



Title	A novel approach, based on BLSOMs (Batch Learning Self-Organizing Maps), to the microbiome analysis of ticks
Author(s)	Nakao, Ryo; Abe, Takashi; Nijhof, Ard M; Yamamoto, Seigo; Jongejan, Frans; Ikemura, Toshimichi; Sugimoto, Chihiro
Citation	The ISME Journal, 7(5), 1003-1015 https://doi.org/10.1038/ismej.2012.171
Issue Date	2013-03
Doc URL	http://hdl.handle.net/2115/53167
Type	article (author version)
File Information	ISME_Nakao.pdf



[Instructions for use](#)

**A novel approach, based on BLSOMs (Batch Learning Self-Organizing Maps), to the
microbiome analysis of ticks**

Ryo Nakao^{1,a}, Takashi Abe^{2,3,a}, Ard M. Nijhof⁴, Seigo Yamamoto⁵, Frans Jongejan^{6,7},
Toshimichi Ikemura², Chihiro Sugimoto¹

¹*Division of Collaboration and Education, Research Center for Zoonosis Control, Hokkaido University, Kita-20, Nishi-10, Kita-ku, Sapporo, Hokkaido 001-0020, Japan*

²*Nagahama Institute of Bio-Science and Technology, Nagahama, Shiga 526-0829, Japan*

³*Graduate School of Science & Technology, Niigata University, 8050, Igarashi 2-no-cho, Nishi-ku, Niigata 950-2181, Japan*

⁴*Institute for Parasitology and Tropical Veterinary Medicine, Freie Universität Berlin, Königsseg 67, 14163 Berlin, Germany*

⁵*Miyazaki Prefectural Institute for Public Health and Environment, 2-3-2 Gakuen Kibanadai Nishi, Miyazaki 889-2155, Japan*

⁶*Utrecht Centre for Tick-borne Diseases (UCTD), Department of Infectious Diseases and Immunology, Faculty of Veterinary Medicine, Utrecht University, Yalelaan 1, 3584 CL Utrecht, The Netherlands*

⁷*Department of Veterinary Tropical Diseases, Faculty of Veterinary Science, University of Pretoria, Private Bag X04, 0110 Onderstepoort, South Africa*

^aThese authors contributed equally to this work.

Keywords: BLSOMs/emerging diseases/metagenomics/microbiomes/symbionts/ticks

Running title: Tick microbiomes revealed by BLSOMs

Subject category: Microbe-microbe and microbe-host interactions

Abstract

Ticks transmit a variety of viral, bacterial and protozoal pathogens, which are often zoonotic. The aim of this study was to identify diverse tick microbiomes, which may contain as-yet unidentified pathogens, using a metagenomic approach. DNA prepared from

5 bacteria/archaea-enriched fractions obtained from seven tick species, namely *Amblyomma testudinarium*, *Amblyomma variegatum*, *Haemaphysalis formosensis*, *Haemaphysalis longicornis*, *Ixodes ovatus*, *Ixodes persulcatus* and *Ixodes ricinus*, was subjected to pyrosequencing after whole-genome amplification. The resulting sequence reads were phylotyped using a Batch Learning Self-Organizing Map (BLSOM) program, which allowed phylogenetic estimation based

10 on similarity of oligonucleotide frequencies, and functional annotation by BLASTX similarity searches. In addition to bacteria previously associated with human/animal diseases, such as *Anaplasma*, *Bartonella*, *Borrelia*, *Ehrlichia*, *Francisella* and *Rickettsia*, BLSOM analysis detected microorganisms belonging to the phylum Chlamydiae in some tick species. This was confirmed by pan-Chlamydia PCR and sequencing analysis. Gene sequences associated with

15 bacterial pathogenesis were also identified, some of which were suspected to originate from horizontal gene transfer. These efforts to construct a database of tick microbes may lead to the ability to predict emerging tick-borne diseases. Furthermore, a comprehensive understanding of

tick microbiomes will be useful for understanding tick biology, including vector competency and interactions with pathogens and symbionts.

20

Introduction

Ticks (Acari: Ixodida) are amongst the most well-known blood-sucking arthropods. They can transmit a variety of pathogens, including viruses, bacteria and protozoa, to humans and animals. Opportunities for human/animal contact in natural foci of tick-borne infection are increasing, along with expanded distributions of vector ticks due to environmental changes (Beugnet and Marié 2009; Parola 2004). Meanwhile, reports of novel tick-borne pathogens are also increasing (Spolidorio *et al.*, 2010; Yu *et al.*, 2011). Above all, many of the emerging rickettsial pathogens identified in recent decades were first discovered in ticks and subsequently recognized as the causative agents of animal or human diseases (Parola *et al.*, 2005). This indicates that active surveillance for potential pathogens in ticks may be a feasible approach to predict and combat emerging tick-borne diseases. Given the diversity of tick species and their wide geographic distribution (Bowman and Nuttall, 2008), it is likely that as-yet unidentified pathogens are harboured by ticks.

35 Morphological detection by electron microscopy has shown that *Rickettsia*- and
Wolbachia-like bacteria infect a variety of tick species (Lewis, 1979; Roshdy, 1968), although
these bacteria are not necessarily pathogenic in either mammalian hosts or ticks. The availability
of molecular genetic tools such as PCR and real-time PCR has increased the number of reports on
a variety of microorganisms in ticks (Benson *et al.*, 2004; Noda *et al.*, 1997; van Overbeek *et al.*,
40 2008); however, their biological roles and potential to cause disease in mammals remain poorly
understood. Interactions between arthropods and their symbionts have been relatively well
studied in aphids. A variety of symbionts are reported to have detrimental, neutral or beneficial
effects on the aphid hosts; for example, through pathogenicity (Grenier *et al.*, 2006), providing
nutrition (Nakabachi and Ishikawa, 1999), altering gene expression associated with body colour
45 (Tsuchida *et al.*, 2010) and, most intriguingly, displaying protective functions against pathogens
and predators (Scarborough *et al.*, 2005). It has been suggested that certain symbiotic microbes
might also affect the vector competency of their arthropod host (Weiss and Aksoy, 2011). A full
characterization of microorganisms harboured by ticks will provide valuable information towards
a better understanding of tick biology in terms of the relationship between ticks and their
50 microbiomes, including microbes that are potential pathogens in mammals.

Recent progress in second generation sequencing technologies has led to the discovery of
previously unknown organisms in a variety of samples, including arthropods (Cox-Foster *et al.*,

2007; Suen *et al.*, 2010; Warnecke *et al.*, 2007), since uncultured microorganisms represent the vast majority of symbiotic, commensal and pathogenic microorganisms. Attempts to assess the microbial diversity in two tick species, *Rhipicephalus microplus* and *Ixodes ricinus*, using novel sequencing technologies, resulted in the identification of microbes (including bacterial species) previously unreported in ticks (Andreotti *et al.*, 2011; Carpi *et al.*, 2011). However, since the data processing methods used in these two studies were based on conventional homology-based sequence searches, the discoveries were limited to microorganisms for which genomic information was already available in the database used. An alternative approach to characterizing DNA fragments might offer improvements for investigating microbial diversity in ticks that are expected to harbour poorly characterized microorganisms.

Batch Learning Self-Organizing Map (BLSOM) is a bioinformatics method that follows a learning process and generates a map independently of the order of data input (Abe *et al.*, 2003; Kanaya *et al.*, 2001). This method was designed to separate and cluster genomic sequence fragments based on the similarity of oligonucleotide frequencies without any other taxonomical information. Since BLSOM does not require orthologous sequence sets and sequence alignments, it could provide a new systematic strategy for revealing the microbial diversity and relative abundance of different phylotype members of uncultured microorganisms in a wide variety of environmental metagenomic samples (Abe *et al.*, 2005; Uehara *et al.*, 2011). Of note, similar

SOM-based methods were recently used to obtain clear phylotype-specific classification of metagenomic sequences (Chan *et al.*, 2008; Dick *et al.*, 2009; Martin *et al.*, 2008; Weber *et al.*, 2011).

In the present study, we used BLSOM to characterize diverse microbiomes in seven tick
75 species. The ticks harboured a variety of bacteria, including those previously reported to be associated with human and animal diseases, and others not previously reported from ticks. This study was intended as a first step towards the discovery of emerging tick-borne pathogens and a systematic understanding of the relationship between ticks and their microbiomes.

80

Materials and methods

Tick species.

Ticks were collected from the field by dragging flannel sheets over the vegetation. Three tick species, namely *Amblyomma testudinarium* (AT), *Haemaphysalis formosensis* (HF) and
85 *Haemaphysalis longicornis* (HL) were collected in March 2009 in Miyazaki, Japan. This is an area where *Rickettsia japonica*, the causative agent of Japanese spotted fever and *Rickettsia tamurae*, the spotted fever group rickettsia recently associated with human disease (Imaoka *et al.*,

2011), are known to be endemic (Morita *et al.*, 1990; S Yamamoto, personal communication).

Two *Ixodes* species, *Ixodes ovatus* (IO) and *Ixodes persulcatus* (IP), were collected in June 2010
90 in Hokkaido, Japan, where *Borrelia burgdorferi sensu lato* (*s.l.*) and *Anaplasma*
phagocytophilum, the causative agents of Lyme disease and human granulocytic anaplasmosis,
respectively, are endemic (Miyamoto *et al.*, 1992; Murase *et al.*, 2011). *Ixodes ricinus* (IR) was
collected in August 2010 in Soesterberg, The Netherlands, where Lyme disease is endemic
(Hofhuis *et al.*, 2006). *Amblyomma variegatum* (AV) ticks were collected in The Gambia in 2005
95 and subsequently maintained under laboratory conditions at the Utrecht Centre for Tick-borne
Diseases, The Netherlands. Adult AV ticks from the second laboratory generation were used in
this study. The collected ticks were pooled according to species, life stage and sex (except for HF
and HL) and specific pool IDs (AVf, AVm, IOf, IOm, IPf, IPm, IRf, ATn, HFfmm and HLfmm)
were provided as indicated in Table 1 (“f”, “m” and “n” represent female, male and nymph,
100 respectively).

DNA preparation.

The live ticks were washed twice with 70% ethanol supplemented with 1% povidone-iodine (Meiji Seika, Tokyo, Japan) solution and rinsed three times with distilled water to decontaminate the tick body surface. Ticks were then ground with 4.8 mm stainless steel beads

105 (TOMY, Tokyo, Japan) using the Micro Smash MS-100R (TOMY). Tick homogenates were treated with a membrane lysis buffer [10 mM Tris-HCl (pH 8.0), 150 mM NaCl, 10 mM MgCl₂, 0.1% IGEPAL CA-630 (Sigma Chemical Co., St. Louis, MO)] for 10 min on ice and centrifuged at 400 × g for 25 min. The supernatant was centrifuged at 20,000 × g for 30 min and the pellet containing the bacterial/archaeal cells was resuspended in a DNase buffer [25 mM Tris-HCl (pH 110 8.0), 10 mM MgCl₂, 0.8 U/μL DNase I (Takara, Shiga, Japan)] to remove any contaminating host-cell DNA. After 60 min at 37°C, the reaction was stopped by adding 25 mM EDTA (pH 8.0). The solution was filtered through a 5.0 μm pore-size membrane filter (Millipore, Bedford, MA) by centrifugation at 1,000 × g. The filtrate was centrifuged at 20,000 × g for 30 min and the resulting pellet was washed twice with PBS. All centrifugation steps were performed at 4°C. To 115 lyse a wide range of bacterial/archaeal species, this study utilized an achromopeptidase, which is reported to have broad-spectrum bacteriolytic activity (Ezaki and Suzuki, 1982). The bacterial/archaeal pellet was treated with an achromopeptidase buffer [10 mM Tris-HCl (pH 8.0), 10 mM NaCl, 1 U/μL achromopeptidase (WAKO, Osaka, Japan)] for 60 min at 37°C. Bacterial/archaeal DNA was then extracted using the NucleoSpin Tissue XS Kit (Macherey- 120 Nagel, Düren, Germany) according to the manufacturer's instructions. The genomic DNA was then subjected to whole-genome amplification using the GenomiPhi™ V2 DNA Amplification Kit (GE Healthcare, Chalfont St Giles, UK) according to the manufacturer's instructions. The

DNA concentration was measured using the Quant-iT dsDNA BR assay with a Qubit
Fluorometer (Invitrogen, Carlsbad, CA). All procedures were conducted under sterile conditions
125 in a flow cabinet. A schematic flow diagram showing the extraction process of bacterial/archaeal
DNA from ticks is presented in Figure 1.

Pyrosequencing.

Genomic DNA from each tick pool was subjected to pyrosequencing on a Roche/454
Genome Sequencer FLX Titanium (Roche Applied Science/454 Life Science, Brandford, CT) at
130 Hokkaido System Science (Sapporo, Japan). Two GS FLX shotgun libraries were prepared using
a standard protocol from the manufacturer and analysed in four independent sequencing runs: the
first library contained three pools (ATn, HFfmm and HLfmm) with different multiplex identifiers
and was analysed on 1/8 of the PicoTiterPlate (Roche Applied Science/454 Life Science), while
the second library contained seven pools (AVf, AVm, IOf, IOm, IPf, IPm and IRf) and was
135 analysed on 1/4 of the plate. The metagenomic sequences were deposited in the DNA Data Bank
of Japan (DDBJ) (<http://www.ddbj.nig.ac.jp>) Sequence Read Archive under the accession no.
DRA000590.

BLSOM analysis.

Self-Organizing Map (SOM) is an unsupervised neural network algorithm that
140 implements a characteristic nonlinear projection from the high-dimensional space of input data
onto a two-dimensional array of weight vectors (Kohonen, 1990). We used the ‘Batch Learning
SOM’ (BLSOM), which is a modified version of the conventional SOM for genome informatics
that makes the learning process and creation of the resulting map independent of the order of data
input (Abe *et al.*, 2003; Kanaya *et al.*, 2001). The initial weight vectors were defined by principal
145 component analysis instead of random values. BLSOM learning was conducted as described
previously (Abe *et al.*, 2003), and the BLSOM program was obtained from UNTROD Inc., Japan
(y_wada@nagahama-i-bio.ac.jp).

To estimate phylotypes of the metagenomic sequences, three types of large-scale
BLSOMs, namely Kingdom-, Bacteria/Archaea- and Genus group-BLSOM, were constructed in
150 advance, using sequences deposited in DDBJ/EMBL/GenBank as previously described (Abe *et al.*
et al., 2003). The clustering conditions were adopted from the previously optimised parameters to
minimise the computation time without loss of clustering power (Abe *et al.*, 2005). Kingdom-
BLSOM was constructed with tetranucleotide frequencies for 5-kb sequences from the whole-
genome sequences of 111 eukaryotes, 2,813 bacteria/archaea, 1,728 mitochondria, 110
155 chloroplasts and 31,486 viruses. To obtain more detailed phylotype information for
bacterial/archaeal sequences, Bacteria/Archaea- and Genus group-BLSOM were constructed with

a total of 3,500,000 5-kb sequences from 3,157 species, for which at least 10 kb of sequence was available from DDBJ/EMBL/GenBank.

Mapping of metagenomic sequences longer than 300 bp on Kingdom-BLSOMs, after
160 normalization of the sequence length, was conducted by finding the lattice point with the
minimum Euclidean distance in the multidimensional space. To identify further detailed
phylogenies of the metagenomic sequences that had been mapped to the bacterial/archaeal
territories on Kingdom-BLSOM, these sequences were successively mapped on
Bacteria/Archaea-BLSOM. Similar stepwise mappings of metagenomic sequences on the Genus
165 group-BLSOM were subsequently conducted to obtain further detailed phylogenetic information.
In order to evaluate the accuracy of BLSOMs for taxonomic classification of metagenomic
sequences, the BLSOM analysis was applied to three previously published simulated datasets of
varying complexities (Mavromatis *et al.*, 2007). When sequences longer than 300 bp were
mapped, approximately 90%, 70% and 40%-50% were correctly classified to the kingdom,
170 phylum and genus levels, respectively (Supplementary Figure S1).

Pan-Chlamydia PCR, sequencing and data analysis.

A pan-Chlamydia PCR was conducted on the tick species in which bacteria belonging to
the phylum Chlamydiae had been detected by BLSOM analysis. PCR amplifications were

performed using three different sets of primers targeting the 16S rRNA gene (Supplementary
175 Table 1). The successfully amplified products were cloned into a pGEM-T vector (Promega,
Madison, WI). Each plasmid clone was sequenced using the BigDye Terminator version 3.1
Cycle Sequencing Kit (Applied Biosystems, Foster City, CA) and an ABI Prism 3130x genetic
analyzer (Applied Biosystems) according to the manufacturer's instructions. The DNA sequences
obtained were submitted to the DDBJ under accession nos. AB725685 to AB725705.

180 The obtained sequences were aligned using ClustalW2 (Larkin *et al.*, 2007) and the
resulting alignments were used to generate an uncorrected pairwise distance matrix using the
MOTHUR program version 1.25.0 (Schloss *et al.*, 2009). An average neighbour clustering
algorithm in the MOTHUR program was used to assign sequences into operational taxonomic
units (OTUs), with a 97% sequence similarity cut-off. A representative sequence from each OTU
185 was taxonomically identified using the Classifier tool available on the Ribosomal Database
Project (RDP) database (Cole *et al.*, 2009) and was compared against the NCBI GenBank
database using a BLASTn search.

Functional annotation and characterization.

The metagenomic sequences were annotated using BLASTX (Altschul *et al.*, 1990)
190 against the Clusters of Orthologous Groups (COG) database (Tatusov *et al.*, 1997) with an E-

value threshold of 1×10^{-5} and were classified according to COG functional categories. We also identified putative virulence-associated factors using the MvirDB database (Zhou *et al.*, 2007), which consists of sequence information from multiple microbial databases of protein toxins, virulence factors and antibiotic resistance genes. Sequences associated with “antibiotic
195 resistance”, “pathogenicity islands” and “virulence protein” were further phylotyped using BLSOM as described above. All sequence reads were used to identify putative virulence-associated factors in BLASTX-based analyses, while only sequences longer than 300 bp were phylotyped.

200

Results

Kingdom classification.

Bacterial/archaeal cells in tick homogenates were purified using centrifugation, DNase treatment and filtration. DNA was extracted from bacteria/archaea-enriched fractions and was
205 subjected to pyrosequencing after whole-genome amplification. The total sequence reads for each species ranged from 21,213 (HFfmm) to 42,258 (AVf), with average lengths between 143.5 bp (IPm) and 391.4 bp (AVm) (Table 2). The sequence reads longer than 300 bp were classified as

bacterial/archaeal, eukaryotic, virus, mitochondrial or chloroplast using Kingdom-BLSOM (Figure 2). The percentage of the reads categorized as bacteria/archaea ranged from 37.8% (AVm) to 72.3% (HLfmn), indicating that the present protocol could efficiently enrich bacterial/archaeal cells from tick homogenates. An average of 35% of the eukaryotic reads were clustered together with those of other arthropods, which were included in the Kingdom-BLSOM construction, while the remaining reads were with a variety of eukaryotes including protists and mammals (data not shown).

215

Phylum classification.

The bacterial/archaeal sequence reads were further classified at the phylum level using Bacteria/Archaea-BLSOM. Between 97.2% (AVf) and 99.4% (HLfmn) of the reads were successfully assigned to specific phyla (Table 2). The number of different phyla obtained from a single library ranged from 20 (HLfmn) to 28 (AVm, IOf and IPf). The composition of bacterial/archaeal phyla in each tick pool is shown in Figure 3. Each tick species had a different microbial composition. Differences were also found between males and females within the same tick species. Sequences classified into the phyla Firmicutes and Gammaproteobacteria constituted nearly half of the total sequences in most tick species. Sequences classified into the phylum

225 Alphaproteobacteria were also observed in all tick species, while those classified into the phylum Chlamydiae were observed at high abundance only in IPf, IPm and HFfmm.

Genus classification.

Between 85.4% (AVm) and 92.3% (ATn) of the sequences characterized into certain bacterial/archaeal phyla were successfully characterized to the genus level using Genus group-
230 BLSOM (Table 2). The number of different genera obtained from a single library ranged from 133 (IRf) to 223 (AVm). The twenty most prevalent bacterial/archaeal genera within each tick species are listed in Table 3, while a complete list of bacterial/archaeal genera recovered from each tick pool is shown in Supplementary Table 2. A total of 64 different genera were detected and some of the dominant bacterial/archaeal genera were commonly observed between tick
235 species. Some of the sequences could not be assigned to known taxa. For example, 4.7% of the sequences in AVm were characterized to the phylum Euryarchaeota, but could not be characterized at the genus level (Table 3). In accordance with the results of phylum classification, many of the dominant bacteria belonged to genera within the phyla Firmicutes and Gammaproteobacteria. The sequences of known tick-borne pathogens such as *Anaplasma*,
240 *Bartonella*, *Borrelia*, *Ehrlichia*, *Francisella* and *Rickettsia*, as well as those of known tick-symbionts such as *Coxiella*, *Rickettsiella* and *Wolbachia*, were also listed as top 20 genera (Table

3). Only three pools (IPf, IPm and HFfmm) contained bacterial sequences classified into the family *Chlamydiaceae* (genera *Chlamydia* and *Chlamydophila*) as dominant taxa/populations.

Pan-Chlamydia PCR and sequencing analysis.

245 DNA extracted from five tick pools, IPf, IPm, HFf, HFm and HFn, was tested by pan-Chlamydia PCR using three different primer sets. Only one primer set (16S FOR2 and 16S REV2) produced PCR products of the expected size (approximately 260 bp) in all tested samples (data not shown). A total of 14 different clones of partial 16S rRNA gene sequences were recovered, four of which were identified in both tick species (Table 4). The MOTHUR program
250 assigned these 14 sequences into five OTUs, the representative sequences of which were predicted to belong to the genera *Neochlamydia*, *Parachlamydia* or *Simkania* based on RDP database classification. BLASTn analysis of a representative sequence for each OTU showed highest similarity (88%–98%) to Chlamydial 16S rRNA gene sequences recovered from various animal sources including cockroach, sea bass, koala, cat and leafy seadragon (Table 4).

255 **Functional annotation.**

Functional annotation of the bacterial/archaeal sequences using BLASTX against the COG database indicated similar functional gene compositions between pools, except for AVm (Supplementary Figure 2). More than 80% of the AVm sequences were annotated as “Replication, recombination and repair”-related genes. This result may indicate that the sequences in this

260 library were not randomly amplified during the whole-genome amplification process. Sequences associated with putative virulence-associated factors were identified using BLASTX against the MvirDB database. The number of identified sequences differed greatly between libraries: IOm showed the highest number (6,715) and ATn the lowest (365) (Figure 4). Most of these were categorized as either “antibiotic resistance”, “pathogenicity island” or “virulence protein”-related
265 genes. Both IOm and IPf contained a large number of sequences associated with “virulence protein” ($n = 6,515$ and $2,189$, respectively), in comparison to other pools. Sequences associated with “pathogenicity island” were observed to some extent in all pools, with the highest number obtained from IPf ($n = 2,072$), followed by IOf, AVm and HLfmm. The origins of these virulence-associated sequences were predicted using Bacteria/Archaea-BLSOM, and five of the most
270 dominant origins of the sequences associated with “antibiotic resistance”, “pathogenicity island” and “virulence protein” in each tick pool are shown in Table 5. The most common phyla were Firmicutes, Alphaproteobacteria, Gammaproteobacteria and Chlamydiae, while many of the sequences could not be phylotyped by BLSOM and remained unclassified.

275

Discussion

The present study investigated the diversity of microbial communities in ticks using
280 pyrosequencing technology coupled with a composition-based data processing approach called
BLSOM. In order to provide an overview of microbiomes associated with different tick species,
ticks were pooled according to species and subjected to the analyses. We showed that ticks
harboured a variety of bacteria, including those previously associated with human/animal
diseases such as *Anaplasma*, *Bartonella*, *Borrelia*, *Ehrlichia*, *Francisella* and *Rickettsia*, as well
285 as potential pathogens such as *Chlamydia*. This is a first attempt to apply BLSOM to the
detection of potential pathogens. Since this approach can be directly applied to other vector
arthropods of medical and veterinary importance, it has great potential for the detection of diverse
microbes, including as-yet unidentified microorganisms, and to pre-empt emerging infectious
diseases.

290 Unlike the conventional sequence homology searches used in most metagenomic studies,
BLSOM does not require orthologous sequences for phylogenetic classification of metagenomic
sequences. This is a great advantage, especially when applied to a group composed of poorly
characterized microorganisms (Abe *et al.*, 2005). In addition, ticks have highly variable genome
sizes, ranging from nearly one third to over two times the size of the human genome, but only
295 one whole-genome sequencing project is underway (Nene, 2009). This indicates that alignment-

based techniques may encounter difficulties in sorting out tick-derived sequences in the metagenomic libraries. On the other hand, BLSOM can separate bacterial/archaeal sequences from those of eukaryotes with high accuracy (Supplementary Figure S1) and thus have great potential for detecting symbiotic and commensal microorganisms in eukaryotic hosts with scarce or no genomic information. Despite the fact that Kingdom-BLSOM did not include tick genome sequences as these were not available at the time of construction, an average of 35% of the eukaryotic sequences were clustered together with those of other arthropods and thus likely to be derived from tick cells. Once it becomes available, the inclusion of a tick genome in BLSOM construction would facilitate the identification of eukaryotic reads as being tick-derived.

305 In a previous pyrosequencing study, in which the whole bodies of *I. ricinus* were directly subjected to DNA and RNA extraction, only 0.095% of the genomic DNA and 0.30% of the cDNA-derived sequences were assigned to known bacterial taxa (Carpi *et al.*, 2011). Although this shotgun strategy seems ideal for capturing a snapshot of a microbial community with minimum bias introduced during sample preparation, there is a concern that genomic data for 310 bacteria/archaea present in very low numbers may be masked by the huge amount of host-genome sequences. To address this, our strategy included a process of microbial enrichment using centrifugation, DNase treatment and filtration, to enrich bacterial/archaeal cells from tick homogenates prior to pyrosequencing. The microbial profiles obtained in this study, therefore,

may not completely reflect the whole picture regarding microbial composition in the tested ticks
315 due to the biases introduced during the microbial enrichment and whole-genome amplification
processes, especially in the case of AVm, where the sequences in the library were not likely to be
randomly amplified (Supplementary Figure 2). Nonetheless, the fact that nearly half of the
sequence reads were identified as bacterial/archaeal, and were associated with a diverse range of
bacterial/archaeal phyla, indicated sufficient coverage by our method to obtain an overview of
320 tick microbiomes.

Sequences of genera containing several tick-borne pathogens, such as *Anaplasma*,
Borrelia and *Rickettsia*, were present in the tested samples (Table 3 and Supplementary Table 2).
Further identification to the species level can be achieved by the use of traditional methods such
as conventional species-specific PCR. This can be done to confirm the presence of tick-borne
325 bacterial species such as *A. phagocytophilum*, *B. burgdorferi s.l.*, *R. japonica* and *R. tamurae*,
which have previously been detected in the tested tick species (Hofhuis *et al.*, 2006; Miyamoto *et al.*,
et al., 1992; Morita *et al.*, 1990; Murase *et al.*, 2011). The data obtained in the present study,
however, are useful for broadening our understanding of phylogenetic diversity among bacterial
genera, and the extent of their habitats. For example, the genera *Chlamydia* and *Chlamydophila*
330 were predicted to be harboured by IPf, IPm and HFfmn in high abundance (Table 3). Further pan-
Chlamydia PCR and sequencing analysis supported the existence of microorganisms possibly

related to the genera *Neochlamydia*, *Parachlamydia* and *Simkania* (Table 4). Some representatives of these genera are implicated in diseases affecting their host animals (Corsaro *et al.*, 2007; Meijer *et al.*, 2006; von Bomhard *et al.*, 2003). Since only the sequences of the genera *Chlamydia* and *Chlamydophila* were available from the phylum Chlamydiae at the time of BLSOM construction, those of the genera *Neochlamydia*, *Parachlamydia* and *Simkania* were not included in Genus group-BLSOM, so sequences from the phylum Chlamydiae were allocated to either genus *Chlamydia* or *Chlamydophila*. This technical limitation will be overcome as sequence data from a wider range of microorganisms become available. Although possible transmission of Chlamydial bacteria to cattle and humans through ticks was reported in several previous studies (Caldwell and Belden, 1973; Facco *et al.*, 1992; McKercher *et al.*, 1980), there has been no molecular genetic evidence to support these observational studies. Thus, this is the first molecular genetic evidence supporting the presence of Chlamydial organisms in ticks. Given the wide genetic diversity and host range, including several arthropods, of members of this phylum of bacteria (Horn, 2008), further research is warranted to improve the understanding of their biology, especially with respect to their pathogenic potential for humans and animals.

Genus group-BLSOM identified 133 to 223 different genera from each single tick pool (Table 2). The presence of a variety of bacterial/archaeal genera in ticks was demonstrated in previous metagenomic studies; in *R. microplus* and *I. ricinus*, 120 and 108 different genera were

350 identified, respectively, using a rDNA-based amplicon pyrosequencing approach (Andreotti *et al.*,
2011; Carpi *et al.*, 2011). This discrepancy may be explained by the fact that these two
approaches rely on different analytic principles. One possible explanation is that BLSOM might
have misclassified genome segments introduced through horizontal gene transfer (HGT) (Abe *et al.*
al., 2005). Tamames and Moya (2008) reported evidence of HGT in 0.8%–1.5% and 2%–8% of
355 the sequences in environmental metagenomic libraries using phylogenetic and composition-based
methods, respectively, indicating that the present approach may require further refinement; for
example, by combining BLSOM with rDNA-based phylogenetic classification as proposed
elsewhere (Abe *et al.*, 2005). It is also possible that the Genus group-BLSOM has overestimated
the microbial diversity in the tested samples as a result of a decreased classification specificity on
360 the genus level (Supplementary Figure S1). This drawback can be remedied by using only longer
sequences to improve the specificity of the analysis or by using more genome sequences
associated with the habitat under study during BLSOM construction (Abe *et al.*, 2005; Weber *et al.*,
al., 2011).

The present study included the ticks maintained under laboratory conditions. The
365 microbial diversity of those ticks was comparable to that of those collected in the field in terms of
the number of different genus identified (Table 2). Environmental factors and developmental
stages were reported to affect microbial communities in ticks (Andreotti *et al.*, 2011; Carpi *et al.*,

2011; Clay *et al.*, 2008; van Overbeek *et al.*, 2008). Therefore, a further comparison of microbial diversity between laboratory colony and field ticks of the same species may be useful to elucidate environmental factors playing key roles on microbial compositions. IO and IP were sampled at the same site and time and have a common host preference for sika deer (*Cervus nippon*) in the sampling area (Isogai *et al.*, 1996), allowing a comparison of microbial compositions between these two tick species. IO harboured *Rickettsia* and *Ehrlichia*, but IP harboured *Bartonella*, at high abundance, despite the presence of *Borrelia* and *Francisella* in both tick species at similar abundance (Table 3). These findings may be useful for improving the understanding of the vector potential of different tick species. In addition to these species differences, microbial communities are expected to vary between individual ticks as reported elsewhere (Andreotti *et al.*, 2011; Carpi *et al.*, 2011). One aspect to be addressed, which may help to explain these differences, is the interaction between microbial lineages. Indeed, antagonistic interactions between microbial lineages in ticks, including pathogenic *Rickettsia*, were demonstrated previously (de la Fuente *et al.*, 2003; Macaluso *et al.*, 2002). Thus, one of the most interesting challenges for future studies on tick microbiomes might be the association between microbial composition in ticks and the potential for pathogen transmission. The pictures of whole microbiomes harboured by each tick species generated in the present study provide the foundation for a statistical correlation analysis between pathogens and certain microbial lineages in individual ticks. Such studies could lead to

the identification of microbiological factors affecting the prevalence and persistence of pathogenic lineages in ticks.

The sequences associated with putative virulence-associated factors were identified from metagenomic libraries using BLASTX-based methods (Figure 4). It should be noted, however, that some of the factors, such as type III secretion systems, are responsible for both bacterial pathogenesis and symbiosis, including mutualism and commensalism (Coombes, 2009), and thus not all factors are directly or necessarily related to human/animal pathogens. The present study predominantly identified sequences related to antibiotic resistance (Figure 4). It has been recognized for some time that antibiotic resistance genes are not only associated with pathogens but also with a wide variety of environmental microbes (Allen *et al.*, 2010). Several lines of evidence indicate that some of the clinically relevant resistance mechanisms in pathogenic organisms originated from such environmental bacteria through HGT (Wright, 2010). Direct contact between pathogens and resistance genes in commensals harboured by ticks could therefore represent a mechanism for the emergence of resistant pathogens.

Several studies described mutualistic relationships between microbes and host arthropods (Currie *et al.*, 1999; Kaltenpoth *et al.*, 2005; Scott *et al.*, 2008). It is plausible that the tick microbiome also affects aspects of tick physiology. This was for example demonstrated by a reduction in reproductive fitness of *Amblyomma americanum* females following antibiotic

treatment, suggesting that the microbe plays a role in tick growth and development (Zhong *et al.*,
405 2007). Further studies should consider the potential contributions of tick microbiomes to the
survival of ticks in the environment; this could provide novel targets for the control of ticks and
tick-borne pathogens.

BLSOM was used to estimate the phylogenetic origins of the sequences associated with
putative virulence-associated factors; however, many of these remained unclassified (Table 5).
410 This may indicate that those sequences had retained their functions, but lost sequence
characteristics such as oligonucleotide frequency. One possible explanation is that the
unclassified sequences were introduced through ancient HGT, and lost their oligonucleotide
frequency in the recipient genome (either tick or other microbes). Indeed, the oligonucleotide
composition of sequences introduced through ancient HGT drift over time (Brown, 2003), makes
415 it difficult to correctly phylotype those sequences using composition-based methods such as
BLSOM. Furthermore, numerous reports have associated bacterial virulence factors with HGT
(Ho Sui *et al.*, 2009; Juhas *et al.*, 2009; Pallen and Wren, 2007). Another explanation is that these
sequences simply originated from as-yet unclassified microorganisms. Regardless, the detection
of virulence-associated factors in ticks may suggest either the existence or possible emergence of
420 tick-borne pathogens and, thus, warrants further investigation to assess potential risks to human
and animal health.

In conclusion, these results provide a foundation for the construction of a database of tick microbes, which may aid in the prediction of emerging tick-borne diseases. The present approach can be extended to detect or predict emerging pathogens in other vector arthropods such as mosquitoes and flies. A comprehensive understanding of tick microbiomes will also be useful for understanding tick biology, including vector competence and interactions with pathogens.

430

Acknowledgements

We thank everyone who helped with tick collection, especially Nariaki Nonaka, Jacqueline Schipper, Joseph W. Magona, Joseph M. Kamau, Noboru Sasaki and Noriko Nakao. The first author was supported by a research grant fellowship from the Japanese Society for the Promotion of Science (JSPS) for young scientists. This work was supported by a Grant-in-Aid for JSPS fellows and for Scientific Research from the Ministry of Education, Culture, Sports, Science and Technology of Japan (MEXT), the program of Funding Research Center for Emerging and Re-emerging Infectious Disease, MEXT. The BLSOM calculation was done, in part, by the Earth Simulator of the Japan Agency for Marine-Earth Science and Technology.

- Abe T, Kanaya S, Kinouchi M, Ichiba Y, Kozuki T, Ikemura T. (2003). Informatics for unveiling hidden genome signatures. *Genome Res* **13**: 693-702.
- Abe T, Sugawara H, Kinouchi M, Kanaya S, Ikemura T. (2005). Novel phylogenetic studies of genomic sequence fragments derived from uncultured microbe mixtures in environmental and
445 clinical samples. *DNA Res* **12**: 281-290.
- Allen HK, Donato J, Wang HH, Cloud-Hansen KA, Davies J, Handelsman J. (2010). Call of the wild: antibiotic resistance genes in natural environments. *Nat Rev Microbiol* **8**: 251-259.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. (1990). Basic local alignment search tool. *J Mol Biol* **215**: 403-410.
- 450 Andreotti R, Pérez de León AA, Dowd SE, Guerrero FD, Bendele KG, Scoles GA. (2011). Assessment of bacterial diversity in the cattle tick *Rhipicephalus (Boophilus) microplus* through tag-encoded pyrosequencing. *BMC Microbiol* **11**: 6.
- Bowman AS, Nuttall, PA. (2008). *Ticks: Biology, Disease and Control*. Cambridge University Press: Cambridge.

- 455 Benson MJ, Gawronski JD, Eveleigh DE, Benson DR. (2004). Intracellular symbionts and other bacteria associated with deer ticks (*Ixodes scapularis*) from Nantucket and Wellfleet, Cape Cod, Massachusetts. *Appl Environ Microbiol* **70**: 616-620.
- Beugnet F, Marié JL. (2009). Emerging arthropod-borne diseases of companion animals in Europe. *Vet Parasitol* **163**: 298-305.
- 460 Brown JR. (2003). Ancient horizontal gene transfer. *Nat Rev Genet* **4**: 121-132.
- Caldwell HD, Belden EL. (1973). Studies of the role of *Dermacentor occidentalis* in the transmission of bovine chlamydial abortion. *Infect Immun* **7**: 147-151.
- Carpi G, Cagnacci F, Wittekindt NE, Zhao F, Qi J, Tomsho LP *et al.* (2011). Metagenomic profile of the bacterial communities associated with *Ixodes ricinus* ticks. *PLoS One* **6**: e25604.
- 465 Chan C-KK, Hsu AL, Tang S-L, Halgamuge SK. (2008). Using growing self-organising maps to improve the binning process in environmental whole-genome shotgun sequencing. *J Biomed Biotechnol* 2008, doi:10.1155/2008/513701
- Clay K, Klyachko O, Grindle N, Civitello D, Oleske D, Fuqua C. (2008). Microbial communities and interactions in the lone star tick, *Amblyomma americanum*. *Mol Ecol* **17**: 4371-4381.

- 470 Cole JR, Wang Q, Cardenas E, Fish J, Chai B, Farris RJ *et al.* (2009). The Ribosomal Database Project: improved alignments and new tools for rRNA analysis. *Nucleic Acids Res* **37**: D141-145.
- Coombes BK. (2009). Type III secretion systems in symbiotic adaptation of pathogenic and non-pathogenic bacteria. *Trends Microbiol* **17**: 89-94.
- 475 Corsaro D, Thomas V, Goy G, Venditti D, Radek R, Greub G. (2007). '*Candidatus Rhabdochlamydia crassificans*', an intracellular bacterial pathogen of the cockroach *Blatta orientalis* (Insecta: Blattodea). *Syst Appl Microbiol* **30**: 221-228.
- Cox-Foster DL, Conlan S, Holmes EC, Palacios G, Evans JD, Moran NA *et al.* (2007). A metagenomic survey of microbes in honey bee colony collapse disorder. *Science* **318**: 283-287.
- 480 Currie CR, Mueller UG, Malloch D. (1999). The agricultural pathology of ant fungus gardens. *Proc Natl Acad Sci U S A* **96**: 7998-8002.
- de la Fuente J, Blouin EF, Kocan KM. (2003). Infection exclusion of the rickettsial pathogen *Anaplasma marginale* in the tick vector *Dermacentor variabilis*. *Clin Diagn Lab Immunol* **10**: 182-184.

- 485 Dick GJ, Andersson AF, Baker BJ, Simmons SL, Thomas BC, Yelton AP *et al.* (2009).
Community-wide analysis of microbial genome sequence signatures. *Genome Biol* **10**: R85.
- Ezaki T, Suzuki S. (1982). Achromopeptidase for lysis of anaerobic gram-positive cocci. *J Clin Microbiol* **16**: 844-846.
- Facco F, Grazi G, Bonassi S, Magnani M, Di Pietro P. (1992). Chlamydial and rickettsial
490 transmission through tick bite in children. *Lancet* **339**: 992-993.
- Grenier AM, Duport G, Pagès S, Condemine G, Rahbé Y. (2006). The phytopathogen *Dickeya dadantii* (*Erwinia chrysanthemi* 3937) is a pathogen of the pea aphid. *Appl Environ Microbiol* **72**: 1956-1965.
- Ho Sui SJ, Fedynak A, Hsiao WW, Langille MG, Brinkman FS. (2009). The association of
495 virulence factors with genomic islands. *PLoS One* **4**: e8094.
- Hofhuis A, van der Giessen JWB, Borgsteede FHM, Wielinga PR, Notermans DW, van Pelt W.
(2006). Lyme borreliosis in the Netherlands: strong increase in GP consultations and hospital admissions in past 10 years. *Euro Surveillance: European Communicable Disease Bulletin* **11**.
- Horn M. (2008). Chlamydiae as symbionts in eukaryotes. *Annu Rev Microbiol* **62**: 113-131.

- 500 Imaoka K, Kaneko S, Tabara K, Kusatake K, Morita E. (2011). The First Human Case of *Rickettsia tamurae* Infection in Japan. *Case Rep Dermatol* **3**: 68-73.
- Isogai E, Isogai H, Masuzawa T, Postic D, Baranton G, Kamewaka Y *et al.* (1996). *Borrelia burgdorferi sensu lato* in an endemic environment: wild sika deer (*Cervus nippon yesoensis*) with infected ticks and antibodies. *Microbiol Immunol* **40**: 13-19.
- 505 Juhas M, van der Meer JR, Gaillard M, Harding RM, Hood DW, Crook DW. (2009). Genomic islands: tools of bacterial horizontal gene transfer and evolution. *FEMS Microbiol Rev* **33**: 376-393.
- Kaltenpoth M, Göttler W, Herzner G, Strohm E. (2005). Symbiotic bacteria protect wasp larvae from fungal infestation. *Curr Biol* **15**: 475-479.
- 510 Kanaya S, Kinouchi M, Abe T, Kudo Y, Yamada Y, Nishi T *et al.* (2001). Analysis of codon usage diversity of bacterial genes with a self-organizing map (SOM): characterization of horizontally transferred genes with emphasis on the *E. coli* O157 genome. *Gene* **276**: 89-99.
- Kohonen T. (1990). The self-organizing map. *Proceedings of the IEEE*. pp 1464-1480.
- Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H *et al.* (2007).
- 515 Clustal W and Clustal X version 2.0. *Bioinformatics* **23**: 2947-2948.

Lewis D. (1979). The detection of rickettsia-like microorganisms within the ovaries of female *Ixodes ricinus* ticks. *Z Parasitenkd* **59**: 295-298.

Macaluso KR, Sonenshine DE, Ceraul SM, Azad AF. (2002). Rickettsial infection in *Dermacentor variabilis* (Acari: Ixodidae) inhibits transovarial transmission of a second *Rickettsia*. *J Med Entomol* **39**: 809-813.

Martin C, Diaz NN, Ontrup J, Nattkemper TW. (2008). Hyperbolic SOM-based clustering of DNA fragment features for taxonomic visualization and classification. *Bioinformatics* **24**: 1568–1574.

Mavromatis K, Ivanova N, Barry K, Shapio H, Goltsman E, McHardy AC, *et al.* (2007). Use of simulated data sets to evaluate the fidelity of metagenomic processing methods. *Nat Methods* **4**: 495-500.

McKercher DG, Wada EM, Ault SK, Theis JH. (1980). Preliminary studies on transmission of *Chlamydia* to cattle by ticks (*Ornithodoros coriaceus*). *Am J Vet Res* **41**: 922-924.

Meijer A, Roholl PJ, Ossewaarde JM, Jones B, Nowak BF. (2006). Molecular evidence for association of *chlamydiales* bacteria with epitheliocystis in leafy seadragon (*Phycodurus*

eques), silver perch (*Bidyanus bidyanus*), and barramundi (*Lates calcarifer*). *Appl Environ Microbiol* **72**: 284-290.

Miyamoto K, Nakao M, Uchikawa K, Fujita H. (1992). Prevalence of Lyme borreliosis spirochetes in ixodid ticks of Japan, with special reference to a new potential vector, *Ixodes ovatus* (Acari: Ixodidae). *J Med Entomol* **29**: 216-220.

Morita C, Yamamoto S, Tsuchiya K, Yoshida Y, Yabe T, Kawabata N *et al.* (1990). Prevalence of spotted fever group rickettsia antibody in *Apodemus speciosus* captured in an endemic focus in Miyazaki Prefecture, Japan. *Jpn J Med Sci Biol* **43**: 15-18.

Murase Y, Konnai S, Hidano A, Githaka NW, Ito T, Takano A *et al.* (2011). Molecular detection of *Anaplasma phagocytophilum* in cattle and *Ixodes persulcatus* ticks. *Vet Microbiol* **149**: 504-507.

Nakabachi A, Ishikawa H. (1999). Provision of riboflavin to the host aphid, *Acyrtosiphon pisum*, by endosymbiotic bacteria, *Buchnera*. *J Insect Physiol* **45**: 1-6.

Nene V. (2009). Tick genomics--coming of age. *Front Biosci* **14**: 2666-2673.

- 545 Noda H, Munderloh UG, Kurtti TJ. (1997). Endosymbionts of ticks and their relationship to *Wolbachia* spp. and tick-borne pathogens of humans and animals. *Appl Environ Microbiol* **63**: 3926-3932.
- Pallen MJ, Wren BW. (2007). Bacterial pathogenomics. *Nature* **449**: 835-842.
- Parola P. (2004). Tick-borne rickettsial diseases: emerging risks in Europe. *Comp Immunol*
- 550 *Microbiol Infect Dis* **27**: 297-304.
- Parola P, Davoust B, Raoult D. (2005). Tick- and flea-borne rickettsial emerging zoonoses. *Vet Res* **36**: 469-492.
- Roshdy MA. (1968). A rickettsialike microorganism in the tick *Ornithodoros savignyi*: Observations on its structure and distribution in the tissues of the tick. *J Invertebr Pathol* **11**:
- 555 155-169.
- Scarborough CL, Ferrari J, Godfray HC. (2005). Aphid protected from pathogen by endosymbiont. *Science* **310**: 1781.
- Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB *et al.* (2009). Introducing mothur: open-source, platform-independent, community-supported software for
- 560 describing and comparing microbial communities. *Appl Environ Microbiol* **75**: 7537-7541.

- Scott JJ, Oh DC, Yuceer MC, Klepzig KD, Clardy J, Currie CR. (2008). Bacterial protection of beetle-fungus mutualism. *Science* **322**: 63.
- Spolidorio MG, Labruna MB, Mantovani E, Brandao PE, Richtzenhain LJ, Yoshinari NH. (2010). Novel spotted fever group rickettsiosis, Brazil. *Emerg Infect Dis* **16**: 521-523.
- 565 Suen G, Scott JJ, Aylward FO, Adams SM, Tringe SG, Pinto-Tomás AA *et al.* (2010). An insect herbivore microbiome with high plant biomass-degrading capacity. *PLoS Genet* **6**: e1001129.
- Tamames J, Moya A. (2008). Estimating the extent of horizontal gene transfer in metagenomic sequences. *BMC Genomics* **9**: 136.
- Tatusov RL, Koonin EV, Lipman DJ. (1997). A genomic perspective on protein families. *Science* **278**: 631-637.
- 570
- Tsuchida T, Koga R, Horikawa M, Tsunoda T, Maoka T, Matsumoto S *et al.* (2010). Symbiotic bacterium modifies aphid body color. *Science* **330**: 1102-1104.
- Uehara H, Iwasaki Y, Wada C, Ikemura T, Abe T. (2011). A novel bioinformatics strategy for searching industrially useful genome resources from metagenomic sequence libraries. *Genes Genet Syst* **86**: 53-66.
- 575

- van Overbeek L, Gassner F, van der Plas CL, Kastelein P, Nunes-da Rocha U, Takken W. (2008).
Diversity of *Ixodes ricinus* tick-associated bacterial communities from different forests. *FEMS
Microbiol Ecol* **66**: 72-84.
- von Bomhard W, Polkinghorne A, Lu ZH, Vaughan L, Vöggtlin A, Zimmermann DR *et al.* (2003).
580 Detection of novel chlamydiae in cats with ocular disease. *Am J Vet Res* **64**: 1421-1428.
- Warnecke F, Luginbühl P, Ivanova N, Ghassemian M, Richardson TH, Stege JT *et al.* (2007).
Metagenomic and functional analysis of hindgut microbiota of a wood-feeding higher termite.
Nature **450**: 560-565.
- Weber M, Teeling H, Huang S, Waldmann J, Kassabgy M, Fuchs BM *et al.* (2011). Practical
585 application of self-organizing maps to interrelate biodiversity and functional data in NGS-
based metagenomics. *ISME J* **5**: 918-928.
- Weiss B, Aksoy S. (2011). Microbiome influences on insect host vector competence. *Trends
Parasitol* **27**: 514-522.
- Wright GD. (2010). Antibiotic resistance in the environment: a link to the clinic? *Curr Opin
590 Microbiol* **13**: 589-594.

Yu XJ, Liang MF, Zhang SY, Liu Y, Li JD, Sun YL *et al.* (2011). Fever with thrombocytopenia associated with a novel bunyavirus in China. *N Engl J Med* **364**: 1523-1532.

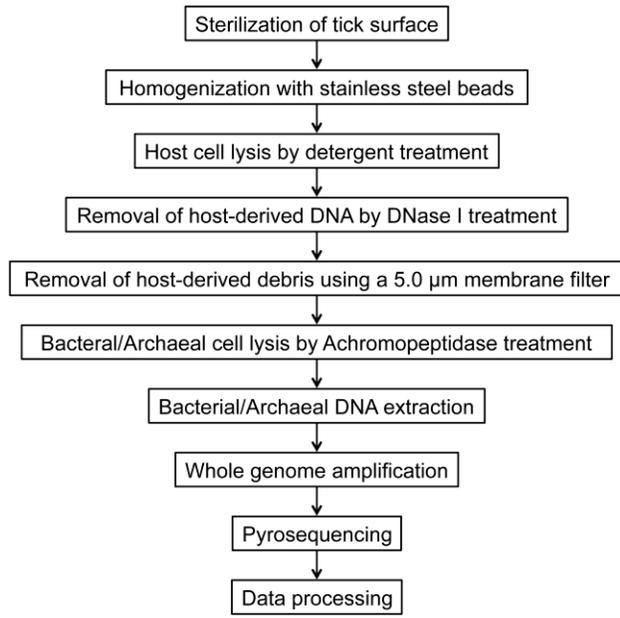
Zhong J, Jasinskas A, Barbour AG. (2007). Antibiotic treatment of the tick vector *Amblyomma americanum* reduced reproductive fitness. *PLoS One* **2**: e405.

595 Zhou CE, Smith J, Lam M, Zemla A, Dyer MD, Slezak T. (2007). MvirDB--a microbial database of protein toxins, virulence factors and antibiotic resistance genes for bio-defence applications. *Nucleic Acids Res* **35**: D391-394.

600

605

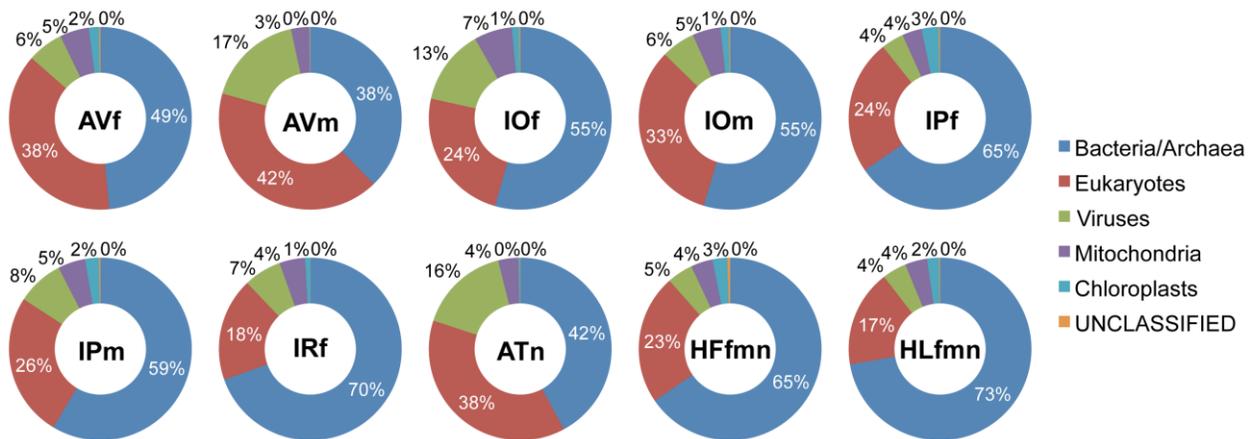
Figure 1. Workflow for bacterial/archaeal purification and metagenomic analysis.



610 The present strategy employed a process for purifying of bacteria/archaea, which comprised centrifugation, DNase treatment and filtration, to enrich bacterial/archaeal cells from tick homogenates prior to whole-genome amplification and pyrosequencing.

615

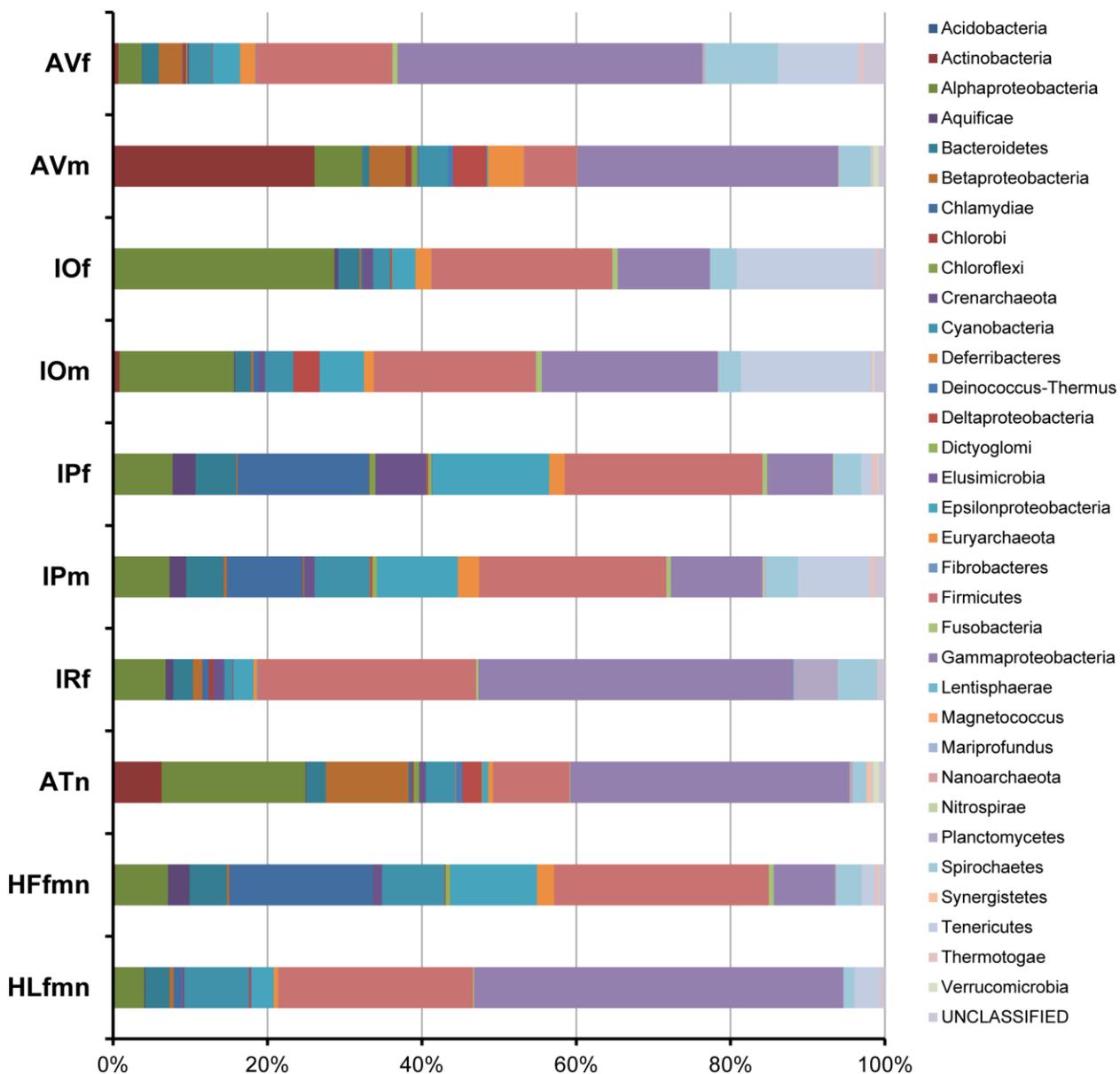
Figure 2. Kingdom classification of the metagenomic sequences from each tick pool.



620 Sequence reads longer than 300 bp were classified into bacteria/archaea, eukaryotes, viruses, mitochondria or chloroplasts using Kingdom-BLSOM. The percentage of sequences in each category is provided. The pool ID is shown in the centre of each pie chart.

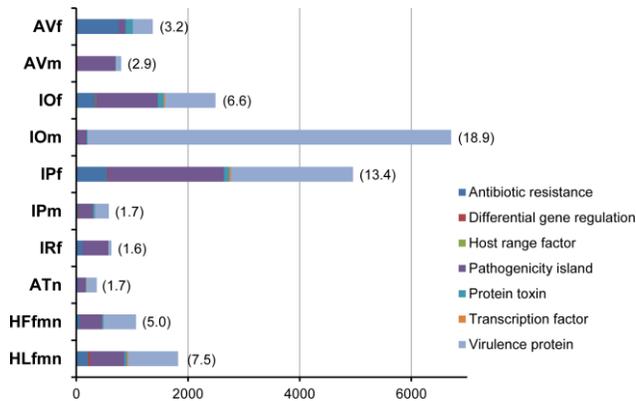
625

630 **Figure 3. Phylum classification of the metagenomic sequences from each tick pool.**



635 Sequence reads classified as bacterial/archaeal using Kingdom-BLSOM were further classified at the phylum level using Bacteria/Archaea-BLSOM. Thirty-three phyla are indicated by different colours. The pool ID is shown on the left of the graph.

Figure 4. Sequences associated with virulence-associated factors.



640 Sequences associated with putative virulence-associated factors were identified by BLASTX against the MvirDB database (Zhou *et al.*, 2007), applying a cut-off value of $1e^{-5}$. The resulting data are shown in seven categories denoted by different colours. The numbers at the bottom of the graph indicate the number of sequences. The numbers in parentheses indicate the percentage of total sequence reads. The pool ID is shown on the left of the graph.

645

650 **Table 1. Collection details for the ticks used in the present study.**

Pool ID	Tick species	Collection site	Adult		Nymph	Total
			Female	Male		
AVf	<i>Amblyomma variegatum</i>	Laboratory colony ^a	20			20
AVm	<i>Amblyomma variegatum</i>	Laboratory colony ^a		20		20
IOf	<i>Ixodes ovatus</i>	Hokkaido, Japan	59			59
IOm	<i>Ixodes ovatus</i>	Hokkaido, Japan		110		110
IPf	<i>Ixodes persulcatus</i>	Hokkaido, Japan	56			56
IPm	<i>Ixodes persulcatus</i>	Hokkaido, Japan		75		75
IRf	<i>Ixodes ricinus</i>	Soesterberg, The Netherlands	55			55
ATn	<i>Amblyomma testudinarium</i>	Miyazaki, Japan			57	57
HFfmn	<i>Haemaphysalis formosensis</i>	Miyazaki, Japan	32	24	513	569
HLfmn	<i>Haemaphysalis longicornis</i>	Miyazaki, Japan	16	20	84	120

^aThis laboratory colony was originally obtained from The Gambia and has been maintained at Utrecht Centre for Tick-borne Diseases, Utrecht University, The Netherlands.

Table 2. Summary of pyrosequencing and BLSOM classifications.

	AVf	AVm	IOf	IOm	IPf	IPm	IRf	ATn	HFfmn	HLfmn
Total no. of reads	42,258	27,582	37,667	35,544	37,136	33,789	39,810	21,639	21,213	24,249
Average read length (nt)	163.6	391.4	379.8	154.6	371.7	143.5	149.9	149.0	185.0	187.5
No. of long reads (≥ 300 bp)	8,160	22,201	29,804	2,962	28,618	4,210	5,438	2,906	4,502	5,186
No. of reads assigned to bacteria/archaea (%) ^a	3,951 (48.4)	8,389 (37.8)	16,239 (54.5)	1,620 (54.7)	18,711 (65.4)	2,468 (58.6)	3,782 (69.5)	1,222 (42.1)	2,946 (65.4)	3,750 (72.3)
No. of reads assigned to a bacterial/archaeal phyla (%) ^b	3,841 (97.2)	8,272 (98.6)	16,088 (99.1)	1,598 (98.6)	18,556 (99.2)	2,439 (98.8)	3,748 (99.1)	1,213 (99.3)	2,925 (99.3)	3,728 (99.4)
No. of phyla obtained	27	28	28	23	28	25	22	27	24	20
No. of reads assigned to known bacterial/archaeal genera (%) ^c	3,419 (89.0)	7,068 (85.4)	14,410 (89.6)	1,410 (88.2)	16,566 (89.3)	2,131 (87.4)	3,398 (90.7)	1,119 (92.3)	2,605 (89.1)	3,319 (89.0)
No. of reads assigned to candidate bacterial/archaeal genera (%)	71 (1.8)	62 (0.7)	700 (4.4)	66 (4.1)	669 (3.6)	98 (4.0)	76 (2.0)	13 (1.1)	89 (3.0)	75 (2.0)
No. of bacterial/archaeal genera obtained	158	223	215	143	185	162	133	201	138	156

^aCalculated using the formula: (No. of reads assigned to bacteria/archaea)/{No. of long reads (≥ 300 bp)} $\times 100$. ^bCalculated using the
660 formula: (No. of reads assigned to any phyla)/(No. of reads assigned to bacteria/archaea) $\times 100$. ^cCalculated using the formula: (No. of
reads assigned to any genera)/(No. of reads assigned to any phyla) $\times 100$.

Table 3. Top 20 genera associated with the metagenomic sequences from each tick pool.

	AVf	AVm	IOf	IOm	IPf	IPm	IRf	ATn	HFfmn	HLfmn	
1	<i>Borrelia</i>	(8.6) <i>Clavibacter</i>	(13.4) <i>Rickettsia</i>	(19.7) <i>Ureaplasma</i>	(11.0) <i>Chlamydomphila</i>	(8.8) <i>Ureaplasma</i>	(7.6) <i>Borrelia</i>	(5.3) <i>Pseudoalteromonas</i>	(17.2) <i>Chlamydomphila</i>	(10.5) <i>Lactobacillus</i>	(7.1)
2	<i>Mycoplasma</i>	(6.2) <i>Pseudoalteromonas</i>	(5.1) <i>Ureaplasma</i>	(15.4) <i>Rickettsia</i>	(7.3) <i>Streptococcus</i>	(8.5) <i>Chlamydomphila</i>	(5.4) <i>Lactobacillus</i>	(5.3) <i>Rickettsia</i>	(10.0) <i>Streptococcus</i>	(8.5) U (Gammaproteobacteria)	(5.1)
3	<i>Staphylococcus</i>	(4.0) U (Euryarchaeota)	(4.7) <i>Mycoplasma</i>	(8.2) <i>Mycoplasma</i>	(7.1) <i>Helicobacter</i>	(6.9) <i>Streptococcus</i>	(5.0) <i>Streptococcus</i>	(4.7) <i>Synechococcus</i>	(2.1) <i>Chlamydia</i>	(6.7) <i>Salmonella</i>	(4.6)
4	<i>Actinobacillus</i>	(4.0) <i>Treponema</i>	(3.6) <i>Clostridium</i>	(5.8) <i>Clostridium</i>	(3.1) <i>Chlamydia</i>	(5.9) <i>Campylobacter</i>	(4.7) <i>Ureaplasma</i>	(4.3) <i>Wigglesworthia</i>	(2.1) <i>Helicobacter</i>	(5.5) <i>Grimontia</i>	(4.3)
5	<i>Coxiella</i>	(4.0) <i>Synechococcus</i>	(3.2) <i>Ehrlichia</i>	(4.0) <i>Helicobacter</i>	(3.1) <i>Campylobacter</i>	(5.5) <i>Helicobacter</i>	(4.4) U (Gammaproteobacteria)	(4.1) <i>Clostridium</i>	(1.8) <i>Prochlorococcus</i>	(4.2) <i>Providencia</i>	(3.7)
6	<i>Grimontia</i>	(3.7) U (Gammaproteobacteria)	(3.1) <i>Francisella</i>	(3.4) <i>Desulfovibrio</i>	(3.1) <i>Prochlorococcus</i>	(3.5) <i>Prochlorococcus</i>	(3.9) <i>Grimontia</i>	(3.9) U (Alphaproteobacteria)	(1.6) <i>Campylobacter</i>	(3.5) <i>Coxiella</i>	(3.4)
7	<i>Thiocystis</i>	(3.7) <i>Arcanobacterium</i>	(2.9) <i>Borrelia</i>	(2.9) U (Gammaproteobacteria)	(3.1) <i>Borrelia</i>	(3.0) <i>Mycoplasma</i>	(3.4) <i>Bacillus</i>	(3.8) <i>Orientia</i>	(1.6) <i>Bacillus</i>	(3.3) <i>Cyanothecce</i>	(3.1)
8	<i>Luteimonas</i>	(3.4) <i>Actinobacillus</i>	(2.7) U (Euryarchaeota)	(2.1) <i>Streptococcus</i>	(2.7) <i>Clostridium</i>	(2.7) <i>Borrelia</i>	(3.1) <i>Luteimonas</i>	(3.0) <i>Bordetella</i>	(1.5) <i>Clostridium</i>	(3.0) <i>Streptococcus</i>	(3.0)
9	<i>Dichelobacter</i>	(2.7) <i>Micromonospora</i>	(2.6) <i>Campylobacter</i>	(2.0) <i>Phytoplasma</i>	(2.6) <i>Bacillus</i>	(2.7) <i>Chlamydia</i>	(3.0) <i>Vibrio</i>	(2.9) <i>Bacillus</i>	(1.4) <i>Thermoanaerobacter</i>	(2.8) <i>Staphylococcus</i>	(2.7)
10	<i>Phytoplasma</i>	(2.5) <i>Salmonella</i>	(2.6) <i>Staphylococcus</i>	(1.9) <i>Francisella</i>	(2.4) <i>Thermoanaerobacter</i>	(2.6) U (Euryarchaeota)	(2.8) <i>Rickettsia</i>	(2.7) U (Gammaproteobacteria)	(1.4) <i>Borrelia</i>	(2.5) <i>Bacillus</i>	(2.7)
11	<i>Campylobacter</i>	(2.4) U (Actinobacteria)	(2.6) <i>Orientia</i>	(1.9) <i>Borrelia</i>	(2.3) <i>Sulfurihydrogenibium</i>	(2.4) <i>Bacillus</i>	(2.4) <i>Thiothrix</i>	(2.5) <i>Amycolatopsis</i>	(1.2) <i>Lactobacillus</i>	(2.3) <i>Acinetobacter</i>	(2.5)
12	U (Tenericutes)	(2.4) <i>Hydrogenovibrio</i>	(2.5) C (Gammaproteobacteria)	(1.5) <i>Thiothrix</i>	(2.2) <i>Francisella</i>	(2.3) <i>Clostridium</i>	(2.3) <i>Actinobacillus</i>	(2.4) <i>Micromonospora</i>	(1.2) U (Euryarchaeota)	(2.2) <i>Mycoplasma</i>	(2.5)
13	<i>Lactobacillus</i>	(2.4) <i>Aeromonas</i>	(2.3) <i>Streptococcus</i>	(1.4) <i>Ehrlichia</i>	(1.9) U (Euryarchaeota)	(2.0) <i>Thermoanaerobacter</i>	(2.0) <i>Coxiella</i>	(2.4) <i>Francisella</i>	(1.2) <i>Sulfurihydrogenibium</i>	(2.1) <i>Pseudoxanthomonas</i>	(2.4)
14	<i>Thiothrix</i>	(2.1) U (Alphaproteobacteria)	(1.9) <i>Wolbachia</i>	(1.3) <i>Staphylococcus</i>	(1.6) C (Chlamydiae)	(2.0) <i>Bartonella</i>	(1.6) <i>Enterococcus</i>	(2.3) <i>Streptococcus</i>	(1.2) <i>Bartonella</i>	(1.5) <i>Xenorhabdus</i>	(2.2)
15	<i>Salmonella</i>	(2.1) <i>Streptococcus</i>	(1.7) <i>Bacillus</i>	(1.2) <i>Campylobacter</i>	(1.4) <i>Wolbachia</i>	(1.9) <i>Lactobacillus</i>	(1.6) <i>Mycoplasma</i>	(2.3) <i>Xanthomonas</i>	(1.2) C (Chlamydiae)	(1.4) <i>Actinobacillus</i>	(1.9)
16	U (Euryarchaeota)	(2.0) <i>Hafnia</i>	(1.6) <i>Buchnera</i>	(1.2) <i>Lactobacillus</i>	(1.4) <i>Bartonella</i>	(1.8) U (Gammaproteobacteria)	(1.6) <i>Cobwellia</i>	(2.3) <i>Aeromonas</i>	(1.0) <i>Wolbachia</i>	(1.3) <i>Rickettsia</i>	(1.9)
17	<i>Neisseria</i>	(2.0) <i>Geobacter</i>	(1.4) U (Gammaproteobacteria)	(1.1) <i>Listeria</i>	(1.4) <i>Sulfurimonas</i>	(1.2) <i>Francisella</i>	(1.5) <i>Francisella</i>	(1.9) <i>Borrelia</i>	(1.0) U (Firmicutes)	(1.2) <i>Thiothrix</i>	(1.8)
18	U (Gammaproteobacteria)	(1.8) <i>Rickettsiella</i>	(1.4) <i>Phytoplasma</i>	(1.1) C (Alphaproteobacteria)	(1.3) <i>Lactobacillus</i>	(1.2) <i>Sulfurihydrogenibium</i>	(1.4) <i>Salmonella</i>	(1.7) <i>Methylobacillus</i>	(1.0) C (Bacteroidetes)	(1.2) <i>Campylobacter</i>	(1.7)
19	<i>Providencia</i>	(1.8) <i>Methylobacterium</i>	(1.3) C (Bacteroidetes)	(1.1) U (Euryarchaeota)	(1.3) <i>Alkaliphilus</i>	(1.1) C (Chlamydiae)	(1.3) <i>Acinetobacter</i>	(1.7) <i>Shigella</i>	(1.0) <i>Cyanothecce</i>	(1.2) <i>Listeria</i>	(1.5)
20	<i>Vibrio</i>	(1.8) <i>Actinoalloteichus</i>	(1.2) C (Tenericutes)	(1.0) <i>Bacillus</i>	(1.3) <i>Lamprocystis</i>	(1.1) <i>Phytoplasma</i>	(1.3) <i>Clostridium</i>	(1.7) <i>Staphylococcus</i>	(1.0) <i>Sulfurimonas</i>	(1.2) <i>Nostoc</i>	(1.5)

665

Numbers in parentheses indicate the percentage. U, Unclassified; C, Candidatus.

Table 4. Phylogenetic analysis of 16S rRNA gene sequences amplified by pan-Chlamydia PCR.

Operational taxonomic unit (OTU)	Sequence ID	Clones assigned to each OTU					RDP result ^a		BLASTn search result ^a		
		IPf	IPm	HFf	HFm	HFn	Family (%) ^b	Genus (%) ^b	Identity (%) ^c	Source	Accession no.
1	1-1*	Yes	Yes			Yes					
	1-2	Yes									
	1-3	Yes					<i>Simkaniaceae</i> (80%)	<i>Simkania</i> (80%)	<i>Rhabdochlamydia crassificans</i> strain CRIB01 (98%)	Cockroach (<i>Blatta orientalis</i>)	AY928092
	1-4		Yes								
	1-5	Yes									
2	2-1*		Yes	Yes	Yes	Yes			Uncultured Chlamydiae bacterium clone UFC3 (90%)	Sea bass (<i>Dicentrarchus labrax</i>)	FJ376381
	2-2		Yes	Yes			<i>Parachlamydiae</i> (53%)	<i>Neochlamydia</i> (22%)			
	2-3			Yes							
	2-4				Yes						
3	3-1*				Yes		<i>Parachlamydiae</i> (89%)	<i>Neochlamydia</i> (73%)	Uncultured Chlamydiales bacterium isolate UKC7 (90%)	Koala (<i>Phascolarctos cinereus</i>)	AY167120
	3-2		Yes		Yes						
4	4-1*			Yes			<i>Parachlamydiae</i> (95%)	<i>Neochlamydia</i> (67%)	Uncultured clinical <i>Neochlamydia</i> sp. clone WB13 (92%)	Cat (<i>Felis silvestris catus</i>)	AY225593
	4-2			Yes							
5	5-1*				Yes		<i>Parachlamydiae</i> (76%)	<i>Parachlamydia</i> (74%)	Uncultured Chlamydiales CRG20 (88%)	Leafy seadragon (<i>Phycodurus eques</i>)	AY013396

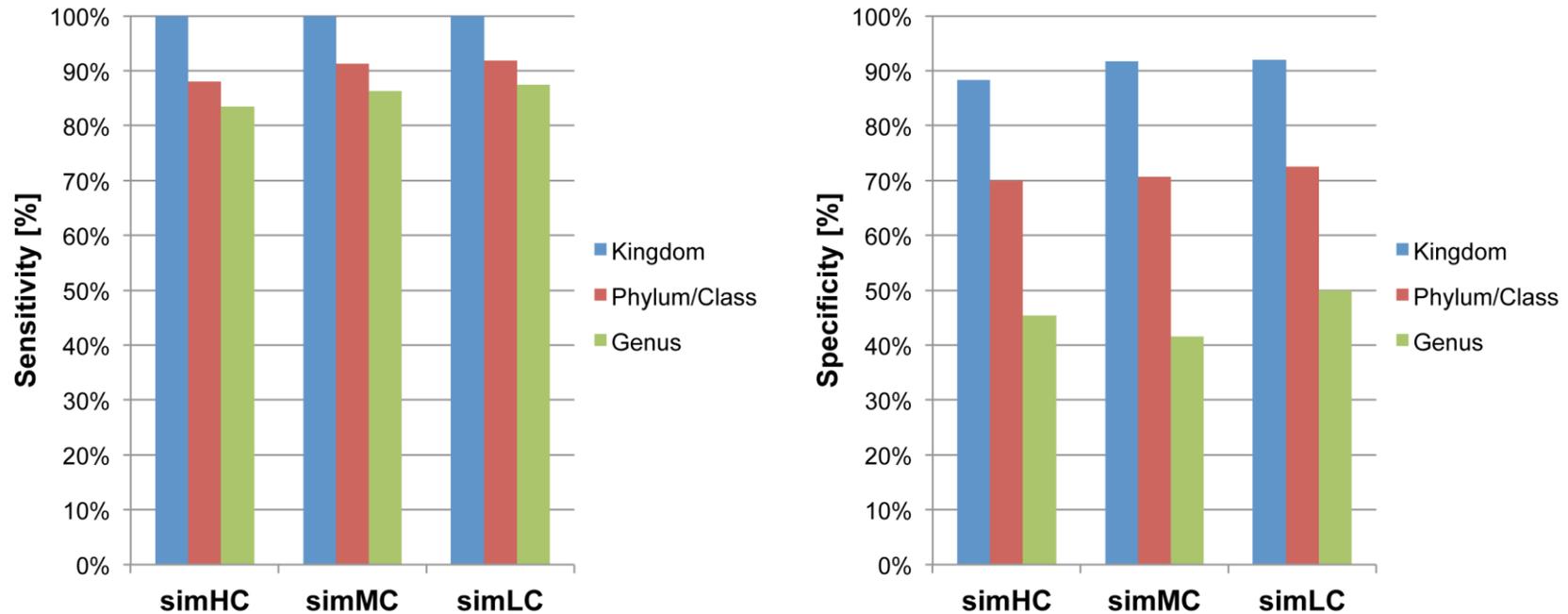
*Representative sequence. ^aAnalysis was conducted using a representative sequence for each OTU. ^bPercentages indicate bootstrap values. ^cPercentages indicate the similarity between two sequences.

Table 5. Origins of the sequences of virulence-associated factors.

AVf	AVm	IOf	IOm	IPf	IPm	IRf	ATn	HfMn	HLfMn	
Antibiotic resistance										
1 U	(424) Actinobacteria	(1) Alphaproteobacteria	(90) U	(6) U	(130) Firmicutes	(8) Gammaproteobacteria	(68) Gammaproteobacteria	(8) U	(15) Gammaproteobacteria	(50)
2 Gammaproteobacteria	(85) Gammaproteobacteria	(1) U	(82) Firmicutes	(5) Firmicutes	(110) Chlamydiae	(6) U	(20) U	(7) Firmicutes	(9) U	(39)
3 Firmicutes	(73)	Firmicutes	(49) Gammaproteobacteria	(4) Gammaproteobacteria	(50) U	(5) Alphaproteobacteria	(6) Alphaproteobacteria	(5) Epsilonproteobacteria	(4) Firmicutes	(16)
4 Spirochaetes	(38)	Gammaproteobacteria	(44) Alphaproteobacteria	(2) Chlamydiae	(49) Gammaproteobacteria	(4) Firmicutes	(5) Actinobacteria	(1) Gammaproteobacteria	(4) Cyanobacteria	(5)
5 Euryarchaeota	(26)	Crenarchaeota	(17) Cyanobacteria	(1) Cyanobacteria	(44) Cyanobacteria	(3) Cyanobacteria	(4) Betaproteobacteria	(1) Chlamydiae	(3) Tenericutes	(2)
Pathogenicity island										
1 Gammaproteobacteria	(46) U	(505) U	(329) U	(86) U	(546) U	(61) Gammaproteobacteria	(164) Gammaproteobacteria	(30) U	(69) Gammaproteobacteria	(183)
2 Firmicutes	(28) Gammaproteobacteria	(98) Alphaproteobacteria	(293) Firmicutes	(29) Firmicutes	(346) Firmicutes	(37) Firmicutes	(149) U	(24) Firmicutes	(38) Firmicutes	(88)
3 U	(16) Actinobacteria	(29) Firmicutes	(152) Gammaproteobacteria	(12) Chlamydiae	(286) Chlamydiae	(36) U	(59) Actinobacteria	(15) Chlamydiae	(36) U	(58)
4 Betaproteobacteria	(9) Alphaproteobacteria	(18) Gammaproteobacteria	(80) Tenericutes	(9) Gammaproteobacteria	(165) Bacteroidetes	(19) Alphaproteobacteria	(12) Alphaproteobacteria	(13) Alphaproteobacteria	(23) Alphaproteobacteria	(13)
5 Alphaproteobacteria	(3) Cyanobacteria	(12) Bacteroidetes	(32) Alphaproteobacteria	(8) Epsilonproteobacteria	(140) Epsilonproteobacteria	(17) Cyanobacteria	(11) Betaproteobacteria	(7) Epsilonproteobacteria	(19) Cyanobacteria	(13)
Virulence protein										
1 Gammaproteobacteria	(192) Gammaproteobacteria	(51) Alphaproteobacteria	(291) Gammaproteobacteria	(81) U	(625) U	(60) Gammaproteobacteria	(14) U	(35) U	(96) Gammaproteobacteria	(197)
2 U	(55) U	(14) U	(238) U	(34) Firmicutes	(411) Firmicutes	(44) U	(14) Gammaproteobacteria	(28) Firmicutes	(68) U	(116)
3 Firmicutes	(34) Firmicutes	(12) Firmicutes	(92) Alphaproteobacteria	(23) Chlamydiae	(285) Chlamydiae	(27) Firmicutes	(9) Alphaproteobacteria	(8) Chlamydiae	(46) Firmicutes	(86)
4 Betaproteobacteria	(22) Actinobacteria	(3) Gammaproteobacteria	(90) Deltaproteobacteria	(21) Epsilonproteobacteria	(159) Gammaproteobacteria	(24) Alphaproteobacteria	(4) Firmicutes	(8) Cyanobacteria	(26) Cyanobacteria	(35)
5 Bacteroidetes	(7) Betaproteobacteria	(2) Tenericutes	(25) Firmicutes	(15) Gammaproteobacteria	(117) Epsilonproteobacteria	(19) Cyanobacteria	(1) Betaproteobacteria	(6) Gammaproteobacteria	(18) Alphaproteobacteria	(15)

675 Numbers in parentheses indicate the total number of sequences. U, Unclassified.

Supplementary Figure S1. Evaluation of BLSOM-based classification using simulated datasets.



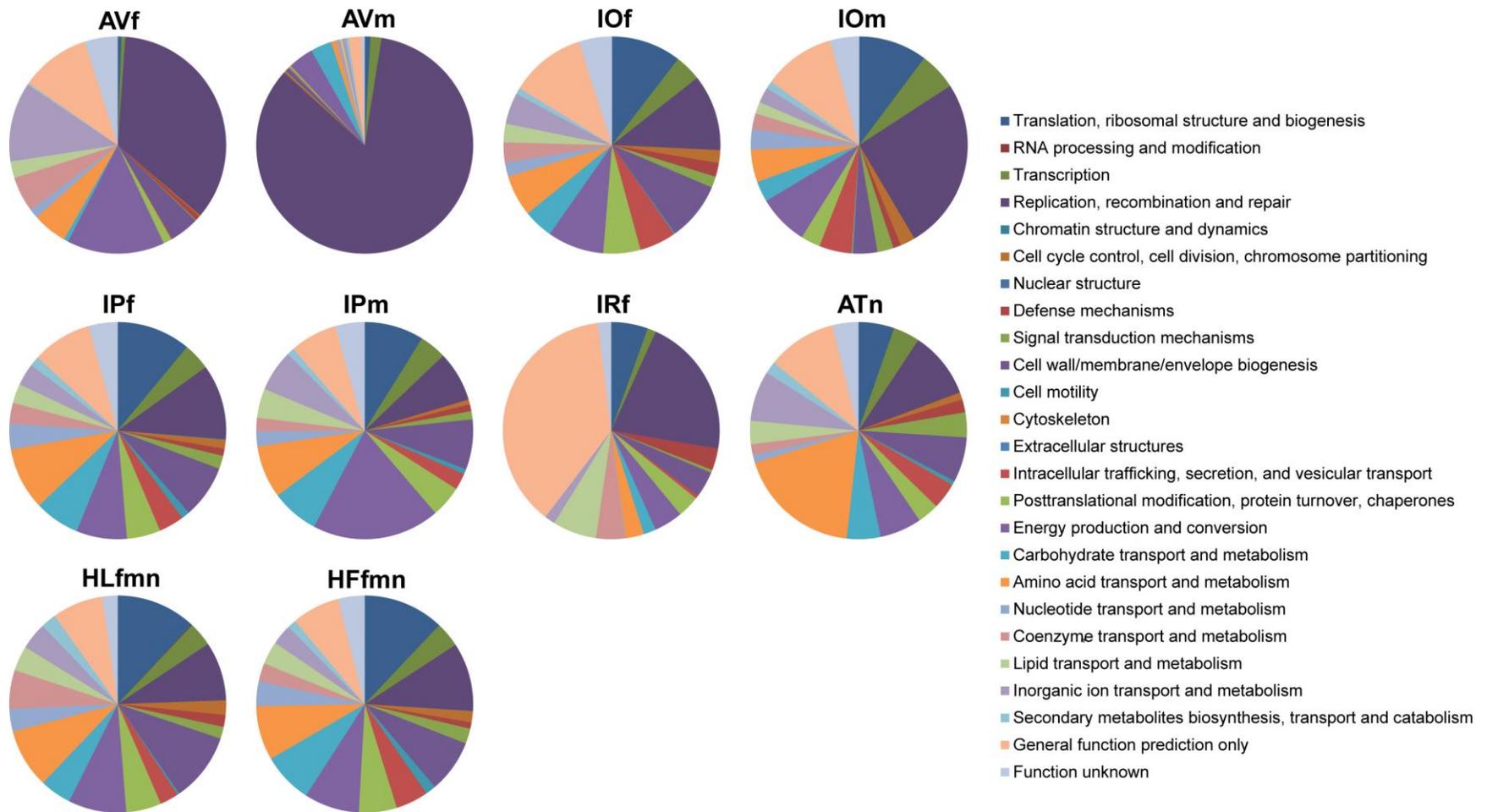
680

Evaluation of Kingdom-, Bacteria/Archaea- and Genus group BLSOM-based classification using simulated datasets of high (simHC), medium (simMC) and low (simLC) complexities (Mavromatis *et al.*, 2007) on contigs longer than 300 bp after normalization of sequence length. The left graph shows the sensitivity {true positives/(true positives + false negatives)} (%) and the right graph shows the specificity {true positives/(true positives + false positives)} (%) of the BLSOM classification. Different taxonomic levels are shown in different colours.

685

Reference: Mavromatis K, Ivanova N, Barry K, Shapio H, Goltsman E, McHardy AC, et al. (2007). Use of simulated data sets to evaluate the fidelity of metagenomic processing methods. *Nat Methods* 4: 495-500.

Supplementary Figure S2. Functional annotation using the Kyoto Encyclopedia of Genes and Genomes (KEGG) database.



690 Sequence reads were annotated by BLASTX against the Clusters of Orthologous Groups (COG) database with an E-value threshold of 1×10^{-5} and were classified according to COG functional categories. Twenty-four COG categories are indicated by different colours. The pool ID is shown on the top of each pie chart.

Supplementary Table S1. Primers amplifying the 16S rRNA gene of pan-Chlamydial bacteria.

Primer name	Sequence (5' to 3')	Amplicon size (bp)	Reference
CF1	CGTGGATGAGGCATGC(A/G)AGTCG	about 1500	(Corsaro <i>et al.</i> , 2002)
CR6	GTCATC(A/G)GCC(T/C)(T/C)ACCTT(A/C/G)(C/G)(A/G)C(A/G)(T/C)(T/C)TCT		
16SIGF	CGGCGTGGATGAGGCAT	about 1500	(Thomas <i>et al.</i> , 2006)
rP2Chlam	CTACCTTGTTACGACTTCAT		
16S FOR2	CGTGGATGAGGCATGCAAGTCGA	about 260	(Ossewaarde <i>et al.</i> , 1999)
16S REV2	CAATCTCTCAATCCGCCTAGACGTCTTAG		

695

References:

Corsaro D, Venditti D, Valassina M. (2002). New chlamydial lineages from freshwater samples. *Microbiology* **148**: 343-344.

Ossewaarde JM, Meijer A. (1999). Molecular evidence for the existence of additional members of the order *Chlamydiales*. *Microbiology* **145**: 411-417.

700 Thomas V, Casson N, Greub G. (2006). *Criblamydia sequanensis*, a new intracellular *Chlamydiales* isolated from Seine river water using amoebal co-culture. *Environ Microbiol* **8**: 2125-2135.

