

# 16S rRNA amplicon sequencing for epidemiological surveys of bacteria in wildlife: the importance of cleaning post-sequencing data before estimating positivity, prevalence and co-infection

*Maxime Galan*<sup>1</sup>, *Maria Razzauti*<sup>1</sup>, *Emilie Bard*<sup>2</sup>, *Maria Bernard*<sup>3,4</sup>, *Carine Brouat*<sup>5</sup>, *Nathalie Charbonnel*<sup>1</sup>, *Alexandre Dehne-Garcia*<sup>1</sup>, *Anne Loiseau*<sup>1</sup>, *Caroline Tatard*<sup>1</sup>, *Lucie Tamisier*<sup>1</sup>, *Muriel Vayssier-Taussat*<sup>6</sup>, *Hélène Vignes*<sup>7</sup>, *Jean-François Cosson*<sup>1,6</sup>

1: INRA, CBGP, Montferrier sur Lez, France

2: INRA, EpiA, Clermont-Ferrand, France

3: INRA, Sigenae, France

4 : INRA, GABI, AgroParisTech, Université Paris-Saclay, Jouy-en-Josas, France

5: Ird, CBGP, Montferrier sur Lez, France

6: INRA, Bipar, Maisons-Alfort, France

7: CIRAD, AGAP, Montpellier, France

**Corresponding authors:** [galan@supagro.inra.fr](mailto:galan@supagro.inra.fr); [cosson@supagro.inra.fr](mailto:cosson@supagro.inra.fr)

## Importance

Several recent public health crises have shown that the surveillance of zoonotic agents in wildlife is important to prevent pandemic risks. Rodents are intermediate hosts for numerous zoonotic bacteria. High-throughput sequencing (HTS) technologies are very useful for the detection and surveillance of zoonotic bacteria, but rigorous experimental processes are required for the use of these cheap and effective tools in such epidemiological contexts. In particular, HTS introduces biases into the raw dataset that might lead to incorrect interpretations. We describe here a procedure for cleaning data before estimating reliable biological parameters, such as bacterial positivity, prevalence and coinfection, by 16S rRNA amplicon sequencing on the MiSeq platform. This procedure, applied to 711 commensal rodents collected from 24 villages in Senegal, Africa, detected several emerging bacterial genera, some in high prevalence, while never before reported for West Africa. This study constitutes a step towards the use of HTS to improve our understanding of the risk of zoonotic disease transmission posed by wildlife, by providing a new strategy for the

33 use of HTS platforms to monitor both bacterial diversity and infection dynamics in  
34 wildlife. In the future, this approach could be adapted for the monitoring of other  
35 microbes such as protists, fungi, and even viruses.

36

## 37 Summary

38 Human impact on natural habitats is increasing the complexity of human-wildlife  
39 interfaces and leading to the emergence of infectious diseases worldwide. Highly  
40 successful synanthropic wildlife species, such as rodents, will undoubtedly play an  
41 increasingly important role in transmitting zoonotic diseases. We investigated the  
42 potential of recent developments in 16S rRNA amplicon sequencing to facilitate the  
43 multiplexing of large numbers of samples, to improve our understanding of the risk of  
44 zoonotic disease transmission posed by urban rodents in West Africa. In addition to  
45 listing pathogenic bacteria in wild populations, as in other high-throughput  
46 sequencing (HTS) studies, our approach can estimate essential parameters for  
47 studies of zoonotic risk, such as prevalence and patterns of coinfection within  
48 individual hosts. However, the estimation of these parameters requires cleaning of  
49 the raw data to eliminate the biases generated by HTS methods. We present here an  
50 extensive review of these biases and of their consequences, and we propose a  
51 trimming strategy for managing them and cleaning the dataset. We also analyzed  
52 711 commensal rodents collected from 24 villages in Senegal, including 208 *Mus*  
53 *musculus domesticus*, 189 *Rattus rattus*, 93 *Mastomys natalensis* and 221 *Mastomys*  
54 *erythroleucus*. Seven major genera of pathogenic bacteria were detected: *Borrelia*,  
55 *Bartonella*, *Mycoplasma*, *Ehrlichia*, *Rickettsia*, *Streptobacillus* and *Orientia*. The last  
56 five of these genera have never before been detected in West African rodents.  
57 Bacterial prevalence ranged from 0% to 90%, depending on the bacterial taxon,  
58 rodent species and site considered, and a mean of 26% of rodents displayed  
59 coinfection. The 16S rRNA amplicon sequencing strategy presented here has the  
60 advantage over other molecular surveillance tools of dealing with a large spectrum of  
61 bacterial pathogens without requiring assumptions about their presence in the  
62 samples. This approach is, thus, particularly suitable for continuous pathogen  
63 surveillance in the framework of disease monitoring programs

64

## 65 Introduction

66 Pathogen monitoring in wildlife is a key method for preventing the emergence of  
67 infectious diseases in humans and domestic animals. More than half the pathogens  
68 causing disease in humans originate from animal species [1]. The early identification  
69 of zoonotic agents in animal populations is therefore of considerable human health  
70 interest. Wildlife species may also act as a reservoir for pathogens capable of  
71 infecting livestock, with significant economic consequences [2]. The monitoring of  
72 emerging diseases in natural populations is also important for preserving biodiversity,  
73 because pathogens carried by invasive species may cause the decline of endemic  
74 species [3]. There is, therefore, a need to develop screening tools for identifying a  
75 broad range of pathogens in samples consisting of large numbers of individual hosts  
76 or vectors.

77 High-throughput sequencing (HTS) approaches require no prior assumptions about  
78 the bacterial communities present in samples of diverse natures, including non-  
79 cultivable bacteria. Such metagenomics approaches are based on the sequencing of  
80 all (WGS: whole-genome sequencing) or some (RNAseq or 16S rRNA amplicon  
81 sequencing) of the bacterial DNA or RNA in a sample, with the sequences obtained  
82 then compared with those in a reference sequence database [4]. Metagenomics has  
83 made a major contribution to the generation of comprehensive inventories of the  
84 bacteria, including pathogens, present in humans [5]. Such approaches are now  
85 being extended to the characterization of bacteria in wildlife [6-13, 90]. However,  
86 improvements in the estimation of infectious risks will require more than just the  
87 detection of bacterial pathogens. Indeed, we will also need to estimate the  
88 prevalence of these pathogens by host taxon and/or environmental features, together  
89 with coinfection rates [14,15] and pathogen interactions [16,17].

90 Razzauti *et al.* [8] recently used 16S rRNA amplicon sequencing with the dual-index  
91 sequencing strategy of Kozich *et al.* [18] to detect bacterial pathogens in very large  
92 numbers of rodent samples (up to several hundred samples in a single run) on the  
93 MiSeq Illumina sequencing platform. The 16S rRNA amplicon sequencing technique  
94 is based on the amplification of small fragments of the hypervariable region of the  
95 16S rRNA gene. The sequences of these fragments are then obtained and compared  
96 with those in a dedicated database, for taxonomic identification [4,19]. Multiplexed

97 approaches of this kind include short indices (or tags) specific to a PCR product. This  
98 makes it possible to assign the sequences generated by the HTS run to a particular  
99 sample following bioinformatic analysis of the dataset [18]. Razzauti *et al.* [8]  
100 demonstrated the considerable potential of this approach for determining the  
101 prevalence of bacteria within populations and for analyzing bacterial interactions  
102 within hosts and vectors, based on the good characterization of bacterial diversity  
103 within each individual samples it provides. However, the various sources of error  
104 during the generation and processing of HTS data [20] may make it difficult to  
105 determine which samples are really positive or negative for a given bacterium. The  
106 detection of one or a few sequences assigned to a given taxon in a sample does not  
107 necessary mean that the bacterium is effectively present in that sample. We carried  
108 out an extensive literature review, from which we identified several potential sources  
109 of error involving all stages of a 16S rRNA amplicon sequencing experiment — from  
110 the collection of samples to the bioinformatic analysis — that might lead to false-  
111 negative or false-positive screening results (Table 1, [18,19,21-40]). These error  
112 sources have now been documented, and recent initiatives have called for the  
113 promotion of open sharing of standard operating procedures and best practices in  
114 microbiome research [41]. However, no experimental designs minimizing the impact  
115 of these sources of error on HTS data interpretation have yet been reported.

116 We describe here a rigorous experimental design for the direct estimation of biases  
117 from the data produced by 16S rRNA amplicon sequencing. We used these bias  
118 estimates to control and filter out potential false-positive and false-negative samples  
119 during screening for bacterial pathogens. We applied this strategy to 711 commensal  
120 rodents collected from 24 villages in Senegal, Western Africa: 208 *Mus musculus*  
121 *domesticus*, 189 *Rattus rattus*, 93 *Mastomys natalensis* and 221 *Mastomys*  
122 *erythroleucus*. Rodents were screened for bacteria as described by Kozich *et al.* [18],  
123 in a protocol based on MiSeq sequencing (Illumina) of the V4 hypervariable region of  
124 the 16SrRNA gene. We considered the common pitfalls listed in Table 1 during the  
125 various stages of the experiment (see details in the workflow procedure, Figure 1).  
126 Biases in assessments of the presence or absence of bacteria in rodents were  
127 estimated directly from the dataset, by including and analyzing negative controls  
128 (NC) and positive controls (PC) at the various stages of the experiment, and  
129 systematically using sample replicates. This strategy delivers realistic and reliable

130 **Table 1.** Sources of bias during the experimental and bioinformatic steps of 16S rRNA  
 131 amplicon sequencing, consequences for data interpretation and solutions for  
 132 decreasing these biases.

Experimental steps	Sources of errors	Consequences	Solutions
<b>Sample collection</b>	Cross-contamination between individuals [21]	False-positive samples	Rigorous processing (decontamination of the instruments, cleaning of the autopsy table, use of sterile bacterial-free consumables, gloves, masks) Negative controls during sampling (e.g., organs of healthy mice during dissection)
	Collection and storage conditions [21]	False-positive & negative samples	Use of appropriate storage conditions/buffers. Use of unambiguously identified samples. Double checking of tube labeling during sample collection.
<b>DNA extraction</b>	Cross-contamination between samples [22]	False-positive samples	Rigorous processing (separation of pre- and post-PCR steps, use of a sterile hood, filter tips and sterile bacterial-free consumables)
	Reagent contamination with bacterial DNA [21,23]	False-positive samples	Negative controls for extraction (extraction without sample)
	Small amounts of DNA [21, 24]	False-negative samples	Use of an appropriate DNA extraction protocol. Discarding of samples with a low DNA concentration
<b>Target DNA region and primer design</b>	Target DNA region efficacy [19,25]	False-negative due to poor taxonomic identification	Selection of an appropriate target region and design of effective primers for the desired taxonomic resolution
	Primer design [21,26]	False-negative samples due to biases in PCR amplification for some taxa	Checking of the universality of the primers with reference sequences
<b>Tag/Index design and preparation</b>	False-assignments of sequences due to cross-contamination between tags/indices [27,28]	False-positive samples	Rigorous processing (use of sterile hood, filter tips and sterile bacterial-free consumables, brief centrifugation before the opening of index storage tubes, separation of pre- and post-PCR steps) Negative controls for tags/indices (empty wells without PCR reagents for particular tags or index combinations) Positive controls for alien DNA, i.e. a bacteria strain highly unlikely to infect the samples studied (e.g., a host-specific bacterium unable to persist in the environment) to estimate false assignment rate
	False-assignments of sequences due to inappropriate tag/index design [29]	False-positive samples	Fixing of a minimum number of substitutions between tags or indices. Each nucleotide position in the sets of tags or indices should display about 25% occupation by each base for Illumina sequencing
<b>PCR amplification</b>	Cross-contamination between PCRs [28]	False-positive samples	Rigorous processing (brief centrifugation before opening the index storage tubes, separation of pre- and post-PCR steps) Negative controls for PCR (PCR without template) with microtubes left open during sample processing
	Reagent contamination with bacterial DNA [21,23]	False-positive samples	Rigorous processing (use of sterile hood, filter tips and sterile bacterial-free consumables) Negative controls for PCRs (PCR without template), with microtubes closed during sample processing
	Chimeric recombinations by jumping PCR [27,30,31,32,33]	False-positive samples due to artifactual chimeric sequences	Increasing the elongation time. Use of a bioinformatic strategy to remove the chimeric sequences (e.g., Uchime program)
	Poor or biased amplification [45]	False-negative samples	Increasing the amount of template DNA; Optimizing the PCR conditions (reagents and program) Use of technical replicates to validate sample positivity Positive controls for PCR (extraction from infected tissue and/or bacterial isolates)
<b>Library preparation</b>	Cross-contamination between PCRs/libraries [22]	False-positive samples	Rigorous processing (use of a sterile hood, filter tips and sterile bacterial-free consumables, electrophoresis and gel excision with clean consumables, separation of pre and post-PCR steps) Use of a protocol with an indexing step during target amplification Negative controls for indices (changing well positions between library preparation sessions)
	Chimeric recombinations by jumping PCR [27]	False-positive samples due to inter-individual recombinations	Avoiding PCR library enrichment of pooled samples. Positive controls for alien DNA, i.e. a bacterial strain that should not be identified in the sample (e.g. a host-specific bacterium unable to persist in the environment)
<b>MiSeq sequencing (Illumina)</b>	Sample sheet errors [21]	False-positive and negative samples	Negative controls (wells without PCR reagents for a particular index combination) Washing of the MiSeq with dilute sodium hypochlorite solution
	Run-to-run carryover (Illumina Technical Support Note No. 770-2013-046)	False-positive samples	qPCR quantification of the library before sequencing.
	Poor quality of reads due to flowcell overloading [34]	False-negative due to low quality of sequences	Decreasing cluster density. Creation of artificial sequence diversity at the flowcell surface (e.g., by adding 5 to 10% PhiX DNA control library)
	Poor quality of reads due to low-diversity libraries (Illumina Technical Support Note No. 770-2013-013)	False-negative due to low depth of sequencing	Decreasing the level of multiplexing Discard the sample with a low number of reads
	Small number of reads per sample [35,36]	False-negative due to low quality of sequences	Increasing paired-end sequence length or decreasing the length of the target sequence
	Too short overlapping read pairs [18]	False-positive due to false index-pairing	Use of a single barcode sequence for both the i5 and i7 indices for each sample (when possible, e.g. small number of samples) Positive controls for alien DNA, i.e., a bacterial strain highly unlikely to be found in the rodents studied (e.g., a host-specific bacterium unable to persist in the environment)
<b>Bioinformatics and taxonomic classification</b>	Mixed clusters on the flowcell [27]	False-positive due to false index-pairing	Use of a single barcode sequence for both the i5 and i7 indices for each sample (when possible, e.g. small number of samples) Positive controls for alien DNA, i.e., a bacterial strain highly unlikely to be found in the rodents studied (e.g., a host-specific bacterium unable to persist in the environment)
	Poor quality of reads	False-negative samples due to poor taxonomic resolution	Removal of low-quality reads
	Errors during processing (sequence trimming, alignment) [18,37,38]	False-positive and negative samples	Use of standardized protocols and reproducible workflows
	Incomplete reference sequence databases [39]	False-negative samples	Selection of an appropriate database for the selected target region and testing of the database for bacteria of particular interest
Error of taxonomic classification [40]	False-positive samples	Positive controls for PCRs (extraction from infected tissue and/or bacterial isolates and/or mock communities) Checking of taxonomic assignments by other methods (e.g., Blast analyses on different databases)	

134 estimates of bacterial prevalence in wildlife populations, and could be used to  
135 analyze the co-occurrence of different bacterial species within individuals.

136

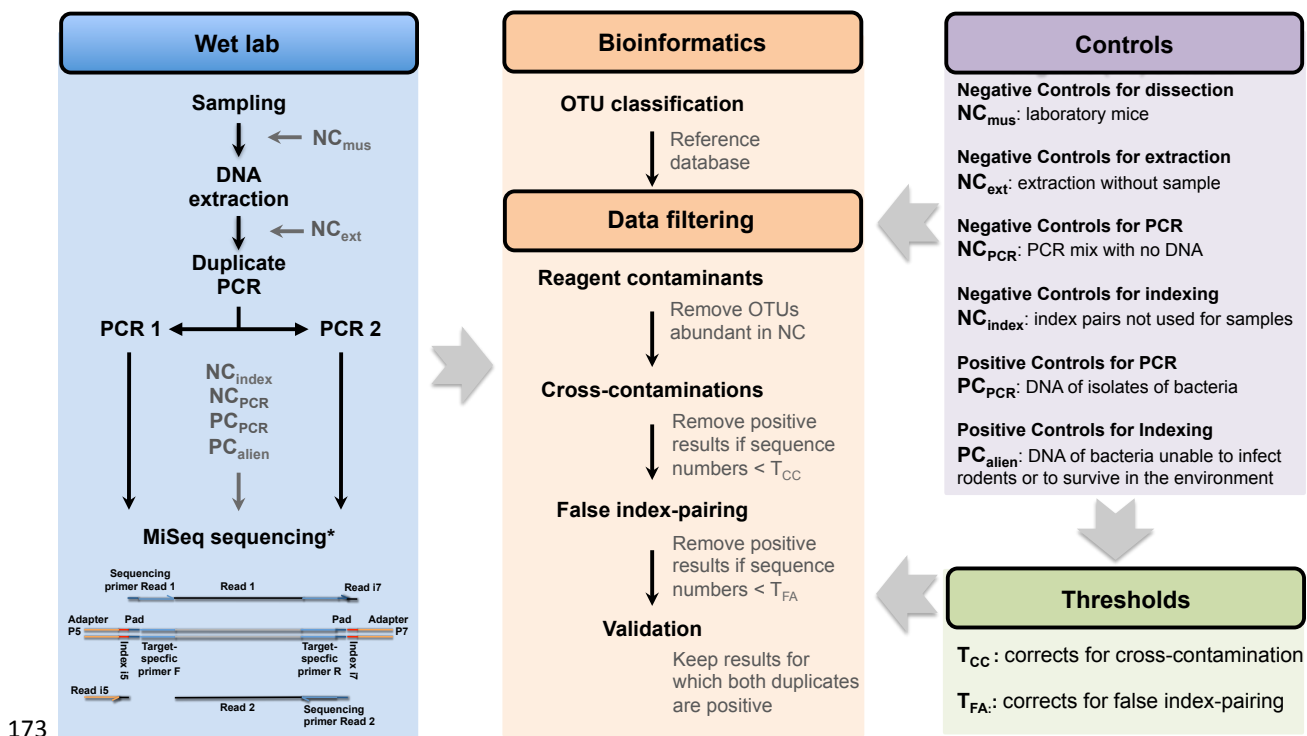
## 137 **Results & Discussion**

138 **Raw sequencing results.** The sequencing of 1569 PCR products (from rodent  
139 samples and controls, see details in Table S1) in two MiSeq runs generated a total of  
140 23,698,561 raw paired-end sequence reads (251-bp) of the V4 region of the  
141 16SrRNA gene. Overall, 99% of wild rodent PCRs generated more than 3,000 raw  
142 reads (mean: 11,908 reads; standard deviation: 6,062). The raw sequence read files  
143 in FASTQ format are available for each PCR and each MiSeq run on request to the  
144 corresponding author. Using mothur v1.34 [42] and the MiSeq standard operating  
145 procedure ([http://www.mothur.org/wiki/MiSeq\\_SOP](http://www.mothur.org/wiki/MiSeq_SOP)), we removed 20.1% of paired-  
146 end reads because they were misassembled, 1.5% of sequences because they were  
147 misaligned, 2.6% because they were chimeric and 0.2% because they were non-  
148 bacterial. The remaining reads were grouped into operational taxonomic units  
149 (OTUs) with a divergence threshold of 3%. Bioinformatics analysis identified 13,296  
150 OTUs, corresponding to a total of 7,960,533 sequences in run 1 and 6,687,060  
151 sequences in run 2.

152 **Taxonomic assignment of sequences.** We used the Bayesian classifier  
153 (bootstrap cutoff = 80%) implemented in mothur with the Silva SSU Ref database  
154 v119 [43] as a reference, for the taxonomic assignment of OTUs. The 50 most  
155 abundant OTUs accounted for 89% (min: 15,284 sequences; max: 2,206,731  
156 sequences) of the total sequence dataset (Table S2). The accuracy of taxonomic  
157 assignment (to genus level) was assessed with positive controls for PCR,  
158 corresponding to DNA extracts from laboratory isolates of *Bartonella taylorii*, *Borrelia*  
159 *burgdorferi* and *Mycoplasma mycoides* (PC<sub>Bartonella\_t</sub>, PC<sub>Borrelia\_b</sub> and PC<sub>Mycoplasma\_m</sub>,  
160 respectively), which were correctly assigned to a single OTU corresponding to the  
161 appropriate genuine sequences (Table 2). Note that the sequences of PC<sub>Mycoplasma\_m</sub>  
162 were assigned to Entomoplasmataceae rather than Mycoplasmataceae because of a  
163 frequent taxonomic error reflected in most databases, including Silva [44]. This  
164 problem might also affect other taxa. We therefore recommend systematically  
165 carrying out a blast analysis against the sequences of taxa of interest in GenBank to



166 confirm the taxonomic assignment obtained with the 16S databases. Finally, we  
 167 assumed that the small number of sequences per sample might limit the  
 168 completeness of bacterial detection [36]. For this reason, we discarded seven rodent  
 169 samples (2 *M. erythroleucus* and 5 *M. domesticus*) yielding fewer than 500  
 170 sequences for at least one of the two PCR replicates. This threshold corresponds to  
 171 99% of the distribution of the numbers of sequences between PCR products.  
 172



174 **Figure 1. Workflow of the wet laboratory, and for bioinformatics and data filtering**  
 175 **procedures, and a list of controls and thresholds included in the process of data**  
 176 **filtering for the elimination of false-positive results for 16S rRNA amplicon sequencing.**  
 177 Reagent contaminants were detected by analyzing the sequences in the  $NC_{ext}$  and  $NC_{PCR}$ ;  $T_{CC}$ : sequence  
 178 number threshold for correcting for cross-contamination.  $T_{CC}$  values are OTU- and run-dependent and  
 179 were estimated by analyzing the sequences in the controls,  $NC_{mus}$ ,  $NC_{ext}$ ,  $NC_{PCR}$  and  $PC_{index}$ ;  $T_{FA}$ :  
 180 sequence number threshold for correcting for false index-pairing.  $T_{FA}$  values are OTU- and run-dependent  
 181 and were estimated by analyzing the sequences in the  $NC_{index}$  and  $PC_{alien}$ . A result was considered  
 182 positive if the number of sequences was  $> T_{CC}$  and  $> T_{FA}$ . Samples were considered positive if a positive  
 183 result was obtained for both PCR replicates. \*see Kozich et al 2013 for details on the sequencing.

184

185 **Filtering for reagent contaminants.** Metagenomics data may be affected by  
 186 the contamination of reagents [23]. We therefore filtered the data, using negative  
 187 controls for extraction ( $NC_{ext}$ ), corresponding to extraction without the addition of a  
 188 tissue sample, and negative controls for PCR ( $NC_{PCR}$ ), corresponding to PCR

189 mixtures to which no DNA was added. This made it possible to identify the most  
190 abundant contaminants, including *Pseudomonas*, *Acinetobacter*, *Herbaspirillum*,  
191 *Streptococcus*, *Pelomonas*, *Brevibacterium*, *Brachybacterium*, *Dietzia*,  
192 *Brevundimonas*, *Delftia*, *Comamonas*, *Corynebacterium*, and *Geodermatophilus*,  
193 which accounted for 29% of the sequences in the dataset (Table S3). The bacterial  
194 contaminants detected differed between MiSeq runs: *Pseudomonas*, *Pelomonas* and  
195 *Herbaspirillum* predominated in run 1, whereas *Brevibacterium*, *Brachybacterium* and  
196 *Dietzia* predominated in run 2. This difference probably reflects the use of two  
197 different PCR kits manufactured at several months apart (Qiagen technical service,  
198 pers. com.). Other taxa, such as *Streptococcus*, most originated from the DNA  
199 extraction kits used, as they were detected in abundance in the negative controls for  
200 extraction (NC<sub>ext</sub>). These results highlight the importance of carrying out systematic  
201 negative controls to filter the taxa concerned, to prevent inappropriate data  
202 interpretation, particularly for the *Streptococcus* genus, which contains a number of  
203 important pathogenic species. The use of DNA-free reagents would improve the  
204 quality of sequencing data without affecting the depth of sequencing of the samples.

205 After filtering for the above reagent contaminants, the seven most relevant  
206 pathogenic bacterial genera, *Bartonella*, *Borrelia*, *Ehrlichia*, *Mycoplasma*, *Orientia*,  
207 *Rickettsia* and *Streptobacillus*, accounted for 66% of the sequences identified in wild  
208 rodent samples. Six different OTUs were obtained for *Mycoplasma*  
209 (*Mycoplasma\_OTU\_1* to *Mycoplasma\_OTU\_6*), with one OTU each for the other  
210 genera (Table 2). The other 34% of sequences probably corresponded to commensal  
211 bacteria (Bacteroidales, Bacteroides, Enterobacteriaceae, *Helicobacter*,  
212 *Lactobacillus*), undetected contaminants and rare taxa of unknown function.

213 **Filtering for false-positive results.** Mothur analysis produced a table of  
214 abundance, giving the number of sequences for each OTU in each PCR product  
215 (data available on request to the corresponding author). The multiple biases during  
216 experimental steps and data processing listed in Table 1 made it impossible to infer  
217 prevalence and co-occurrence directly from the table of sequence presence/absence  
218 in the PCR products. We suggest filtering the data with data-based estimates of the  
219 different biases calculated from the multiple controls introduced during the process.  
220 This strategy involves calculating sequence number thresholds from our bias  
221 estimates. Two different thresholds were set for each of the 12 OTUs and two MiSeq



222 runs. We then discarded positive results associated with numbers of sequences  
223 below the thresholds (Figure 1).

224 **Threshold  $T_{CC}$ : Filtering for cross-contamination.** One source of false positives is  
225 cross-contamination between samples processed in parallel (Table 1). Negative  
226 controls for dissection ( $NC_{mus}$ ), consisting of the spleens of healthy laboratory mice  
227 manipulated during sessions of wild rodent dissection, and negative controls for  
228 extraction ( $NC_{ext}$ ) and PCR ( $NC_{PCR}$ ) were used, together with positive controls for  
229 PCR ( $PC_{Bartonella_t}$ ,  $PC_{Borrelia_b}$  and  $PC_{Mycoplasma_m}$ ), to estimate cross-contamination.  
230 For each sequencing run, we calculated the maximal number of sequences for the 12  
231 pathogenic OTUs in the negative and positive controls. These numbers ranged from  
232 0 to 115 sequences, depending on the OTU and the run considered (Table 2), and  
233 we used them to establish OTU-specific thresholds ( $T_{CC}$ ) for each run. The use of  
234 these  $T_{CC}$  led to 0% to 69% of the positive results being discarded, corresponding to  
235 only 0% to 0.14% of the sequences, depending to the OTU considered (Figure 2,  
236 Table S4). A PCR product may be positive for several bacteria in cases of  
237 coinfection. In such cases, the use of a  $T_{CC}$  makes it possible to discard the positive  
238 result for one bacterium whilst retaining positive results for other bacteria.

239 **Threshold  $T_{FA}$ : Filtering out incorrectly assigned sequences.** Another source of  
240 false positives is the incorrect assignment of sequences to a PCR product (Table 1).  
241 This phenomenon is essentially due to mixed clusters during the sequencing [27].  
242 We used two kinds of controls to detect incorrect assignments (Figure 1).

243 First, negative control index pairs ( $NC_{index}$ ), corresponding to particular index pairs  
244 not used to identify our samples, were included to check for cross-contamination  
245 between indices or for errors during completion of the Illumina sample sheet.  $NC_{index}$   
246 returned very few read numbers (1 to 12), suggesting that there was little or no cross-  
247 contamination between indices in our experiment.

248 Second, we used “alien” positive controls ( $PC_{alien}$ ) in the PCR amplification step:  
249  $PC_{Mycoplasma_m}$ , corresponding to the DNA of *Mycoplasma mycoides*, which cannot  
250 infect rodents, and  $PC_{Borrelia_b}$ , corresponding to the DNA of *Borrelia burgdorferi*,  
251 which is not present in Africa. Neither of these bacteria can survive in abiotic  
252 environments, so the presence of their sequences in African rodent PCR products  
253 indicates a misassignment of sequences due to false index-pairing [27]. Using  
254  $PC_{Mycoplasma_m}$ , we obtained an estimate of the global false index-pairing rate of

255 0.14% (i.e. 398 of 280,151 sequences of the *Mycoplasma mycoides* OTU were  
 256 assigned to samples other than PC<sub>Mycoplasma\_m</sub>). Using PC<sub>Borrelia\_b</sub>, we obtained an  
 257 estimate of 0.22% (534 of 238,772 sequences of the *Borrelia burgdorferi* OTU were  
 258 assigned to samples other than PC<sub>Borrelia\_b</sub>). These values are very close to the  
 259 estimate of 0.3% obtained by Kircher *et al.* [27]. Close examination of the distribution  
 260 of misassigned sequences within the PCR 96-well microplates showed that all PCR  
 261 products with misassigned sequences had one index in common with either  
 262 PC<sub>Mycoplasma\_m</sub> or PC<sub>Borrelia\_b</sub> (Figure S1).

263 **Table 2. Number of sequences for 12 pathogenic OTUs observed in wild rodents, in**  
 264 **negative controls and in positive controls, together with T<sub>CC</sub> and T<sub>FA</sub> threshold**  
 265 **values.** Data are given for the two MiSeq runs separately. NC<sub>PCR</sub>: negative controls for PCR; NC<sub>ext</sub>:  
 266 negative controls for extraction; NC<sub>mus</sub>: negative controls for dissection; PC<sub>Bartonella\_t</sub>: positive controls for  
 267 PCR; PC<sub>Borrelia\_b</sub> and PC<sub>Mycoplasma\_m</sub>: positive controls for PCR and positive controls for indexing; T<sub>CC</sub> and  
 268 T<sub>FA</sub>: thresholds for positivity for a particular bacterium according to bacterial OTU and MiSeq run (see also  
 269 Figure 1).

OTUs	Total	Wild rodents (n=711)												Negative controls			Positive controls			Thresholds	
		Wild rodents (n=711)			NC <sub>PCR</sub>		NC <sub>ext</sub>		NC <sub>mus</sub>		PC <sub>Bartonella_t</sub>		PC <sub>Borrelia_b</sub>		PC <sub>Mycoplasma_m</sub>		T <sub>CC</sub> <sup>*</sup>	T <sub>FA</sub> <sup>**</sup>			
		Total no. of sequences	Maximum no. of sequences in one PCR	Maximum no. of sequences in one PCR	Total no. of sequences	Maximum no. of sequences in one PCR	Total no. of sequences	Maximum no. of sequences in one PCR	Total no. of sequences	Maximum no. of sequences in one PCR	Total no. of sequences	Maximum no. of sequences in one PCR	Total no. of sequences	Maximum no. of sequences in one PCR	Total no. of sequences	Maximum no. of sequences in one PCR					
<b>Whole dataset</b>	7960533	7149444	64722	45900	8002	39308	8741	68350	26211	137424	73134	239465	120552	280642	82933	/	/				
<i>Mycoplasma</i> _OTU_1	1410218	1410189	61807	2	1	3	2	9	5	3	3	8	6	4	3	6	282				
<i>Mycoplasma</i> _OTU_3	507376	507369	36335	2	1	0	0	0	0	2	2	1	1	2	2	2	101				
<i>Ehrlichia</i> _OTU	649451	649423	63137	4	2	3	2	7	4	1	1	1	1	12	6	6	130				
<i>Borrelia</i> _OTU	345873	345845	28528	4	4	7	4	9	4	1	1	0	0	7	3	4	69				
<i>Orientia</i> _OTU	279965	279957	29503	1	1	4	1	0	0	2	2	0	0	1	1	2	56				
<i>Bartonella</i> _OTU	202127	67973	16145	1	1	1	1	1	1	134124	71163	7	4	20	9	9	40				
<i>M. mycoides</i> ***	280151	338	28	0	0	0	0	2	2	34	20	24	18	279753	82767	/	/				
<i>B. burgdorferi</i> ***	238772	420	43	0	0	0	0	0	0	38	21	238238	119586	76	23	/	/				
<b>Whole dataset</b>	6687060	6525107	42326	61231	9145	53334	7669	/	/	12142	7518	13378	7164	21868	6520	/	/				
<i>Mycoplasma</i> _OTU_1	155486	155486	7703	0	0	0	0	/	/	0	0	0	0	0	0	0	31				
<i>Mycoplasma</i> _OTU_2	1036084	1035890	23588	1	1	192	115	/	/	0	0	0	0	1	1	115	207				
<i>Mycoplasma</i> _OTU_3	127591	127590	5072	1	1	0	0	/	/	0	0	0	0	0	0	1	26				
<i>Mycoplasma</i> _OTU_4	85596	85583	20146	0	0	13	13	/	/	0	0	0	0	0	0	13	17				
<i>Mycoplasma</i> _OTU_5	56324	56324	10760	0	0	0	0	/	/	0	0	0	0	0	0	0	11				
<i>Mycoplasma</i> _OTU_6	13356	13356	1482	0	0	0	0	/	/	0	0	0	0	0	0	0	3				
<i>Ehrlichia</i> _OTU	74017	74017	19651	0	0	0	0	/	/	0	0	0	0	0	0	0	15				
<i>Borrelia</i> _OTU	21636	21636	3085	0	0	0	0	/	/	0	0	0	0	0	0	0	4				
<i>Orientia</i> _OTU	307	307	181	0	0	0	0	/	/	0	0	0	0	0	0	0	0				
<i>Bartonella</i> _OTU	1559028	1547652	14515	1	1	2	2	/	/	11297	6714	2	2	74	59	59	312				
<i>Streptobacillus</i> _OTU	32399	32399	6245	0	0	0	0	/	/	0	0	0	0	0	0	0	6				
<i>Rickettsia</i> _OTU	589	589	329	0	0	0	0	/	/	0	0	0	0	0	0	0	0				
<i>M. mycoides</i> ***	16854	2	1	0	0	0	0	/	/	0	0	0	0	16852	5766	/	/				
<i>B. burgdorferi</i> ***	12197	0	0	0	0	0	0	/	/	0	0	12197	6426	0	0	/	/				

\*: Threshold T<sub>CC</sub> is based on the maximum number of sequences observed in a negative or positive control for a particular OTU in each run

\*\* : Threshold T<sub>FA</sub> is based on the false assignment rate (0.02%) weighted by the total number of sequences of each OTU in each run

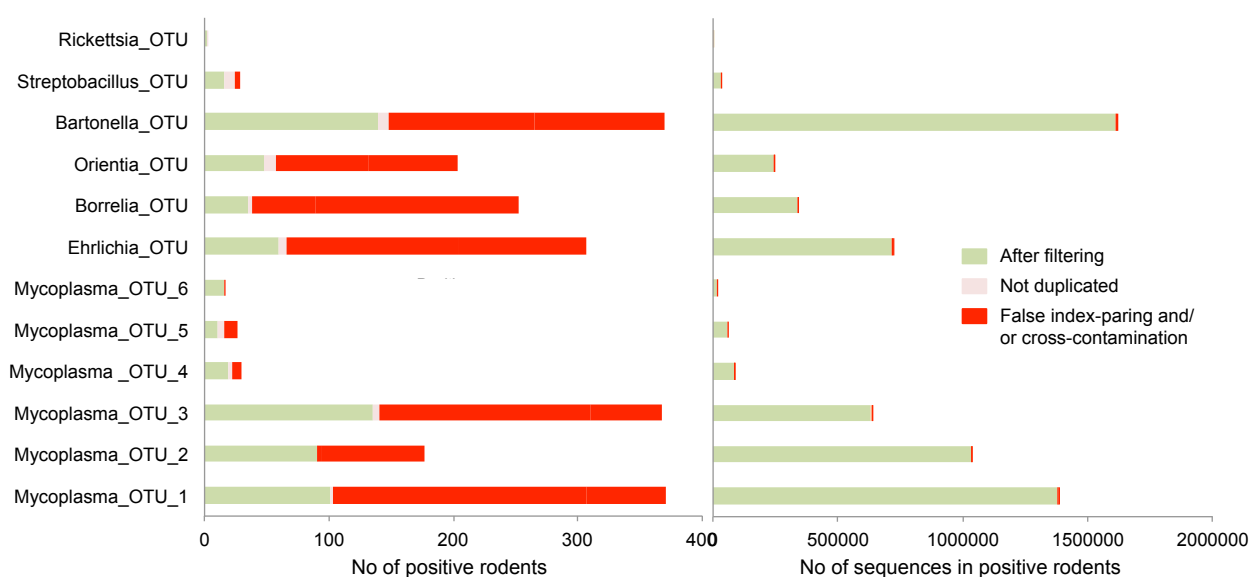
\*\*\*: *Mycoplasma mycoides* and *Borrelia burgdorferi* bacterial isolates added as positive controls for PCR and indexing (i.e., PC<sub>alien</sub>, see Figure 1)

270

271 We then estimated the impact of false index-pairing for each PCR product, by  
 272 calculating the maximal number of sequences of “alien” bacteria assigned to PCR  
 273 products other than the corresponding PC. These numbers varied from 28 to 43,  
 274 depending on the positive control for run 1 (Table 2) — run 2 was discarded because  
 275 of the low values of the numbers of sequences, which is likely due to the fact that  
 276 DNAs of PC of were hundred-fold diluted in run 2 (Table S1) —. We then estimated a  
 277 false-assignment rate for each PCR product (R<sub>fa</sub>), by dividing the above numbers by  
 278 the total number of sequences from “alien” bacteria in the sequencing run 1. R<sub>fa</sub> was

279 estimated for  $PC_{Mycoplasma\_m}$  and  $PC_{Borrelia\_b}$  separately.  $R_{fa}$  reached 0.010% and  
280 0.018% for  $PC_{Mycoplasma\_m}$  and  $PC_{Borrelia\_b}$ , respectively. We adopted a conservative  
281 approach, by fixing the  $R_{fa}$  value to 0.020%. This number signifies that each PCR  
282 product may receive a maximum 0.020% of the total number of sequences of an  
283 OTU present in a run due to false index-pairing. Moreover, the number of  
284 misassigned sequences for a specific OTU into a PCR product should increase with  
285 the total number of sequences of the OTU in the MiSeq run. We therefore defined the  
286 second threshold ( $T_{FA}$ ) as the total number of sequences in the run for an OTU  
287 multiplied by  $R_{fa}$ .  $T_{FA}$  values varied with the abundance of each OTU in the  
288 sequencing run (Table 2). Because the abundance of each OTU varied from one  
289 sequencing run to another,  $T_{FA}$  also varied according to the sequencing run. The use  
290 of the  $T_{FA}$  led to 0% to 87% of positive results being discarded. This corresponded to  
291 0% to 0.71% of the sequences, depending on the OTU (Figure 2, Table S4).

292



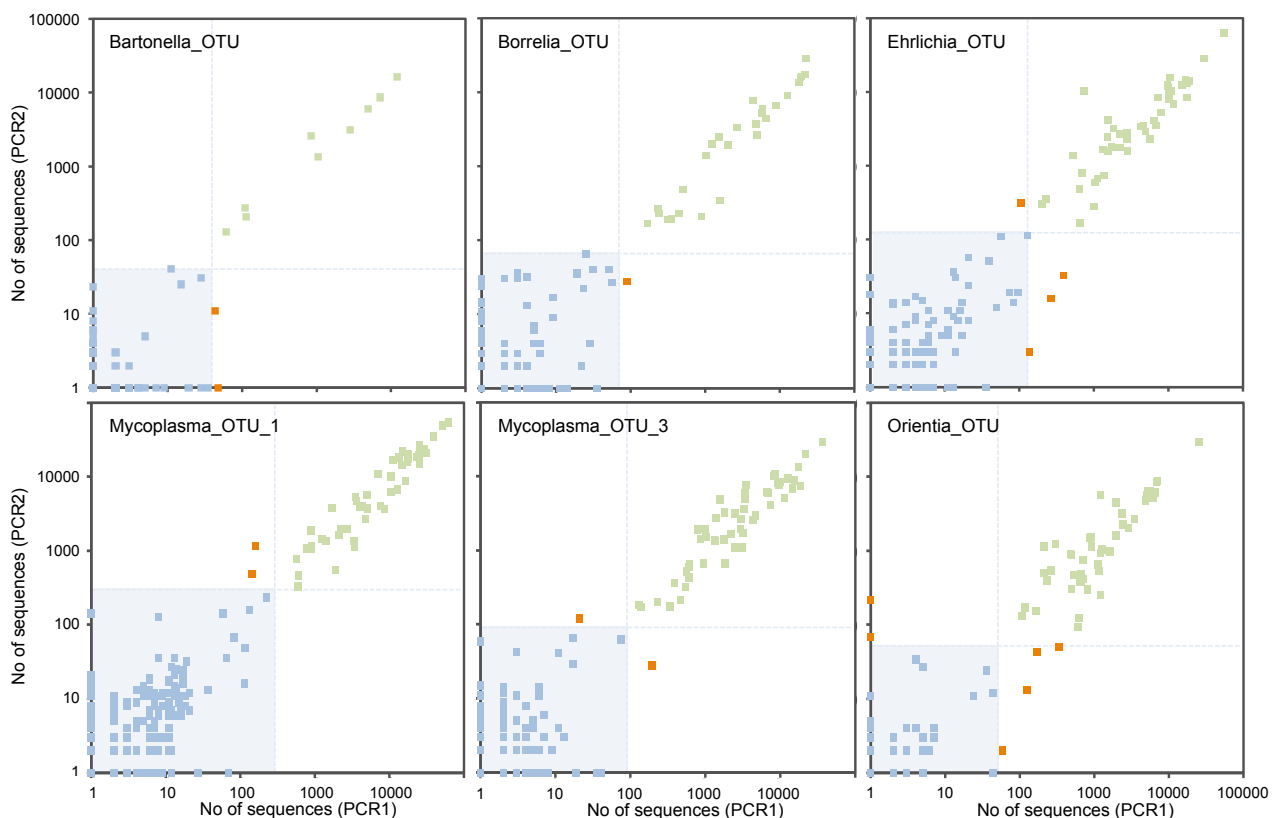
293

294 **Figure 2. Numbers of positive rodents, and of sequences in positive rodents,**  
295 **removed for each OTU at each step in data filtering.** These findings demonstrate  
296 that the positive rodents filtered out corresponded to only a very small number of  
297 sequences. (A) The histogram shows the number of positive rodents discarded because of likely  
298 cross-contamination, false index-pairing and for a negative result in a replicate PCR, and, finally the  
299 positive results retained at the end of data filtering in green. (B) The histogram shows the number of  
300 sequences corresponding to the same class of positive rodents. Note that several positive results may  
301 be recorded for the same rodent in cases of co-infection.

302 **Validation with PCR replicates.** Random contamination may occur during the  
303 preparation of PCR 96-well microplates. These contaminants may affect some of the

304 wells, but not those for the negative controls, leading to the generation of false-  
305 positive results. We thus adopted a conservative approach, in which we considered  
306 rodents to be positive for a given OTU only if both PCR replicates were considered  
307 positive after the filtering steps described above. The relevance of this strategy was  
308 supported by the strong correlation between the numbers of sequences for the two  
309 PCR replicates for each rodent ( $R^2 > 0.90$ , Figure 3 and Figure S2). At this stage, 673  
310 positive results for 419 rodents were validated for both replicates (note that a rodent  
311 may be positive for several bacteria, and may thus be counted several times),  
312 whereas only 52 positive results were discarded because the result for the other  
313 replicate was negative.

314



316 **Figure 3. Plots of the number of sequences (log (x+1) scale) from bacterial**  
317 **OTUs in both PCR replicates (PCR1 & PCR2) of the 348 wild rodents analyzed**  
318 **in the first MiSeq run.** Note that each rodent was tested with two replicate PCRs. Green points  
319 correspond to rodents with two positive results after filtering; red points correspond to rodents with one  
320 positive result and one negative result; and blue points correspond to rodents with two negative  
321 results. The light blue area and lines correspond to the threshold values used for the data filtering:  
322 samples below the lines are filtered out. See Figure S2 for plots corresponding to the second MiSeq  
323 run.

324 At this final validation step, 0% to 60% of the positive results for a given OTU were  
325 discarded, corresponding to only 0% to 7.17% of the sequences (Figure 2, Table S4  
326 and Table S5). Note that the number of replicates may be increased, as described in  
327 the strategy of Gómez-Díaz *et al* [45].

328 **Post-filtering results.** Finally, the proportion of rodents positive for a given OTU  
329 filtered out by the complete filtering approach varied from 6% to 86%, depending on  
330 the OTU, corresponding to only 1% of the total sequences (Figure 2). Indeed, our  
331 filtering strategy mostly excluded rodents with a small number of sequences for the  
332 OTU concerned. These rodents were considered to be false-positive.

333 **Refining bacterial taxonomic identification.** We refined the taxonomic  
334 identification of the 12 bacterial OTUs through phylogenetic and blast analyses. We  
335 were able to identify the bacteria present down to genus level and, in some cases,  
336 we could even identify the most likely species (Table 3 and Figure S3). For instance,  
337 the sequences of the six *Mycoplasma* OTUs were consistent with three different  
338 species — *M. haemomuris* for OTU\_1 and 3, *M. coccoides* for OTU\_4, 5 and 6, and  
339 *M. species novo* [46] for OTU\_2 — with high percentages of sequence identity  
340 ( $\geq 93\%$ ) and strong bootstrap support ( $\geq 80\%$ ). All three of these species belong to the  
341 Hemoplasma group, which is known to infect mice, rats and other mammals [47,48],  
342 and is thought to cause anemia in humans [49,50]. The *Borrelia* sequences grouped  
343 with three different species of the relapsing fever group (*crocidurae*, *duttonii* and  
344 *recurrentis*) with a high percentage of identity (100%) and a reasonably high  
345 bootstrap value (71%). In West Africa, *B. crocidurae* causes severe borreliosis, a  
346 rodent-borne disease transmitted by ticks and lice [51]. The *Ehrlichia* sequences  
347 were 100% identical to and clustered with the recently described Candidatus  
348 *Ehrlichia khabarensis* isolated from voles and shrews in the Far East of Russia [52].  
349 The *Rickettsia* sequences were 100% identical to the sequence of *R. typhi*, a species  
350 of the typhus group responsible for murine typhus [53], but this clade was only  
351 weakly differentiated from many other *Rickettsia* species and had only moderate  
352 bootstrap support (61%). The most likely species corresponding to the sequences of  
353 the *Streptobacillus* OTU was *S. moniliformis*, with a high percentage of identity  
354 (100%) and a high bootstrap value (100%). This bacterium is common in rats and  
355 mice and causes a form of rat-bite fever, Haverhill fever [54]. The *Orientia* sequences  
356 corresponded to *O. chuto*, with a high percentage of identity (100%) and a high

357 **Table 3. Detection of 12 bacterial OTUs in the four wild rodent species (n=704)**  
 358 **sampled in Senegal, and the biology and pathogenicity of the corresponding**  
 359 **bacterial genus.** n= number of rodents analyzed.

OTUs of interest (genus level)	Closest species* (% identity in GenBank)	Number of positive wild rodents					Biology & epidemiology
		<i>Mastomys erythroleucus</i> (n=219)	<i>Mastomys natalensis</i> (n=93)	<i>Mus musculus</i> (n=203)	<i>Rattus rattus</i> (n=189)		
<i>Bartonella</i>	undetermined	60	68	1	6	<i>Bartonella</i> spp. are intracellular fastidious hemotropic gram-negative organisms identified in a wide range of domestic and wild mammals and transmitted by arthropods. Several rodent-borne <i>Bartonella</i> species have emerged as zoonotic agents, and various clinical manifestations are reported, including fever, bacteremia and neurological symptoms [83].	
<i>Borrelia</i>	<i>crocidurae</i> (100%) <i>duttonii</i> (100%) <i>recurrentis</i> (100%)	21	0	8	6	<i>Borrelia</i> is a genus of spiral gram-negative bacteria of the spirochete phylum. These bacteria are obligate parasites of animals and are responsible for relapsing fever borreliosis, a zoonotic disease transmitted by arthropods (tick and lice). This disease is the most frequent human bacterial disease in Africa. <i>B. crocidurae</i> is endemic to West Africa, including Senegal, and <i>B. duttonii</i> and <i>B. recurrentis</i> have been reported in Central, southern and East Africa [51].	
<i>Ehrlichia</i>	<i>khabarensis</i> (100%)	40	0	12	8	The genus <i>Ehrlichia</i> includes five species of small gram-negative obligate intracellular bacteria. The life cycle includes the reproduction stages taking place in both ixodid ticks, acting as vectors, and vertebrates. <i>Ehrlichia</i> spp. can cause a persistent infection in the vertebrate hosts, which thus become reservoirs of infection. A number of new genetic variants of <i>Ehrlichia</i> have been recently detected in rodent species (e.g., <i>Candidatus Ehrlichia khabarensis</i> [52]).	
<i>Mycoplasma</i> OTU_1	<i>haemomuris</i> (96%)	28	41	30	1	<i>Mycoplasma</i> is a genus including over 100 species of bacteria that lack of a cell wall around their cell membrane. <i>Mycoplasma coccoides</i> and <i>Mycoplasma haemomuris</i> are blood parasites of wild and laboratory rodents. A new closely related species was recently isolated from brown rats (AB752303 [46]). These species are commonly referred as "hemoplasmas". Hemoplasmas have been detected within the erythrocytes of cats, dogs, pigs, rodents and cattle, in which they may cause anaemia. There have been sporadic reports of similar infections in humans, but these infections have been poorly characterized [50].	
<i>Mycoplasma</i> OTU_2	<i>sp. novo</i> (100%) GenBank AB752303	0	0	0	90		
<i>Mycoplasma</i> OTU_3	<i>haemomuris</i> (93%)	93	23	1	1		
<i>Mycoplasma</i> OTU_4	<i>coccoides</i> (96%)	0	0	0	18		
<i>Mycoplasma</i> OTU_5	<i>coccoides</i> (95%)	3	7	0	0		
<i>Mycoplasma</i> OTU_6	<i>coccoides</i> (97%)	3	14	0	0		
<i>Orientia</i>	<i>chuto</i> (100%) <i>tsutsugamushi</i> (98%)	0	3	46	0	<i>Orientia</i> is a genus of obligate intracellular gram-negative bacteria found in mites and rodents. <i>Orientia tsutsugamushi</i> is the agent of scrub typhus in humans. This disease, one of the most underdiagnosed and underreported febrile illnesses requiring hospitalization, has an estimated 10% fatality rate unless treated appropriately. A new species, <i>Orientia chuto</i> , was recently characterized in sick patients from the Arabian Peninsula, and new <i>Orientia</i> haplotypes have been identified in France and Senegal [9].	
<i>Rickettsia</i>	<i>typhi</i> (100%)	1	0	0	1	<i>Rickettsia</i> is a genus of obligate intracellular gram-negative bacteria found in arthropods and vertebrates. <i>Rickettsia</i> spp. are symbiotic species transmitted vertically in invertebrates, and some are pathogenic invertebrates. <i>Rickettsia</i> species of the typhus group cause many human diseases, including murine typhus, which is caused by <i>Rickettsia typhi</i> and transmitted by fleas [53].	
<i>Streptobacillus</i>	<i>moniliformis</i> (100%)	10	1	0	5	<i>Streptobacillus</i> is a genus of aerobic, gram-negative facultative anaerobe bacteria, which grow in culture as rods in chains. <i>Streptobacillus moniliformis</i> is common in rats and mice and is responsible of the Streptobacillosis form of rat-bite fever, the Haverhill fever. This zoonosis begins with high prostrating fevers, rigors (shivering), headache and polyarthralgia (joint pain). Untreated, rat-bite fever has a mortality rate of approximately 10% [54].	

\*based on phylogenetic analysis, see Figure S3  
 n: number of rodents screened



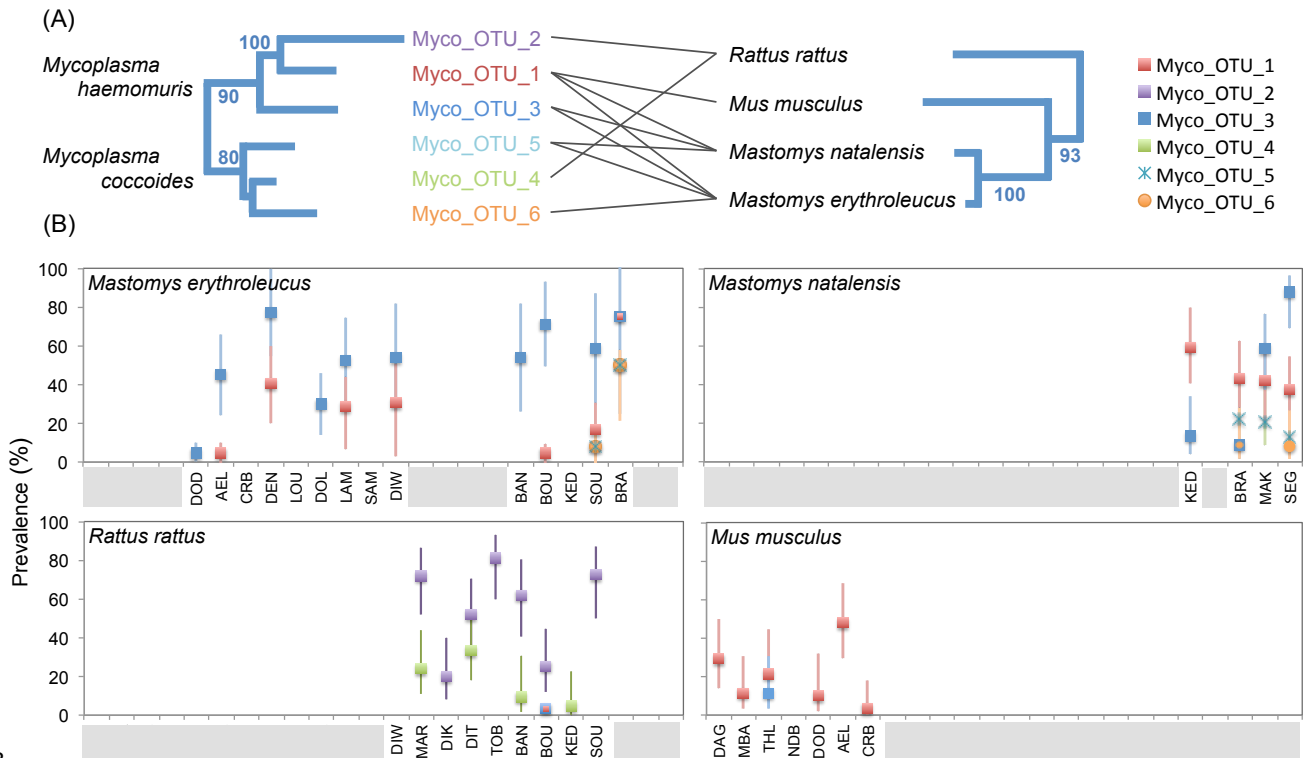
361 bootstrap value (77%). This species was recently isolated from a patient infected in  
362 Dubai [55]. Finally, accurate species determination was not possible for *Bartonella*,  
363 as the 16S rRNA gene does not resolve the species of this genus well [56]. Indeed,  
364 the sequences from the *Bartonella* OTU detected in our rodents corresponded to at  
365 least seven different species (*elizabethae*, *japonica*, *pachyuromydis*, *queenslandis*,  
366 *rattaaustraliani*, *tribocorum*, *vinsonii*) and a putative new species recently identified in  
367 Senegalese rodents [57].

368 These findings demonstrate the considerable potential of 16S rRNA amplicon  
369 sequencing for the rapid identification of zoonotic agents in wildlife, provided that the  
370 post-sequencing data are cleaned beforehand. *Borrelia* [51] and *Bartonella* [57] were  
371 the only ones of the seven pathogenic bacterial genera detected here in Senegalese  
372 rodents to have been reported as present in rodents from West Africa before. The  
373 other bacterial genera identified here have previously been reported to be presented  
374 in rodents only in other parts of Africa or on other continents. *S. moniliformis* has  
375 recently been detected in rodents from South Africa [58] and there have been a few  
376 reports of human streptobacillosis in Kenya [59] and Nigeria [60]. *R. typhi* was  
377 recently detected in rats from Congo, in Central Africa [61], and human seropositivity  
378 for this bacterium has been reported in coastal regions of West Africa [62]. With the  
379 exception of one report in Egypt some time ago [63], *Mycoplasma* has never before  
380 been reported in African rodents. Several species of *Ehrlichia* (from the *E. canis*  
381 group: *E. chaffeensis*, *E. ruminantium*, *E. muris*, *E. ewingii*) have been characterized  
382 in West Africa, but only in ticks from cattle [89] together with previous reports of  
383 possible cases of human ehrlichioses in this region [64]. Finally, this study reports the  
384 first identification of *Orientia* in African rodents [9]. There have already been a few  
385 reports of suspected human infection with this bacterium in Congo, Cameroon,  
386 Kenya and Tanzania [65].

387 ***Estimating prevalence and coinfection.*** After data filtering, we were able to  
388 estimate the prevalence in rodent populations and to assess coinfection in individual  
389 rodents, for the 12 bacterial OTUs. Bacterial prevalence varied considerably between  
390 rodent species (Table 3). *Bartonella* was highly prevalent in the two multimammate  
391 rats *M. natalensis* (73%) and *M. erythroleucus* (27%); *Orientia* was prevalent in the  
392 house mouse *M. musculus* (22%) and *Ehrlichia* occurred frequently in only one on  
393 the two multimammate rats *M. erythroleucus* (18%). By contrast, the prevalence of

394 *Streptobacillus* and *Rickettsia* was low in all rodent species (<5%). Coinfection was  
 395 common, as 184 rodents (26%) were found to be coinfecting with bacteria from two  
 396 (19%), three (5%), four (2%) or five (0.1%) different bacterial pathogens.

397



398

399

400 **Figure 4. Prevalence of *Mycoplasma* lineages in Senegalese rodents, by site, and**  
 401 **phylogenetic associations between *Mycoplasma* lineages and rodent species.** (A)  
 402 Comparison of phylogenetic trees based on the 16S rRNA V4-sequences of *Mycoplasma*, and on the  
 403 mitochondrial cytochrome *b* gene and the two nuclear gene fragments (IRBP exon 1 and GHR) for  
 404 rodents (rodent tree redrawn from [91]). Lines link the *Mycoplasma* lineages detected in the various  
 405 rodent species (for a minimum site prevalence exceeding 10%). The numbers next to branches are  
 406 bootstrap values (only shown if >70%). (B) Plots of OTU prevalence with 95% confidence intervals  
 407 calculated by Sterne's exact method [92] by rodent species and site (see [67] for more information about  
 408 site codes and their geographic locations). The gray bars in the X-legend indicate sites from which the  
 409 rodent species concerned is absent.

410

411 Interestingly, several *Mycoplasma* OTUs appeared to be specific to a rodent genus  
 412 or species (Table 3, Figure 4). OTU\_2, putatively identified as a recently described  
 413 lineage isolated from brown rat, *Rattus norvegicus* [46], was specifically associated  
 414 with *R. rattus* in this study. Of the OTUs related to *M. coccoides*, OTU\_4 was found  
 415 exclusively in *R. rattus*, whereas OTUs\_5 and 6 seemed to be specific to the two  
 416 multimammate rats (*M. erythroleucus* and *M. natalensis*). Comparative phylogenies of  
 417 *Mycoplasma* OTUs and rodents showed that *R. rattus*, which is phylogenetically

418 more distantly related to the other three rodents, contained a *Mycoplasma*  
419 community different from that in the *Mus-Mastomys* rodent clade (Figure 4).  
420 Pathogen prevalence also varied considerably between sites, as shown for the six  
421 *Mycoplasma* OTUs (Figure 4). This suggests that the infection risks for animals and  
422 humans vary greatly according to environmental characteristics and/or biotic features  
423 potentially related to recent changes in the distribution of rodent species in Senegal  
424 [66,67]

425

## 426 **Perspectives**

427 **Improving HTS for epidemiological surveillance.** The screening strategy  
428 described here has the considerable advantage of being non-specific, making it  
429 possible to detect unanticipated or novel bacteria. Razzauti *et al.* [8] recently showed  
430 that the sensitivity of 16S rRNA amplicon sequencing on the MiSeq platform was  
431 equivalent to that of whole RNA sequencing (RNAseq) on the HiSeq platform for  
432 detecting bacteria in rodent samples. However, little is known about the comparative  
433 sensitivity of HTS approaches relative to qPCR with specific primers, the current gold  
434 standard for bacterial detection within biological samples. Additional studies are  
435 required to address this question. Moreover, as 16S rRNA amplicon sequencing is  
436 based on a short sequence, it does not yield a high enough resolution to distinguish  
437 between species in some bacterial genera, such as *Bartonella*. Whole-genome  
438 shotgun or RNAseq techniques provide longer sequences, through the production of  
439 longer reads or the assembly of contigs, and they might therefore increase the  
440 accuracy of species detection [68]. However, these techniques would be harder to  
441 adapt for the extensive multiplexing of samples [8]. Other methods could be used to  
442 assign sequences to bacterial species for individuals found positive for a bacterial  
443 genera following the 16S rRNA screening. For example, positive PCR assays could  
444 be carried out with bacterial genus-specific primers, followed by amplicon  
445 sequencing, as commonly used in MLSA (multilocus sequence analysis) strategies  
446 [69] or high-throughput microfluidic qPCR assays based on bacterial species-specific  
447 primers could be used [70]. High-throughput amplicon sequencing approaches could  
448 be fine-tuned to amplify several genes for species-level assignment, such as the *gltA*

449 gene used by Gutierrez *et al.* [71] for the *Bartonella* genus, in parallel with the 16S  
450 rRNA-V4 region.

451 This strategy could also easily be adapted for other microbes, such as protists, fungi  
452 and even viruses, provided that universal primers are available for their detection  
453 (see [72,73] for protists and fungi, and [74] for degenerate virus family-level primers  
454 for viruses). Finally, our filtering method could also be translated to any other post-  
455 sequencing dataset of indexed or tagged amplicons in the framework of  
456 environmental studies (e.g. metabarcoding for diet analysis and biodiversity  
457 monitoring [75], the detection of rare somatic mutations [76] or the genotyping of  
458 highly polymorphic genes (e.g. MHC or HLA typing, [77,78]).

459 **Monitoring the risk of zoonotic diseases.** Highly successful synanthropic  
460 wildlife species, such as the rodents studied here, will probably play an increasingly  
461 important role in the transmission of zoonotic diseases [79]. Many rodent-borne  
462 pathogens cause only mild or undifferentiated disease in healthy people, and these  
463 illnesses are often misdiagnosed and underreported [54,80-83]. The information  
464 about pathogen circulation and transmission risks in West Africa provided by this  
465 study is important in terms of human health policy. We show that rodents carry seven  
466 major pathogenic bacterial genera: *Borrelia*, *Bartonella*, *Mycoplasma*, *Ehrlichia*,  
467 *Rickettsia*, *Streptobacillus* and *Orientia*. The last five of these genera have never  
468 before been reported in West African rodents. The data generated with our HTS  
469 approach could also be used to assess zoonotic risks and to formulate appropriate  
470 public health strategies involving the focusing of continued pathogen surveillance and  
471 disease monitoring programs on specific geographic areas or rodent species likely to  
472 be involved in zoonotic pathogen circulation, for example.

473

## 474 **Materials & Methods**

475 **Ethics statement.** Animals were treated in accordance with European Union  
476 guidelines and legislation (Directive 86/609/EEC). The CBGP laboratory received  
477 approval (no. B 34-169-003) from the Departmental Direction of Population  
478 Protection (DDPP, Hérault, France), for the sampling of rodents and the storage and  
479 use of their tissues. None of the rodent species investigated in this study has  
480 protected status (see IUCN and CITES lists).

481 **Sample collection.** Rodents were killed by cervical dislocation, as recommended  
482 by Mills *et al.* [84] and dissected as described in Herbreteau *et al.* [85]. Rodent  
483 species were identified by morphological and/or molecular techniques [67]. Cross-  
484 contamination during dissection was prevented by washing the tools used  
485 successively in bleach, water and alcohol between rodents. We used the spleen for  
486 bacterial detection, because this organ is a crucial site of early exposure to bacteria  
487 [86]. Spleens were placed in RNAlater (Sigma) and stored at 4°C for 24 hours and  
488 then at -20°C until their use for genetic analyses.

489 **Target DNA region and primer design.** We used primers with sequences  
490 slightly modified from those of the universal primers of Kozich *et al.* [18] to amplify a  
491 251-bp portion of the V4 region of the 16S rRNA gene (16S-V4F:  
492 GTGCCAGCMGCCGCGGTAA; 16S-V4R: GGACTACHVGGGTWTCTAATCC). The  
493 ability of these primers to hybridize to the DNA of bacterial zoonotic pathogens was  
494 assessed by checking that there were low numbers of mismatched bases over an  
495 alignment of 41,113 sequences from 79 zoonotic genera inventoried by Taylor *et al.*  
496 [1], extracted from the Silva SSU database v119 [43] (Table S6). The FASTA file is  
497 available on request to the corresponding author. We used a slightly modified version  
498 of the dual-index method of Kozich *et al.* [18] to multiplex our samples. The V4  
499 primers included different 8-bp indices (i5 in the forward and i7 in the reverse  
500 position) and Illumina adapters (i.e. P5 in the forward and P7 in the reverse position)  
501 in the 5' position. The combinations of 24 i5-indexed primers and 36 i7-indexed  
502 primers made it possible to identify 864 different PCR products loaded onto the same  
503 MiSeq flowcell. Each index sequence differed from the others by at least two  
504 nucleotides, and each nucleotide position in the sets of indices contained  
505 approximately 25% of each base, to prevent problems due to Illumina low-diversity  
506 libraries (Table 1).

507 **DNA extraction and PCRs.** All laboratory manipulations were conducted with  
508 filter tips, under a sterile hood, in a DNA-free room. DNA was extracted with the  
509 DNeasy 96 Tissue Kit (Qiagen) with final elution in 200 µl of elution buffer. One  
510 extraction blank (NC<sub>ext</sub>), corresponding to an extraction without sample tissue, was  
511 systematically added to each of the eight DNA extraction microplates. DNA was  
512 quantified with a NanoDrop 8000 spectrophotometer (Thermo Scientific), to confirm

513 the presence of a minimum of 10 ng/ $\mu$ L of DNA in each sample. DNA amplification  
514 was performed in 5  $\mu$ L of Multiplex PCR Kit (Qiagen) Master Mix, with 4  $\mu$ L of  
515 combined i5 and i7 primers (3.5 $\mu$ M) and 2  $\mu$ L of genomic DNA. PCR began with an  
516 initial denaturation at 95°C for 15 minutes, followed by 40 cycles of denaturation at  
517 95°C for 20 s, annealing at 55°C for 15 s and extension at 72°C for 5 minutes,  
518 followed by a final extension step at 72°C for 10 minutes. PCR products (3  $\mu$ L) were  
519 verified by electrophoresis in a 1.5% agarose gel. One PCR blank (NC<sub>PCR</sub>),  
520 corresponding to the PCR mix with no DNA, was systematically added to each of the  
521 18 PCR microplates. DNA was amplified in replicate for all wild rodent samples  
522 ( $n=711$ ) (Table S1).

523 ***Library preparation and MiSeq sequencing.*** Two MiSeq (Illumina) runs  
524 were conducted, including PCR products from wild rodents and the positive and  
525 negative controls detailed in Figure 1 and Table S1. The MiSeq platform was chosen  
526 because it generates lower error rates than other HTS platforms [87]. The number of  
527 PCR products multiplexed was 823 for the first MiSeq run and 746 for the second  
528 MiSeq run (Table S1). Additional PCR products from other projects were added to  
529 give a total of 864 PCR products per run. PCR products were pooled by volume for  
530 each 96-well PCR microplate: 4  $\mu$ L for rodents and controls, and 1.5  $\mu$ L for bacterial  
531 isolates. Mixes were checked by electrophoresis in 1.5% agarose gels before their  
532 use to generate a “super-pool” of 864 PCR products for each MiSeq run. We  
533 subjected 100  $\mu$ L of each “super-pool” to size selection for the full-length amplicon  
534 (expected size: 375 bp including primers, indexes and adaptors), by excision in a  
535 low-melting agarose gel (1.25%) to discard non-specific amplicons and primer  
536 dimers. The PCR Clean-up Gel Extraction kit (Macherey-Nagel) was used to purify  
537 the excised bands. DNA was quantified by using the KAPA library quantification kit  
538 (KAPA Biosystems) on the final library before loading on a MiSeq (Illumina) flow cell  
539 (expected cluster density: 700-800 K/ $\text{mm}^2$ ) with a 500-cycle Reagent Kit v2  
540 (Illumina). We performed runs of 2 x 251 bp paired-end sequencing, which yielded  
541 high-quality sequencing through the reading of each nucleotide of the V4 fragments  
542 twice after the assembly of reads 1 and reads 2. The raw sequence reads (.fastq  
543 format) are available on request to the corresponding author.

544 ***Bioinformatic and taxonomic classification.*** MiSeq datasets were  
545 processed with mothur v1.34 [42] and with the MiSeq standard operating procedure



546 (SOP) [18]. We used the Silva SSU Reference database v119 [43] and the Silva  
547 taxonomy file for taxonomic assignment. The abundance table generated by mothur  
548 for each PCR product and each OTU was filtered as described in the Results section.  
549 The most abundant sequence for each OTU in each sample was extracted from the  
550 sequence dataset with a custom-written Perl script. The most abundant sequences  
551 for the 12 OTUs are available from GenBank (Accession Number KU697337 to  
552 KU697350). The sequences were aligned with reference sequences from bacteria of  
553 the same genus available from the SILVA SSU Ref NR database v119, using  
554 SeaView v4 [88]. The FASTA files used are available on request to the  
555 corresponding author. Phylogenetic trees were generated from the K2P distance with  
556 SeaView and species were identified on the basis of the “closest phylogenetic  
557 species”. We also used our sequences for blast analyses of GenBank, to identify the  
558 reference sequences to which they displayed the highest percentage identity.

559

## 560 **Acknowledgments**

561 This study was funded by the French National Institute for Agricultural Research  
562 (INRA) Meta-omics and microbial ecosystems metaprogram (Patho-ID project:  
563 Rodent and tick pathobiomes) and the ANR ENEMI (ANR-11-JSV7-0006). We would  
564 like to thank Virginie Dupuy for extracting DNA from bacterial cultures and Julie  
565 Sappa from Alex Edelman & Associates for improving the English writing. Analyses  
566 were performed on the CBGP HPC computational platform. The funders had no role  
567 in study design, data collection and analysis, the decision to publish, or preparation  
568 of the manuscript.

569

## 570 **Authors' contributions**

571 The study was conceived and designed by MG and JFC. MG, AL, CT, LT, HV and  
572 MR carried out the molecular biology procedures and validated the MiSeq data. MG,  
573 EB, MB and ADG contributed to the development of bioinformatics methods and  
574 validated taxonomic assignments. JFC and MTV coordinated the Patho-ID project  
575 and CB and NC coordinated the ENEMI project. MG, JFC, LT, CB and NC analyzed  
576 the data. MG and JFC wrote the manuscript. CB, NC, MR and MVT helped to draft

577 and to improve the manuscript. All the authors have read and approved the final  
578 manuscript.

579

## 580 **Supplementary materials**

581 **Table S1. Numbers of samples and numbers of PCRs for wild rodents and**  
582 **controls.** Negative Controls for dissection, NC<sub>mus</sub> ; Negative Controls for extraction, NC<sub>ext</sub> ; Negative  
583 Controls for PCR, NC<sub>PCR</sub> ; Negative Controls for indexing, NC<sub>index</sub> ; Positive Controls for PCR, PC<sub>PCR</sub> ;  
584 Positive Controls for Indexing, PC<sub>alien</sub>. See also Figure 1 for more details concerning negative controls  
585 (NC) and positive controls (PC).

586 **Table S2. The 50 most abundant OTUs in wild rodents and controls.**

587 **Table S3. Bacterial contaminants observed in negative and positive controls.**  
588 They were identified as contaminants on the basis of negative controls for extraction and PCR. Taxa  
589 in bold correspond to the sequences of DNA extracted from laboratory isolates.

590 **Table S4. Proportion of sequences and proportion of positive results removed**  
591 **at each step in data filtering.** Note that several positive results may be recorded for the same  
592 rodent in cases of co-infection.

593 **Table S5. Proportion of positive results for both PCR products at each step in**  
594 **data filtering.** Note that several positive results may be recorded for the same rodent in cases of  
595 co-infection.

596 **Table S6. Number of mismatches between PCR forward and reverse primers**  
597 **and 41,113 bacterial 16S rRNA V4 sequences of 79 zoonotic genera.** Data [1] was  
598 extracted from the Silva SSU database v119. Numbers of mismatches > 3 correspond to sequences  
599 of bad quality from different taxon. The number of mismatches in the 3' side of primers was always <2.

600 **Figure S1. Numbers of sequences of the positive controls for indexing**  
601 **PC<sub>Borrelia\_b</sub> (in blue) and PC<sub>Mycoplasma\_m</sub> (in red) in the various PCR products, with**  
602 **a dual-indexing design, for MiSeq runs 1 (a) and 2 (b).** The two PCRs for PC<sub>Borrelia\_b</sub>  
603 were performed with plate 9, positions A1 and E1 for run 1 and B1 and F1 for run 2, and the four  
604 PCRs for PC<sub>Mycoplasma\_m</sub> were performed with plate 9, positions C1, D1, G1 and H1 for the two runs.  
605 The numbers of sequences for the other wells correspond to indexing mistakes due to false index-  
606 pairing due to mixed clusters during the sequencing (see Table 1).

607 **Figure S2. Plots of the number of sequences (log (x+1) scale) from bacterial**  
608 **OTUs in both PCR replicates (PCR1 & PCR2) for the 356 wild rodents analyzed**  
609 **in the second MiSeq run.** Note that each rodent was tested with two duplicate PCRs. Green  
610 points correspond to rodents with two positive results after the filtering process; orange points  
611 correspond to rodents with one positive result and one negative result; and blue points correspond to  
612 rodents with two negative results. The light blue area and lines correspond to the threshold values  
613 used for the data filtering: samples below the lines are filtered out. See Figure S2 for plots  
614 corresponding to the second MiSeq run. See Figure 3 for plots corresponding to the first MiSeq run.

615 **Figure S3. Phylogenetic trees of the 16S rRNA V4 sequences for 12 pathogenic**  
616 **bacterial OTUs detected in wild rodents from Senegal.** Sequences boxed with an  
617 orange line were retrieved from African rodents and/or corresponds to positive controls (PC) for

618 *Borellia burgdorferi*, *Mycoplasma mycoides* and *Bartonella taylorii*. The other sequences were  
619 extracted from the SILVA database and GenBank. Trees include all lineages collected for *Rickettsia*,  
620 *Bartonella*, *Ehrlichia* and *Orientia*, but only lineages of the Spotted Fever Group for *Borrelia*, and  
621 lineages of the pneumonia group for *Mycoplasma*. The numbers indicated are the bootstrap values  
622 >55%. The Fasta files used are available on request to the corresponding author.  
623

## 624 References

- 625 1. **Taylor LH, Latham SM, Woolhouse ME.** 2001. Risk factors for human disease  
626 emergence. *Philos Trans R Soc Lond B Biol Sci* 356: 983-989.
- 627 2. **King DA, Peckham C, Waage JK, Brownlie J, Woolhouse ME.** 2006.  
628 Epidemiology. Infectious diseases: preparing for the future. *Science* 313:  
629 1392-1393.
- 630 3. **Grogan LF, Berger L, Rose K, Grillo V, Cashins SD, et al.** 2014. Surveillance  
631 for emerging biodiversity diseases of wildlife. *PLoS Pathog* 10: e1004015.
- 632 4. **Petrosino JF, Highlander S, Luna RA, Gibbs RA, Versalovic J.** 2009.  
633 Metagenomic pyrosequencing and microbial identification. *Clin Chem* 55: 856-  
634 866.
- 635 5. **Hugon P, Dufour JC, Colson P, Fournier PE, Sallah K, et al.** 2015. A  
636 comprehensive repertoire of prokaryotic species identified in human beings.  
637 *Lancet Infect Dis* 15: 1211-1219.
- 638 6. **Rynkiewicz EC, Hemmerich C, Rusch DB, Fuqua C, Clay K.** 2015.  
639 Concordance of bacterial communities of two tick species and blood of their  
640 shared rodent host. *Mol Ecol* 24: 2566-2579.
- 641 7. **Goffton AW, Doggett S, Ratchford A, Oskam CL, Papparini A, et al.** 2015.  
642 Bacterial Profiling Reveals Novel "Ca. Neoehrlichia", Ehrlichia, and  
643 Anaplasma Species in Australian Human-Biting Ticks. *PLoS One* 10:  
644 e0145449.
- 645 8. **Razzauti M, Galan M, Bernard M, Maman S, Klopp C, et al.** 2015. A  
646 Comparison between Transcriptome Sequencing and 16S Metagenomics for  
647 Detection of Bacterial Pathogens in Wildlife. *PLoS Negl Trop Dis* 9: e0003929.
- 648 9. **Cosson JF, Galan M, Bard E, Razzauti M, Bernard M, et al.** 2015. Detection of  
649 *Orientia* sp. DNA in rodents from Asia, West Africa and Europe. *Parasit*  
650 *Vectors* 8: 172.
- 651 10. **Vayssier-Taussat M, Moutailler S, Michelet L, Devillers E, Bonnet S, et al.**  
652 2013. Next generation sequencing uncovers unexpected bacterial pathogens  
653 in ticks in western Europe. *PLoS One* 8: e81439.
- 654 11. **Williams-Newkirk AJ, Rowe LA, Mixson-Hayden TR, Dasch GA.** 2014.  
655 Characterization of the bacterial communities of life stages of free living lone  
656 star ticks. (*Amblyomma americanum*). *PLoS One* 9: e102130.
- 657 12. **Williams-Newkirk AJ, Rowe LA, Mixson-Hayden TR, Dasch GA.** 2012.  
658 Presence, genetic variability, and potential significance of "Candidatus  
659 *Midichloria mitochondrii*" in the lone star tick *Amblyomma americanum*. *Exp*  
660 *Appl Acarol* 58: 291-300.

- 661 13. **Carpi G, Cagnacci F, Wittekindt NE, Zhao F, Qi J, et al.** 2011. Metagenomic  
662 profile of the bacterial communities associated with *Ixodes ricinus* ticks. *PLoS*  
663 *One* 6: e25604.
- 664 14. **Vaumourin E, Vourc'h G, Gasqui P, Vayssier-Taussat M.** 2015. The  
665 importance of multiparasitism: examining the consequences of co-infections  
666 for human and animal health. *Parasit Vectors* 8: 545.
- 667 15. **Tollenaere C, Susi H, Laine AL.** 2016. Evolutionary and Epidemiological  
668 Implications of Multiple Infection in Plants. *Trends Plant Sci* 21: 80-90.
- 669 16. **Vayssier-Taussat M, Albina E, Citti C, Cosson JF, Jacques MA, et al.** 2014.  
670 Shifting the paradigm from pathogens to pathobiome: new concepts in the  
671 light of meta-omics. *Front Cell Infect Microbiol* 4: 29.
- 672 17. **Vayssier-Taussat M, Kazimirova M, Hubalek Z, Hornok S, Farkas R, et al.**  
673 2015. Emerging horizons for tick-borne pathogens: from the 'one pathogen-  
674 one disease' vision to the pathobiome paradigm. *Future Microbiol* 10: 2033-  
675 2043.
- 676 18. **Kozich JJ, Westcott SL, Baxter NT, Highlander SK, Schloss PD.** 2013.  
677 Development of a dual-index sequencing strategy and curation pipeline for  
678 analyzing amplicon sequence data on the MiSeq Illumina sequencing  
679 platform. *Appl Environ Microbiol* 79: 5112-5120.
- 680 19. **Claesson MJ, Wang Q, O'Sullivan O, Greene-Diniz R, Cole JR, et al.** 2010.  
681 Comparison of two next-generation sequencing technologies for resolving  
682 highly complex microbiota composition using tandem variable 16S rRNA gene  
683 regions. *Nucleic Acids Res* 38: e200.
- 684 20. **Kuczynski J, Lauber CL, Walters WA, Parfrey LW, Clemente JC, et al.** 2012.  
685 Experimental and analytical tools for studying the human microbiome. *Nat Rev*  
686 *Genet* 13: 47-58.
- 687 21. **Sinha R, Chen J, Amir A, Vogtmann E, Shi J, et al.** 2015. Collecting Fecal  
688 Samples for Microbiome Analyses in Epidemiology Studies. *Cancer Epidemiol*  
689 *Biomarkers Prev*.
- 690 22. **Kircher M, Heyn P, Kelso J.** 2011. Addressing challenges in the production and  
691 analysis of illumina sequencing data. *BMC Genomics* 12: 382.
- 692 23. **Salter SJ, Cox MJ, Turek EM, Calus ST, Cookson WO, et al.** 2014. Reagent  
693 and laboratory contamination can critically impact sequence-based  
694 microbiome analyses. *BMC Biol* 12: 87.
- 695 24. **Horvath A, Peto Z, Urban E, Vagvolgyi C, Somogyvari F.** 2013. A novel,  
696 multiplex, real-time PCR-based approach for the detection of the commonly  
697 occurring pathogenic fungi and bacteria. *BMC Microbiol* 13: 300.
- 698 25. **Caro-Quintero A, Ochman H.** 2015. Assessing the Unseen Bacterial Diversity in  
699 Microbial Communities. *Genome Biol Evol* 7: 3416-3425.
- 700 26. **Walters W, Hyde ER, Berg-Lyons D, Ackermann G, Humphrey G, Parada A,**  
701 **Gilbert JA, Jansson JK, Caporaso JG, Fuhrman JA, Apprill A, Knight B.**  
702 2015. Improved bacterial 16S rRNA gene. V4 and V4-5. and fungal internal  
703 transcribed spacer marker gene primers for microbial community surveys.  
704 *mSystems* 1:e00009-15

- 705 27. **Kircher M, Sawyer S, Meyer M.** 2012. Double indexing overcomes inaccuracies  
706 in multiplex sequencing on the Illumina platform. *Nucleic Acids Res* 40: e3.
- 707 28. **Esling P, Lejzerowicz F, Pawlowski J.** 2015. Accurate multiplexing and filtering  
708 for high-throughput amplicon-sequencing. *Nucleic Acids Res* 43: 2513-2524.
- 709 29. **Bystrykh LV.** 2012. Generalized DNA barcode design based on Hamming  
710 codes. *PLoS One* 7: e36852.
- 711 30. **Meyerhans A, Vartanian JP, Wain-Hobson S.** 1990. DNA recombination during  
712 PCR. *Nucleic Acids Res* 18: 1687-1691.
- 713 31. **Paabo S, Irwin DM, Wilson AC.** 1990. DNA damage promotes jumping between  
714 templates during enzymatic amplification. *J Biol Chem* 265: 4718-4721.
- 715 32. **Odelberg SJ, Weiss RB, Hata A, White R.** 1995. Template-switching during  
716 DNA synthesis by *Thermus aquaticus* DNA polymerase I. *Nucleic Acids Res*  
717 23: 2049-2057.
- 718 33. **Lahr DJ, Katz LA.** 2009. Reducing the impact of PCR-mediated recombination in  
719 molecular evolution and environmental studies using a new-generation high-  
720 fidelity DNA polymerase. *Biotechniques* 47: 857-866.
- 721 34. **Whiteford N, Skelly T, Curtis C, Ritchie ME, Lohr A, et al.** 2009. Swift: primary  
722 data analysis for the Illumina Solexa sequencing platform. *Bioinformatics* 25:  
723 2194-2199.
- 724 35. **Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, et al.** 2012.  
725 Ultra-high-throughput microbial community analysis on the Illumina HiSeq and  
726 MiSeq platforms. *ISME J* 6: 1621-1624.
- 727 36. **Smith DP, Peay KG.** 2014. Sequence depth, not PCR replication, improves  
728 ecological inference from next generation DNA sequencing. *PLoS One* 9:  
729 e90234.
- 730 37. **Kilianski A, Haas JL, Corriveau EJ, Liem AT, Willis KL, et al.** 2015. Bacterial  
731 and viral identification and differentiation by amplicon sequencing on the  
732 MinION nanopore sequencer. *Gigascience* 4: 12.
- 733 38. **Callahan B, Proctor D, Relman D, Fukuyama J, Holmes S.** 2016.  
734 Reproducible Research Workflow in R for the Analysis of Personalized Human  
735 Microbiome Data. *Pac Symp Biocomput* 21: 183-194.
- 736 39. **Yilmaz P, Parfrey LW, Yarza P, Gerken J, Pruesse E, et al.** 2014. The SILVA  
737 and "All-species Living Tree Project. LTP)" taxonomic frameworks. *Nucleic*  
738 *Acids Res* 42: D643-648.
- 739 40. **Westcott SL, Schloss PD.** 2015. De novo clustering methods outperform  
740 reference-based methods for assigning 16S rRNA gene sequences to  
741 operational taxonomic units. *PeerJ* 3: e1487.
- 742 41. **Sinha R, Abnet CC, White O, Knight R, Huttenhower C.** 2015. The microbiome  
743 quality control project: baseline study design and future directions. *Genome*  
744 *Biol* 16: 276.
- 745 42. **Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, et al.** 2009.  
746 Introducing mothur: open-source, platform-independent, community-supported  
747 software for describing and comparing microbial communities. *Appl Environ*  
748 *Microbiol* 75: 7537-7541.

- 749 43. **Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, et al.** 2013. The SILVA  
750 ribosomal RNA gene database project: improved data processing and web-  
751 based tools. *Nucleic Acids Res* 41: D590-596.
- 752 44. **Gasparich GE, Whitcomb RF, Dodge D, French FE, Glass J, et al.** 2004. The  
753 genus *Spiroplasma* and its non-helical descendants: phylogenetic  
754 classification, correlation with phenotype and roots of the *Mycoplasma*  
755 *mycoides* clade. *Int J Syst Evol Microbiol* 54: 893-918.
- 756 45. **Gomez-Diaz E, Doherty PF, Jr., Duneau D, McCoy KD.** 2010. Cryptic vector  
757 divergence masks vector-specific patterns of infection: an example from the  
758 marine cycle of Lyme borreliosis. *Evol Appl* 3: 391-401.
- 759 46. **Sashida H, Sasaoka F, Suzuki J, Fujihara M, Nagai K, et al.** 2013. Two  
760 clusters among *Mycoplasma haemomuris* strains, defined by the 16S-23S  
761 rRNA intergenic transcribed spacer sequences. *J Vet Med Sci* 75: 643-648.
- 762 47. **Neimark H, Johansson KE, Rikihisa Y, Tully JG.** 2001. Proposal to transfer  
763 some members of the genera *Haemobartonella* and *Eperythrozoon* to the  
764 genus *Mycoplasma* with descriptions of '*Candidatus Mycoplasma haemofelis*',  
765 '*Candidatus Mycoplasma haemomuris*', '*Candidatus Mycoplasma haemosuis*'  
766 and '*Candidatus Mycoplasma wenyonii*'. *Int J Syst Evol Microbiol* 51: 891-899.
- 767 48. **Neimark H, Peters W, Robinson BL, Stewart LB.** 2005. Phylogenetic analysis  
768 and description of *Eperythrozoon coccoides*, proposal to transfer to the genus  
769 *Mycoplasma* as *Mycoplasma coccoides* comb. nov. and Request for an  
770 Opinion. *Int J Syst Evol Microbiol* 55: 1385-1391.
- 771 49. **Steer JA, Tasker S, Barker EN, Jensen J, Mitchell J, et al.** 2011. A novel  
772 hemotropic *Mycoplasma*. *hemoplasma*. in a patient with hemolytic anemia and  
773 pyrexia. *Clin Infect Dis* 53: e147-151.
- 774 50. **Pitcher DG, Nicholas RA.** 2005. *Mycoplasma* host specificity: fact or fiction? *Vet*  
775 *J* 170: 300-306.
- 776 51. **Trape JF, Diatta G, Arnathau C, Bitam I, Sarih M, et al.** 2013. The  
777 epidemiology and geographic distribution of relapsing fever borreliosis in West  
778 and North Africa, with a review of the *Ornithodoros erraticus* complex. *Acari:*  
779 *Ixodida*). *PLoS One* 8: e78473.
- 780 52. **Rar VA, Pukhovskaya NM, Ryabchikova EI, Vysochina NP, Bakhmetyeva**  
781 **SV, et al.** 2015. Molecular-genetic and ultrastructural characteristics of  
782 '*Candidatus Ehrlichia khabarensis*', a new member of the *Ehrlichia* genus.  
783 *Ticks Tick Borne Dis* 6: 658-667.
- 784 53. **Perlman SJ, Hunter MS, Zchori-Fein E.** 2006. The emerging diversity of  
785 *Rickettsia*. *Proc Biol Sci* 273: 2097-2106.
- 786 54. **Elliott SP.** 2007. Rat bite fever and *Streptobacillus moniliformis*. *Clin Microbiol*  
787 *Rev* 20: 13-22.
- 788 55. **Izzard L, Fuller A, Blacksell SD, Paris DH, Richards AL, et al.** 2010. Isolation  
789 of a novel *Orientia* species. *O. chuto* sp. nov.. from a patient infected in Dubai.  
790 *J Clin Microbiol* 48: 4404-4409.
- 791 56. **Buffet JP, Kosoy M, Vayssier-Taussat M.** 2013. Natural history of *Bartonella*-  
792 infecting rodents in light of new knowledge on genomics, diversity and  
793 evolution. *Future Microbiol* 8: 1117-1128.



- 794 57. **Mediannikov O, Aubadie M, Bassene H, Diatta G, Granjon L, Fenollar F.** 2014.  
795 Three new *Bartonella* species from rodents in Senegal. *Int J Infect Dis*  
796 21S:335.
- 797 58. **Julius R, Bastos A, Brettschneider H, Chimimba C.** 2012. Dynamics of  
798 Rodent-borne zoonotic diseases and their reservoir hosts: invasive *Rattus* in  
799 South Africa. Proc 25th Vertebrate Pest Conference, Monterey, California,  
800 USA.
- 801 59. **Bhatt KM, Mirza NB.** 1992. Rat bite fever: a case report of a Kenyan. *East Afr*  
802 *Med J* 69: 542-543.
- 803 60. **Gray HH.** 1967. Squirrel bite fever. *Trans R Soc Trop Med Hyg* 61:857.
- 804 61. **Laudisoit A, Falay D, Amundala N, Akaibe D, de Bellocq JG, et al.** 2014. High  
805 prevalence of *Rickettsia typhi* and *Bartonella* species in rats and fleas,  
806 Kisangani, Democratic Republic of the Congo. *Am J Trop Med Hyg* 90: 463-  
807 468.
- 808 62. **Dupont HT, Brouqui P, Faugere B, Raoult D.** 1995. Prevalence of antibodies to  
809 *Coxiella burnetti*, *Rickettsia conorii*, and *Rickettsia typhi* in seven African  
810 countries. *Clin Infect Dis* 21: 1126-1133.
- 811 63. **Ammar AM, Sabry MZ, Kirchhoff H.** 1980. Distribution of mycoplasmas in field  
812 and laboratory rodents in Egypt. *Z Versuchstierkd* 22: 216-223.
- 813 64. **Ndip LM, Labruna M, Ndip RN, Walker DH, McBride JW.** 2009. Molecular and  
814 clinical evidence of *Ehrlichia chaffeensis* infection in Cameroonian patients  
815 with undifferentiated febrile illness. *Ann Trop Med Parasitol* 103: 719-725.
- 816 65. **Kelly DJ, Foley DH, Richards AL.** 2015. A Spatiotemporal Database to Track  
817 Human Scrub Typhus Using the VectorMap Application. *PLoS Negl Trop Dis*  
818 9: e0004161.
- 819 66. **Konecny A, Estoup A, Duplantier JM, Bryja J, Ba K, et al.** 2013. Invasion  
820 genetics of the introduced black rat. *Rattus rattus*. in Senegal, West Africa.  
821 *Mol Ecol* 22: 286-300.
- 822 67. **Dalecky A, Ba K, Piry S, Lippens C, Diagne CA, Kane M, Sow A, Diallo M,**  
823 **Niang Y, Konecny A, Sarr N, Artige E, Charbonnel N, Granjon L,**  
824 **Duplantier JM, Brouat C.** 2015. Range expansion of the invasive house  
825 mouse *Mus musculus domesticus* in Senegal, West Africa: a synthesis of  
826 trapping data over three decades, 1983–2014. *Mammal Review* 45: 176-190
- 827 68. **Ranjan R, Rani A, Metwally A, McGee HS, Perkins DL.** 2016. Analysis of the  
828 microbiome: Advantages of whole genome shotgun versus 16S amplicon  
829 sequencing. *Biochem Biophys Res Commun* 469: 967-977.
- 830 69. **Glaeser SP, Kampfer P.** 2015. Multilocus sequence analysis. MLSA. in  
831 prokaryotic taxonomy. *Syst Appl Microbiol* 38: 237-245.
- 832 70. **Michelet L, Delannoy S, Devillers E, Umhang G, Aspan A, et al.** 2014. High-  
833 throughput screening of tick-borne pathogens in Europe. *Front Cell Infect*  
834 *Microbiol* 4: 103.
- 835 71. **Gutierrez R, Morick D, Cohen C, Hawlena H, Harrus S.** 2014. The effect of  
836 ecological and temporal factors on the composition of *Bartonella* infection in  
837 rodents and their fleas. *ISME J* 8: 1598-1608.

- 838 72. **Amaral-Zettler LA, McCliment EA, Ducklow HW, Huse SM.** 2009. A method  
839 for studying protistan diversity using massively parallel sequencing of V9  
840 hypervariable regions of small-subunit ribosomal RNA genes. *PLoS One* 4:  
841 e6372.
- 842 73. **Mueller RC, Gallegos-Graves LV, Kuske CR.** 2016. A new fungal large subunit  
843 ribosomal RNA primer for high-throughput sequencing surveys. *FEMS*  
844 *Microbiol Ecol* 92.
- 845 74. **Anthony SJ, Epstein JH, Murray KA, Navarrete-Macias I, Zambrana-Torrelío**  
846 **CM, et al.** 2013. A strategy to estimate unknown viral diversity in mammals.  
847 *MBio* 4: e00598-00513.
- 848 75. **Galan M, Pages M, Cosson JF.** 2012. Next-generation sequencing for rodent  
849 barcoding: species identification from fresh, degraded and environmental  
850 samples. *PLoS One* 7: e48374.
- 851 76. **Robasky K, Lewis NE, Church GM.** 2014. The role of replicates for error  
852 mitigation in next-generation sequencing. *Nat Rev Genet* 15: 56-62.
- 853 77. **Galan M, Guivier E, Caraux G, Charbonnel N, Cosson JF.** 2010. A 454  
854 multiplex sequencing method for rapid and reliable genotyping of highly  
855 polymorphic genes in large-scale studies. *BMC Genomics* 11: 296.
- 856 78. **Lange V, Bohme I, Hofmann J, Lang K, Sauter J, et al.** 2014. Cost-efficient  
857 high-throughput HLA typing by MiSeq amplicon sequencing. *BMC Genomics*  
858 15: 63.
- 859 79. **McFarlane R, Sleight A, McMichael T.** 2012. Synanthropy of wild mammals as a  
860 determinant of emerging infectious diseases in the Asian-Australasian region.  
861 *Ecohealth* 9: 24-35.
- 862 80. **Meerburg BG, Singleton GR, Kijlstra A.** 2009. Rodent-borne diseases and their  
863 risks for public health. *Crit Rev Microbiol* 35: 221-270.
- 864 81. **Civen R, Ngo V.** 2008. Murine typhus: an unrecognized suburban vectorborne  
865 disease. *Clin Infect Dis* 46: 913-918.
- 866 82. **Watt G, Parola P.** 2003. Scrub typhus and tropical rickettsioses. *Curr Opin Infect*  
867 *Dis* 16: 429-436.
- 868 83. **Vayssier-Taussat M, Moutailler S, Féménia F, Raymond P, Croce O, La**  
869 **Scola B, Fournier PE, Raoult D.** 2016. Identification of new zoonotic  
870 *Bartonella* species responsible for bacteremia in humans bitten by ticks.  
871 *Emerg Inf Dis* 22:
- 872 84. **Mills JN, Childs J, Ksiazek TG, Peters CJ, Velleca WM.** 1995. Methods for  
873 trapping and sampling small mammals for virologic testing. CDC, Atlanta.
- 874 85. **Herbreteau V, Rerkamnuaychoke W, Jittapalapong S, Chaval Y, Cosson JF,**  
875 **Morand S.** 2011. Field and laboratory protocols for rodent studies. Kasetsart  
876 University Press 46 p. <http://www.ceropath.org/research/protocols>
- 877 86. **Mebius RE, Kraal G.** 2005. Structure and function of the spleen. *Nat Rev*  
878 *Immunol* 5: 606-616.
- 879 87. **D'Amore R, Ijaz UZ, Schirmer M, Kenny JG, Gregory R, et al.** 2016. A  
880 comprehensive benchmarking study of protocols and sequencing platforms for  
881 16S rRNA community profiling. *BMC Genomics* 17: 55.

- 882 88. **Gouy M, Guindon S, Gascuel O.** 2010. SeaView version 4: A multiplatform  
883 graphical user interface for sequence alignment and phylogenetic tree  
884 building. *Mol Biol Evol* 27: 221-224.
- 885 89. **Parola P, Inokuma H, Camicas JL, Brouqui P, Raoult D.** 2001. Detection and  
886 identification of spotted fever group Rickettsiae and Ehrlichiae in African ticks.  
887 *Emerg Infect Dis* 7: 1014-1017.
- 888 90. **Qiu Y, Nakao R, Ohnuma A, Kawamori F, Sugimoto C.** 2014. Microbial  
889 population analysis of the salivary  
890  $\text{L}^{\text{M:P}}$  Phylogeny and biogeography of African Murinae based on  
891 mitochondrial and nuclear gene sequences, with a new tribal classification of  
892 the subfamily. *BMC Evol Biol* 8:199
- 893 92. **Reiczigel J.** 2003. Confidence intervals for the binomial parameter: some new  
894 considerations. *Statistics in Medicine* 22:611-621
- 895