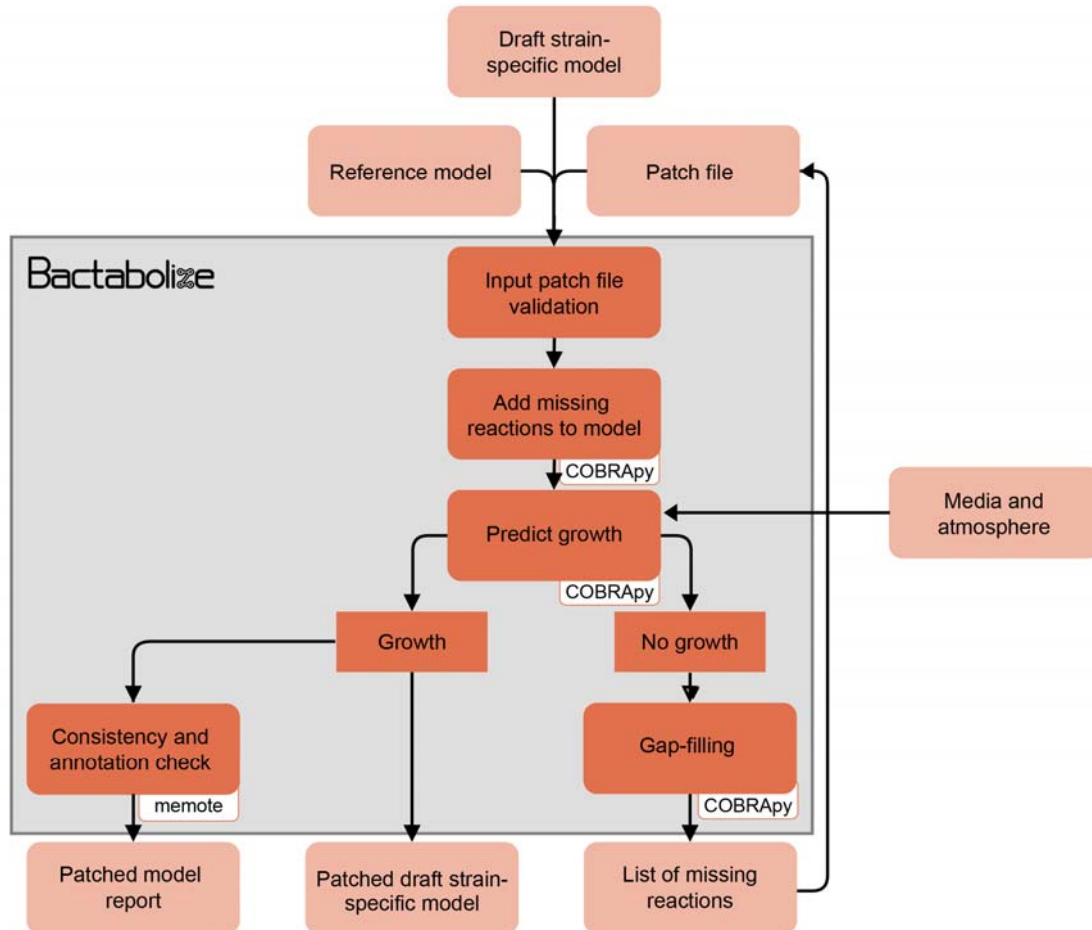
781

782

783 **Figure S2**

784 Flow diagram showing the overview of the patch_model module from Bactabolize, which patches
785 metabolic models that do not simulate growth. Input and output files are shown in light orange
786 while Bactabolize processes are shown in dark orange. Third-party dependencies are indicated
787 within the white boxes.

788



789

790

791

792

793

794

795

796

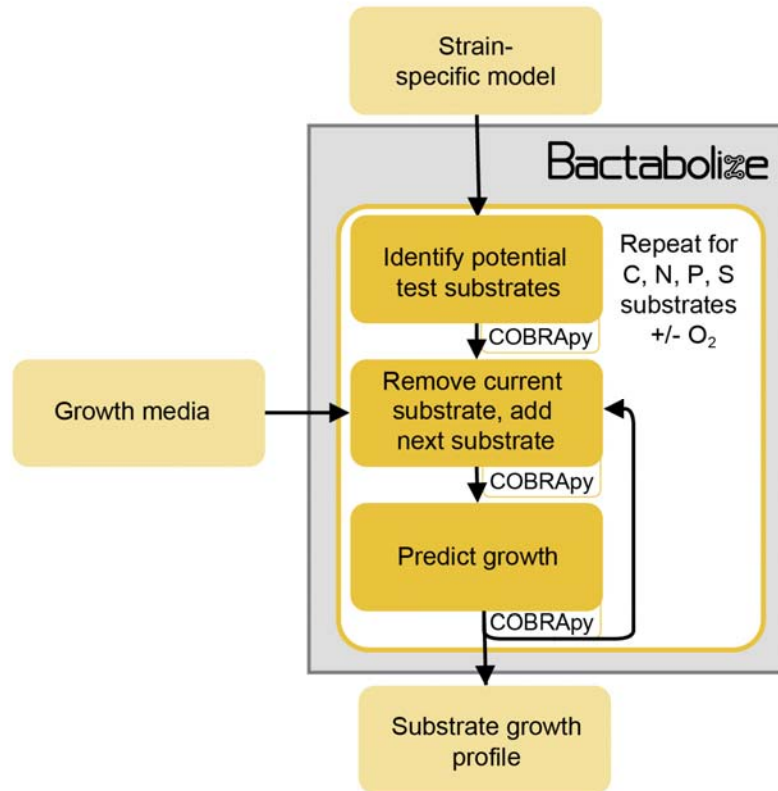
797

798

799 **Figure S3**

300 Flow diagram showing the overview of the fba module from Bactabolize, which performs growth
301 simulations using Flux Balance Analysis. Input and output files are shown in light yellow while
302 Bactabolize processes are shown in dark yellow. Third-party dependencies are indicated within
303 white boxes. C, carbon; N, nitrogen, P, phosphorus; S, sulphur; O₂, oxygen.

304



305

306

307

308

309

310

311

312

313

314

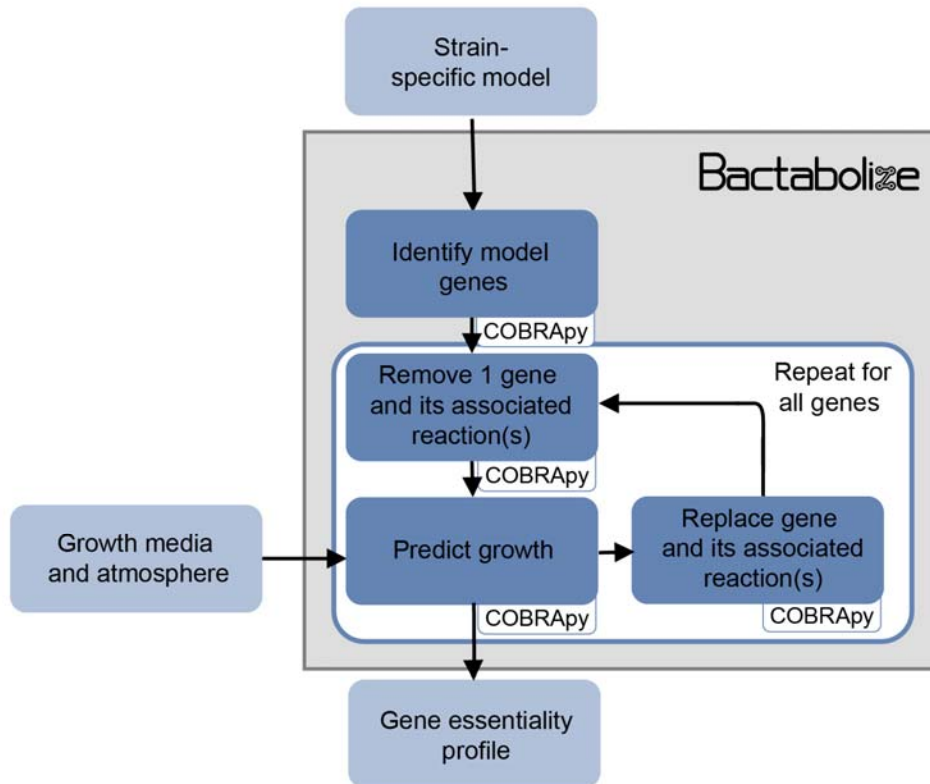
315

316

317 **Figure S4**

318 Flow diagram showing the overview of the *sgk* module from Bactabolize, which performs Single
319 Gene Knockout analysis. Input and output files are shown in light blue while Bactabolize processes
320 are shown in dark blue. Third-party dependencies are indicated within white boxes.

321



322

323

324

325

326

327

328

329

330

331

332

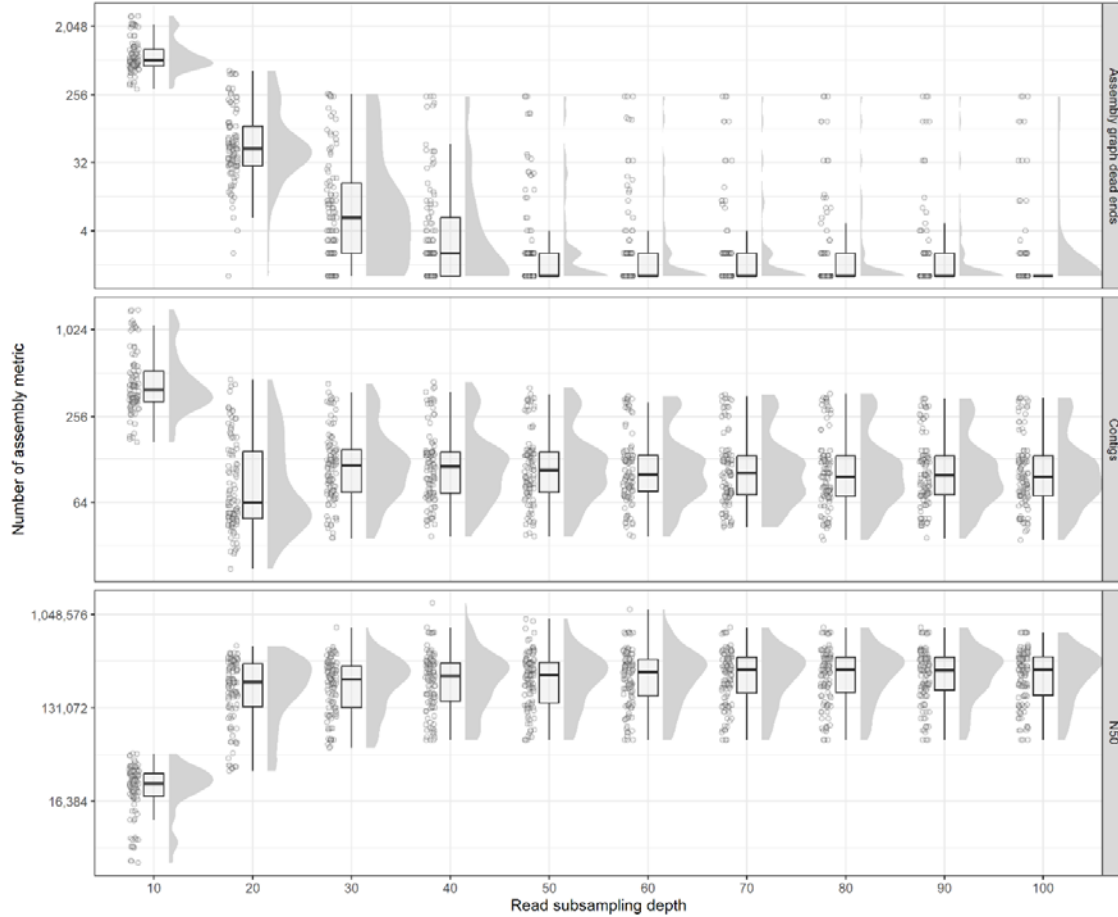
333

334

335 **Figure S5**

336 Raincloud plot showing distributions of assembly metrics across various read subsampling depths
337 (10x increments).

338



339

340

341

342

343

344

345

346

347

348

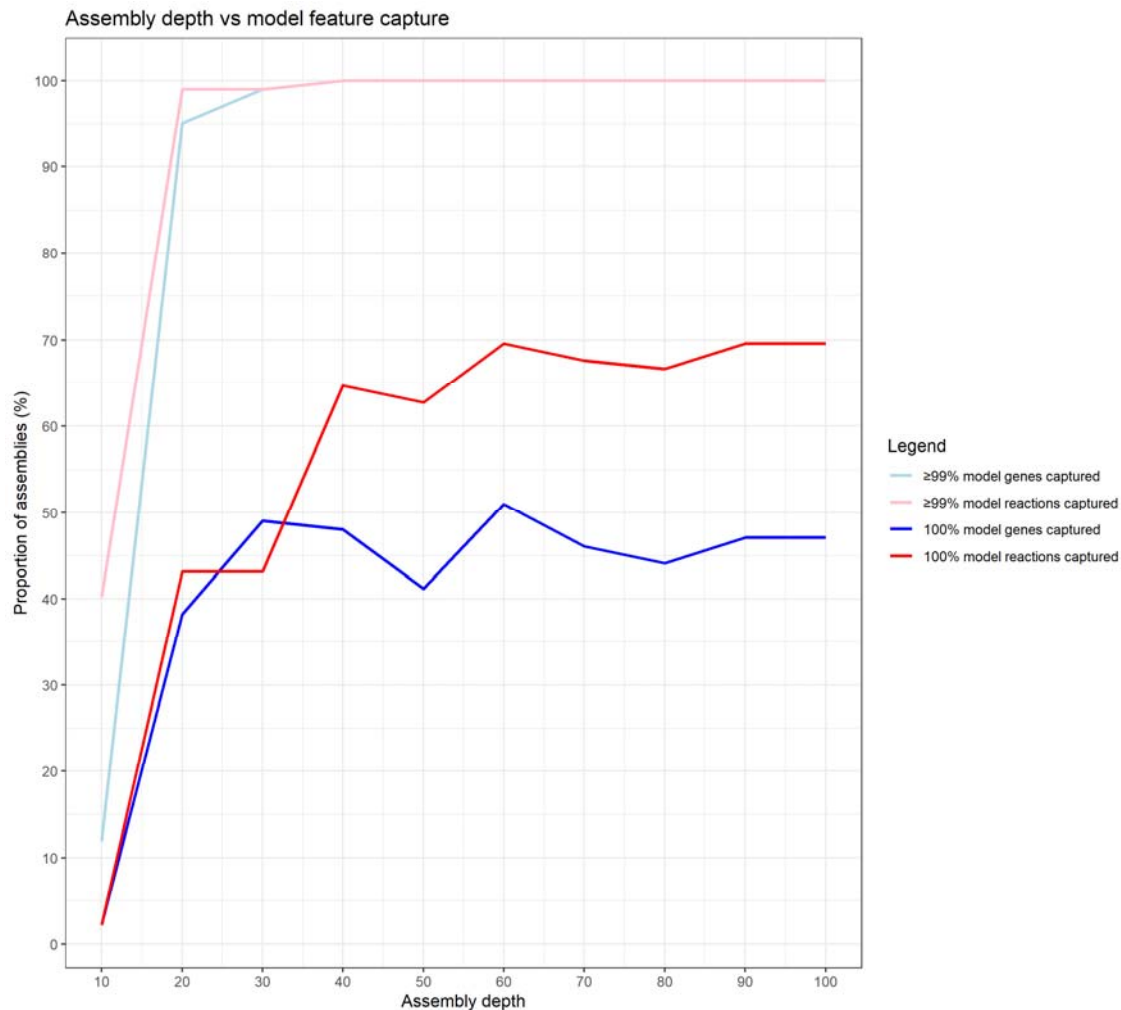
349

350

351 **Figure S6**

352 Line graph showing the capture of model features of draft assemblies (short read only) at various
353 depths, compared to the corresponding completed genome (long-read + short read assemblies).

354



355

356

357

358

359

360

361

362

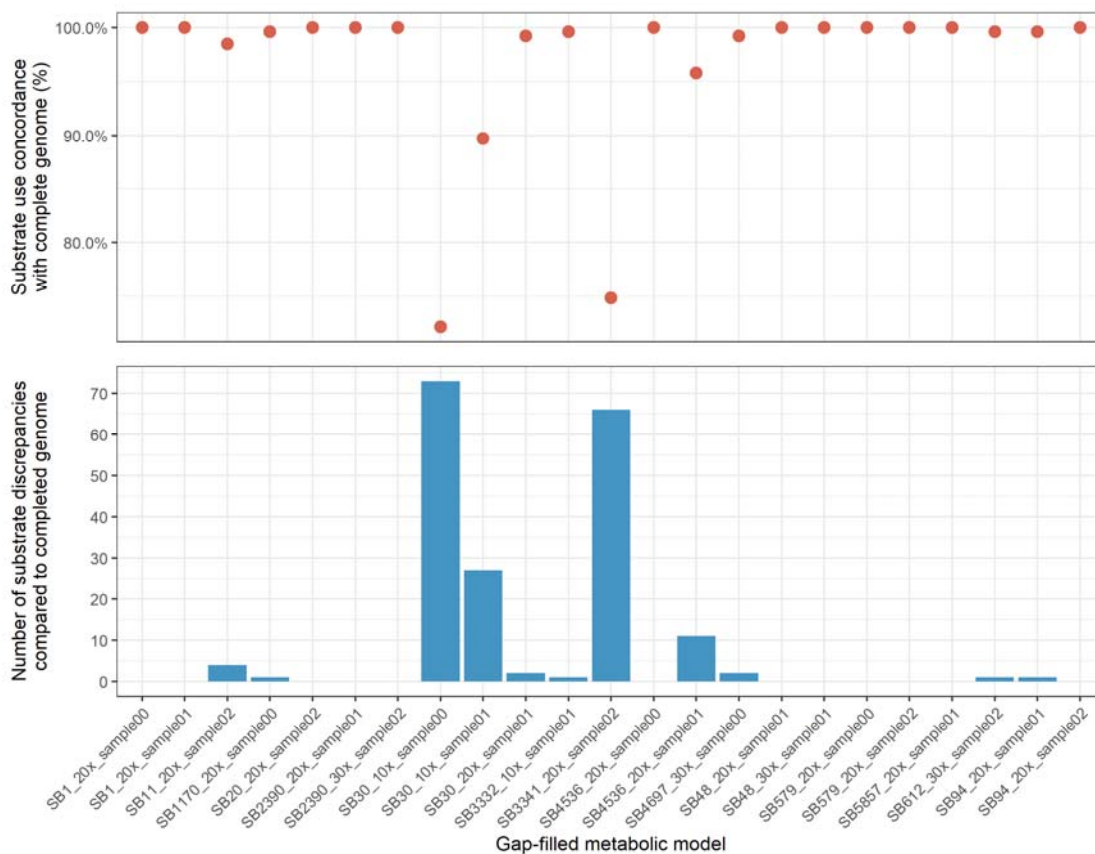
363

364

365 **Figure S7**

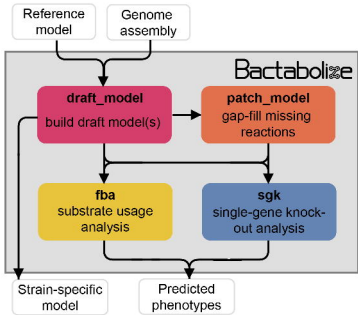
366 Faceted graphs showing the number of substrate usage (fba module) discrepancies of gap-filled
367 models (patch_model module) which initially did not produce biomass (models which failed to
368 simulate growth). The dots indicate percentage concordance with the completed genome model,
369 while the columns indicate number of substrates with discrepancies (no simulated growth in
370 patched model, but growth in completed genome model).

371

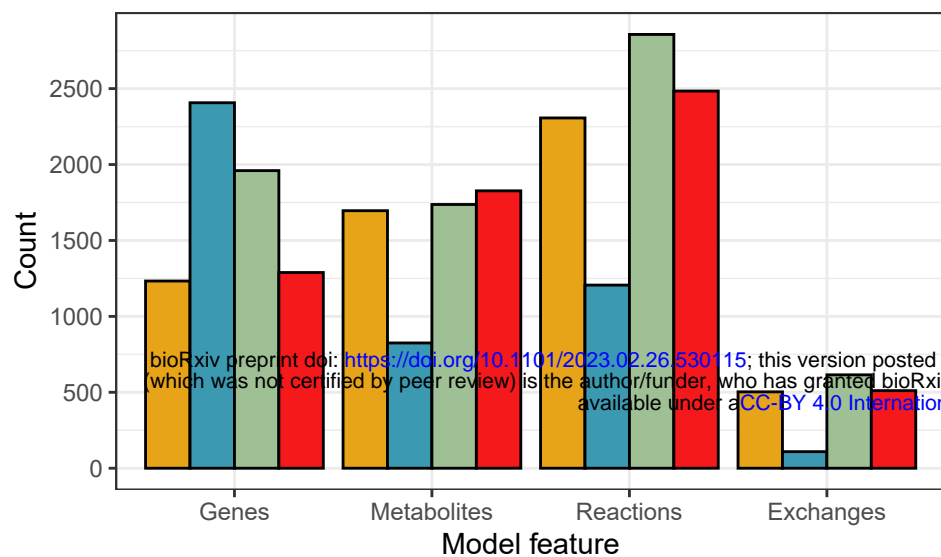


372

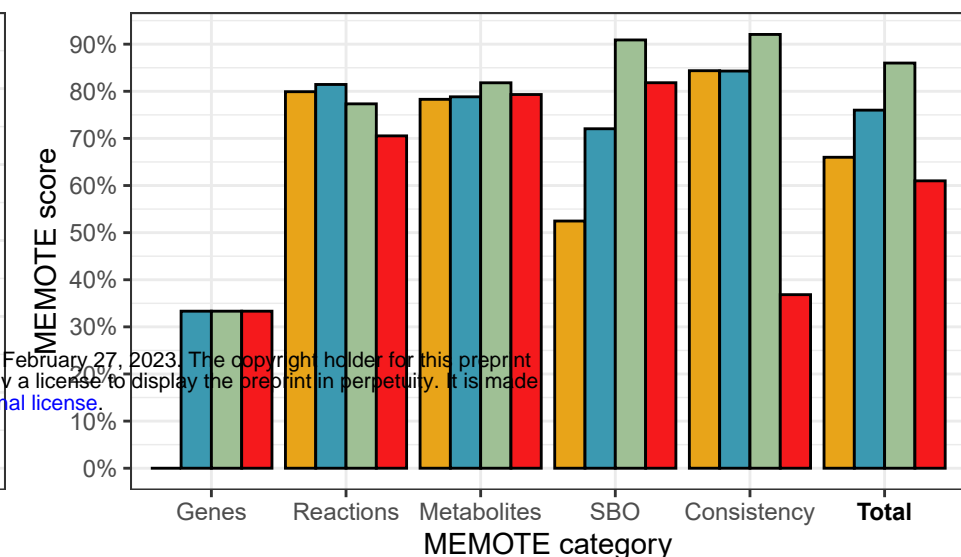
373



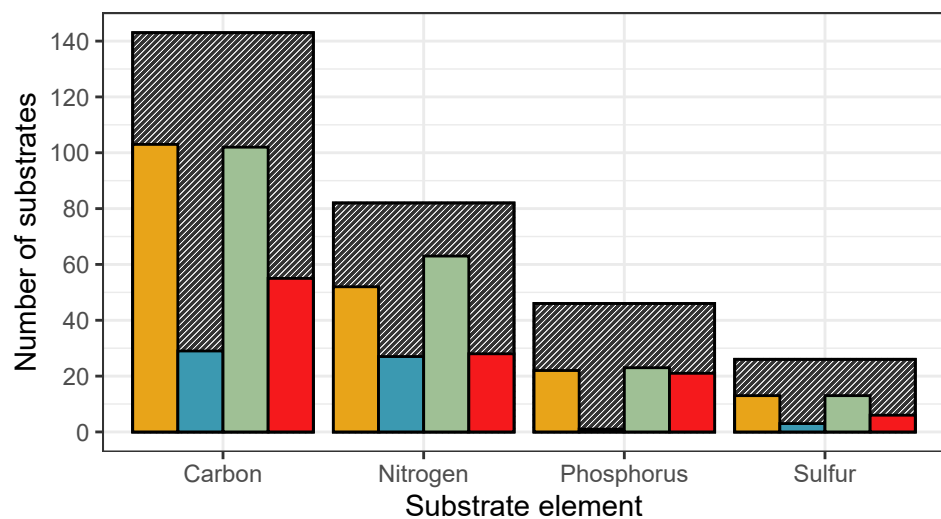
A) Model features



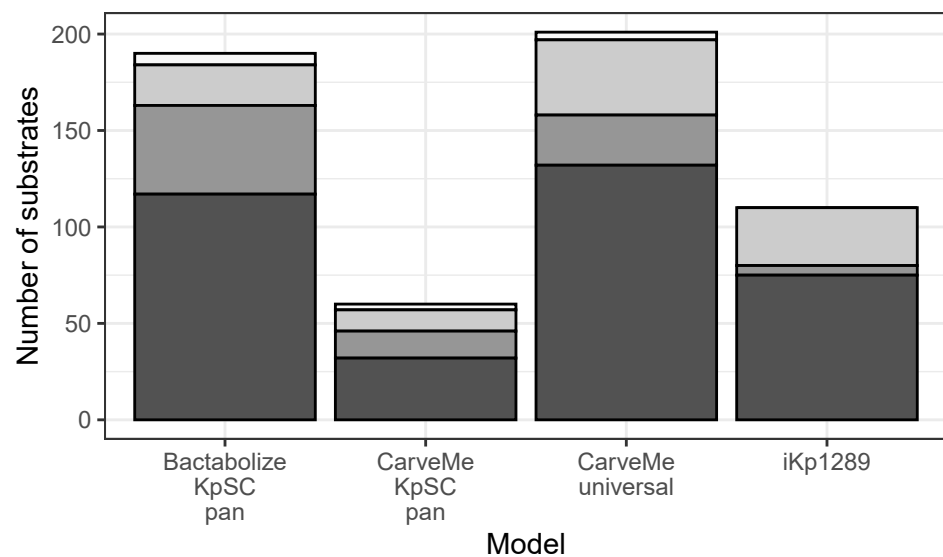
B) MEMOTE scores



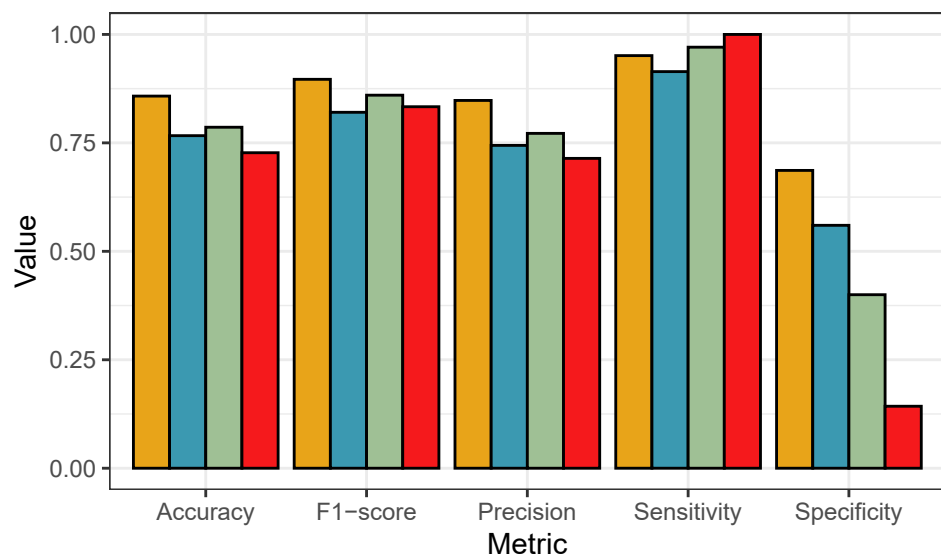
C) Substrates for which growth can be simulated among those with matched phenotypes



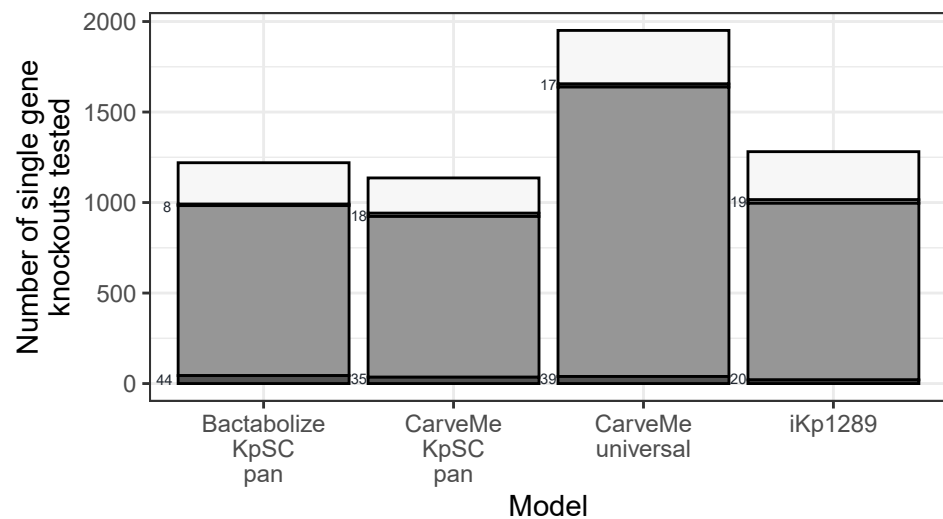
D) Comparison of predicted to true growth phenotypes



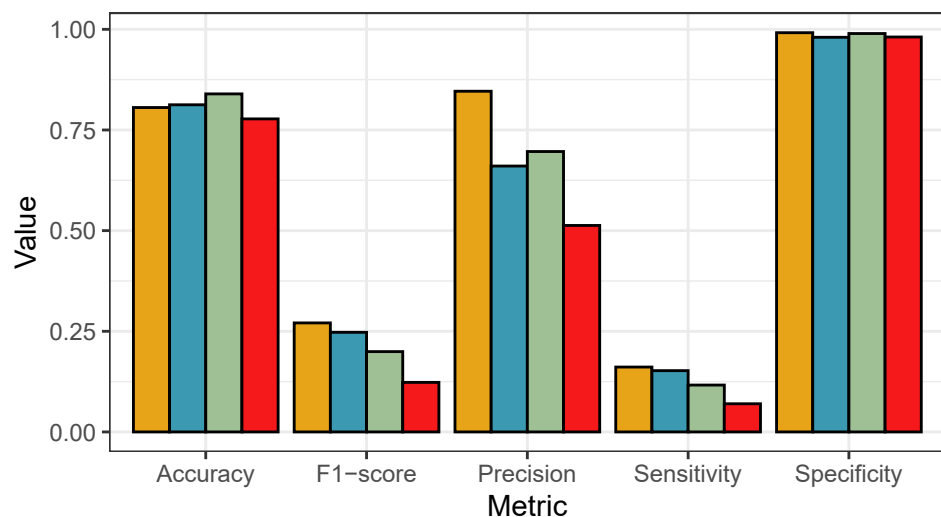
E) Accuracy metrics for predicted growth phenotypes



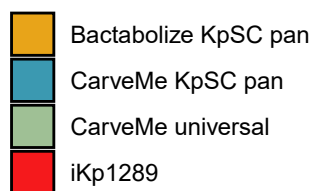
F) Comparison of predicted to true single-gene knockouts



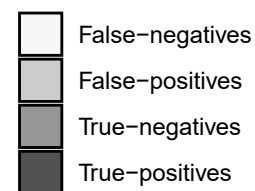
G) Accuracy metrics for predicted single-gene knockout phenotypes



Metabolic model

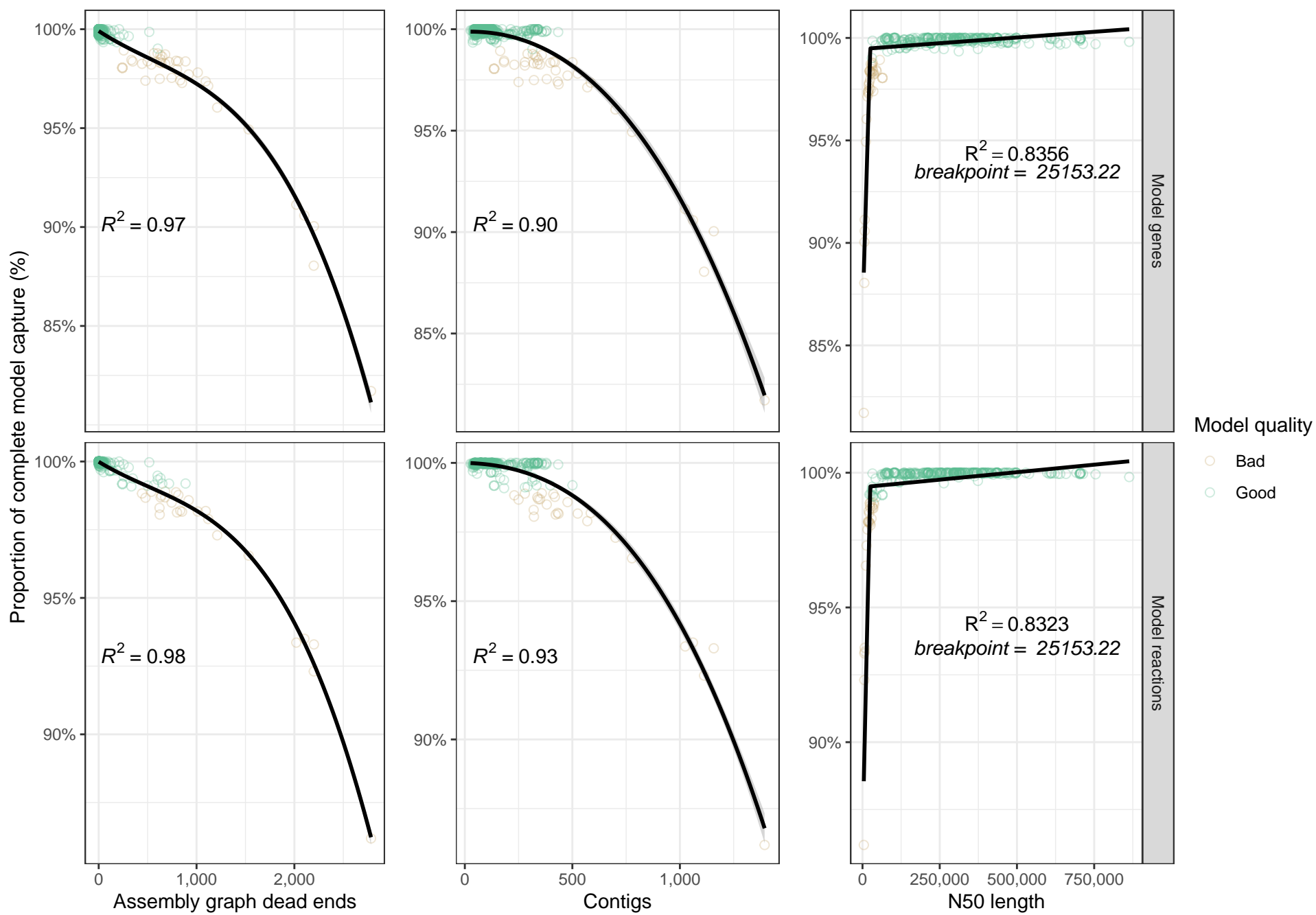


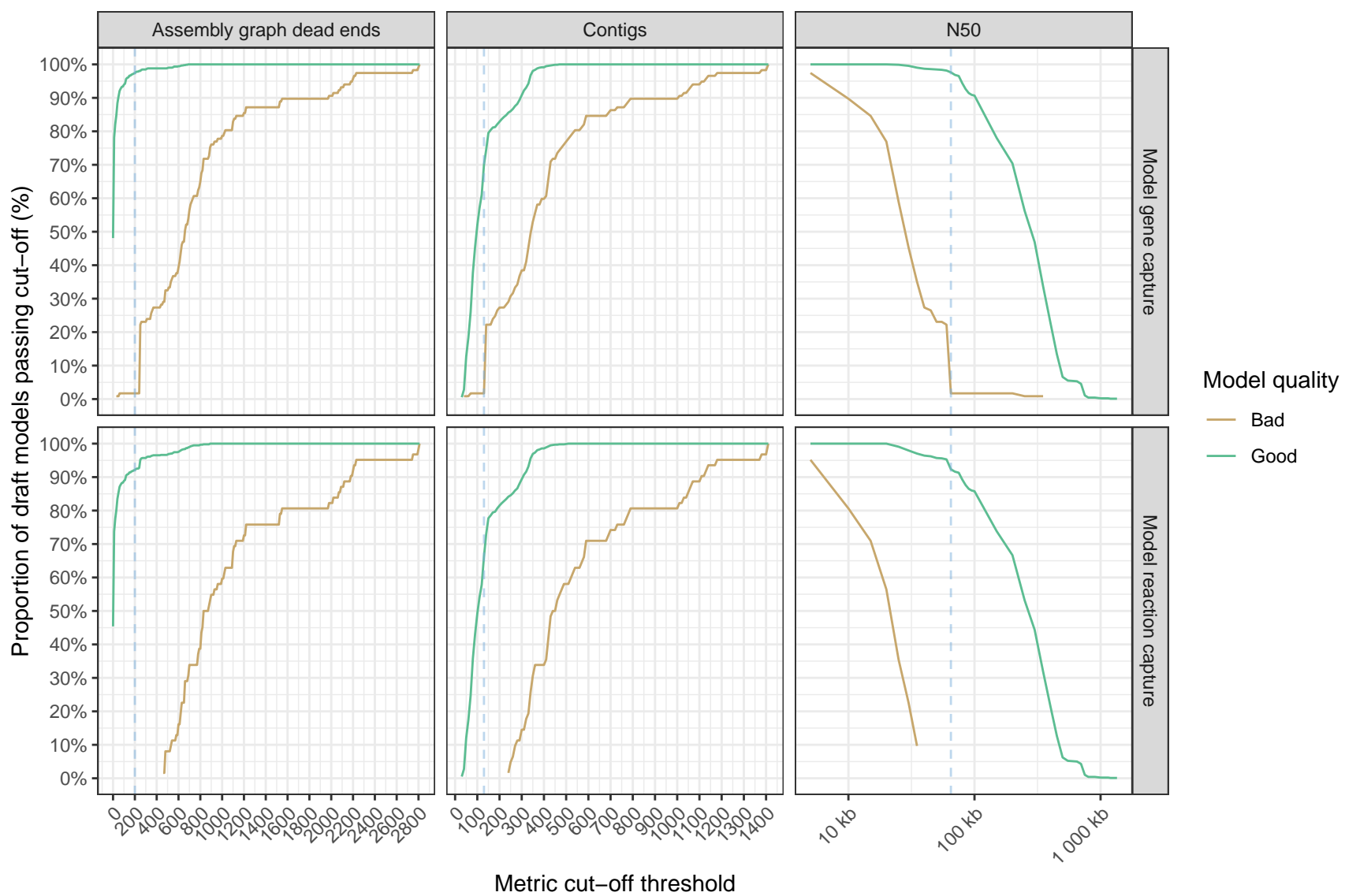
Statistic



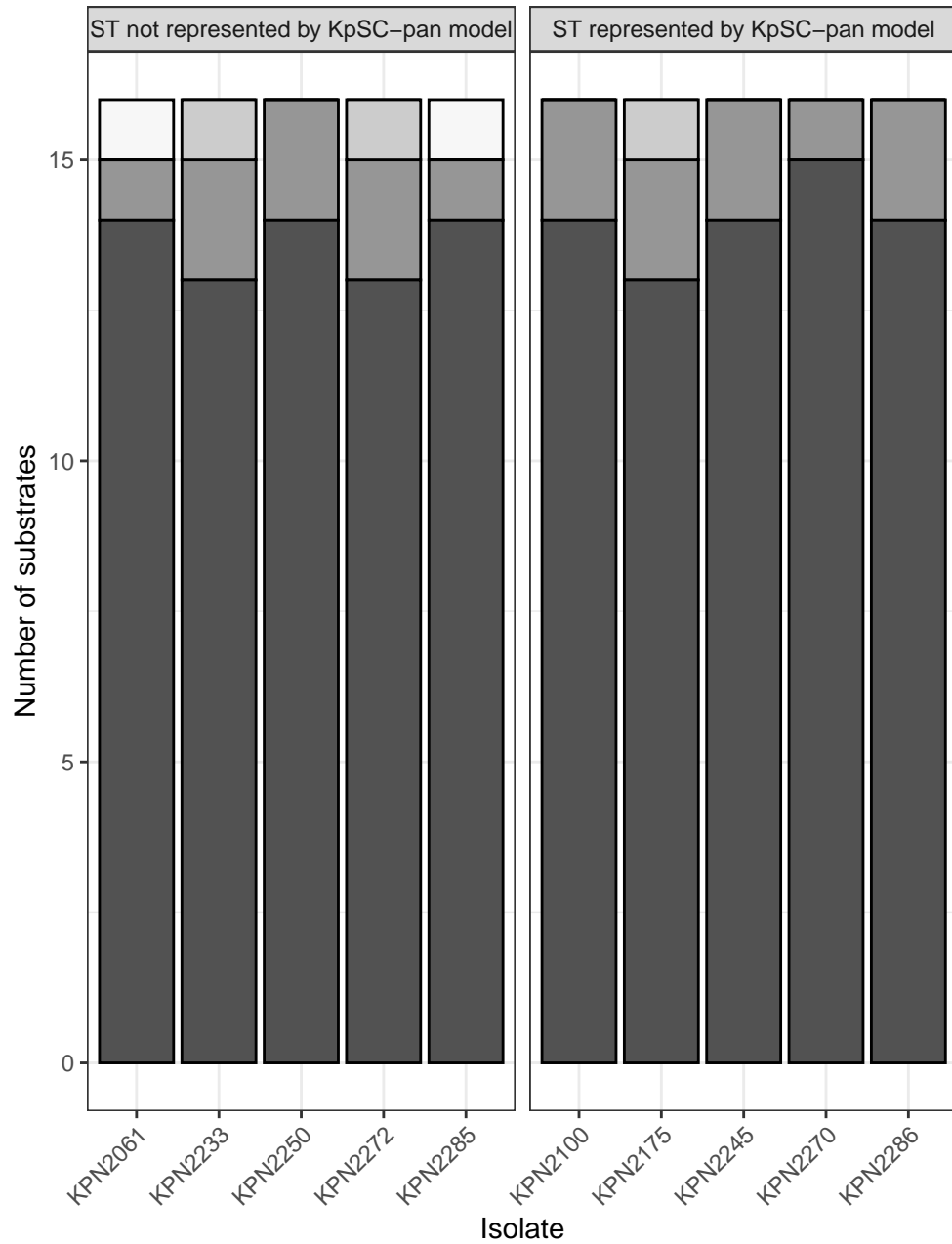
Phenotype data available







A) Comparison of predicted to true growth phenotypes



B) Accuracy metrics for predicted growth phenotypes

