

In-Depth Resistome Analysis by Targeted Metagenomics

**Val F. Lanza^{1,2,3,4}, Fernando Baquero^{1,2,3}, José Luís Martínez^{2,4}, Ricardo Ramos-Ruíz⁵,
Bruno González-Zorn⁶, Antoine Andremont⁷, Antonio Sánchez-Valenzuela¹,
Dusko Ehrlich^{8,9}, Sean Kennedy^{8*}, Etienne Ruppe^{8**}, Willem van Schaik¹⁰,
Rob J. Willems¹⁰, Fernando de la Cruz¹¹, and Teresa M. Coque^{1,2,3}**

¹ Servicio de Microbiología, Hospital Universitario Ramón y Cajal University Hospital, Ramón y Cajal Health Research Institute (IRYCIS), Madrid, Spain,

² Joint Unit of Antibiotic Resistance and Bacterial Virulence associated with the Spanish National Research Council (CSIC),

³ Network Research Centre in Epidemiology and Public Health (CIBER-ESP), Madrid, Spain,

⁴ Centro Nacional de Biotecnología. CSIC. Darwin 3. 28049-Madrid. Spain

⁵ Genomics Unit, Science Park, Madrid.

⁶ Facultad de Veterinaria Universidad Complutense de Madrid, Spain.

⁷ IAME, UMR 1137, INSERM, Paris Diderot University, Sorbonne Paris Cité, Bactériologie Laboratory, Hôpital Bichat, AP-HP, Paris, France.

⁸ MGP MetaGénoPolis, INRA, Université Paris-Saclay, 78350 Jouy-en-Josas, France

⁹ Centre of Host Microbiome Interactions, King's College, London, United Kingdom

¹⁰ Department of Medical Microbiology, University Medical Center Utrecht, The Netherlands.

¹¹ Departamento de Biología Molecular, ⁴Universidad de Cantabria and Instituto de Biomedicina y Biotecnología de Cantabria IBBTEC (UC- CSIC), Santander, Spain

*Current affiliation: Hub Bioinformatique et Biostatistique, C3BI & Biomics Pole, CITECH Institut Pasteur, Centre François Jacob, 28 rue du Docteur Roux, 75015 Paris

** Current affiliation: Genomic Research Laboratory, Service of Infectious Diseases, Geneva University Hospitals, rue Gabrielle-Perret-Gentil 4, 1205 Geneva, Switzerland.

32 **ABSTRACT**

33 We developed ResCap, a targeted sequence capture platform based on SeqCapEZ technology, to
34 analyse resistomes and other genes related to antimicrobial resistance (heavy metals,
35 biocides and plasmids). ResCap includes probes for 8,667 canonical resistance genes (7,963
36 antibiotic resistance genes and 704 genes conferring resistance to metals or biocides),
37 plus 2,517 relaxase genes (plasmid markers). Besides, it includes 78,600 genes homologous to
38 the previous ones (47,806 for antibiotics and 30,794 for biocide or metals). ResCap
39 enriched 279-fold the targeted sequences detected by metagenomic shotgun sequencing
40 and improves their identification. Novel bioinformatic approaches allow quantifying “gene
41 abundance” and “gene diversity”. ResCap, the first targeted sequence capture specifically
42 developed to analyse resistomes, enhances the sensitivity and specificity of available
43 metagenomic methods to analyse antibiotic resistance in complex populations, enables the
44 analysis of other genes related to antimicrobial resistance and opens the possibility to
45 accurately study other complex microbial systems.

46 INTRODUCTION

47

48 Antimicrobial resistance is considered a major Global Health challenge recently included in the
49 political agendas of international bodies such as G8 and IMF¹. The adoption of measures to face
50 the “antibiotic resistance crisis”² is impaired by the controversy about “what” is resistance and
51 “how” and “where” should be detected and analysed³⁻⁵

52 Metagenomic methods are increasingly used to analyse the ensemble of genes encoding
53 antibiotic resistance in different microbial ecosystems which has recently been defined as the
54 “resistome”⁶⁻¹⁶. An important hurdle of the available resistome analyses is the low
55 discrimination in the detection of minority populations harbouring resistance genes (often
56 present at concentrations below the detection level of the methods used)¹⁷ and/or the
57 identification of allelic variants that might confer different resistance phenotypes.

58 A sensitive and specific identification of antibiotic resistance genes in a metagenome
59 background is required for assessing the associated risks in terms of Public Health^{18a,18b}. Such
60 methodological challenge parallels the difficulties of analysing sets of orthologous genes of
61 many individuals for the diagnosis of human multifactorial inherited diseases¹⁹. In this case,
62 the use of “capture-based” or “targeted” sequencing strategies, was a cost-effective and high-
63 throughput alternative that overcame the limitations of metagenomic shotgun sequencing
64 (MSS)^{20,21}. In-solution targeted capture platforms (TCP) take advantage of Next Generation
65 Sequencing to provide technical improvements over array-based platforms or other genome-
66 partitioning approaches in terms of scalability, cost-effectiveness, and enhanced data quality
67 (lower variance in target coverage, more accurate SNP calling, higher reproducibility and
68 longer assembled contigs)²². Currently, TCPs offer a tremendous potential for boosting
69 advances in environmental and ecological studies, particularly involving micro-biodiversity
70 research, which requires the isolation of sequences of interest from a mixture of DNAs of a
71 complex multiplicity of organisms²³.

72 Our work reports the development and validation of the first TCP for the analysis of bacterial
73 resistomes, which we designate as ResCap (for Resistome Capture). We show that ResCap
74 results in a significant improvement in sensitivity and specificity over previous metagenomic
75 analysis of antimicrobial resistance. ResCap also allows the analysis of the presence and
76 diversity of genes conferring resistance to other antimicrobials (heavy metals and biocides),
77 which are frequently co-selected with antibiotic resistance genes and also genes from replicons
78 of the mobilome (as plasmids). An *ad-hoc* advanced bioinformatics pipeline, developed in
79 parallel, exploits the capabilities of Rescap comparative metagenomic analysis. The
80 metagenomic approach described here opens the way for a series of applications in the

81 identification, epidemiological surveillance, ecology, and study of evolutionary trajectories of
82 resistance genes.
83

84 RESULTS

85

86 Targeted metagenomics, a tool for high-resolution analysis of resistome

87 ResCap was designed to establish a standardized framework that would allow performing both
88 quantitative and qualitative analysis of resistomes. Also, to facilitate the analysis of novel genes
89 potentially involved in the resistance to antibiotics, metals, biocide or any combination of
90 them. As a proof of concept, we compared the performance of ResCap with metagenomic
91 shotgun sequencing (MSS) by analysing the resistome in 17 fecal samples, 9 from humans and
92 8 from swine.

93 ResCap exhibits a target capacity (total amount of targeted sequences) of 88.13 Mb, and
94 includes probes for 78,600 non-redundant genes (81,117 redundant genes), including 7,963
95 functionally validated antibiotic resistance genes, 704 functionally validated metal & biocide
96 resistance genes, and 2,517 relaxase genes (genes used for plasmid identification and
97 classification)²⁴. Besides the 8,667 genes that confer functionally proved resistance to
98 antimicrobials (canonical genes), the platform also includes targets for 78,600 homologous of
99 resistance genes (47,806 for antibiotics and 30,794 for biocide and metals resistance). The
100 criteria used to select the targeted genes are explained in the section Material and Methods.

101 ResCap performance was compared with MSS in two ways. First, by applying a reference-based
102 approach that maps metagenome reads against specific databases (AbR, Metal & Biocides and
103 Relaxases). Second, by applying a reference-free approach that assembles metagenomic reads
104 and performs a functional annotation. The results of both evaluations are detailed below.

105

106 Reference-based evaluation

107 This section addresses how the abundance and diversity of resistance genes (ResCap or those
108 already validated) were calculated.

109

110 *ResCap achieves better recovery of target genes than MSS*

111 An average of 1.9×10^7 paired-reads was obtained from the MSS and ResCap datasets ($0.92-$
112 3.2×10^7). The on-target average (the number of reads mapping on the target genes relative to
113 the total read number) against the selected databases (see Material and Methods) was 0.11%
114 (0.07-0.18) for MSS data and 30.26% (20.27–41.83%) for ResCap data, which represents an
115 enrichment of 279-fold (**Table 1**).

116 The analysis of the gene abundance, expressed in RPKMs (reads per kb per million reads, see
117 Materials and Methods), demonstrates better recovery of genes coding for resistance to
118 antibiotics, heavy metals, biocides and relaxases (plasmid genes) when using ResCap than
119 when using MSS. **Figure 1** represents the RPKMs inferred before (MSS) and after capture

120 (ResCap) for all the samples analysed, while **Figure S1** shows the gain plots for each sample.
121 Most canonical genes (99.3%, 1339/1.348) detected by MSS were also detected with ResCap.
122 Furthermore, almost half of the genes detected by ResCap (42% ,975/2323), were not detected
123 by MSS. The linearity of the system was evaluated by using a linear regression model for the
124 genes detected in each paired-sample (MSS vs ResCap). An R^2 mean of 0.813 (0.85-0.99) shows
125 a good match between both protocols.

126 The enrichment of canonical resistance genes when using ResCap was similar in samples from
127 humans and swine. Nonetheless, the differences in the relative abundance of genes encoding
128 resistance to antimicrobials (antibiotics, heavy metals, biocides) and relaxases in different
129 samples (Figure 2) is not surprising due to the variability of microbiotas of different hosts ^{25,26}

130

131 ***ResCap addresses gene diversity***

132 Allele redundancy of some resistance genes hinders the correct estimation of “gene diversity”
133 and precludes a correct estimation of “gene abundance” in metagenomes when using most
134 available metagenomic tools.

135 To overcome this issue, we define the term Mapping Gene Cluster (MGC) as the group of
136 alleles/genes detected by the same set of reads (see Material and Methods). MGCs, firstly
137 defined in this work, allow an estimation of gene diversity across samples, and are measured as
138 the number of MGCs per million reads (MPM). The number of MPMs increased 1.3 fold in
139 humans (0.7-1.74) and 2.1-fold (2.3-1.9) in swine when using ResCap instead of MSS (**Figure**
140 **2**).

141 An increase in reads per MGC does not imply a homogeneous distribution of the reads.
142 Therefore, we also determined the “gene horizontal alignment coverage”, which was defined as
143 the fraction of a gene that is covered by reads. The probability of identifying an allele-specific
144 mutation will also increase with the number of reads per nucleotide or “gene depth coverage”.

145 **Figure 3** shows the improvement of “gene alignment horizontal coverage” using ResCap and
146 MSS (average= 97.5%, range = 66%-99% vs. average= 73.4%, range = 35.9%-94.8%,
147 respectively). Most genes were almost fully covered by reads and there was also a general
148 increase in “gene depth coverage” (**Figure S2**). As a consequence, the number of genes
149 unequivocally detected by ResCap almost doubled that of MSS (n=26, range 17.1-30.0 genes
150 per sample per million of reads vs. n=14.9, range 12-17.6 genes per sample per million of
151 reads). The number of reads unequivocally mapped increased up to 300 fold (2×10^5 for
152 ResCap vs 8×10^2 for MSS) (**Figure 4**).

153 **Figures S3** shows the abundance (RPKMs) and diversity (MPMs) obtained by ResCap and MSS
154 for individual categories of resistance genes (antibiotics, biocides and metals), which also
155 illustrates the improved sensitivity of ResCap vs. MSS. **Figure S4** reflects that although both

156 ResCap and MSS can track the most abundant gene families as those conferring resistance to
157 beta-lactams, macrolides, aminoglycosides and tetracyclines, followed by those conferring
158 resistance to phenicols and sulphonamides, many canonical resistance genes were only
159 detected by the ResCap platform (e.g. *mecA*, *blaZ* in beta-lactams; *ermA*, *ermC*, *ermD*, *erm33* or
160 *Inu* among macrolides; *fexA*, *catA* and *catB* alleles among phenicols). Genes encoding resistance
161 to fluoroquinolones, glycopeptides, or trimethoprim, first line antibiotics families used to treat
162 community and hospital-based infections, were barely detected using MSS but unequivocally
163 detected with ResCap (e.g. *dfrA16*, *dfrA15*, *dfrG*, *dfrK* among those conferring resistance to
164 trimethoprim, *oqxAB*, *qnrB*, *qnrS* among those producing resistance to quinolones, or *vanB*,
165 *vanA* for glycopeptides-resistance). ResCap also detected more genes conferring resistance to
166 heavy metals (e.g. cadmium, copper, silver or mercury), and relaxases, which are markers of
167 plasmid families that carry antibiotic resistance genes (MOB_C, MOB_F, MOB_{P1}, MOB_{P2}) (**Figures**
168 **S5-S7**).

169 170 **Comparative analysis of resistomes from different samples**

171 A statistical analysis of “gene abundance”, analogous to that used for comparing the abundance
172 of mRNA among samples in differential expression analysis²⁷, was performed to quantify the
173 improvement of ResCap over MSS in samples from different hosts. The need for such
174 comparisons is based on the known differences in microbiotas of different hosts.

175 “Gene abundance” data without normalization were processed as “count data” and used as
176 “input data” for differential analysis of the genes (detection of the genes only present in either
177 human or swine samples) (**Figure 5**). Using MSS, the resistome of the total sample analyzed
178 comprises 88 MGCs differentially detected (60 MGCs from humans and 28 MGCs from swine)
179 with a p-value lower than 0.001. Conversely, ResCap detected 262 MGCs (186 from humans and
180 76 from swine) (**Figure 5, panel a**). Out of these 262 MGCs, 185 were differentially detected by
181 ResCap and not by MSS, 77 were differentially detected by both approaches and 11 MGCs were
182 only differentially detected by MSS (**Figure 5, panel C**). This result means that ResCap detected
183 roughly three times more the MGCs on each resistome than MSS. The 11 MGCs that were only
184 detected by MSS belong to common (“present in both human and swine samples”) MGCs by
185 ResCap, suggesting that these differentially detected MGCs might in fact represent false
186 positives. Meanwhile, the number of common MGCs detected in human and swine sets was 437
187 with MSS and 569 with ResCap, of which 269 MGCs were disclosed by both approaches, 300
188 MGCs being specific for ResCap and 168 for MSS. The 168 MGCs detected as common between
189 human and swine metagenomes with MSS but not with ResCap were identified to be
190 differentially present by ResCap as false negatives. This can be explained because the count of

191 reads by MSS is lower than that of ResCap which makes the statistical analysis confidence values
192 by ResCap better for a given MGC.

194 **Reference-free evaluation**

195 ResCap includes approximately 78,600 genes that are homologous to “known” resistance
196 genes, with different degrees of sequence identity with defined resistance genes, which might
197 be involved in antibiotic resistance.

198 Assembly statistics and coverage show that the information obtained with ResCap only covers
199 a small portion of the metagenome, as intended by design (**Figure 7**). As expected, ResCap
200 increases, with respect to MSS, the number of sequenced genes that are homologous
201 (evolutionary close) to the canonical genes included in Arg-ANNOT, BACMet and ConjDB
202 databases. To perform a comparative analysis, the genes were catalogued as “ResCap”,
203 “UniProt” or “Novel”. The “ResCap” gene set includes genes within the ResCap database of
204 canonical genes. The “UniProt” gene set comprehends those that are already described in
205 UniProtKB database and result in a positive blast against ResCap database. The “Novel” gene
206 set corresponds to those genes not included in UniProtKB but resulting in a positive blast
207 against ResCap canonical database. Only Blast hits with e-values lower than 10^{-100} were
208 considered as positive and included in the analysis.

209 The annotation of the genes shows that ResCap also improves the recovery of genes homology
210 with genes coding for resistance, (UniProtKB 752 ± 237 genes with ResCap vs 237 ± 107 for
211 humans and 441 ± 71 genes vs 82 ± 46 for swine with MSS; Novel genes, 79 ± 38 genes with
212 ResCap vs 20 ± 7 107 for humans and 105 ± 26 genes vs 9 ± 4 for swine with MSS) as presented in
213 **Figure 8**. Although the actual role of these genes in antibiotic resistance will require functional
214 validation that is beyond the scope of the current study, its identification as *bona fide*
215 resistance genes as well as the analysis of their abundance upon antibiotic challenge might
216 have a deep impact in further studies on the evolution of antibiotic resistance. **Figure S8**
217 shows the better resolution of ResCap expressed by number of blast hits per gene per
218 megabase.

DISCUSSION

This work reports the development of a novel resistance gene capture platform ResCap and on its comparative evaluation with MSS in resistance gene identification in a collection of human and swine faecal samples. Our study shows that ResCap is ideally suited for high-resolution analysis of resistomes and also offers the possibility to detect genes homologous to “known” resistance genes to further analyse the evolution of antibiotic resistance.

ResCap also provides several technical advantages to study resistomes in comparison with current metagenomic methods. First, the enrichment of ResCap resides in its targeted metagenomics approach, which significantly increases the recovery of sequences of resistance genes. Thus, ResCap reduces the sequencing depth needed to comprehensively detect the targeted genes and, consequently, contributes to lower sequencing costs. More importantly, it can significantly lower the limit of detection of resistance genes in complex microbiomes. Our results indicate that the resistome represents barely 0.2% of the gut metagenome. As a consequence, MSS needs at least 3.75×10^9 reads per sample to reach a similar coverage to that obtained by using ResCap (average of 1.9×10^7 paired reads that represents a relative enrichment of 279x). Second, the tiling of capture probes greatly facilitates the higher level of “gene horizontal alignment” coverage of ResCap relative to MSS (**Figure 2, Figure S2**), thus increasing specificity (**Figure S9**). Third, ResCap ability to detect previously unrecognized DNA fragments with homology to canonical resistance genes will facilitate the discovery of novel genes potentially involved in antibiotic resistance. In case they are antibiotic-selectable, the novel genes will be enriched in the presence of antibiotics. In addition, ResCap will be of interest in Public Health, because it allows a more accurate “ranking risk analysis”¹⁸ of the genes within the resistomes of different microbiotas. Finally, the substantial capacity of the platform (200Mb) makes ResCap extensible up to two fold of its current capacity, thus making possible its updating with new sequences published or added to resistance gene databases. ResCap updates will be publicly available through the GitHub repository and the Nimblegene webpage. Nonetheless, the threshold of detection of ResCap remains unknown due to the lack of a negative control that demonstrate the ability of ResCap to pick antibiotic resistance genes from quantified minority populations (e.g. mock genomic populations). Although appropriate, the complexity and variability of the metagenomic samples makes difficult to use a good negative control to this kind of studies.

The definition of parameters that accurately express antibiotic resistance “gene abundance” and antibiotic resistance “gene diversity” constitute a requirement to comparatively analyse the resistomes of different samples. Relative abundance parameters are widely used in computational analysis of MSS datasets, but require specialized statistics, as these

257 compositional parameters are influenced by the variability in metagenomes of different
258 samples. The application of the novel concept of MGCs (Mapping Gene Clusters, groups of
259 alleles detected by the same set of reads) provides a set of normalized variables that can be
260 measured in abundance and diversity among samples. Furthermore, the MGC-based system
261 permits to evaluate the diversity within and between different functional groups (in our case,
262 families of antibiotics, groups of genes conferring resistance to heavy metals or biocides and
263 plasmid relaxases). To date, only a very few quantitative metagenomic approaches to analyze
264 resistomes are available but they do not achieve this level of accuracy^{14,16}.

265 Because of its sensitivity, specificity, and the possibility to accurately compare results between
266 samples, ResCap complies with the needs of public health epidemiology of antibiotic resistance
267 that include: i) the detection of emerging antibiotic resistance risks in different microbial
268 environments²⁸ (<http://www.efsa.europa.eu/en/press/news/140325>); ii) the need for
269 implementation of accurate risk assessment studies based on resistome analysis in healthy
270 humans, hospitalized patients, animal husbandry, food industry, and the environment; iii) the
271 quality control of sewage and water bodies decontamination of antibiotic resistant genes iv)
272 the update and refining of the list of resistance genes to be considered in monitoring the
273 adverse effects of drugs in microbiomes, including pharmacomicrobiomic applications in
274 clinical trials; v) the close monitoring of the efficacy of microbiome reconstitution/re-biosis,
275 whether through targeted probiotic-live culture administration or fecal microbiota
276 transplantation, to alleviate the adverse impact of antibiotic administration, and vi) to analyse
277 the effect of eco-evo drugs and strategies to combat antibiotic resistance²⁹.

278 In summary, ResCap provides an opportunity to meet the challenge of analyzing samples with
279 complex and heterogeneous mix of genes in low and high concentration DNA samples. Thus,
280 ResCap-like approaches might also be used to other complex microbial systems and their
281 minority bacterial populations (e.g. virulence determinants, key-ecological traits involved in
282 biosynthesis or biodegradation, or relevant genes of biotechnological interest).

283

284 **METHODS**

285 **ResCap design**

286 The ResCap capture library was based on a homemade core reference database (it will be
287 available as per request) that comprises both well-known and hypothetical genes encoding
288 resistance to antimicrobials (antibiotics, heavy metals, biocides) and genes coding for plasmid
289 family markers (relaxases). The core reference database was built by downloading sequences
290 associated with non-redundant antimicrobial genes available in curated databases Arg-
291 ANNOT³⁰, CARD³¹, RED-DB (<http://www.fibim.unisi.it/REDDB/Default.asp>), ResFinder³² and
292 Bacmet³³.

293 The putative antibiotic resistance genes dataset was constructed as follows. All antibiotic
294 resistance databases were combined in a non-redundant set. Proteins were clustered in
295 protein families by homology, using CD-HIT with parameters of 80% identity and 80%
296 coverage. First, each protein family was aligned by MUSCLE v. 3.7³⁴ with default parameters.
297 Then, a Hidden Markov Model (HMM) was built for each family with *hmmbuild* function of the
298 HMMER3³⁵ using default parameters. Hmmer search function (*hmmsearch*) was used against
299 UniProtDB for each HMM profile to search homologous proteins for each family of proteins
300 that confer antibiotic resistance. Manual curation of datasets was performed to remove false
301 positives. Final protein data set was translated to DNA sequence using ENA accession numbers
302 associated with each UniProtDB entry.

303 As a result, the final ResCap targeted sequence panel consists of 78,600 non-redundant genes
304 (81,117 redundant genes) that would search a target space of 88.13Mb, not reaching yet the
305 200Mb target capacity offered by the custom SeqCap EZ library format (NimbleGen). Probes
306 targeting the antibiotic resistome include 47,806 putative antibiotic resistance genes and 7,963
307 functionally characterized, canonical, antibiotic resistance genes. Probes targeting the metal
308 and biocide resistome include 30,794 putative resistant genes and 704 canonical resistance
309 genes. The platform also includes probes for 2,517 relaxases of the Conj database.

310 The consolidated list of target sequences was submitted to Roche NimbleGen for capture
311 library design and synthesis and further implemented under the custom NimbleGen SeqCap EZ
312 Developer Library format. Redistribution of probes for better capture uniformity, redundancy,
313 and comprehensive target base coverage relied on NimbleGen, and was based on patented
314 algorithms. ResCap design covers the 98.3% of the 88.13Mb and 99.6% of the genes have more
315 than 50% of their sequence covered. (**Figure S9**).

316
317 **The ResCap workflow**

318 The Rescap workflow consists of: i) whole-metagenome shotgun library construction, ii)
319 hybridization, and iii) capture. All steps were performed according to NimbleGen standard

320 protocols for Illumina platforms. To evaluate ResCap efficiency, samples were sequenced
321 before and after capture.

322 i) Whole-metagenome shotgun library construction. Total nucleic acid was extracted following
323 the standardized Metahit protocol³⁶ (<http://www.metahit.eu/>) and using the FastPrep
324 instrument (MP Biomedicals, USA). Libraries were prepared following the instructions of
325 “Kapa Library Preparation Kit for Illumina platforms” (Kapabiosystems, KR0935-v1.13).
326 Briefly, 1.0 µg input DNA (measure by Picogreen) was fragmented to 500-600 bp insert size by
327 sonication with Bioruptor (FastPrep®-24). After End repair, A-tailing and Adapter ligation, we
328 follow Dual-SPRI size selection adding 0.5 vol in first cut and 0.2 vol in the second cut to get
329 650-750pb libraries.. Library amplification was carried out using LM-PCR of 7 cycles, as
330 indicated in the SeqCap EZ Library SR User’s Guide v4.2. At this level, samples were labelled
331 with specific barcodes for further sample identification. A first aliquot of the resulting
332 amplified libraries were quality checked in a Bioanalyzer 2100 (Agilent) and pooled in
333 equimolecular amounts for sequencing on Illumina HiSeq 2000 instrument, generating 100-
334 150-bp paired-end reads (“pre-capture” samples).

335
336 ii) Hybridization and capture. The second part of each DNA library was subjected to targeted
337 sequence capture with the custom ResCap probes prior to sequencing (“post-capture”
338 samples). Both experiments were made in separate sequencing runs. Targeted sequence
339 capture was carried out according to the manufacturer’s specifications. The captured DNA was
340 checked for quality and integrity in a Bioanalyzer and titrated by quantitative PCR using the
341 “Kapa-SYBR FAST qPCR kit forLightCycler480” and a reference standard for quantification. The
342 captured libraries were denatured prior to be loaded on a flow-cell at a density of 2,2pM,
343 where clusters were formed and sequenced using a HiSeq 2000 in a 2x100 pair-end mode for
344 swine samples and NextSeq 500 in a 2x150 pair-end mode for human samples. Raw sequences
345 were processed using FastX Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/).

346 **Bioinformatic analysis**

347 ***Reference-based workflow***

348
349 Analysis of sequence data from metagenomes constitutes a challenge because of the inherent
350 variability of the samples analysed, and the limitations of current bioinformatics’ methods to
351 unequivocally identify specific alleles from short length reads (100-150 bp). To overcome such
352 limitations, we developed a novel approach to define variables suitable for inferring “gene
353 abundance” and “gene diversity” and, in our case, to perform quantitative analysis of
354 antimicrobial resistance genes. Moreover, we suggest a workflow of variable normalization in
355 relation to the information content of the targeted variable that would make it possible to

356 compare different samples of different hosts. These tools were developed for ResCap but could
357 be implemented for any other metagenomic sequence dataset. Shotgun metagenomic
358 sequencing allowed assembling the sequences into contigs to infer the functionality of the
359 sequenced metagenome. **Figure 9** shows the workflow that illustrates and defines the
360 variables used.

361 *Raw Data Processing*

362 Reads were mapped against our database comprising ARG-ANNOT, BacMET and ConjDB
363 databases independently, using Bowtie2 software³⁷. Bowtie was set up to retrieve all end to
364 end possible alignments and suppress both discordant alignments and mixed alignments. The
365 output SAM file was parsed to get the fields of *Query template NAME*, *Reference sequence NAME*,
366 *1-based leftmost mapping Position*, *MAPping Quality*, *Position of the mate/next read*. Reads with
367 unavailable information (field *Query Template NAME* equal to ‘*’) were suppressed.
368 Subsequently, a homemade perl script (available per request) was used to count matched
369 reads per gene. Using the SAM parsed file and the length of the reference genes, the perl script
370 generated a table with the following fields i) the number of reads per gene mapped (RPG, “gene
371 depth coverage”), ii) the number of reads per kb of gene (RPK), iii) the number of the reads
372 that were mapped unequivocally to a given gene and iv) the percentage of coverage of the gene
373 sequence (“gene horizontal alignment coverage”) of each mapped gene. Table fields Unique,
374 RPG and RPK were normalized by the total amount of reads in each sample, transforming such
375 fields in “read per gene per million reads” and “reads per kb per million reads” (RPKM),
376 respectively, the last one being a common unit of “gene abundance”³⁸. Several ways to
377 normalize abundance data have been applied to different studies (e.g., expression data in RNA-
378 Seq experiment). The aim of our approach was to estimate the proportion of antimicrobial
379 resistance genes among samples that putatively contained the same amount of DNA, so the
380 normalization using the total amount of DNA (i.e reads) among samples fits better with the
381 initial approach.

382
383 The redundancy of mapped reads may be represented as a network where the nodes are the
384 genes (usually alleles of the same gene) and the edges are the reads that map in the different
385 nodes. Because one read can map in different alleles/genes, all the genes mapped by these
386 reads are linked among them. The resulting network that comprises all the nodes and edges in
387 a set of samples is named “allele Network” (**Figure S10**). In our context, the allele network
388 must be unique for all samples of a given assay, so an allele network was built joining all the
389 SAM parsed files of the study.

390 Each cluster of the Allele Network represents the set of genes that are detected by a set of
391 reads. They are defined as a Mapping Gene Clusters (MGCs) and each one may include

392 hundreds of genes or just one gene. A given MGC will be detected when at least one read maps
393 against any of the genes within that MGC (diversity). To quantify the MGCs in each sample, the
394 highest value shown by an allele (node) within a given MGC is taken as the occurrence of such
395 MGC (abundance). The MGCs system builds a set of normalized variables that can be measured
396 in abundance and diversity among samples and thus, allows comparing datasets of different
397 sources, while maximizing the accuracy of the observable information.

398 A homemade perl script was used to build the allele network from the SAM parsed files, taking
399 the mapped genes as nodes and searching the ambiguously mapped reads to create the edges.
400 Perl script calculates the edges-weight as the number of reads that map the linked nodes at the
401 same time. Allele Network was loaded in R environment³⁹ using the *igraph* package⁴⁰. MGCs
402 were defined using *mcl*, from MCL R package⁴¹, with default parameters except allow loops and
403 cluster with only one member on the allele network.

404

405 *Data Analysis*

406 The resistome of a given experiment was analysed in terms of gene abundance and diversity
407 according to the methodology described above. The *abundance* and the *diversity* of genes in a
408 particular resistome are the (dependent) variables that define this resistome and are measured
409 as the number of RPKM per MGC and the number of MGCs, respectively.

410 The number of MGCs was normalized by the total number of sequencing reads per each sample
411 expressed in millions of reads (MPM), this value being considered as a unit of **diversity**. MGCs
412 of the antibiotic resistance gene database were divided according to antibiotic families³⁰. MGCs
413 of the relaxase database were organized in known different relaxase families⁴². The MGCs of
414 biocide and heavy metal resistance gene database were classified according the susceptibility
415 to specific compounds³³. Genes that belong to more than one functional category (e.g. some
416 conferring resistance to different metals) contribute equally for any of them. Figure S9 shows
417 the 839 MGCs determined in our sample (237 for AbR, 283 for Biocide and Metals and 319 for
418 relaxases). Descriptive statistic was performed using *dplyr*⁴³, *tidyr*⁴⁴ and *ggplot2*⁴⁵ packages of
419 R³⁹.

420 Differential analysis was performed using DESeq2 package⁴⁶. Although DESeq2 was originally
421 designed for differential expression analysis, it also works well with abundance data. Tables
422 containing the original abundance data obtained by ResCap and MSS datasets were used
423 separately as input for DESeq2 package to determine the MGCs differentially detected between
424 swine and human hosts. Normalization and statistical analysis were performed with the
425 default parameters of DESeq2. MGCs with p-value lower than 0.001 were classified as
426 differentially detected, rest of the MGCs (p-value above 0.001) were classified as commonly
427 detected.

428

429 ***Reference-free workflow***

430 Assemblies were performed by MegaHit software with default parameters⁴⁷. Prodigal⁴⁸ was
431 used for gene recognition and translation with the specific parameters for metagenomic
432 sequences. Quality assemblies' quantification was performed by Quast software⁴⁹. Predicted
433 genes were first annotated against the ResCap database by Best Blast Hit approach using blastn
434 software⁵⁰. In order to identify only genes belonging to ResCap database or their homologs, and
435 minimize the false positive ratio, Blast hits were filtered by e-value of 10^{-100} and 80% of
436 coverage. Genes with identities higher than 95% and coverage higher than 80% were
437 considered as belonging to ResCap. The remaining genes were translated to proteins. These
438 proteins were classified as non-ResCap and were compared against UniProt by blastp. Again,
439 hits with higher identity than 95%, coverage higher than 80% and e-value lower than 10^{-100}
440 were considered as UniProt known proteins. The set of proteins that did not accomplish this
441 threshold were considered as novel proteins.

442

443 **Samples analyzed**

444 ResCap was validated by analysing fecal samples from 9 human and 8 swine individuals, all
445 collected as part of FP7 European Research Consortium EvoTAR (www.evotar.eu). Swine
446 samples were collected in Spanish farms linked to large companies which supply broilers and
447 swine processed meat in the EU. Antibiotics as growth promoters or with preventive purposes
448 are not used in these farms. Human samples were collected in the Hôpital Bichat, Paris, France,
449 under the protocol approved by its local ethics committee. DNA preparation was accomplished
450 for animal and human samples using standardized protocols (MetaHIT Protocol). Robustness
451 of the platform was tested by comparative analysis of two technical replicates of two swine
452 samples.

453

454 **ACKNOWLEDGEMENTS**

455 This work was supported by the European Commission, Seven Framework Program
456 (EVOTARFP7-HEALTH- 282004 for VFL, FB, JLM, AA, DE, ER, RJLW, WvS,FdIC and TMC), the
457 Join Programming Initiative in Water (JPI Water StARE JPIW2013-089-C02-01 to JLM), the
458 Ministry of Economy and Competitiveness of Spain (BIO2014-54507-R to JLM, and
459 PLASWIRES-612146/FP7-ICT- 2013-10 and BFU2014-55534-C2-1-P for FdIC). Authors also
460 acknowledge the European Development Regional Fund “A way to achieve Europe” (ERDF) for
461 cofounding the Plan Nacional de I+D+ I 2012-2015 (BIO2014-54507-R to JLM, PI12-01581 to
462 TMC and BFU2014-55534- C2-1-P for FdIC) and CIBER (CIBER in Epidemiology and Public
463 Health, CIBERESP; CB06/02/0053 to FB) and the Spanish Network for Research on Infectious
464 Diseases (REIPI RD12/0015 to JLM), and the Regional Government of Madrid (PROMPT-
465 S2010/BMD2414). Val F. Lanza was further funded by a Research Award Grant 2016 of the
466 European Society for Clinical Microbiology and Infectious Diseases (ESCMID).

467 Conflict of Interest: none declared

468

469

470 **COMPETING FINANTIAL INTEREST**

471 The authors declare no competing financial interests

472

473

474 **REFERENCES**

- 475 1. Kesselheim, A. S. & Outterson, K. Fighting antibiotic resistance: Marrying new financial
476 incentives to meeting public health goals. *Health Aff.* **29**, 1689–1696 (2010).
- 477 2. Phillips, I. *et al.* Does the use of antibiotics in food animals pose a risk to human health? A
478 critical review of published data. *J. Antimicrob. Chemother.* **53**, 28–52 (2004).
- 479 3. World Health Organization & Who. WHO Global Strategy for Containment of Antimicrobial
480 Resistance. *World Health WHO/CDS/CS*, 105 (2001).
- 481 4. Bax, R. *et al.* Surveillance of antimicrobial resistance--what, how and whither? *Clin. Microbiol.*
482 *Infect.* **7**, 316–25 (2001).
- 483 5. Laxminarayan, R. *et al.* Antibiotic resistance-the need for global solutions. *Lancet. Infect. Dis.*
484 **13**, 1057–98 (2013).
- 485 6. D'Costa, V. M., McGrann, K. M., Hughes, D. W. & Wright, G. D. Sampling the antibiotic
486 resistome. *Science* **311**, 374–377 (2006).
- 487 7. Wright, G. D. The antibiotic resistome: the nexus of chemical and genetic diversity. *Nat. Rev.*
488 *Microbiol.* **5**, 175–86 (2007).
- 489 8. Wright, G. D. The antibiotic resistome. *Expert Opin. Drug Discov.* **5**, 779–88 (2010).
- 490 9. Forsberg, K. J. *et al.* Bacterial phylogeny structures soil resistomes across habitats. *Nature*
491 **509**, 612–6 (2014).
- 492 10. Forsberg, K. J. *et al.* The Shared Antibiotic Resistome of Soil Bacteria and Human Pathogens.
493 *Science (80-.).* **337**, 1107–1111 (2012).
- 494 11. Sommer, M. O. A., Dantas, G. & Church, G. M. Functional characterization of the antibiotic
495 resistance reservoir in the human microflora. *Science* **325**, 1128–1131 (2009).
- 496 12. Modi, S. R., Lee, H. H., Spina, C. S. & Collins, J. J. Antibiotic treatment expands the resistance
497 reservoir and ecological network of the phage metagenome. *Nature* **499**, 219–22 (2013).
- 498 13. Bengtsson-Palme, J., Boulund, F., Fick, J., Kristiansson, E. & Larsson, D. G. J. Shotgun
499 metagenomics reveals a wide array of antibiotic resistance genes and mobile elements in a
500 polluted lake in India. *Front. Microbiol.* **5**, (2014).
- 501 14. Hu, Y. *et al.* Metagenome-wide analysis of antibiotic resistance genes in a large cohort of
502 human gut microbiota. *Nat. Commun.* **4**, 2151 (2013).
- 503 15. Ghosh, T. S., Gupta, S. Sen, Nair, G. B. & Mande, S. S. In silico analysis of antibiotic resistance
504 genes in the gut microflora of individuals from diverse geographies and age-groups. *PLoS*
505 *One* **8**, e83823 (2013).
- 506 16. Forslund, K. *et al.* Country-specific antibiotic use practices impact the human gut resistome.
507 *Genome Res.* **23**, 1163–9 (2013).
- 508 17. Lynch, M. D. J. & Neufeld, J. D. Ecology and exploration of the rare biosphere. *Nat. Rev.*
509 *Microbiol.* **13**, 217–29 (2015).

- 510 18. Martínez, J. L., Coque, T. M. & Baquero, F. What is a resistance gene? Ranking risk in
511 resistomes. *Nat. Rev. Microbiol.* **13**, 116–123 (2014).
- 512 19. Ng, S. B. *et al.* Targeted capture and massively parallel sequencing of 12 human exomes.
513 *Nature* **461**, 272–276 (2009).
- 514 20. Olson, M. Enrichment of super-sized resequencing targets from the human genome. *Nat.*
515 *Methods* **4**, 891–892 (2007).
- 516 21. Tennessen, J. A. *et al.* Evolution and functional impact of rare coding variation from deep
517 sequencing of human exomes. *Science* **337**, 64–9 (2012).
- 518 22. Mamanova, L. *et al.* Target-enrichment strategies for next-generation sequencing. *Nat.*
519 *Methods* **7**, 111–8 (2010).
- 520 23. Jones, M. R. & Good, J. M. Targeted capture in evolutionary and ecological genomics. *Mol. Ecol.*
521 **25**, 185–202 (2016).
- 522 24. Garcillán-Barcia, M. P., Francia, M. V. & de la Cruz, F. The diversity of conjugative relaxases
523 and its application in plasmid classification. *FEMS Microbiol. Rev.* **33**, 657–87 (2009).
- 524 25. Ley, R. E., Lozupone, C. a, Hamady, M., Knight, R. & Gordon, J. I. Worlds within worlds:
525 evolution of the vertebrate gut microbiota. *Nat. Rev. Microbiol.* **6**, 776–88 (2008).
- 526 26. Dethlefsen, L., McFall-Ngai, M. & Relman, D. A. An ecological and evolutionary perspective on
527 humang-microbe mutualism and disease. *Nature* **449**, 811–818 (2007).
- 528 27. Wang, Z., Gerstein, M. & Snyder, M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat.*
529 *Rev. Genet.* **10**, 57–63 (2009).
- 530 28. Frieden, T. Antibiotic resistance threats in the United States. *Centers Dis. Control Prev.* 114
531 (2013). doi:CS239559-B
- 532 29. Baquero, F., Coque, T. M. & de la Cruz, F. Ecology and evolution as targets: the need for novel
533 eco-evo drugs and strategies to fight antibiotic resistance. *Antimicrob. Agents Chemother.* **55**,
534 3649–60 (2011).
- 535 30. Gupta, S. S. K. *et al.* ARG-ANNOT, a New Bioinformatic Tool To Discover Antibiotic Resistance
536 Genes in Bacterial Genomes. *Antimicrob. Agents Chemother.* **58**, 212–20 (2014).
- 537 31. McArthur, A. G. *et al.* The comprehensive antibiotic resistance database. *Antimicrob. Agents*
538 *Chemother.* **57**, 3348–57 (2013).
- 539 32. Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes. *J. Antimicrob.*
540 *Chemother.* **67**, 2640–2644 (2012).
- 541 33. Pal, C., Bengtsson-Palme, J., Rensing, C., Kristiansson, E. & Larsson, D. G. J. BacMet:
542 antibacterial biocide and metal resistance genes database. *Nucleic Acids Res.* **42**, D737-43
543 (2014).
- 544 34. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
545 *Nucleic Acids Res.* **32**, 1792–7 (2004).

- 546 35. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 547 36. Wesolowska-Andersen, A. *et al.* Choice of bacterial DNA extraction method from fecal
548 material influences community structure as evaluated by metagenomic analysis. *Microbiome*
549 **2**, 19 (2014).
- 550 37. Pongor, L. S., Vera, R. & Ligeti, B. Fast and sensitive alignment of microbial whole genome
551 sequencing reads to large sequence datasets on a desktop PC: Application to metagenomic
552 datasets and pathogen identification. *PLoS One* **9**, (2014).
- 553 38. Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying
554 mammalian transcriptomes by RNA-Seq. *Nat. Methods* **5**, 621–8 (2008).
- 555 39. R Development Core Team. R: A language and environment for statistical computing. R
556 Foundation for Statistical Computing, Vienna, Austria. *R Foundation for Statistical Computing*,
557 *Vienna, Austria.* (2014).
- 558 40. Csardi, G. & Nepusz, T. The igraph software package for complex network research.
559 *InterJournal Complex Sy*, 1695 (2006).
- 560 41. Jäger, M. L. MCL: Markov Cluster Algorithm. (2015).
- 561 42. Guglielmini, J., de la Cruz, F. & Rocha, E. P. C. Evolution of Conjugation and Type IV Secretion
562 Systems. *Mol. Biol. Evol.* **30**, 315–31 (2013).
- 563 43. Wickham, H. & Francois, R. dplyr: A Grammar of Data Manipulation. (2015).
- 564 44. Wickham, H. Tidy data. *J. Stat. Softw.* **59**, 1–23 (2014).
- 565 45. Wickham, H. *ggplot2: elegant graphics for data analysis.* (Springer New York, 2009).
- 566 46. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome Biol.*
567 **11**, R106 (2010).
- 568 47. Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node
569 solution for large and complex metagenomics assembly via succinct de Bruijn graph.
570 *Bioinformatics* **31**, 1674–6 (2015).
- 571 48. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
572 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 573 49. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: Quality assessment tool for genome
574 assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
- 575 50. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database
576 search programs. *Nucleic Acids Res.* **25**, 3389–402 (1997).
- 577
578

579 **TABLE 1**

Sample	Metagenome shotgun sequence MSS					ResCap					Gain
	<i>N° Reads</i>	<i>AbR</i>	<i>BacMet</i>	<i>Relaxases</i>	Total	<i>N° Reads</i>	<i>AbR</i>	<i>BacMet</i>	<i>Relaxases</i>	Total	
Bichat1	14,127,290	0,05%	0,002%	0,04%	0,10%	16,705,789	19%	0,48%	5,22%	24,68%	244,24
Bichat2	15,128,135	0,05%	0,028%	0,04%	0,12%	33,589,838	12%	5,56%	2,98%	20,27%	170,81
Bichat3	14,488,245	0,05%	0,005%	0,03%	0,09%	17,276,637	34%	2,31%	5,96%	41,83%	480,59
Bichat6	17,476,666	0,07%	0,001%	0,05%	0,12%	19,191,320	25%	0,85%	5,84%	32,13%	261,87
Bichat7	16,732,926	0,07%	0,002%	0,05%	0,13%	28,530,922	27%	0,33%	5,99%	33,60%	267,77
Bichat9	17,058,000	0,03%	0,013%	0,05%	0,09%	18,038,257	14%	10,59%	7,55%	31,80%	336,29
Bichat10	15,039,883	0,03%	0,066%	0,06%	0,15%	34,798,281	6%	28,81%	5,48%	40,63%	265,17
Bichat11	13,425,077	0,03%	0,091%	0,06%	0,18%	35,901,508	5%	26,73%	4,72%	35,98%	201,67
Bichat13	17,903,872	0,05%	0,023%	0,06%	0,14%	26,283,052	16%	13,27%	7,39%	36,85%	270,02
F266	19,557,955	0,06%	0,005%	0,02%	0,08%	14,024,345	21%	4,62%	2,57%	28,23%	337,50
PIG20	27,375,311	0,08%	0,028%	0,01%	0,12%	22,485,364	18%	16,20%	1,81%	36,38%	298,23
PIG26	13,831,057	0,07%	0,005%	0,02%	0,10%	15,756,070	19%	3,88%	2,73%	25,93%	271,20
PIG29	18,945,765	0,09%	0,018%	0,02%	0,12%	26,223,850	18%	10,26%	2,52%	30,65%	248,31
PIG31	12,778,294	0,07%	0,010%	0,02%	0,09%	18,055,019	13%	5,76%	2,08%	20,77%	219,12
PIG528	19,689,471	0,06%	0,003%	0,02%	0,08%	13,864,257	21%	2,70%	3,11%	26,83%	323,23
PIG94	15,985,219	0,07%	0,004%	0,02%	0,10	15,351,408	18%	2,57%	3,31%	24,05%	240,20
PIG96	9,290,402	0,06%	0,001%	0,01%	0,07%	12,225,935	21%	1,13%	1,67%	23,84%	320,90

580

581 **Table 2. Summary of metagenomic comparative analysis**

	Human				Swine				Human & Swine			
	AbR	Bac	MOB	Total	AbR	Bac	MOB	Total	AbR	Bac	MOB	Total
MSS	8	3	49	60	24	0	4	28	106	197	134	437
ResCap	37	49	100	186	58	8	10	76	142	223	204	569

582 MGCs were classified on significantly detected on Humans ($\log_2\text{FoldChange} < 0$, $p\text{value} < 1e-3$),
583 significantly detected on Swine ($\log_2\text{FoldChange} > 0$, $p\text{value} < 1e-3$) and commons on Human and
584 Swine ($p\text{value} \geq 1e-3$).

January 23, 2017

Figure Legends

Figure 1

Gain function plot. Representation of the gain in reads per kilobase per million of reads of each detected gene between MSS protocol (abscissa axis) and ResCap (ordinate axis). Genes only detected by ResCap are represented by the dot cluster in the initial values of abscissa axis. The pictures are represented in log-log scale to a better perception of the linearity of the gain function in genes detected by both protocols.

Figure 2

Platform Efficiency by Source Sample and Data Base Group. Data distribution of the platform efficiency evaluating (a) the number of mapped reads per million of sequenced reads against canonical (well-known) genes data set; and (b) the number of detected genes per million of sequenced reads using as reference the well-known genes data set. Fecal samples were differentiated according to the source (9 from humans and 8 from swine). Data distribution of the platform efficiency evaluating (c) the number of mapped reads per million of sequenced reads against the three canonical genes groups and (d) the number of detected genes per million of sequenced reads using as reference the three canonical genes groups.

Figure 3

Longitudinal coverage distribution. The figure shows the comparison of longitudinal coverage distribution between protocols in each sample. Distributions are represented by density parameter and expressed by the number of genes (ordinate axis) and coverage percent (abscissa axis).

Figure 4

Quantification of unequivocally mapping reads. The figure shows the comparative of the quantification of reads that mapping on just one gene (or allele). First the abundance of reads that are unequivocally mapped on one gene (a). On another hand, the number of genes (or Mapping Gene Cluster) that have almost one read that mapping unequivocally (b). Box plots are differentiated for MSS protocol and ResCap protocol.

Figure 5

Differential study plots. Panel a) summarize the number of statistically significant MGCs of humans, swines and the genes in common between them using both approaches: MSS (up) and ResCap (button). Panel b) show the distribution of abundance variation between swine and human AbR resistomes (left), Metal and Biocide resistome (middle) and mobilome (right) in the form of volcano plots (fold change vs p-value) using the different approaches MSS (up) and ResCap (button) Left and right branches in the volcano refers to higher abundance in humans and swine respectively. Panel c) shown the Venn diagrams between approaches of differentially detected MGCs (up) and commonly (in both sets) detected MGCs (bottom)

Figure 6

Reproducibility of ResCap. Reads from replicates are represented in dot plot to illustrate the linearity of the results from ResCap sequencing. Dots represent the genes detected in any of the replicates. Pearson's product-moment correlation was used to estimate the correlation between technical replicates.

January 23, 2017

Figure 7

Assembly statistics. Statistic summary of the main assembly variables. *Total Length* and *Number of genes* were normalized by the total amount of megabases sequenced by each sample. Coverage data was calculated as the total sequenced bases divided by the total length (without normalizing). Assembly statistics was calculated by Quast software.

Figure 8

Functional annotation distribution. Assembled genes are classified as ResCap, UniProt or Novel (see Material & Methods). All assessed genes have a maximum e-value of 10^{-100} with some of the genes included in ResCap database. Figure show the comparative between human and swine samples and between MSS and ResCap approaches.

Figure 9

ResCap analysis workflow. Processed reads are mapped against reference database, SAM files are parsed to extract the reads unequivocally mapped and the ambiguously mapped to determine the Genes unequivocally detected and to perform the Allele Network. Allele Network is build using all SAM files of the study. The MGCs determines from Allele Network are used to perform the statistical analysis of Abundance and Diversity. Finally with the data of Abundance a differential analysis was performed.

January 23, 2017

LEGENDS TO SUPPLEMENTARY MATERIAL

Supplementary Figure 1

Gain function plot for each sample. Representation of the gain in reads per kilobase per million of reads of each detected gene between MSS (abscissa axis) and ResCap (ordinate axis). Genes only identified by ResCap are represented by the dot cluster in the initial values of abscissa axis. The pictures are represented in log-log scale to a better perception of the linearity of the gain function in genes detected by each protocol.

Supplementary Figure 2

Distribution of Reads Abundance: Figure shown the histograms of reads abundance per each gene. Each frame represent a sample, superposing results from MSS protocol and ResCap protocol. A square scale was used for ordinate axis and a logarithmic scale for abscissa axis to optimize the representation of the data.

Supplementary Figure 3

(1) Diversity and Abundance of Antibiotic Resistance: Comparison of ResCap and MSS protocol in Antibiotic Resistance data. Antibiotic resistance genes were divided among nine families by antibiotic family (AGly: Aminoglycosides, Bla: Beta-Lactamases, Flq: Fluoroquinolones, Gly: Glycopeptides, MLS: Macrolides, Phe: Phenicol, Sul: Sulphonamides, Tet: Tetracyclines and Tmt: Trimethoprim). Abundance (a) was measured as Read Per Kilobase per Million of reads that mapping against genes or allele-cluster genes of each antibiotic resistance family. Diversity (b) was measured as a number of detected Genes Per Million reads of each antibiotic resistance family.

(2) Diversity and Abundance of Relaxases: Comparison of ResCap and MSS protocol in Relaxases data. Relaxases were divided among nine protein families (MOB_B, MOB_C, MOB_F, MOB_H, MOB_{P1}, MOB_{P2}, MOB_Q, MOB_T and MOB_V). Abundance (a) was measured as Read Per Kilobase per Million of reads that mapping against genes or allele-cluster genes of each relaxase family. Diversity (b) was measured as a number of detected Genes Per Million reads of each relaxase family.

(3) Diversity and Abundance of Biocide & Metal resistance: Comparison of ResCap and MSS protocol in Biocide & Metal resistance data. Biocide & Metal resistance genes were divided by compound susceptibility. Abundance (a) was measured as Read Per Kilobase per Million of reads that mapping against genes or allele-cluster genes of each compound family. Diversity (b) was measured as a number of detected Genes Per Million reads of each compound family.

Supplementary Figure 4

MGCs abundance comparative of antibiotic resistance between swine and human samples. MGCs corresponding to antibiotic resistance dataset were classified by antibiotic families (AGly: Aminoglycosides, Bla: Betalactams, Flq: Fluoroquinolones, Gly: Glycopeptides, MLS: Macrolides, Phe: Phenicol, Sul: Sulphonamides, Tet: Tetracyclines, Tmt: Trimethoprim). Abundance was measured as Read per Kilobase per Million of reads. Panel right shown the results of MSS and panel left shown the results of ResCap.

January 23, 2017

Supplementary Figure 5

MGCs abundance comparative of biocide resistance between swine and human samples. Gene abundance was extracted from original count data after normalization. Some sets of genes make complex MGCs. In this representation, MGCs quantification was discarded in order to increase the biological information. Genes were classified by compound susceptibility. Due to biocide resistance genes spectrum of activity, genes are not constricted to one category but some genes show resistance to more than one compound. Genetic abundance is expressed as Reads per Kilobase per Million of Reads (RPKM). The panel right shows the results of MSS and the panel left shows the results of ResCap.

Supplementary Figure 6

Gene abundance comparative of metal resistance between swine and human samples. Gene abundance was extracted from original count data after normalization. Some sets of genes make complex MGCs. In this representation, MGCs quantification was discarded in order to increase the biological information. Genes were classified by metal susceptibility. Due to metal resistance genes spectrum of activity, some genes are not constricted to one category but some genes show resistance to more than one metal. Genetic abundance is expressed as Reads per Kilobase per Million of Reads (RPKM). The panel right shows the results of MSS and the panel left shows the results of ResCap.

Supplementary Figure 7

MGCs abundance comparative of Relaxases between swine and human samples. Relaxases were classified by MOB families. MGCs abundance was summarized in MOB families. Each MOB families are composed by several MGCs. Genetic abundance is expressed as Reads per Kilobase per Million of Reads (RPKM). The panel right shows the results of MSS and the panel left shows the results of ResCap.

Supplementary Figure 8

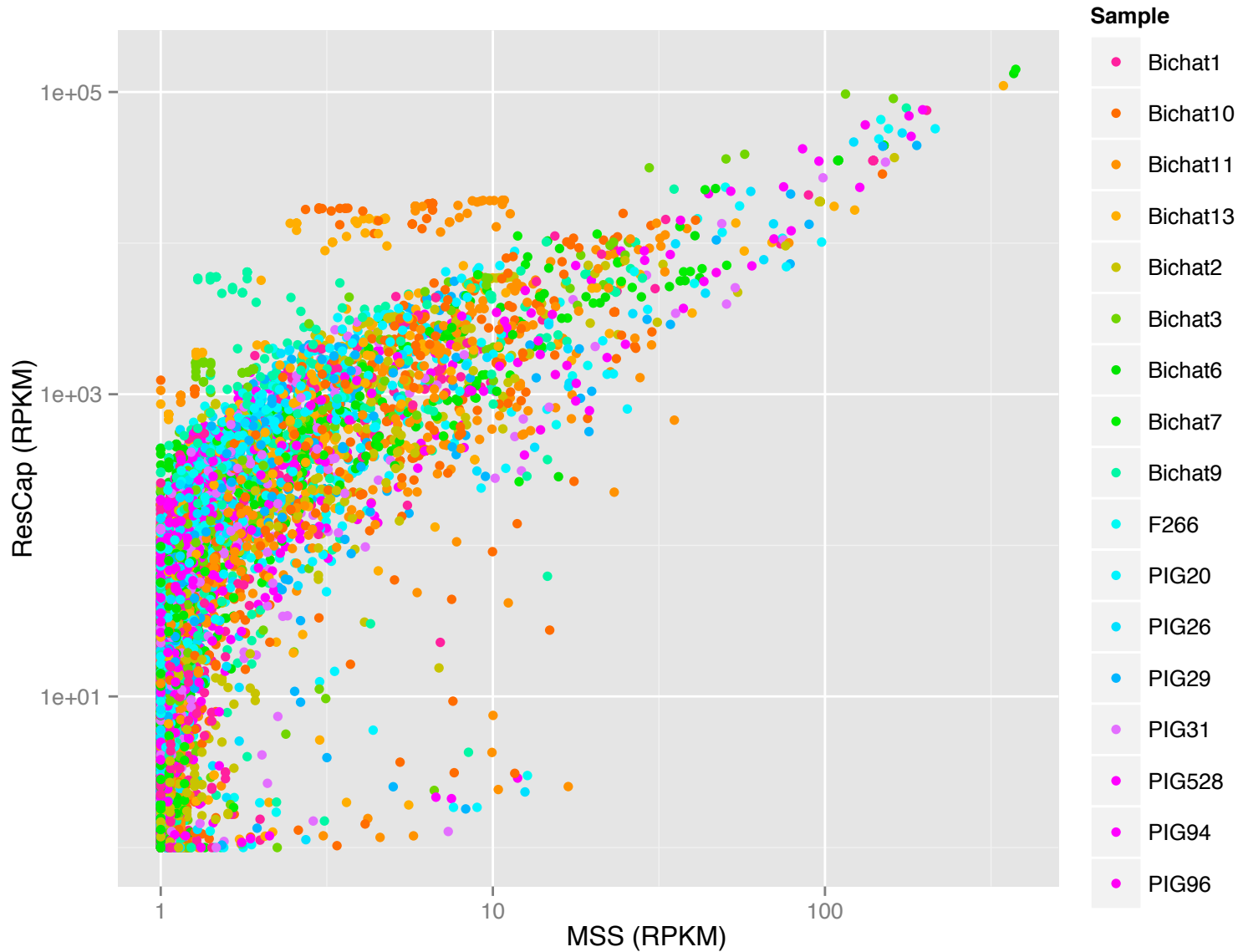
Blast annotations summary. Summary of the classification steps of assembled genes. The sequential annotation comprises a first blastn search for identify resistome homologous genes. Genes with evaluate higher than 10^{-100} were discarded. Filtered genes were split into two groups, genes with identity higher than 95% and genes with identity lower than 95%. The second group were annotated against UniProtKB and were split again into two groups, genes with identity higher than 95% of identity and genes with identity lower of 95%. A number of blast hits were normalized by the number of assembling genes per sequenced megabases.

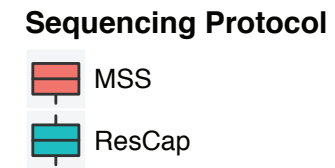
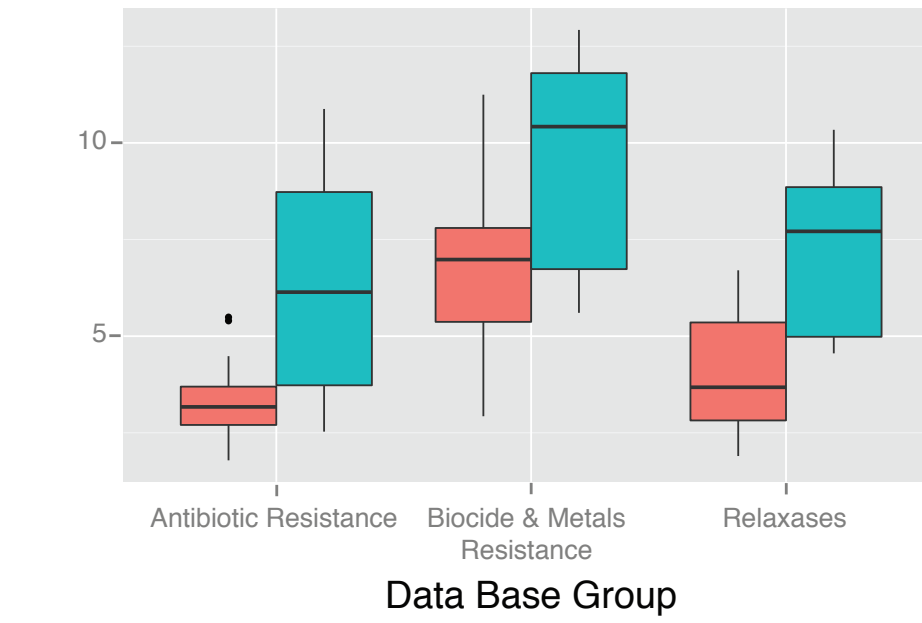
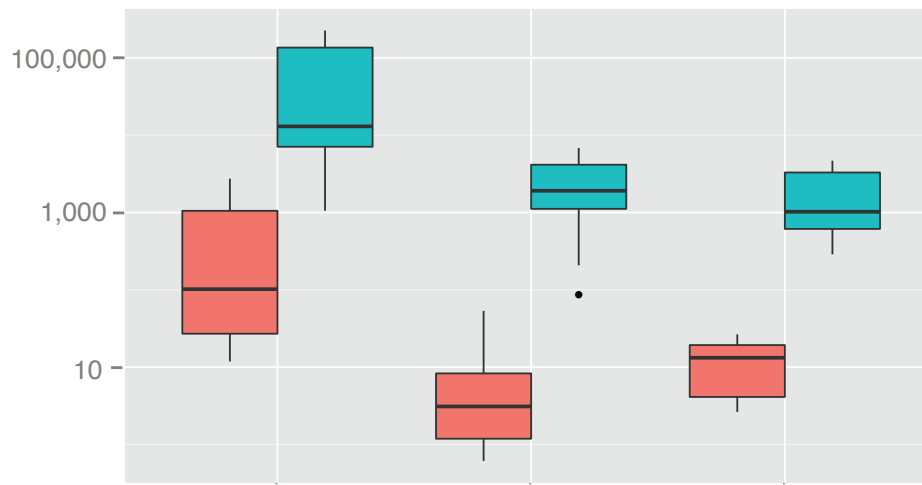
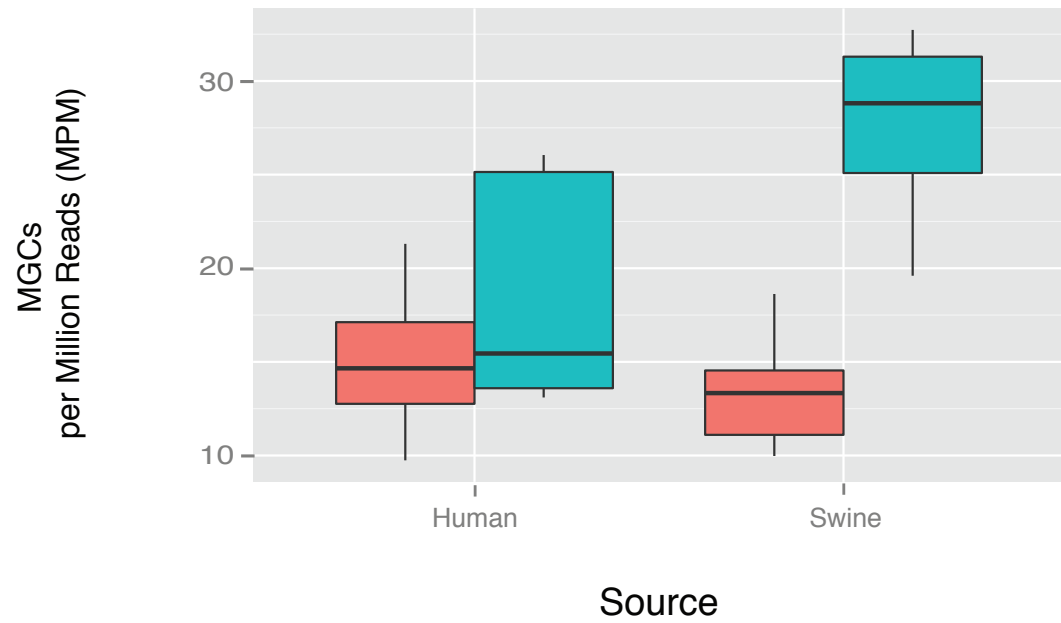
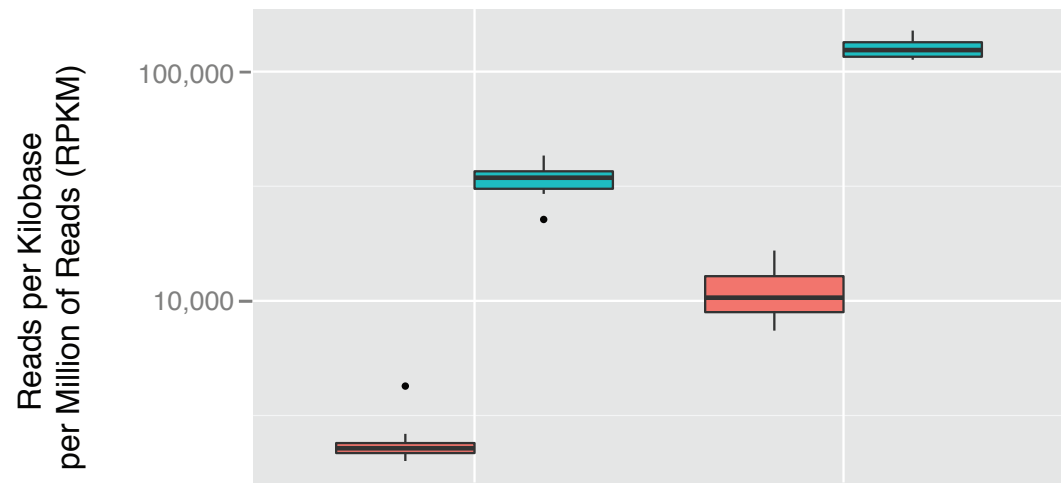
Supplementary Figure 9

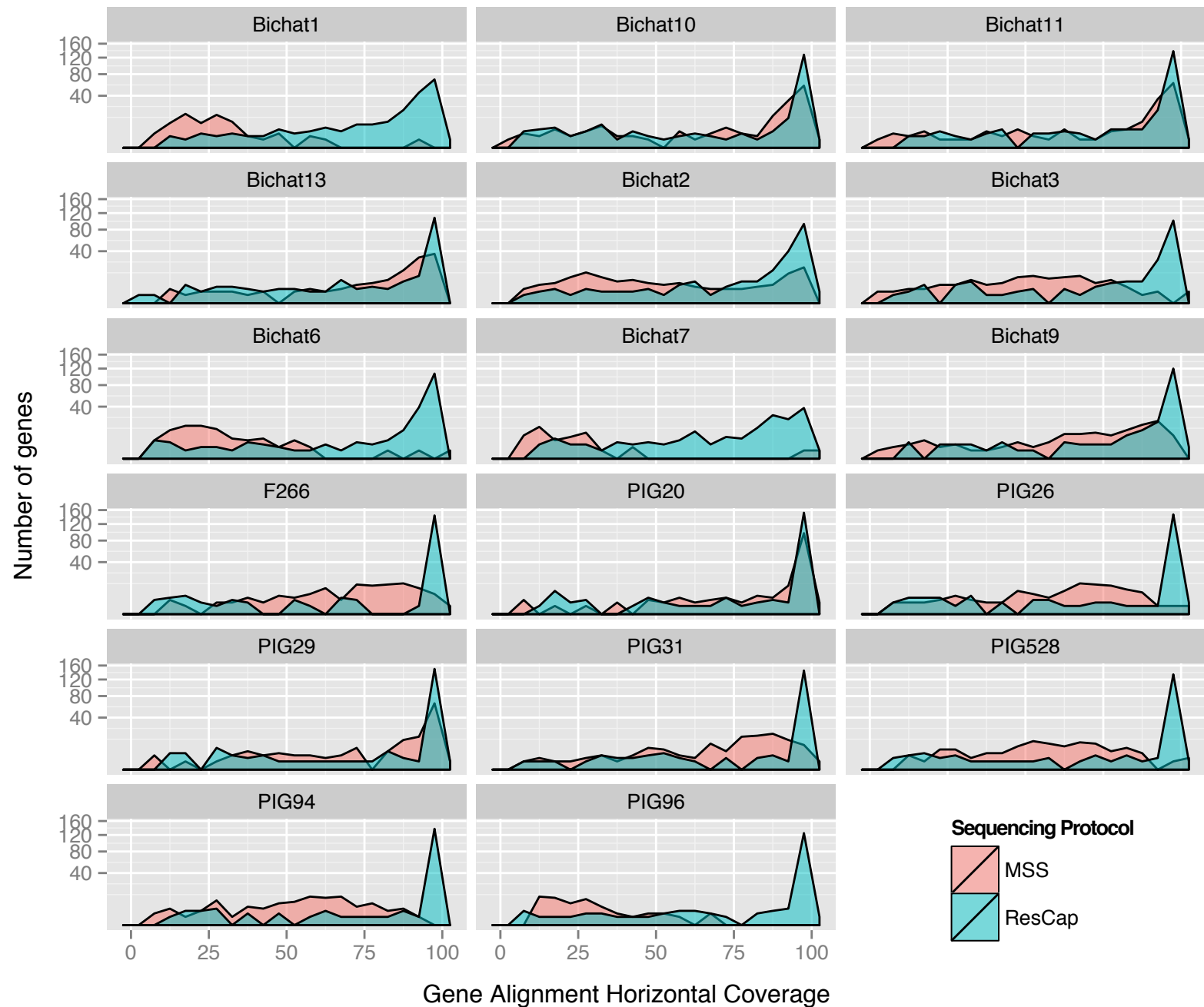
Histogram of gene coverage distribution by hybridizing probes. Two metrics was provided by NimbleGene, Direct Coverage (red bars) and Adjacent Coverage (cyan bars). 90% of the genes are covered at least 96.9% by direct coverage and 90% of the genes are covered at 100% of Adjacent Coverage.

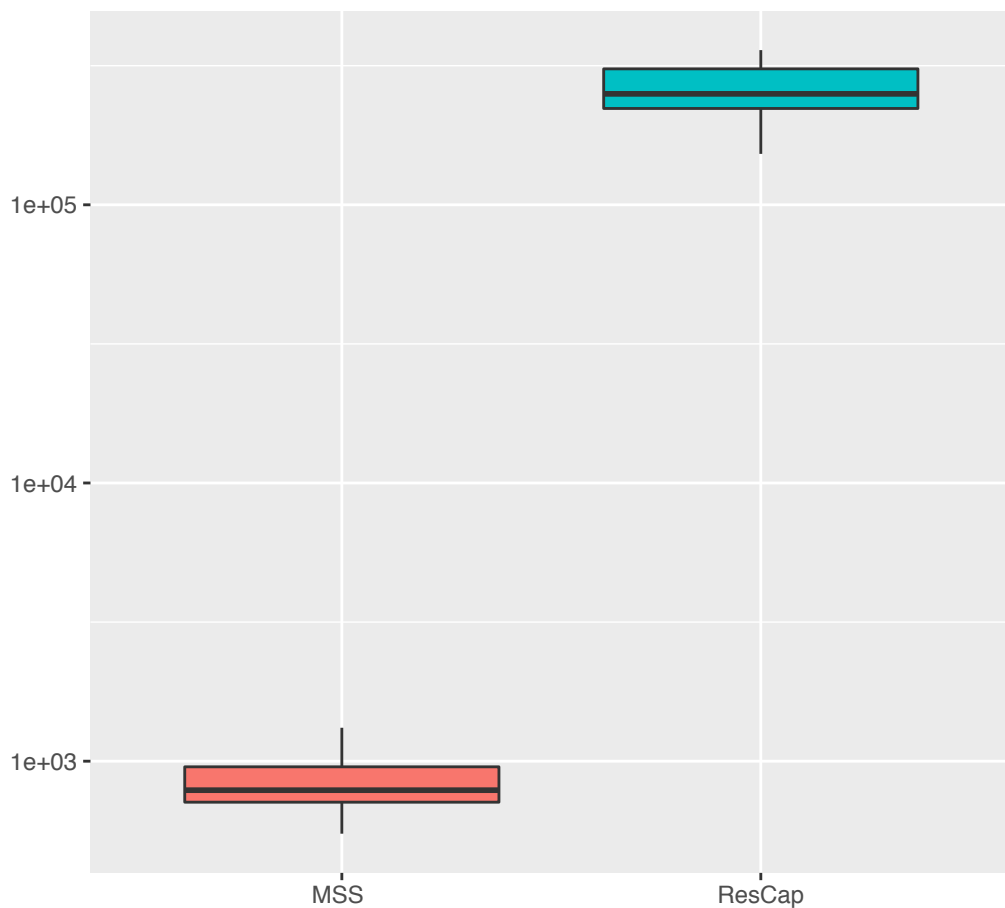
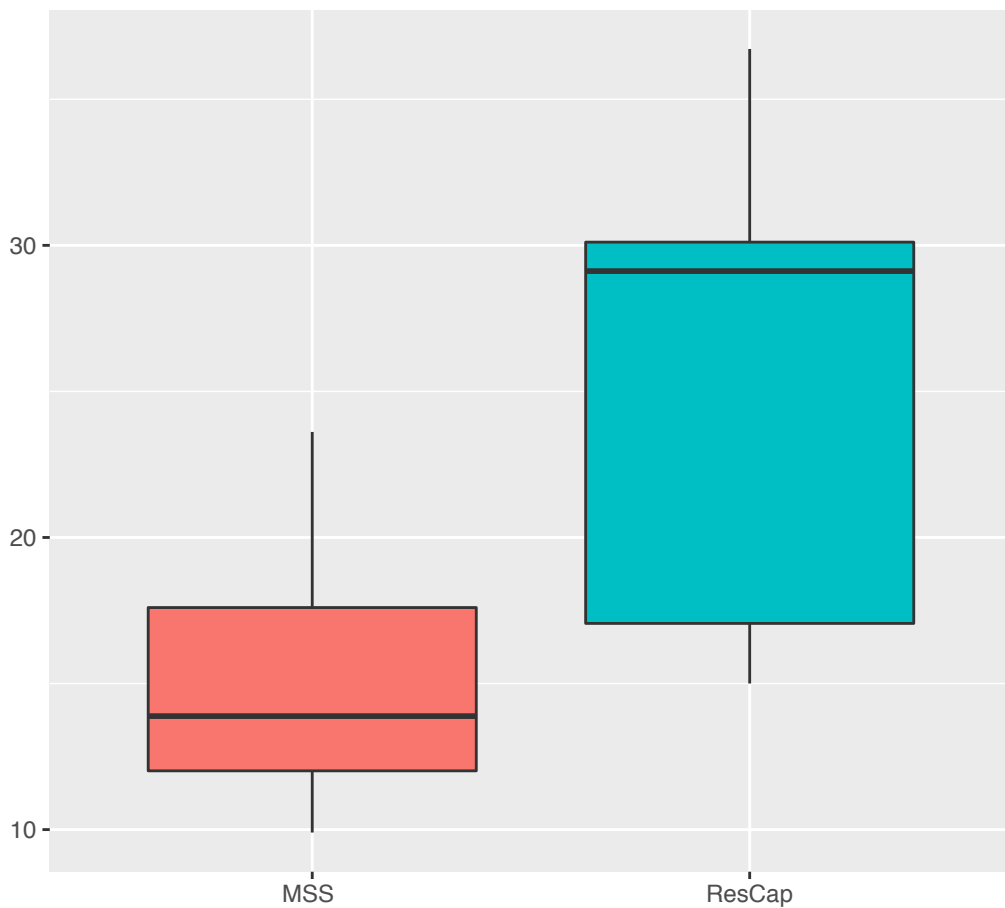
Supplementary Figure 10

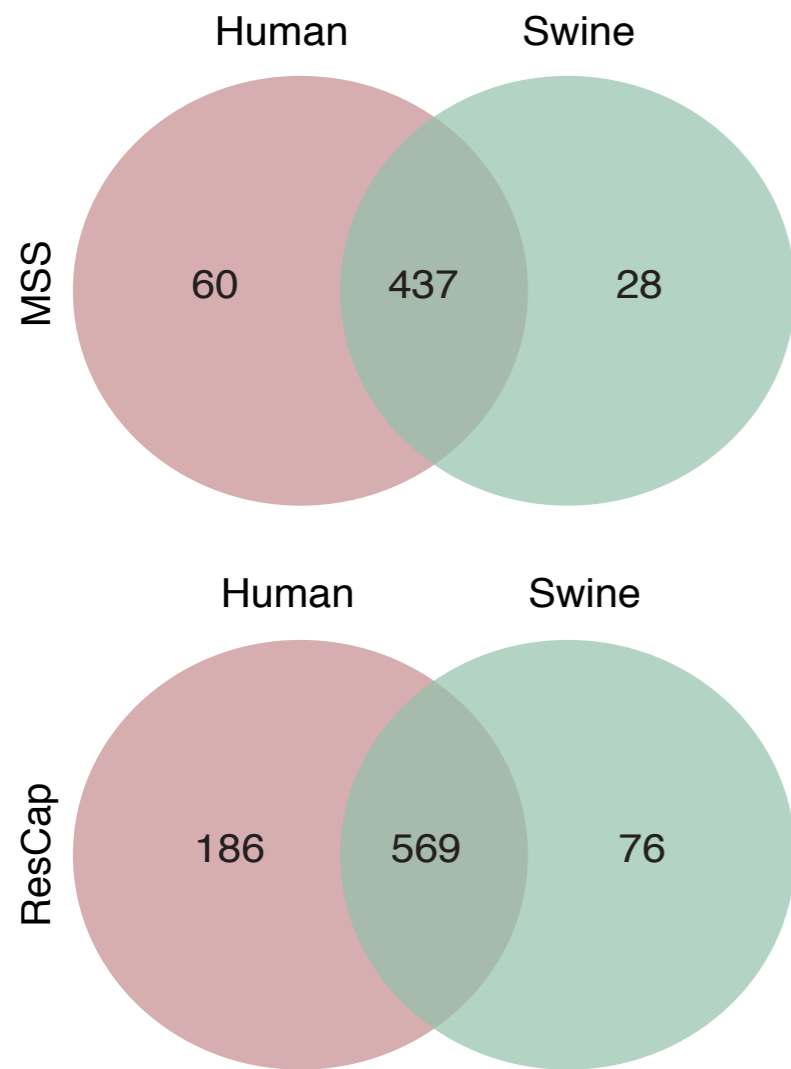
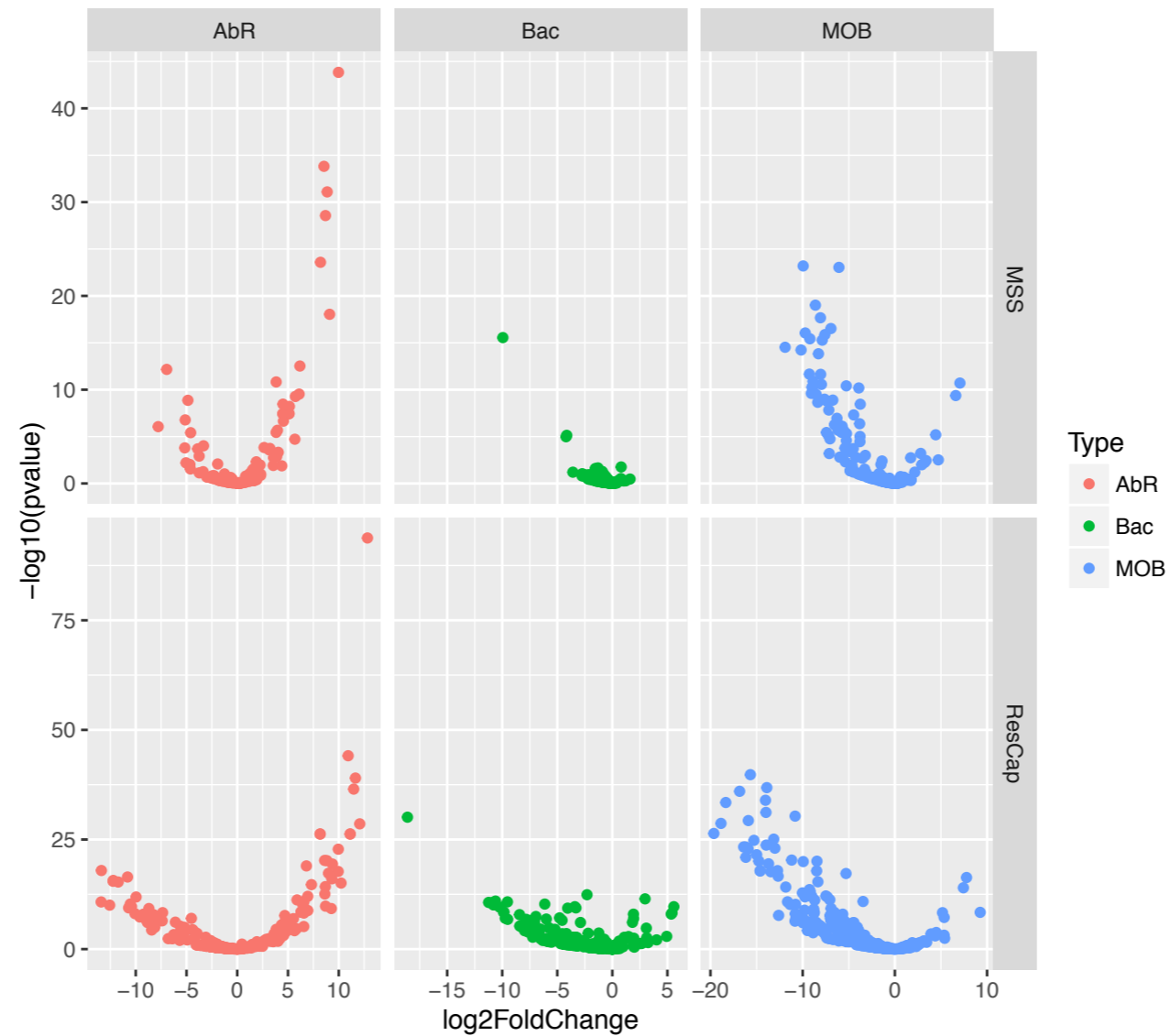
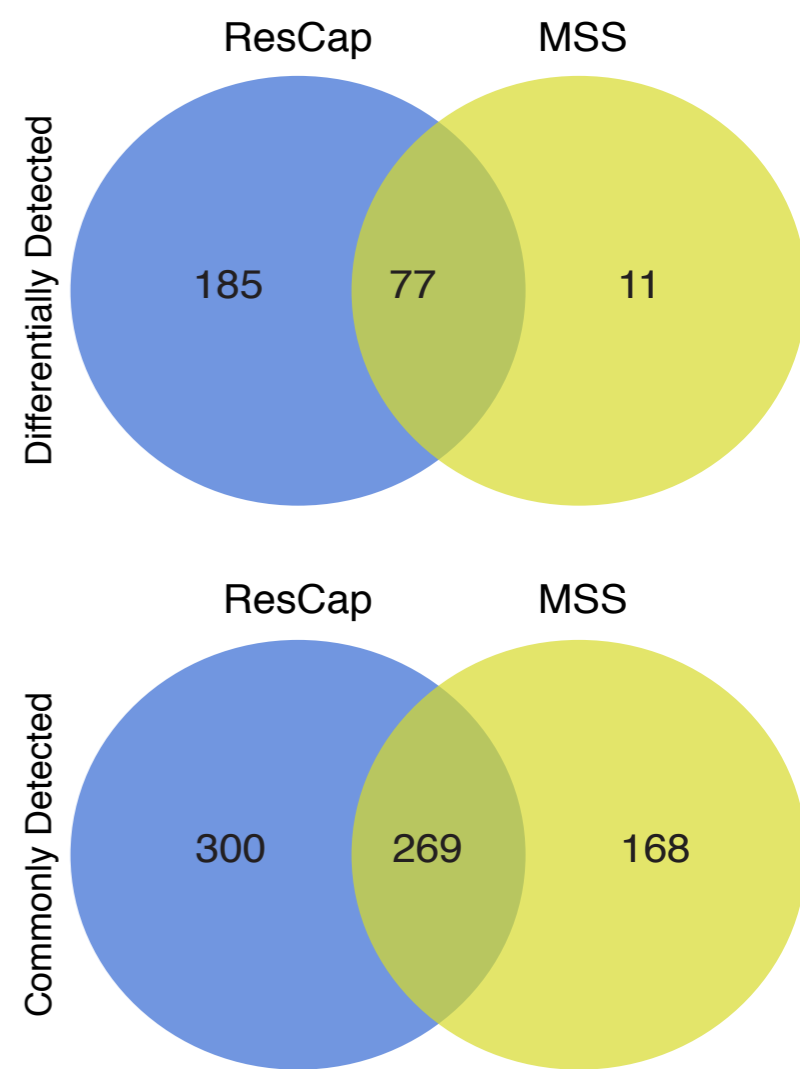
Allele Network: Nodes of the network represents individual genes that are mapped by some read. Edges between nodes represent reads that mapped on both nodes that link. Individual nodes are genes that are unequivocally identified. Gene clusters are mainly composed by different variants of the same gene (alleles). Mapping Gene Cluster (MGC) is defined using Markov cluster algorithm MCL.







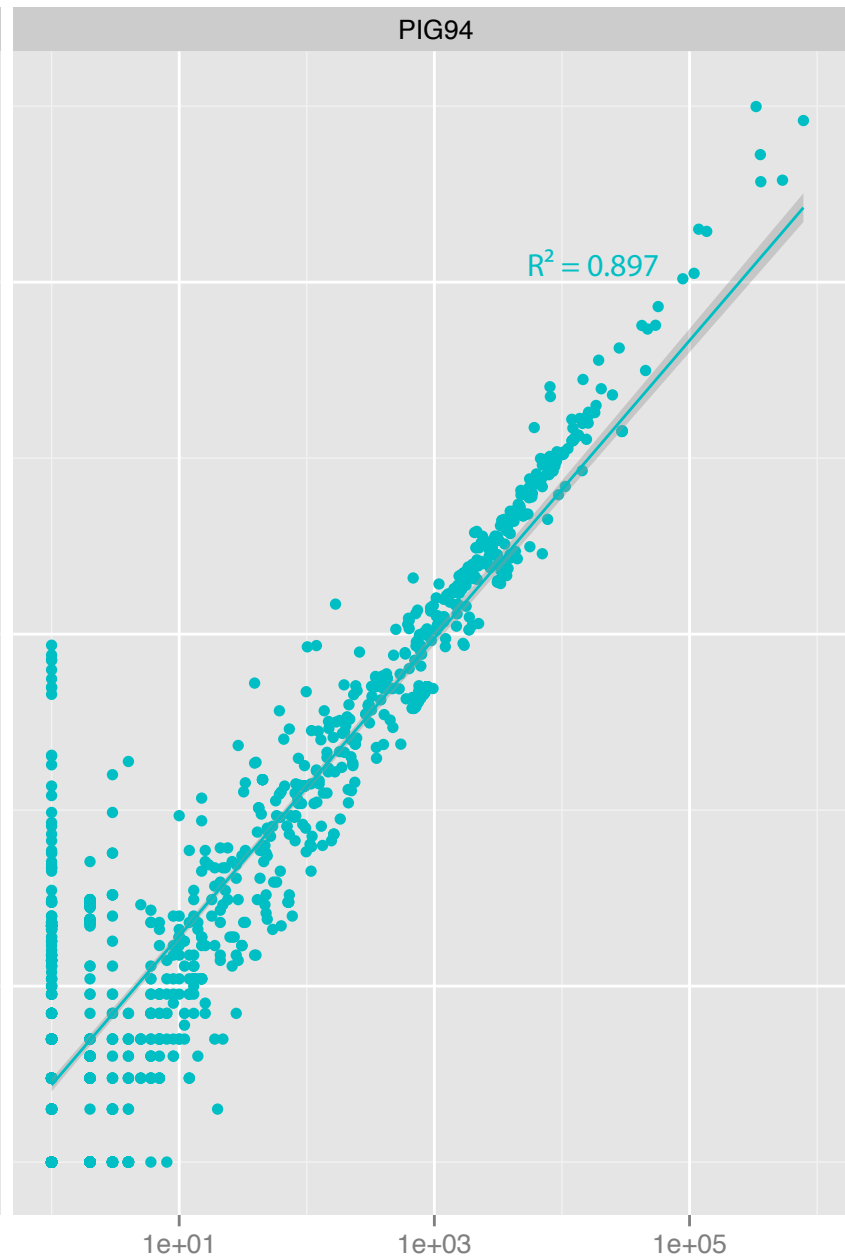
aReads Unequivocally Mapped
per Million of Reads**b**Number of detected genes by unequivocally
mapped reads per Million of Reads**Sequencing Protocol** MSS ResCap

a**b****c**

PIG28



PIG94

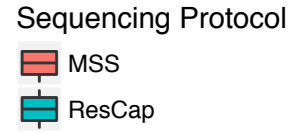
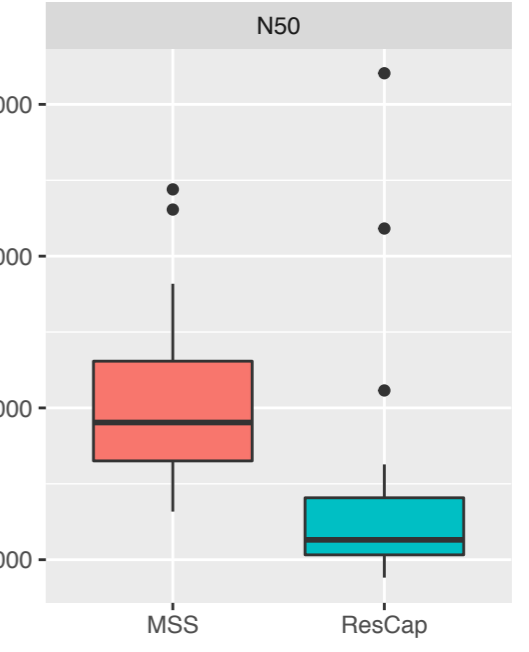
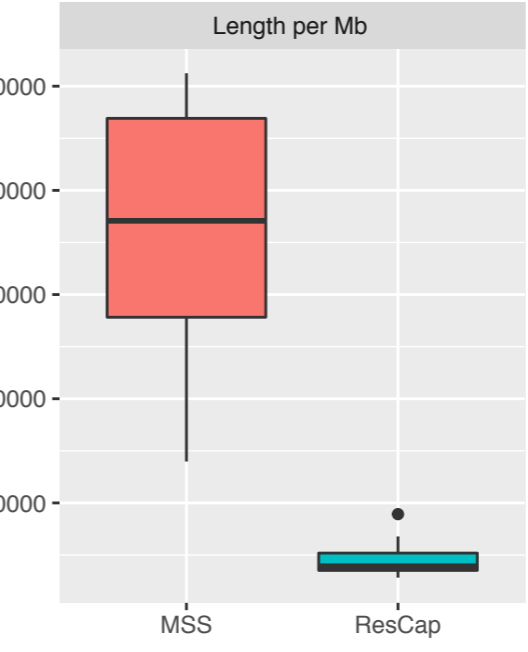
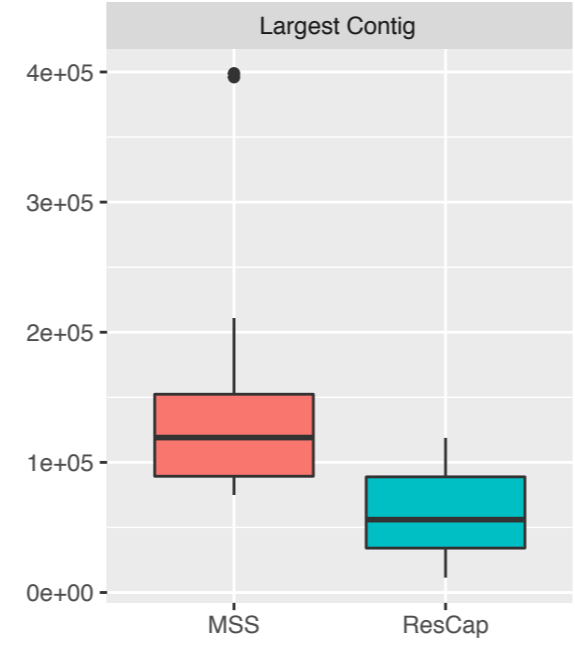
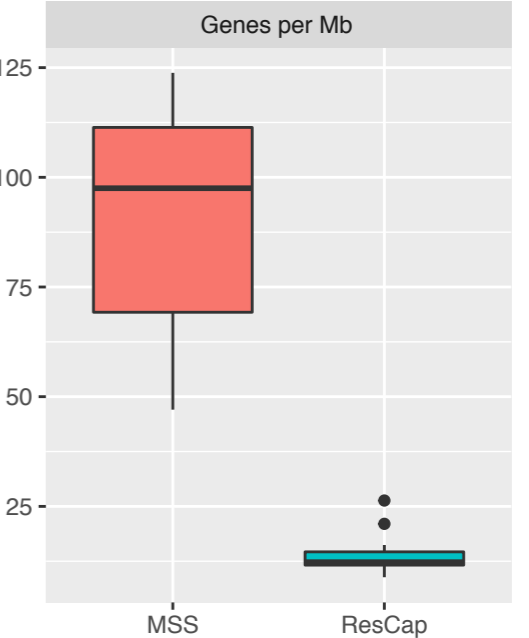
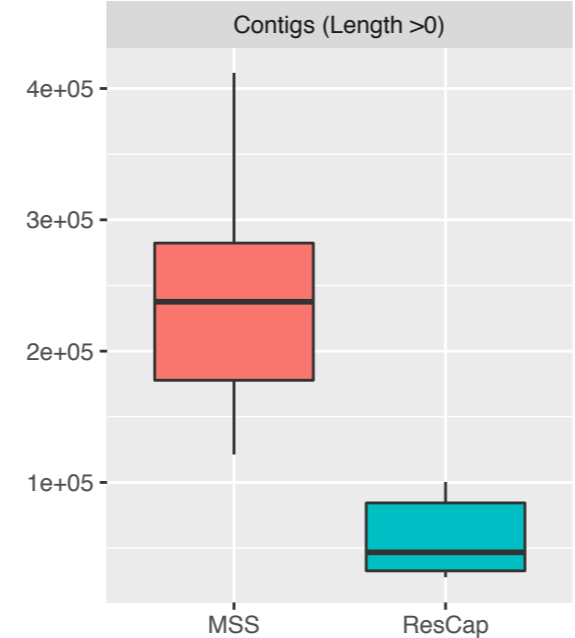
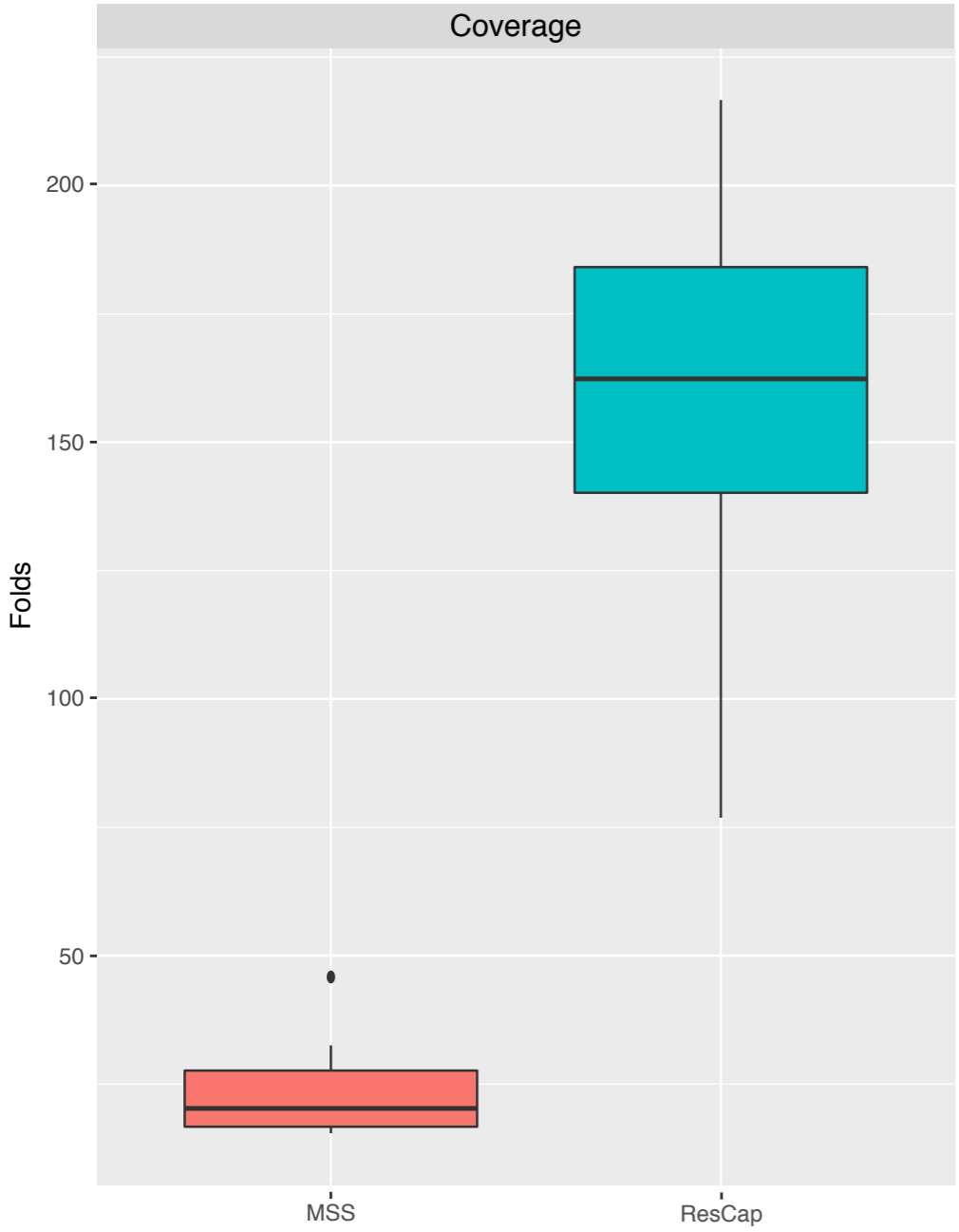


Sample

PIG28

PIG94

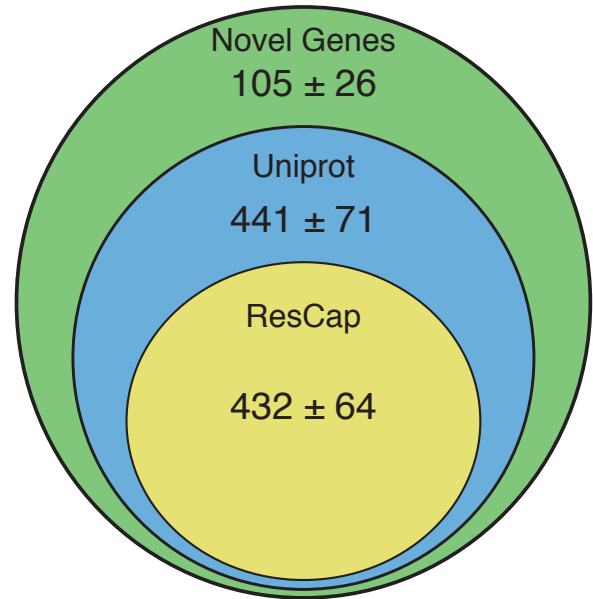
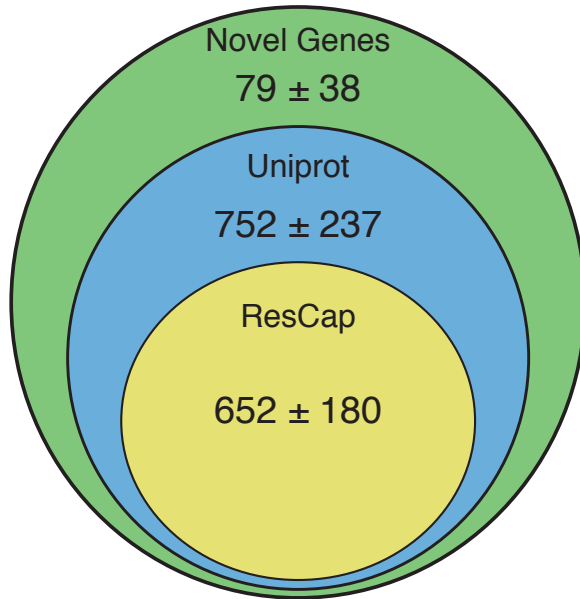
Number of reads in Replicate 2



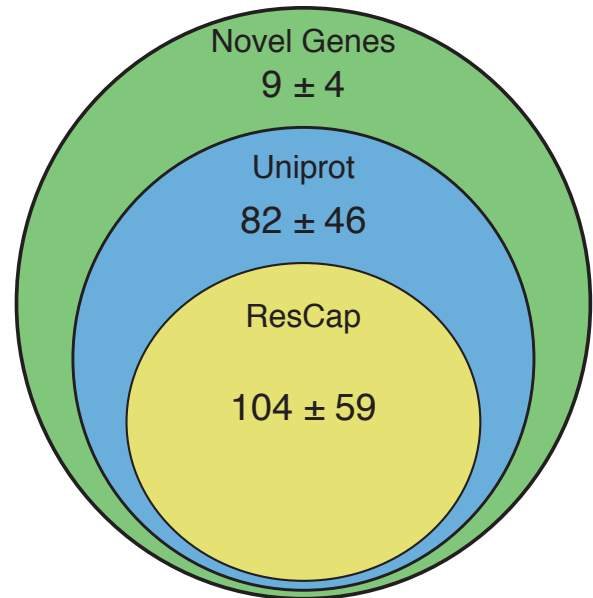
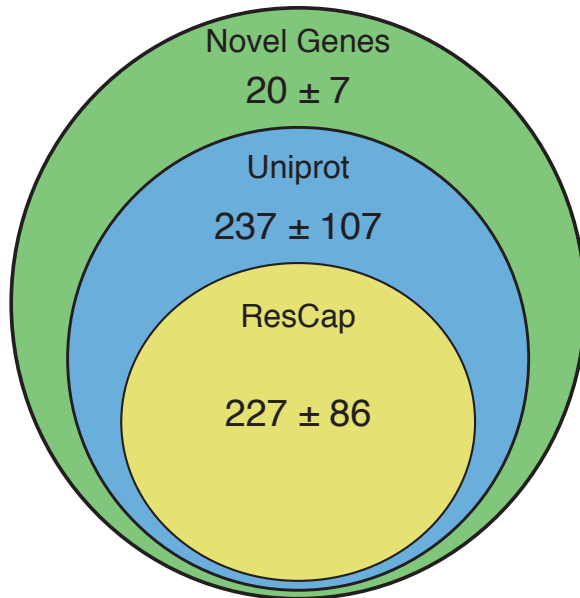
Human

Swine

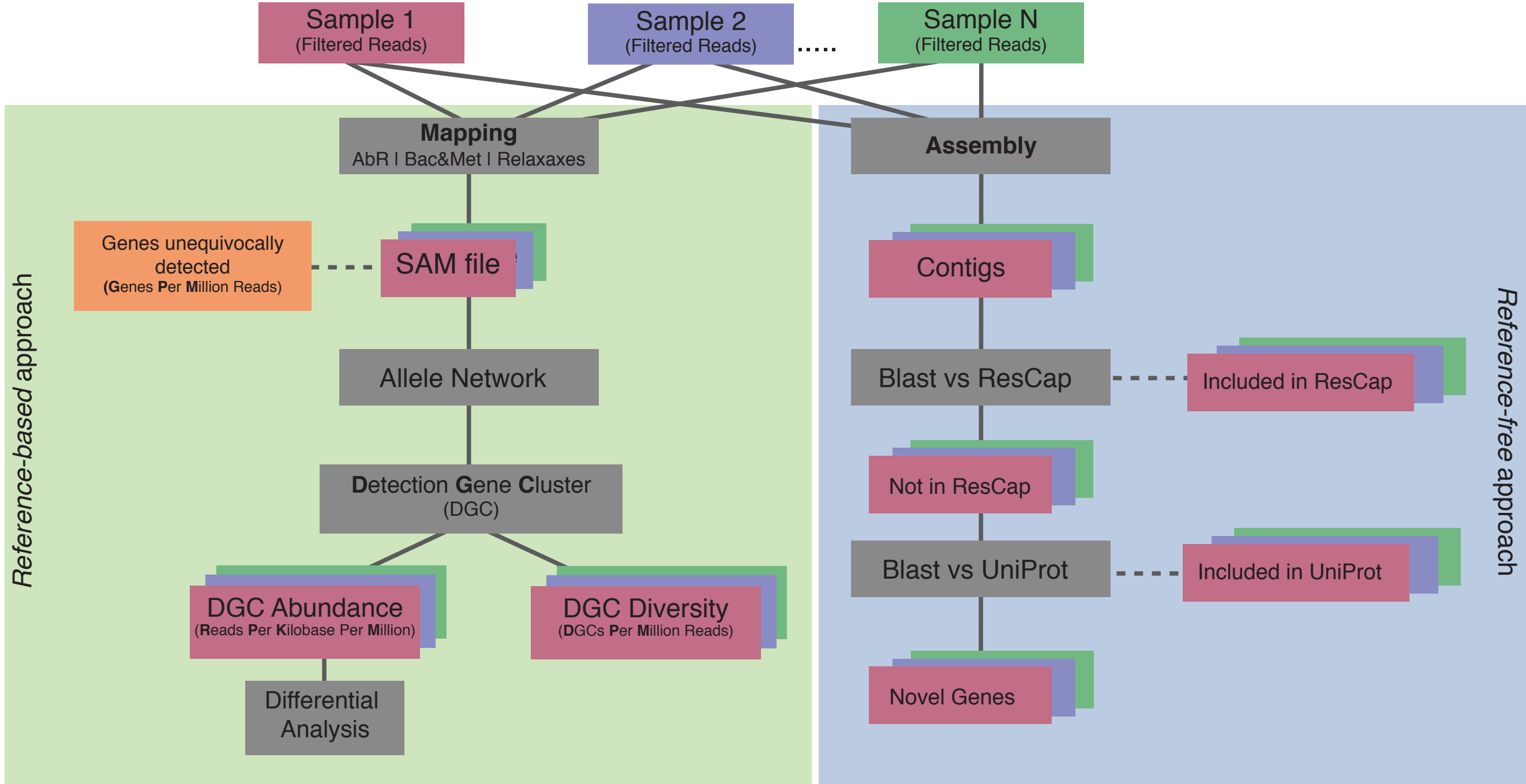
ResCap

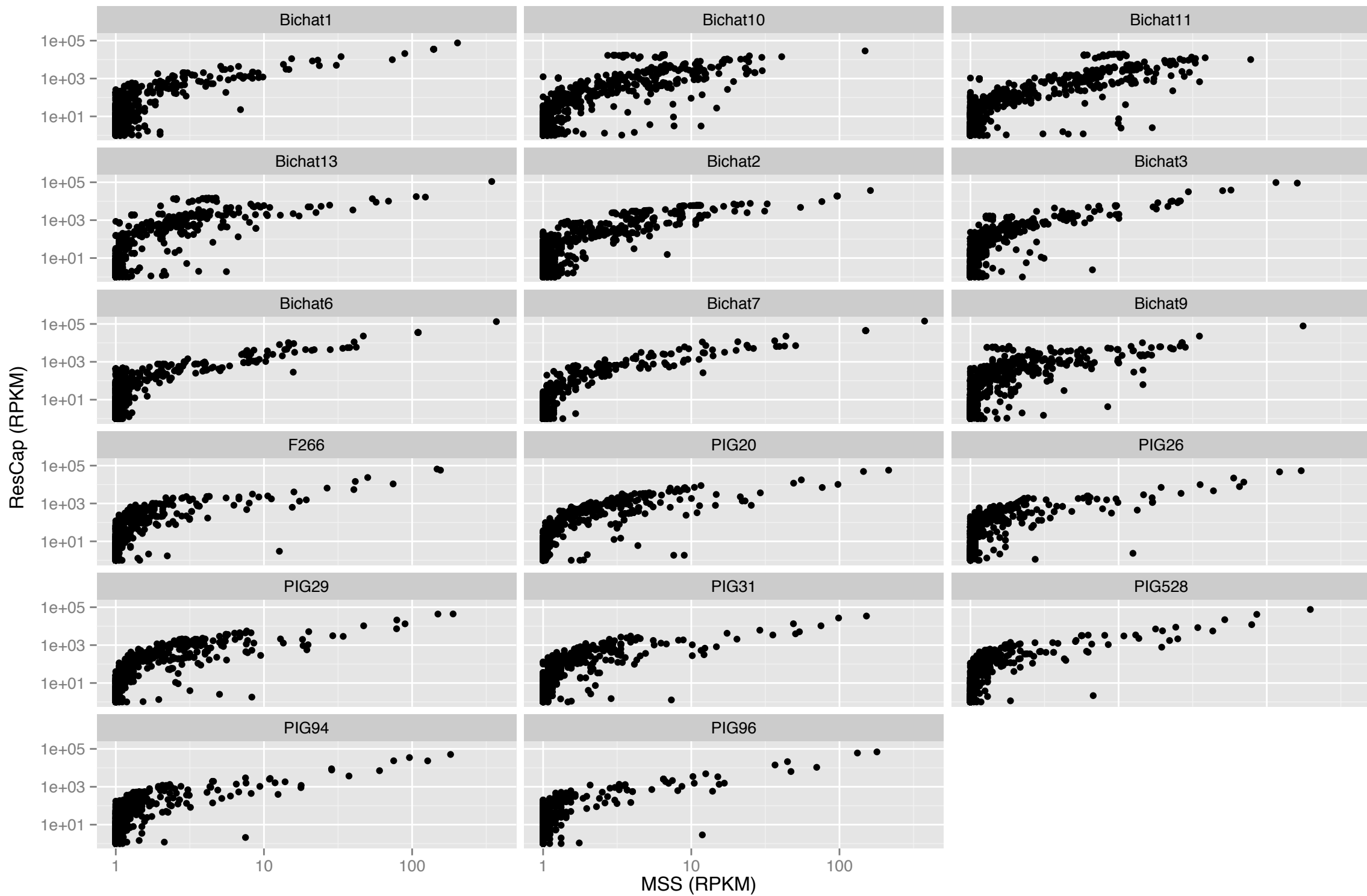


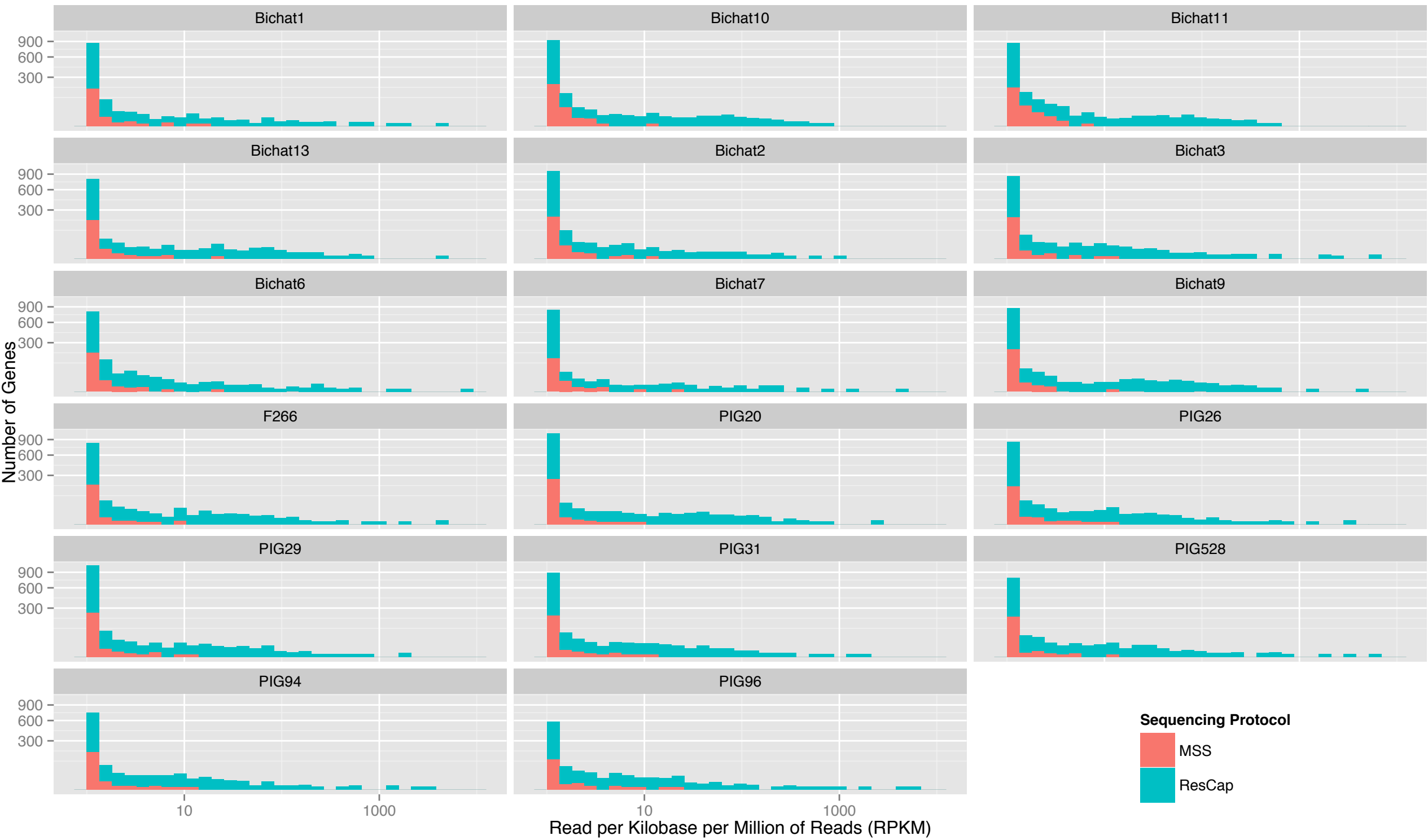
MSS



ResCap Analysis Workflow

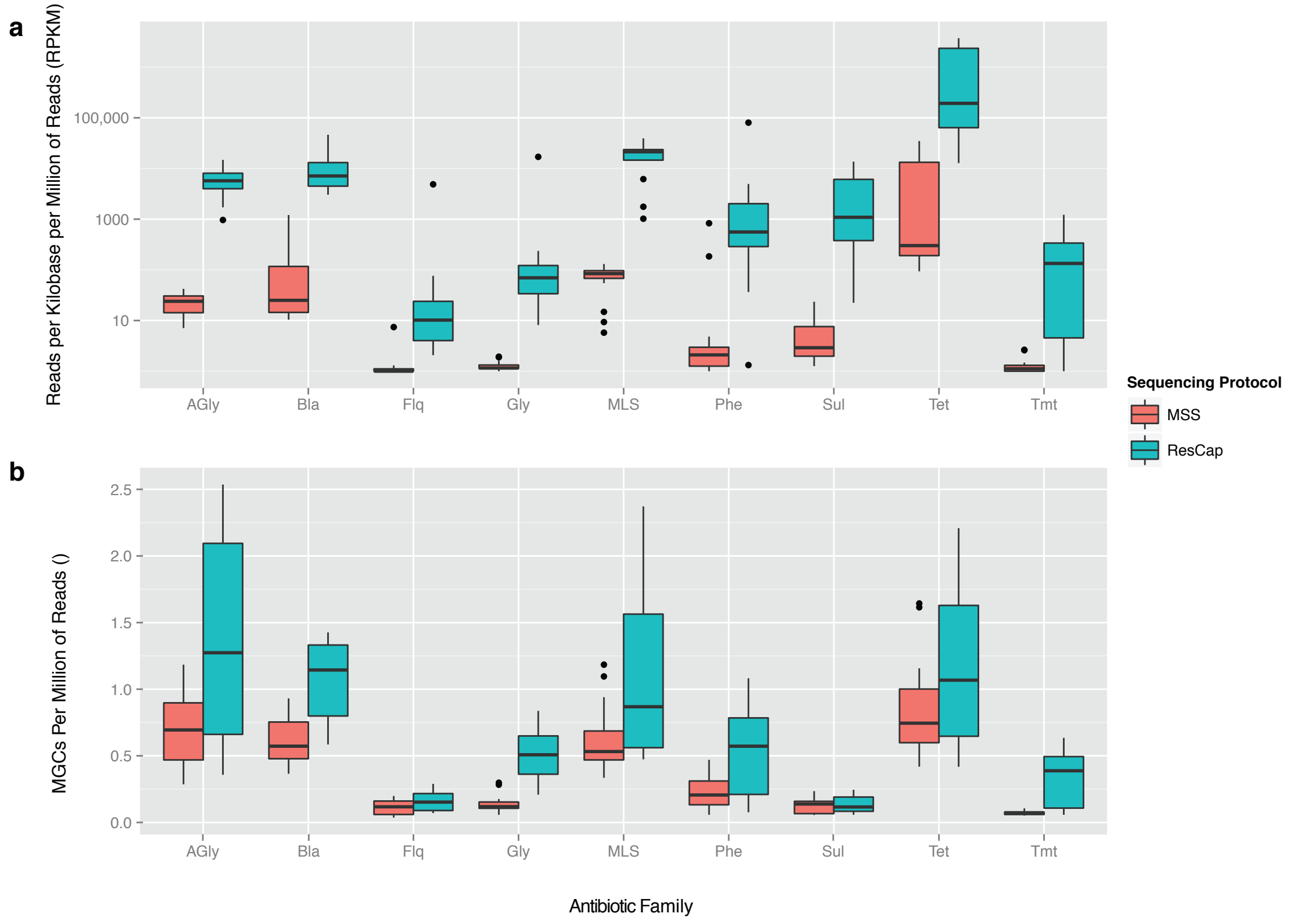






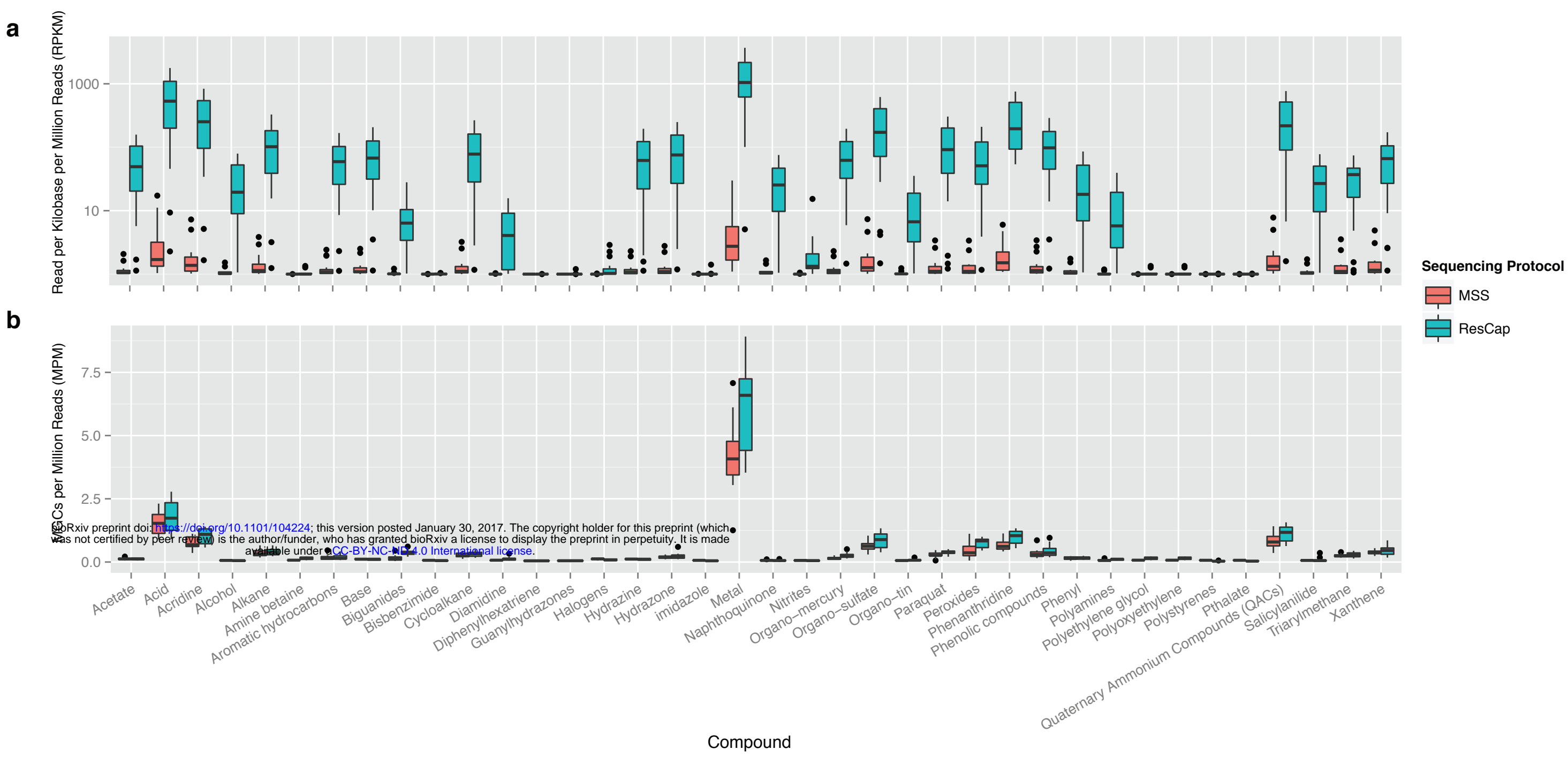
1

Abundance and Diversity of Antibiotic Resistance Genes



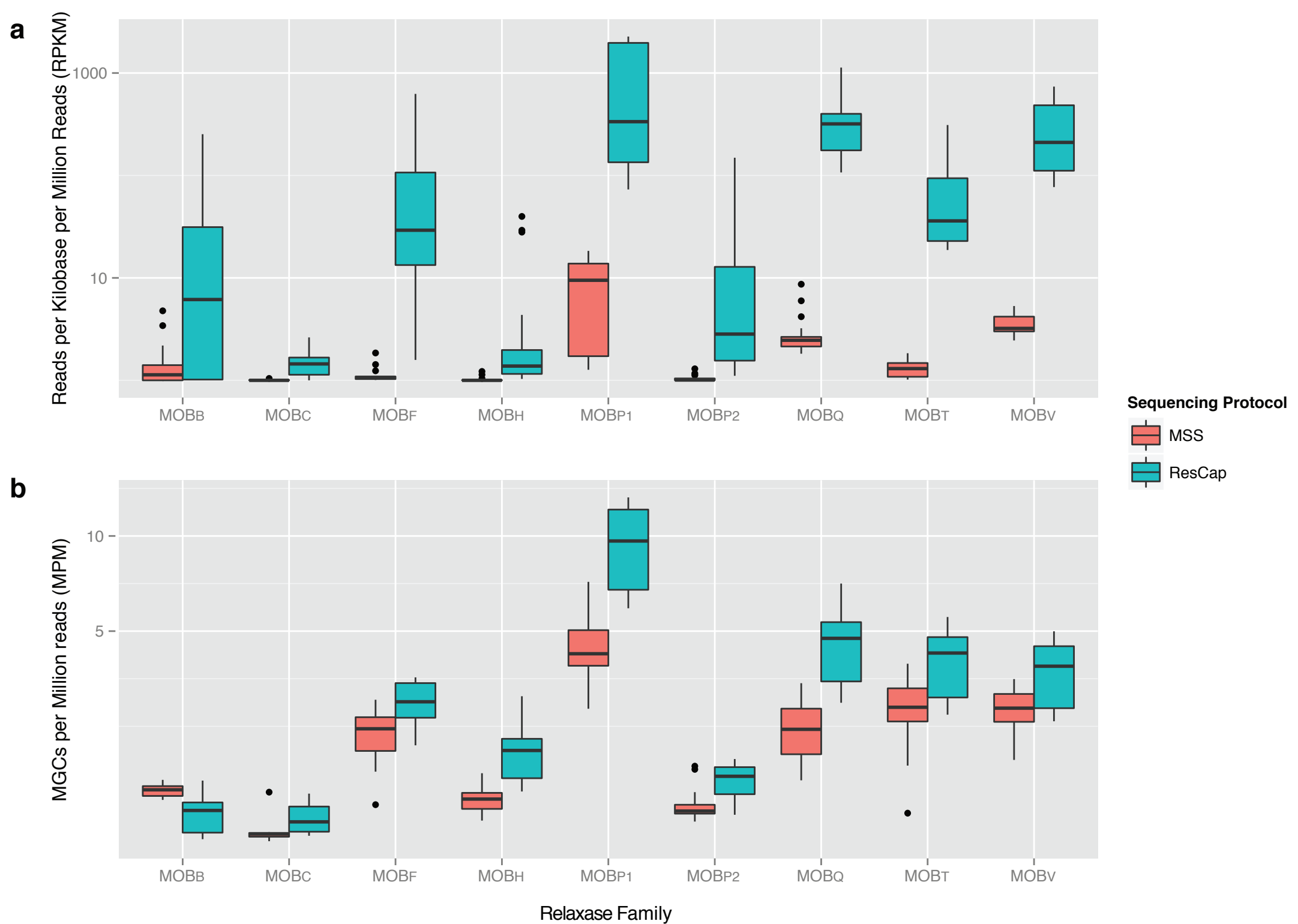
2

Abundance and Diversity of Metal & Biocides Resistance Genes



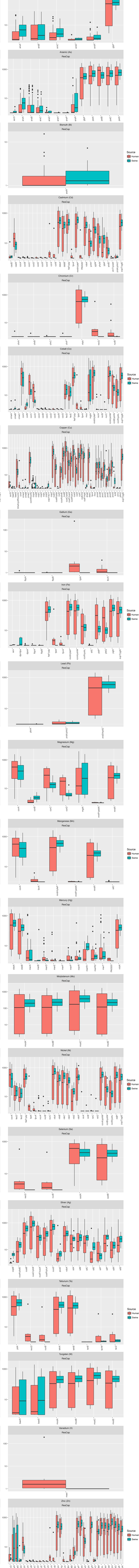
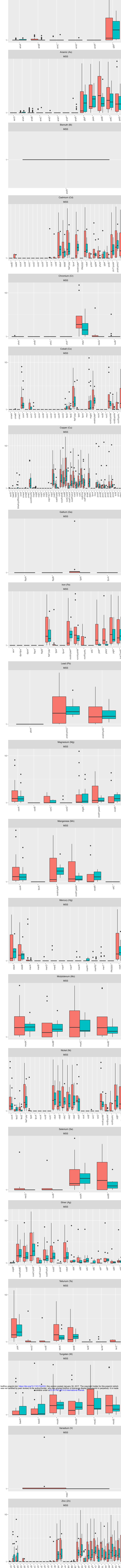
3

Abundance and Diversity of Antibiotic Relaxases





Reads per Kilobase per Million of Reads (RPKM)

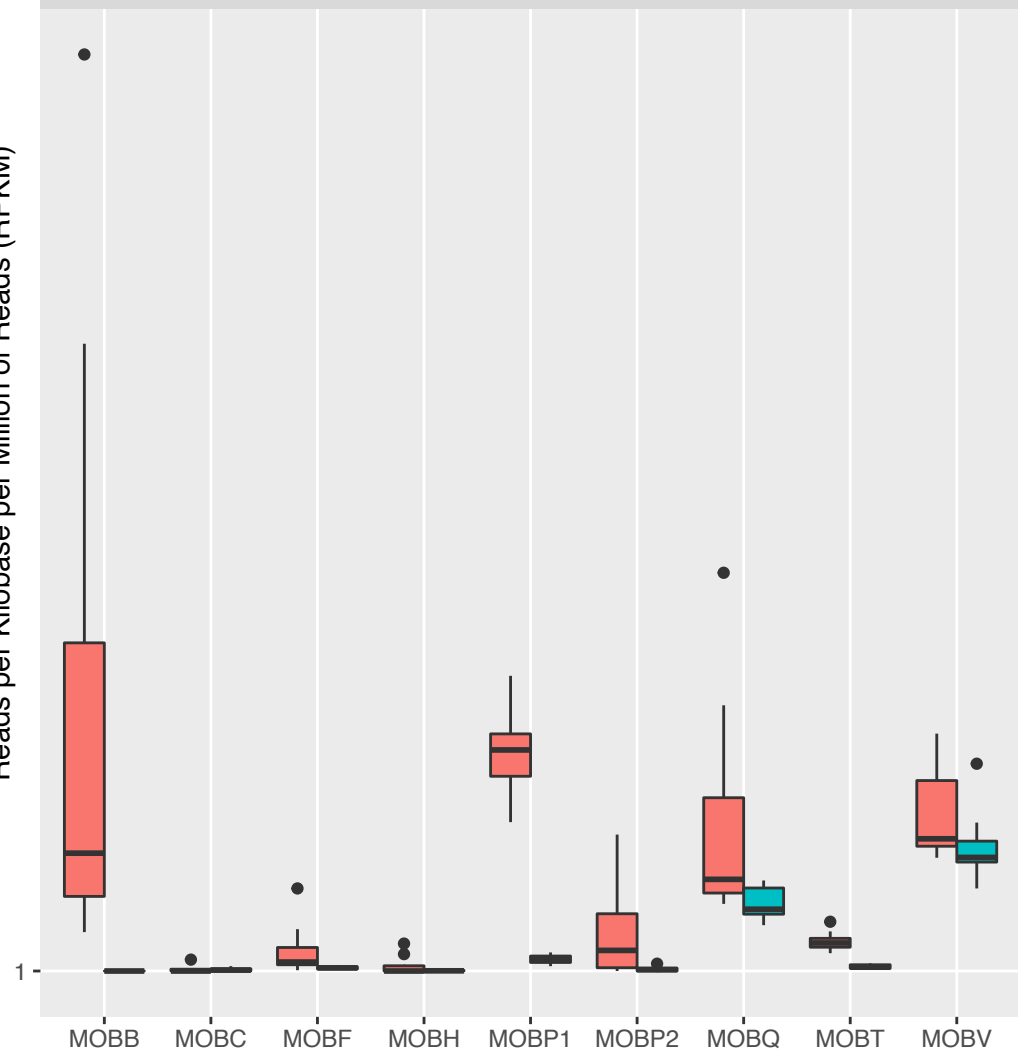


bioRxiv preprint doi: <https://doi.org/10.1101/100224>; this version posted January 30, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

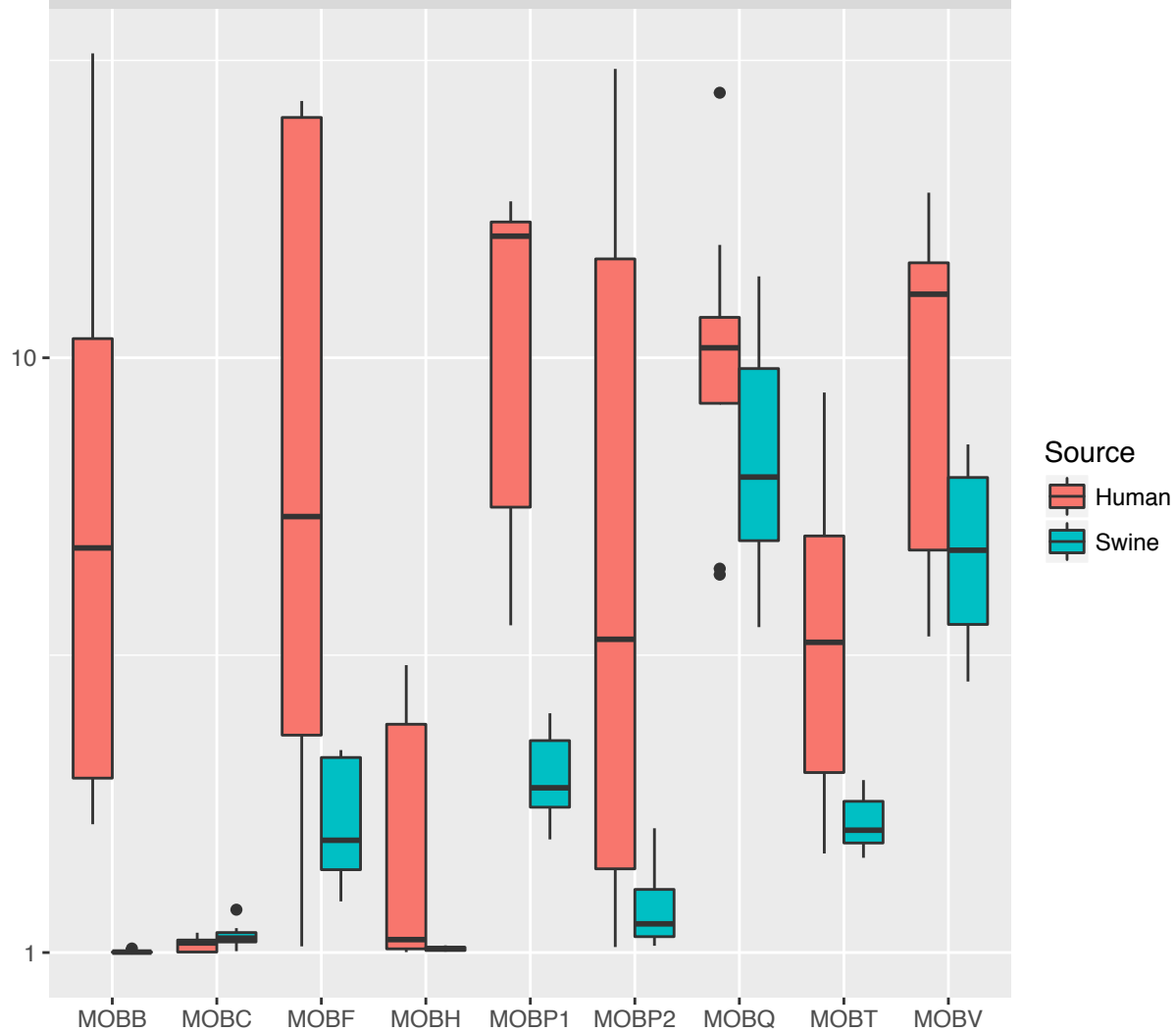
Gene Name

MSS

Reads per Kilobase per Million of Reads (RPKM)



ResCap

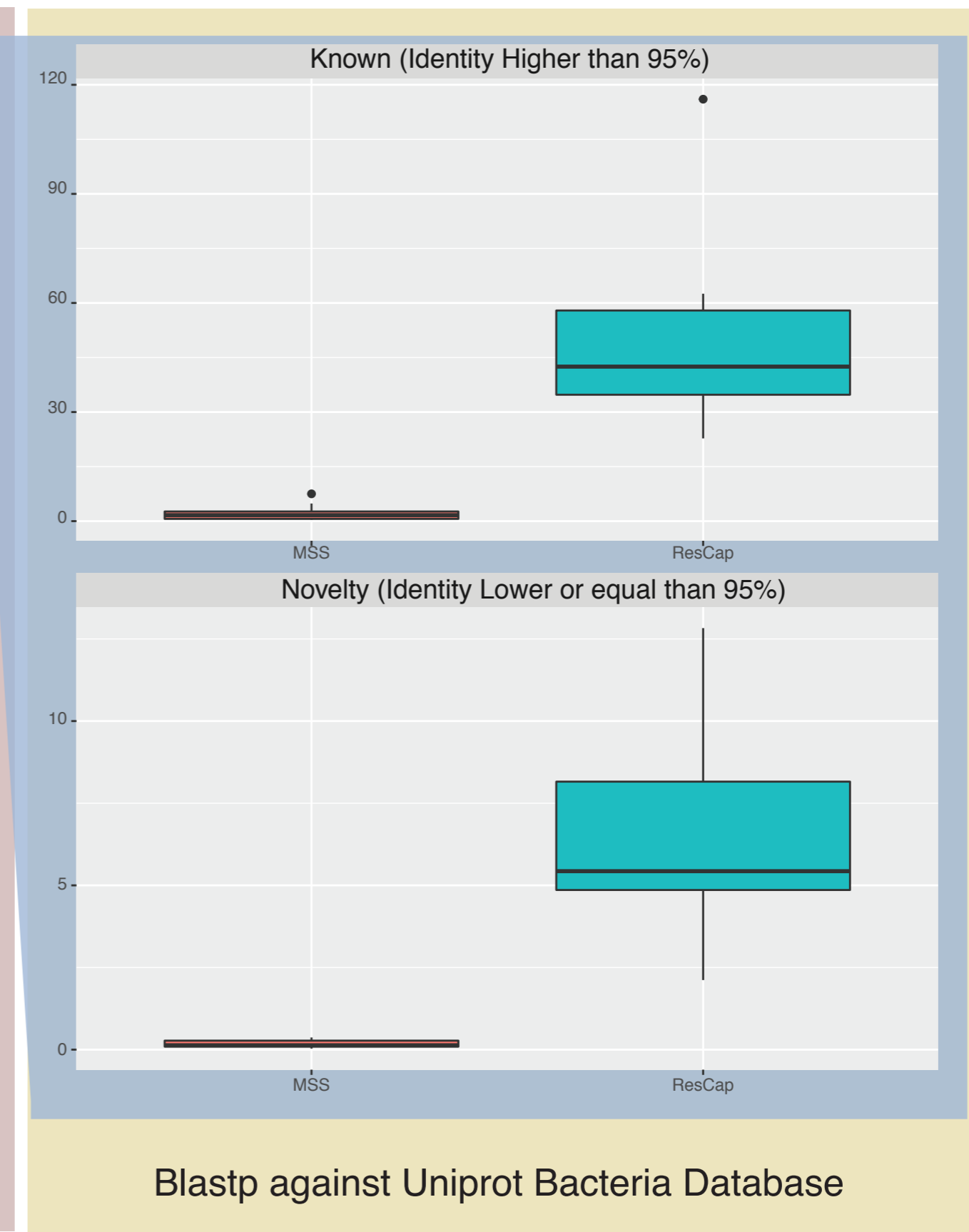
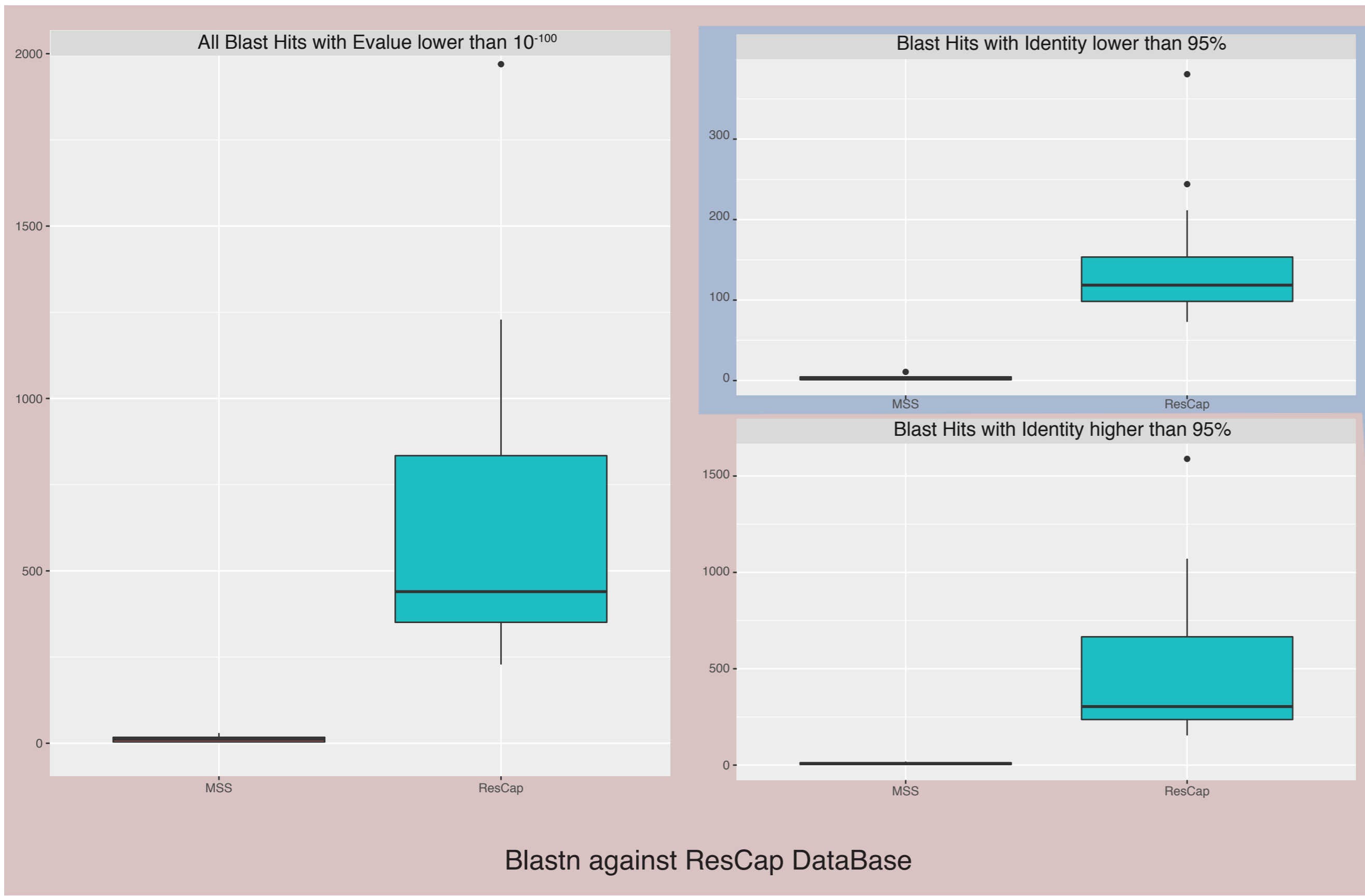


Source

Human

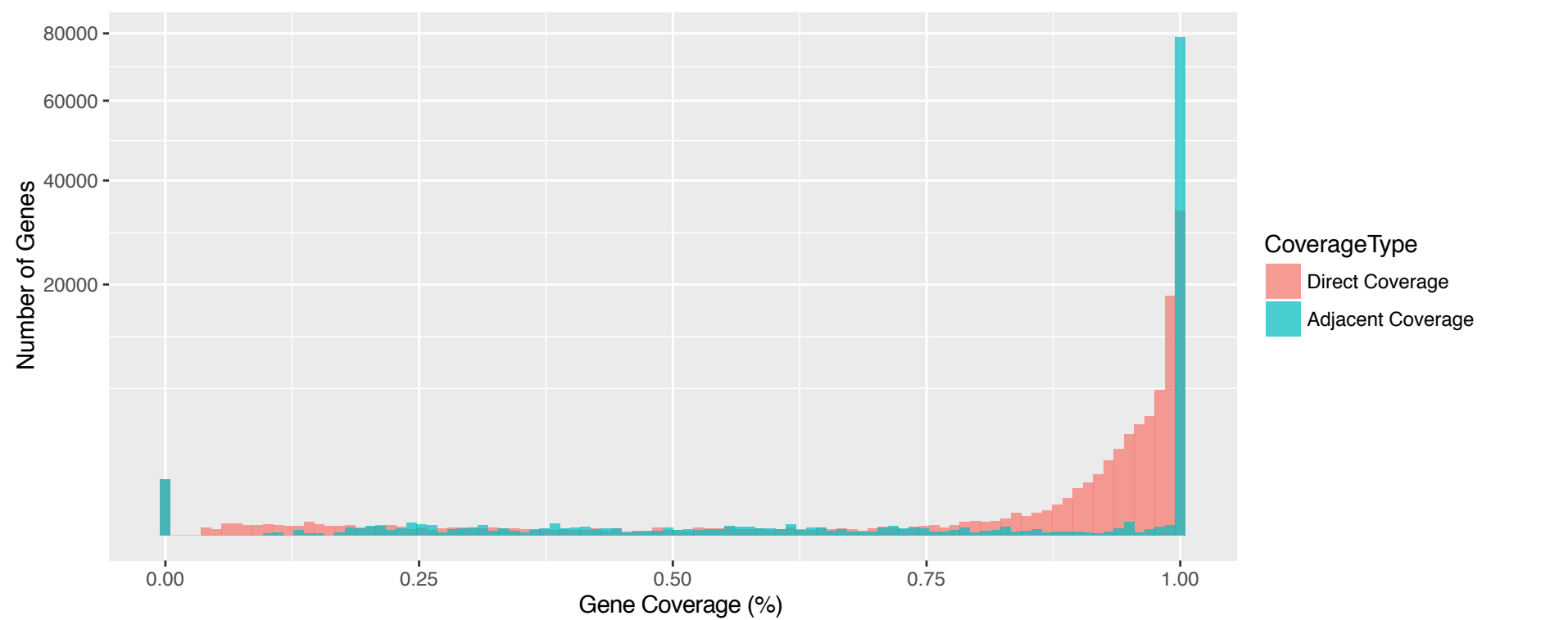
Swine

Number of Blast Hits per Genes per Megabase

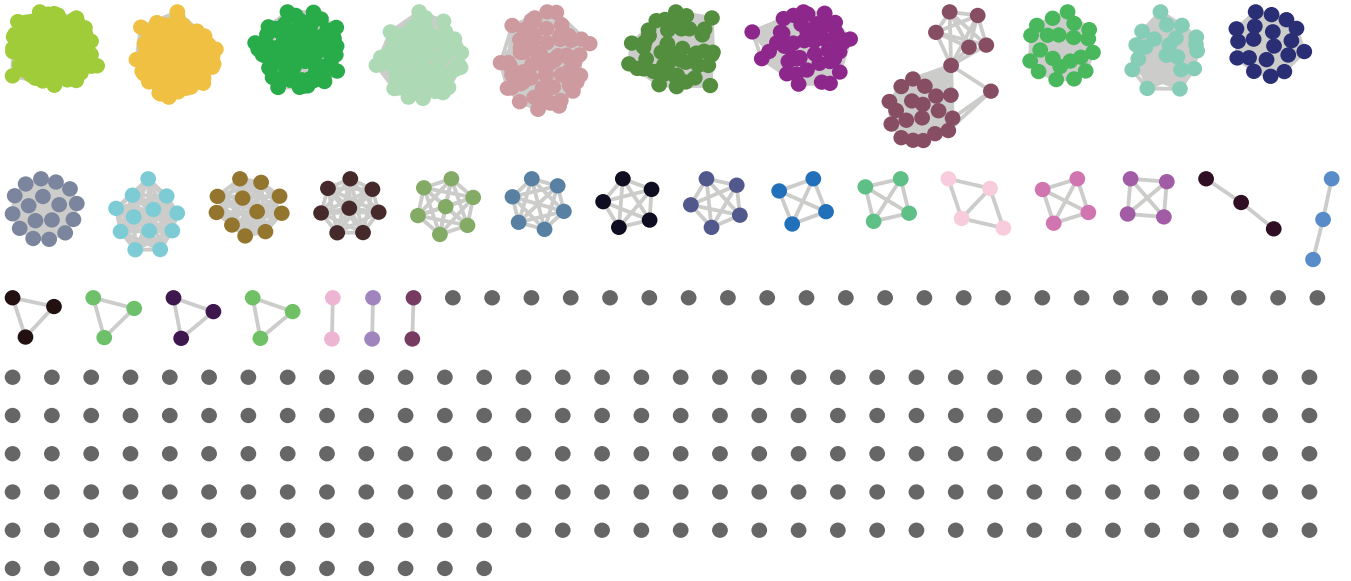


Sequencing Protocol

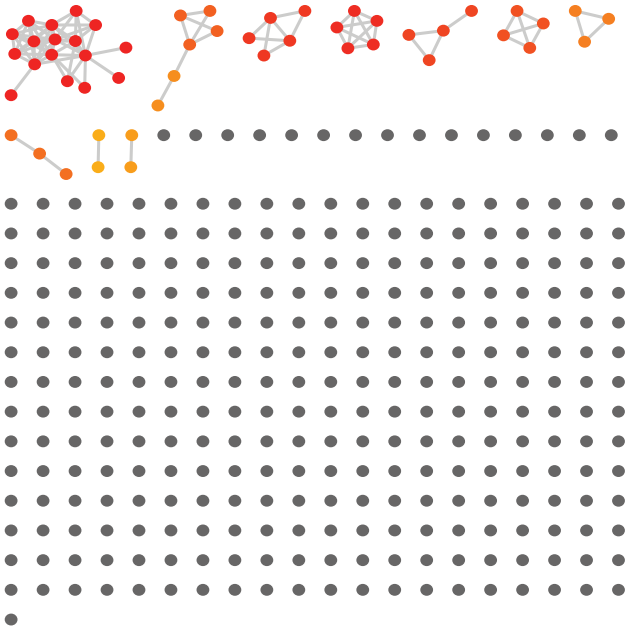
- MSS
- ResCap



Antibiotic Resistance



Biocide & Metal



Relaxases

