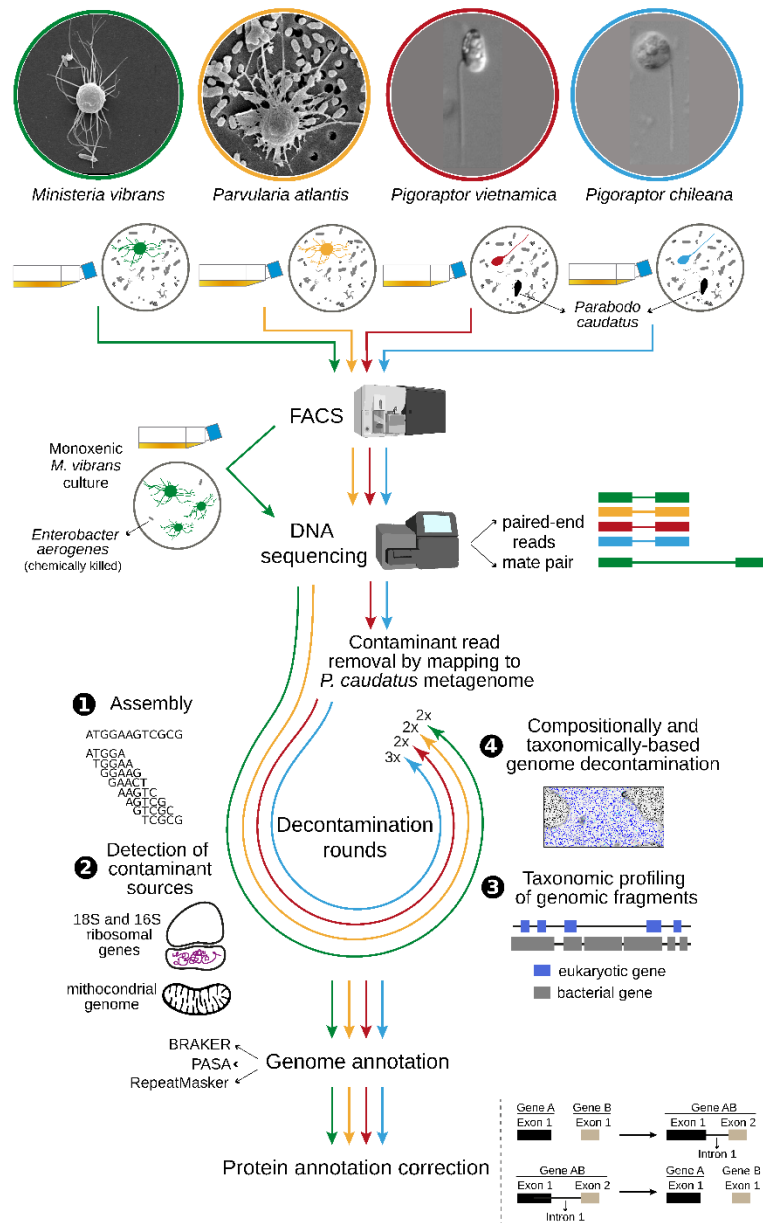


Supplementary Information 1



2

3 **Supplementary Information 1-Fig. 1.** Schematic representation of the methodological approach followed
 4 for the acquisition of genomic data for these four unicellular opisthokonts from polyxenic cultures. Cartoons
 5 from both FACS and Illumina machines were obtained from <https://biorender.com/>.

6

7 **1) *Ministeria vibrans* (Filasterea, Opisthokonta)**

8

9 **1.1) Cultures, Cell cytometry and DNA sequencing**

10 We started from non-axenic cultures of *M. vibrans* ATCC 50519¹ growing in ATCC Medium 1525
 11 and maintained at 23 °C. Fluorescence-activated cell sorting (FACS) was used to isolate *M.*
 12 *vibrans* in a rich medium supplemented with antibiotics and with chemically killed bacteria (see

13 below), also maintained at 23 °C. Flow cytometry analyses and cell sorting of the cultures were
14 performed in a BD FACSAria II cell sorter (Becton Dickinson, San Jose, CA) equipped with 488
15 argon laser. For that, samples were incubated with 5-cyao-2,3-diotolyl tetrazolium chloride and
16 LysoTracker Green DND-26 to differentially label bacterial and eukaryotic cells, respectively. We
17 used the gating strategy based in the following sequence (Supplementary Information 1-Fig. 2):
18 Forward Scatter (FSC) vs. green fluorescence (FITC channel 525/50 nm band-pass filter,
19 LysoTracker-green Fluorescence); and FSC versus Side Scatter (SCC). From these plots, we
20 defined a population including those larger and green fluorescent cells (P2, Supplementary
21 Information 1-Fig. 2A-B). That population was subsequently gated for red (PerCPCy5.5 channel
22 685/35 nm band-pass filter) versus green dot-plot (PerCP-5.5 vs. FITC), to finally sort only green
23 cells (P4, Supplementary Information 1-Fig. 2C) corresponding to the eukaryotic cells. Sorted
24 cells were collected in 48-well plates filled with rich medium. The rich medium was prepared by
25 mixing two-thirds of ATCC® 327-X™ with one-third of Phosphate-buffered saline and adding
26 3.6g/100mL of salts, being sterilized with filters of 0.22 µM. Cultures of *M. vibrans* growing in the
27 rich medium were supplemented with Gentamicin and Chloramphenicol to maintain them free of
28 potential bacterial contamination. They were also fed with chemically killed *Enterobacter*
29 *aerogenes* samples that were prepared by treating them overnight with 0.5% of formaldehyde
30 and then washed two times with PBS.

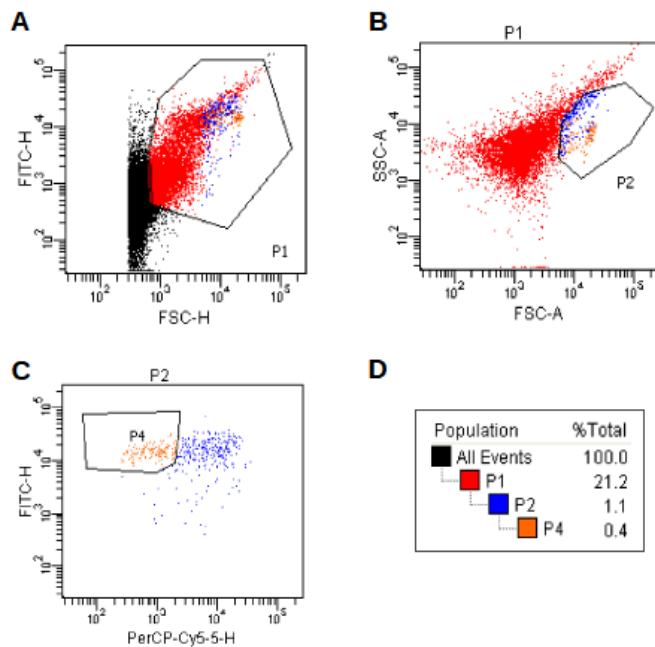
31

32 We used a duplex PCR-based assay for measuring the ratio of 18S/16S ribosomal genes².
33 Results indicated a considerable improvement of the *M. vibrans*/Bacteria ratio in the new cultures
34 compared to the old ones (Supplementary Information 1-Fig. 3A).

35

36

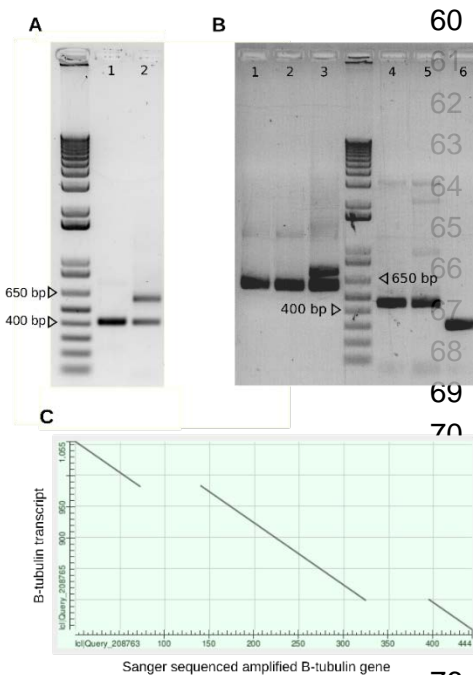
37



Supplementary Information 1-Fig. 2.

Plots showing the strategy to sorter *M. vibrans* cells using fluorescence-activated cell sorting. Samples were incubated with 5-cyao-2,3-diotolyl tetrazolium chloride and LysoTracker Green DND-26 to differentially label bacterial and eukaryotic cells, respectively. Flow cytometry analyses and cell sorting of the cultures were performed in a BD FACSaria II cell sorter (Becton Dickinson, San Jose, CA) equipped with 488 argon laser. We used the gating strategy based in the following sequence: (A) Forward Scatter (FSC) vs. green fluorescence (FITC channel

54 525/50 nm band-pass filter, LysoTracker green Fluorescence); and (B) FSC vs Side Scatter (SSC). From
 55 these plots, we defined a population including those larger and green fluorescent cells (P2). (C) That
 56 population was subsequently gated for red (PerCPCy5.5 channel 685/35 nm band-pass filter) vs green dot-
 57 plot (PerCP-5.5 vs. FITC), to finally sort only orange cells (P4) corresponding to the eukaryotic cells. Sorted
 58 cells were collected in 48-well plates. (D) Percentage of events corresponding to each defined population.
 59



Supplementary Information 1-Fig. 3

(A) Agarose gel showing the amplification products of the duplex PCR-based assay for measuring the ratio of 18S/16S ribosomal genes². Lanes 1 and 2 correspond to DNA extractions from old and new cultures of *M. vibrans*, respectively. The expected band length for the amplified 18S and 16S products are 400 and 600 bp, respectively. The intensity of the bands indicate that the new cultures show better 18S/16S ratios than the old cultures. (B) Agarose gel showing the amplification products of two PCR: (i) using 16S universal primers (lanes 1-3) and (ii) using specific primers for the *P. atlantis* B-tubulin gene (lanes 4-6). Lanes 1, 2, 4, 5 correspond to DNA extractions from the pooling of cells sorted by flow cytometry analyses of *P. atlantis* cultures; whereas lanes 3 and 6 correspond to cDNA from *P. atlantis* cultures. The cDNA from *P. atlantis* was obtained using the protocol described in ³. The B-tubulin was used

77 as a marker for the presence of *P. atlantis* genomic DNA instead of the 18S because we were not able to
 78 amplify this ribosomal gene from previous genomic DNA extractions. The amplified bands suggest the
 79 presence of both 16S and B-tubulin genes in the DNA extractions from the pooling of sorted cells, suggesting
 80 the presence of our organism of interest but also of uncertain bacterial contamination. We confirmed that

81 the bands in the lanes 4 and 5 correspond to the B-tubulin gene by Sanger sequencing. (C) Dot plot
82 alignment between the Sanger sequenced B-tubulin gene and the B-tubulin transcript (obtained from the
83 RNA-seq *de novo* assembly). The dot plot was performed using the online *blastn suite*. The different lengths
84 between the lanes 4 and 5 and the lane 6 is explained by the presence of two intronic sequences in this
85 gene. B-tubulin forward primer sequence: GCAGATGCTTAACGTCCAGAGC. B-tubulin Reverse primer
86 sequence: GATGCCTCCTGGTACTGCTGG.

87

88 DNA was extracted from a pooling of multiple cultures to achieve the required amounts for DNA
89 library preparations. Extractions were done with PureLink® Genomic DNA Mini Kit following the
90 standard protocol. From the extracted DNA, two libraries were prepared for paired-end and mate-
91 pair sequencing (PE and MP, respectively). Each library was sequenced in a 20% Illumina HiSeq
92 2500 lane using the sequencing kit HiSeq v4 chemistry. The insert sizes for PE and MP were 560
93 bp and 3000 bp, respectively, and 125 bp of read length. Library preparation and DNA sequencing
94 of *M. vibrans* and other species were done at the CRG Genomics unit (Barcelona).

95

96 1.2) Read pre-processing

97 We followed different pre-processing strategies for PE and MP reads. PE reads were
98 preprocessed with *trimmomatic*⁴ v0.36 using the following parameters: *SLIDINGWINDOW:12:30*
99 *LEADING:30 ILLUMINACLIP:2:30:10 MINLEN:80*. TruSeq-PE related adapter sequences from
100 *trimmomatic* were used as contaminant database for *ILLUMINACLIP*. To validate that the
101 *ILLUMINACLIP* parameter did not substantially trim false positive adapter reads, we
102 preprocessed a set of simulated PE reads from *Capsaspora owczarzaki* genome, the closest
103 relative to *M. vibrans*. *C. owczarzaki* genome was not sequenced using Illumina chemistry, and
104 hence non TruSeq-PE related adapters are expected among the simulated reads. Reads were
105 simulated using *DWGSIM* v0.1.11 [-e 0 -E 0 -C 40 -1 125 -2 125]
106 (<https://github.com/nh13/DWGSIM>). Only 21 of the 9420724 of the simulated reads were trimmed
107 [*ILLUMINACLIP:2:30:10 MINLEN:80*]. Hence, we expect a negligible false discovery rate from
108 this adapter trimming strategy. We used *FastQC* v0.11.5 during all pre-processing steps for read
109 quality assessment (www.bioinformatics.babraham.ac.uk/projects/fastqc). Before *trimmomatic*
110 [*SLIDINGWINDOW:20:30 LEADING:30 ILLUMINACLIP:2:30:10 MINLEN:50*], MP reads were
111 preprocessed using *nxtrim*⁵ v0.4.1 [--separate --justmp] in order to keep only reads in mate-pair
112 orientation as well as to remove Nextera Transposase sequences. For MP, the contaminant
113 database also included Mate Pair Adapter Sequence Elements (see Table 1 in:
114 [https://www.illumina.com/documents/products/technotes/technote_nextera_matepair_data_proc](https://www.illumina.com/documents/products/technotes/technote_nextera_matepair_data_processing.pdf)
115 [essing.pdf](https://www.illumina.com/documents/products/technotes/technote_nextera_matepair_data_processing.pdf)).

116

117 1.3) First round of read decontamination

118 Because we expected a fraction of the sequenced reads to correspond to bacterial contamination
119 (at least from *E. aerogenes*), prior to a definitive genome assembly, we decided to do a draft
120 assembly in order to identify the contaminant contigs and remove the corresponding reads.

121

122 For the genome assemblies, we decided to use *SPAdes*⁶ v3.10.1 software because (1) it
123 produced the best assemblies in previous studies of our laboratory that included genome data
124 from unicellular relatives to *M. vibrans*⁷, (2) it allows to combine assemblies with different k-mers,
125 (3) it allows to perform both read error and contig miss-match corrections, and (4) because the --
126 meta parameter (i.e., *metaSPAdes*) supports metagenomic data with uneven coverage. For the
127 preliminary assembly (first assembly), we ran *SPAdes* using only the preprocessed paired PE
128 reads and considering the input data as a metagenome [--meta, only non-default software
129 parameters will be specified].

130

131 1.3.1) Evaluation of contaminant sources

132 We first inspected the assembled contigs in order to detect and remove potential vector and/or
133 adapter sequences not trimmed during the pre-processing of the reads. To do so, we used
134 *BLASTn*⁸ with the parameters recommended by *VecScreen* documentation
135 (www.ncbi.nlm.nih.gov/tools/vecscreen) [-task blastn -reward 1 -penalty -5 -gapopen 3 -
136 gapextend 3 -dust yes -soft_masking true -evaluate 700 -searchsp 1750000000000], using UniVec
137 database as reference. We removed from the assembly two contigs that aligned with high identity
138 and query coverage with UniVec sequences.

139

140 We then evaluated the contigs for potential sources of eukaryotic and prokaryotic contamination.
141 For eukaryotic contamination, we searched for 18S ribosomal and mitochondrial sequences. In
142 18S searches, we aligned an in-house curated database of 18S sequences from a variety of
143 eukaryotic groups with the assembled contigs using *BLASTn* [-evaluate 1e-20]. In mitochondrial
144 searches, we aligned with *tBLASTn* [-evaluate 1e-20] the cytochrome c oxidase subunits I and III
145 (COX-I, COX-III) and the cytochrome b (Cyt-b) protein sequences from *Andalucia godoyi*⁹, as
146 these three proteins are found in most of mitochondrial genomes. To check for prokaryotic
147 contamination, we aligned contigs with a local 16S ribosomal database downloaded from NCBI.
148 All contigs found to be potential 18S, mitochondrial or 16S sequences were aligned with the NCBI
149 nt online database and alignment results were manually inspected. We only found 18S and
150 mitochondrial contigs corresponding to *M. vibrans*, suggesting the absence of eukaryotic
151 contamination. However, we found 16S sequences corresponding to *E. aerogenes* but also to
152 *Stenotrophomonas maltophilia*, which suggested an unexpected potential contamination also
153 from this bacterial species.

154

155 1.3.2) Taxonomic classification

156 We used indirect and direct sequence-similarity based approaches to classify contigs into
157 potential contaminant or potential *M. vibrans*. The indirect strategy consisted of classifying contigs
158 according to the average taxonomic signal shown by their preliminary predicted genes. For that,
159 the genome was preliminary annotated with *BRAKER*¹⁰, a RNA-seq-based annotation pipeline
160 that combines *GeneMark-ES/ET*¹¹ v4.33 and *AUGUSTUS*¹² v3.1.0. *M. vibrans* RNA-seq reads
161 were downloaded from NCBI (SRX096925 and SRX096927), corrected with *SEECER*¹³ v.0.1.3
162 and aligned with contigs using *TopHat*¹⁴ v2.1.1. The *accepted_hits.bam* file was used as input for
163 *braker.pl* v1.9. Predicted proteins were then aligned [*BLASTP*: -task blastp-fast, -evaluate 1e-3]
164 with an in-house database including all the prokaryotic Uniprot reference proteomes and 25
165 eukaryotic proteomes [euk_prok_db], mostly from Opisthokonta but also at least one proteome
166 from all the major eukaryotic groups¹⁵. We then classified each gene into Eukaryote (E), Bacteria
167 (B), Putatively Eukaryote (PE), Putatively Bacteria (PB) or Unknown (?) according to the following
168 criteria: (a) If the best hit of a protein was an eukaryotic sequence and the second best hit as well,
169 the corresponding gene was categorized as 'E'. However, if the second best hit was a prokaryotic
170 sequence, and the division between the exponents of the E-values from the second and the first
171 best hits was less than 0.75, the gene was categorized as 'PE'. (b) If the best hit of a protein was
172 a prokaryotic sequence and there were no hits with an eukaryotic sequence, the corresponding
173 gene was categorized as 'B'. In the opposite case, if the division between the exponents of the E-
174 values from the best hit with an eukaryotic sequence and the best hit was higher than 0.75, the
175 corresponding gene was categorized as 'PB'. Genes that did not align were categorized as '?'.
176 Finally, contigs were classified as potentially *M. vibrans* ('Euk profile') or as potentially
177 contaminant ('Bact profile') when the E+PE/B+PB ratio was >1 or <1, respectively. We excluded
178 from this classification contigs with less than 2500 bp.

179

180 For the direct approach, we followed two strategies. The first one (automated strategy), consisted
181 in aligning all the assembled contigs with the NCBI nt database [*BLASTn*: -task megablast -evaluate
182 1e-5], and classify those contigs whose best hit was a prokaryotic sequence as contaminants
183 ('Automated Bact +'). For the second strategy (curated strategy), we separately aligned contigs
184 with the NCBI nt database [*BLASTn*: -evaluate 1e-5]. By using *nucl_gb.accession2taxid*,
185 *nodes.dmp* and *names.dmp* files, we ranked the organisms represented in the database
186 according to the number of best hits received by their sequences. Organisms with more than 5
187 hits were considered as potential contaminants, and their genomes were included into a
188 contaminants database. In particular, the contaminant database included the genomes of *E.*
189 *aerogenes* and *Stenotrophomonas maltophilia*, but also of other bacterial and eukaryotic species
190 that received few hits. We then aligned contigs with this contaminant database [*BLASTn*: -evaluate
191 1e-30], and classified as bacterial contaminants those that aligned with its best scoring database
192 sequence with > 90% of total query coverage and average identity. The total query coverage is
193 the percentage of contig positions that aligned in any hit with the best scoring database sequence.
194 To compute the average identities, we first assigned to each aligned position the highest identity

195 value among the hits in which that position aligned, and then the average was computed for every
196 contig.

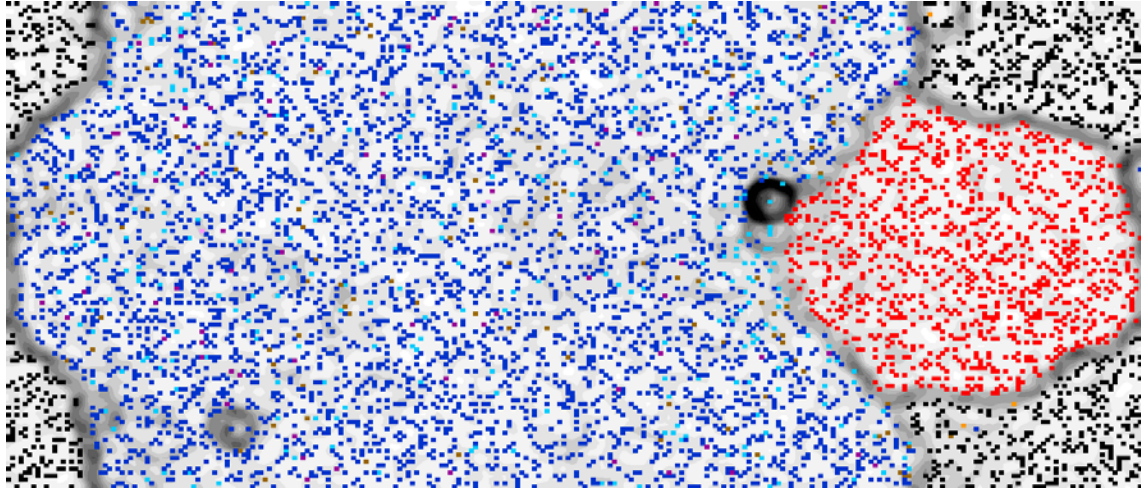
197

198 Contigs that aligned with the contaminant database but that did not satisfy the above-mentioned
199 thresholds were aligned with the NCBI nt online database and alignment results were manually
200 inspected (no evidence of contamination was found on these). All contigs classified as
201 contaminants using this approach ('Curated Bact +') corresponded to *E. aerogenes* and *S.*
202 *maltophilia*, suggesting there are no further contaminant sources. These two contaminant sources
203 were already detected in the fast round of contamination assessment, confirming that the
204 screening of the 18S, 16S and mitochondrial contigs was a valid approach for the identification of
205 contaminant sources in this assembly.

206

207 From the results of the taxonomic classification approaches, all contigs in the assembly can
208 belong to one or more of the following categories: 'Curated Bact +' for the curated strategy,
209 'Automated Bact +' for the automated strategy, and 'Euk profile' or 'Bact profile' for the indirect
210 strategy (the last two are mutually exclusive). We also created a 'No data' category for contigs
211 not classified in any of the four previous categories. The 57.05% of the assembly (in terms of
212 length) was classified as 'Euk profile' (Supplementary Information 1-Fig. 4), suggesting that most
213 of the data is likely to correspond to the *M. vibrans* genome (already expected from the 18S/16S
214 heteroduplex PCR results, see Supplementary Information 1-Fig. 3). A 14.02% of the assembly
215 was in 'No data' category. This includes the *M. vibrans* mitochondrial genome, from which no
216 genes were predicted by *BRAKER1* (probably because of differences in the genetic code). Very
217 few contaminant contigs were expected in 'No data', since the sum of the lengths of those contigs
218 identified as *E. aerogenes* and *S. maltophilia* by the curated strategy is similar to the expected
219 length of both genomes. Thus, most 'No data' contigs likely represent artefactual contigs or non-
220 coding regions of *M. vibrans* genome. The contigs that are neither in 'Euk profile' nor in 'No data',
221 which correspond to the 28.90% of the data, are potential bacterial contaminants because of
222 being in at least one of these categories: 'Curated Bact +', 'Automated Bact +' or 'Bact profile'.
223 The vast majority of them were contaminants according to the three categories (25.08% of the
224 data). All contigs classified as 'Curated Bact +' were also classified as 'Bact profile' and/or
225 'Automated Bact +'. However, these two strategies also identified as contaminants contigs that
226 were not detected by the curated strategy (3.20% of the data).

227



253
 254 **Supplementary Information 1-Fig. 5.** ESOM map of the *M. vibrans* contigs from the first assembly. Each
 255 dot in the map correspond to one contig/contig window, which are colored according to the area where the
 256 corresponding contig was included in the Venn diagram shown in Supplementary Information 1-Fig. 4.
 257 Tetranucleotide frequency distances between neighbour contig/contig windows are represented with a
 258 white/black background gradient for smaller and larger distances, respectively.

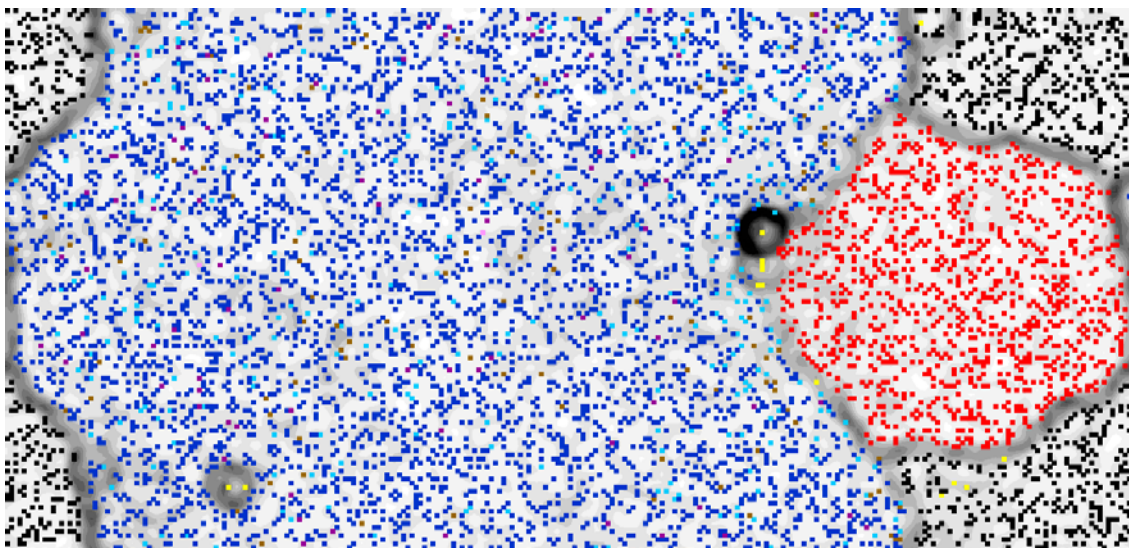
259
 260 The topology of the map shows two big clusters that include all contigs identified as *E. aerogenes*
 261 and *S. maltophilia* by the curated strategy (see red and black dot regions, respectively; note that
 262 the map is continuous from top to bottom and side to side). The black dot region also includes
 263 five dots of distinct color (i.e., fragments that were not detected as *S. maltophilia* but that
 264 apparently share compositional similarity to *S. maltophilia* detected contigs). The vast majority of
 265 contigs classified as 'Euk profile' and 'No data' appear out of the two bacterial regions of the map,
 266 suggesting that the large region in the middle corresponds to *M. vibrans* genome, with the only
 267 exception of one dark blue contig included in the *S. maltophilia* cluster. This suggest a good
 268 precision for the indirect approach in identifying the contigs corresponding to the *M. vibrans*
 269 genome. All 'Euk profile' contigs classified as bacterial by the automated strategy ('Automated
 270 Bact +') are found in the *M. vibrans* region (see dark purple dots). Moreover, the vast majority of
 271 contigs classified as contaminants by the indirect and automated strategies but not by the curated
 272 strategy are also found in *M. vibrans* region (colored in orange, brown and pink). We thus expect
 273 most of these contigs to have been misclassified as contaminants because of false positive
 274 alignments with bacterial sequences (we used relaxed E-value thresholds in the *BLAST*
 275 searches). Also within the *M. vibrans* region, two smaller clusters are observed: a first one in the
 276 left bottom includes two 'Euk profile' contigs; a second one located near the red bacterial region
 277 includes one 'No data' contig (light blue color). Below to the second cluster, there are also some
 278 contigs which we consider in-between the *M. vibrans* and the *E. aerogenes* regions. All these
 279 uncertain contigs (see yellow dots in Supplementary Information 1-Fig. 6) were aligned with the
 280 NCBI nt online database and results were manually inspected. We finally decided to include three

281 of these suspicious contigs, which were classified as contaminants according to the automated
282 or indirect strategies, into the 'Contaminant set', which also included all 'Curated Bact +' contigs.
283 Other contigs evaluated in ESOM were included into the 'Non-contaminant' set.

284

285 Overall, while the results shown by the sequence-similarity and the tetranucleotide distance
286 approaches are highly consistent between them, the combination of both methodologies was
287 necessary to detect the few contigs that were misclassified by either the taxonomic or the ESOM
288 approaches. Among the three strategies used for taxonomic classification, the results from ESOM
289 proved that the curated strategy was the most accurate. However, this strategy was only
290 applicable to *M. vibrans* data because the genomes of the two contaminants were available. The
291 combination of the automated and the indirect strategies also allowed to detect all contaminant
292 contigs, at the expense of some false positives. However, these false positive cases were later
293 corrected with ESOM, this pointing that the combination of the automated and indirect strategies
294 with a tetranucleotide distance analyses are good alternatives for complex metagenomic data.
295 On the one hand, the automated strategy should work for data with undetermined contaminant
296 bacteria because it uses the NCBI nt as database. On the other hand, given that amino acid
297 sequences allow to detect homology at larger evolutionary distances than nucleotides, the indirect
298 strategy should work for cases in which the contaminant genome is not available.

299



300

301 **Supplementary Information 1-Fig. 6.** ESOM map of the *M. vibrans* contigs from the first assembly, as in
302 Supplementary Information 1-Fig. 5, but colored in yellow those contig/contig windows that were further
303 inspected in alignments against NCBI databases.

304

305 1.4) Second round of read decontamination

306 We expected that a second assembly, after having removed at least a major fraction of the
307 contaminant reads, will greatly reduce potential uneven coverage problems, which are typical in
308 metagenomic data and may had limited the quality of the first assembly. The absence of this

309 constraint enables the usage of higher k-mers values and the mismatch correction mode during
310 the assembly (incompatible with the metagenome mode in *SPAdes*).

311

312 We first removed from the PE and MP libraries the reads that aligned with the 'Contaminant' but
313 not with the 'Non-contaminant' post-ESOM sets. For PE reads, we only removed them if both
314 paired reads satisfied this criterion. We used *bowtie2*¹⁷ v2.2.9 for read alignments. Then, the
315 surviving reads were assembled with *SPAdes* [-k 53,75,91,109 --careful --hqmp1-fr --cov-cutoff
316 auto]. The average coverage of the 109-mer assembly was 46.65.

317

318 We next compared the completeness and the contiguity of this second assembly with the first
319 assembly. Completeness and contiguity were estimated by aligning a set of *de novo* assembled
320 *bona fide* and non-redundant *M. vibrans* transcripts with both assemblies [*BLASTn*: -evalue 1e-
321 3]. In particular, we estimated completeness by counting how many transcripts aligned with the
322 genome with an average identity of >95% and with a total query coverage of >95%. Contiguity
323 was estimated as completeness, but only hits with the best scoring target scaffold were
324 considered (i.e., a transcript which sequence is complete but fragmented into distinct scaffolds
325 will sum for completeness but not for contiguity). The set of *bona fide* and non-redundant *M.*
326 *vibrans* transcripts was chosen first by reducing redundancy with *CD-HIT*¹⁸ v4.6 [-c 0.70], and
327 second by keeping only those transcripts without prokaryotes among whose three *BLASTx* best
328 targeting species (a total of 10056 transcript sequences) [-evalue 1e-3, -db euk_prok_db].
329 Completeness and contiguity measures were found to be better for the second assembly (9577
330 and 8889, respectively) than for the first assembly (9520 and 8348, respectively). This supports
331 the strategy of re-assembling the non-contaminant reads identified during the decontamination of
332 the first assembly, and also agrees with our decontamination approach performing well in terms
333 of specificity (i.e., very few *M. vibrans* genomic data was misidentified as bacterial contaminant).
334 Indeed, we only found 13 transcripts present in the first assembly that were not recovered in the
335 second assembly. We hence added to the second assembly the 8 contig fragments (3441 bp) of
336 the first assembly to which those 8 transcripts aligned (their sequence names include the suffix
337 '_fromdraftassembly').

338

339 Scaffolds were next submitted to a second round of decontamination. We first searched for
340 potential remaining vector/adaptor sequences using *BLASTn* with UniVec database (explained in
341 1.3.1). We removed one short scaffold (360 bp) and we also detected two scaffolds likely
342 containing contaminant sequences related to Illumina technology (we chose a score cutoff value
343 of >36.5 to distinguish between true and false positive contaminants, since it was the highest
344 value with which a scaffold from *C. owczarzaki*, the closest relative to *M. vibrans* and whose
345 genome was not sequenced with Illumina technology, aligned to a UniVec target related to
346 Illumina technology). We removed the aligned regions of the other two scaffolds and hence each
347 one was split into two sub-scaffolds. We did that just to ensure that our assembly did not include

348 fragments that were misassembled because reads with contaminant sequences connected during
349 the de Bruijn graph step⁶.

350

351 Before ESOM analyses, scaffolds were taxonomically-classified into eukaryotic/bacterial using
352 the indirect and the automated approaches (both explained in 1.3.2). Scaffolds were split in
353 different categories, according to the results from the taxonomic classification: 'Bact profile', 'Euk
354 profile', 'Automated Bact +', 'No data' (see 1.3.2); and we also considered an extra category to
355 include contigs classified as 'Euk profile' and 'Automated Bact+' ('Euk profile, Automated Bact+').
356 Because in ESOM each scaffold can only be present within one category, those detected as
357 bacterial by the two approaches were included in 'Bact profile' but not in 'Automated Bact+'.
358 Finally, we also incorporated into ESOM analyses all contaminant contigs from the first assembly.

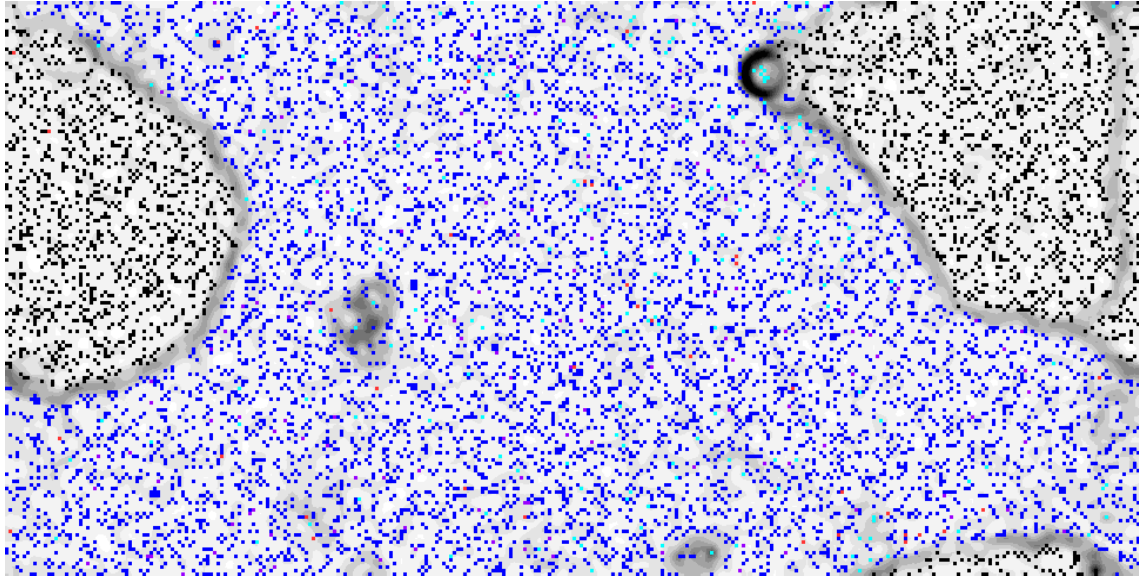
359

360 We did not find scaffolds with windows present inside and outside the region of the map with the
361 contaminant contigs from the first assembly (colored in black in Supplementary Information 1-Fig.
362 7). This suggests the absence of *M. vibrans/Bacteria* chimeric scaffolds. 'Euk profile' and 'Euk
363 profile, Automated Bact+' scaffolds that localized within the bacterial region were further inspected
364 by means of online *BLASTn/BLASTx* searches (colored in yellow in Supplementary Information
365 1-Fig. 8). Scaffolds surrounding grey areas outside the bacterial region (i.e., the region including
366 most of 'Euk profile' scaffolds, filled mostly with dark blue dots in Supplementary Information 1-
367 Fig. 7) were also inspected (colored in pink in Supplementary Information 1-Fig. 8). While the
368 white/black gradient suggests a distinct compositional pattern in these scaffold windows respect
369 to the average genome, we did not remove them because alignment results did not suggest
370 contamination either from *E. aerogenes/S. Maltophilia* scaffolds or other genomes. Moreover,
371 most of them have eukaryotic genes annotated, and we already rejected the possibility of
372 eukaryotic contamination during the decontamination of the first assembly. Still, four of these
373 scaffolds, with a length between and 2093 and 6202 bp, have bacterial but not eukaryotic genes
374 annotated. Because of this, we added the tag "_potentialcontaminant" as a suffix in their scaffold
375 names in the FASTA file.

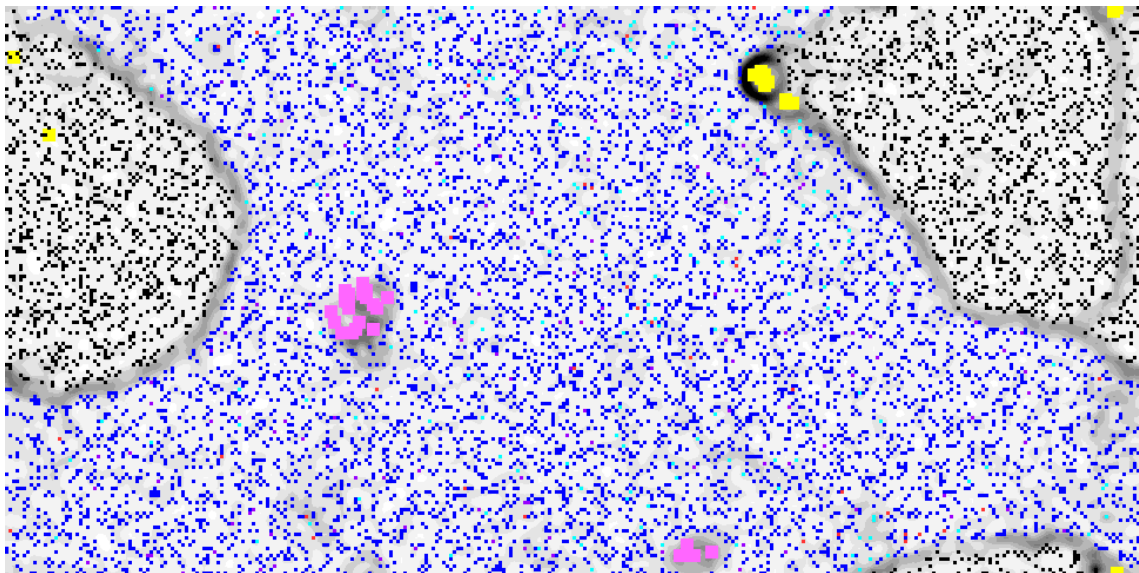
376

377 All scaffolds found outside the bacterial region as well as those labeled as "_potentialcontaminant"
378 were included in the *Mvib.gDNA.clean.v1.fasta* file (all FASTA files produced will be available
379 online as soon as we publish the corresponding manuscript). Because scaffolds with <2000 bp
380 were not considered during ESOM analyses, by default we did not include them in this set, with
381 two exceptions. First, those that were aligned by the set of reliable transcripts. Second, those with
382 eukaryotic genes annotated. In both cases, those scaffolds that also had one bacterial or potential
383 bacterial gene annotated were also labeled with the suffix "_potentialcontaminant". As a final step
384 to ensure that we kept only *bona fide M. vibrans* scaffolds, we aligned them [-evalue 1e-5] with
385 the contaminant database created during the curated strategy (see 1.3.2). Alignment results only
386 suggested one scaffold (141 bp) as potential contaminant, and hence was also labeled as

387 "_potentialcontaminant". In total, 27 of 2860 scaffolds (20516 bp of 29797085) were labeled as
388 "_potentialcontaminant". We also labeled the names of the scaffolds containing the mitochondrial
389 genome (assembled in a single scaffold; 55949 bp) and the 18S genes as "_mitochondrial" and
390 "_ribosomal", respectively.
391



392
393 **Supplementary Information 1-Fig. 7.** ESOM map of the *M. vibrans* scaffolds from the second assembly.
394 Each dot in the map correspond to one scaffold/scaffold window. Tetranucleotide frequency distances
395 between neighbour scaffold/scaffold windows are represented with a white/black background gradient for
396 smaller and larger distances, respectively. Color code: red for 'Bact profile', dark blue for 'Euk profile', green
397 for 'Automated Bact +', light blue 'No data' (see previous sections) purple for 'Euk profile, Automated Bact+',
398 and black for contaminant contigs from the first assembly.
399



400

401 **Supplementary Information 1-Fig. 8.** ESOM map of the *M. vibrans* scaffolds from the second assembly,
402 as in Supplementary Information 1-Fig. 7, but colored in yellow and pink those scaffolds/scaffold windows
403 that were further inspected in alignments against NCBI databases.

404

405 1.5) Benchmarking our supervised binning approach by comparing it to CONCOCT
406 (unsupervised approach)

407 As shown in previous sections, we used a supervised approach to classify (or bin) genomic
408 fragments from the *M. vibrans* metagenome, based on the taxonomic annotation and the
409 tetranucleotide distance-based clustering (i.e., ESOM). At the time of this analysis, unsupervised
410 binning tools such as CONCOCT¹⁹ were already proved to be successful with complex
411 metagenomic data (e.g., ²⁰). We compared the performance of our approach with the results
412 provided by CONCOCT v0.4.1.

413

414 CONCOCT grouped 4525 of the 20548 contigs from the first assembly into 61 bins, this
415 corresponding to 94.95% of the assembly length (excluding vectors). Regarding to the detection
416 of contaminant contigs, in our binning approach of the first assembly, 1008 of the contigs
417 (10331925 bp) were classified as bacterial contaminants after ESOM analyses. The CONCOCT
418 bins 10 and 60 included 83 of these 1008 contigs (10091058 bp). However, these two bins also
419 included one and two contigs, respectively, that were not classified as bacterial by our approach.
420 Based on the analyses of the BLAST results with NCBI nt, two of these three contigs would be
421 *bona fide* bacterial, and hence were misclassified by our approach as non-bacterial. The other
422 contig could have been misclassified as bacterial by CONCOCT. Apart of the bins 10 and 60, the
423 bins 44, 56 and 58 also included in total 7 additional contigs classified as bacterial by our approach
424 (3, 3 and 1 contigs, respectively; 12121 bp). Despite a manual inspection of BLAST results
425 suggest that these 7 contigs would be *bona fide* bacterial, the bins 56 and 58 also included 4 and
426 52 contigs that were not classified as bacterial by our approach. The 4 contigs within the bin 56
427 would be *bona fide* non-bacterial (indeed, the *M. vibrans* 18S ribosomal gene is one of these
428 contigs). In the case of bin 58, only 1 of the 52 contigs would be *bona fide* bacterial.

429

430 Regarding the detection of *M. vibrans* contigs, 4432 of the 19537 contigs that we classified as *M.*
431 *vibrans* were included in CONCOCT bins. These corresponds to the 93.97% of length spanned
432 by the 19537 contigs (29783450 bp, excluding the 3 contigs that our approach misclassified as
433 non-bacterial in the binning of the first assembly); indicating that ~6% of the data would have been
434 directly lost because of not having been included in any bin. 4376 contigs that are non-bacterial
435 according to our approach are in those CONCOCT bins that did not include *bona fide* bacterial
436 contigs (93.43% of 29783450 bp). Among them, the bin 40 includes 20884537 bp (2082 contigs)
437 and none of the *bona fide* bacterial contigs, suggesting that it includes a substantial fraction of
438 the *M. vibrans* genome. The remaining non-bacterial contigs according to our binning approach
439 (7104273 bp) are distributed in 59 other CONCOCT bins. We consider that this extra number of

440 bins does not reflect the presence of further contaminants sources, but instead of different
441 compositional features or uneven coverage between distinct regions of the *M. vibrans* genome.
442 First, because the comprehensive sequence-similarity and ESOM analyses done did not indicate
443 the presence of further contaminants apart of the two bacterial species. Second, because 98.04%
444 of the data that we considered as *M. vibrans* genome were in CONCOCT bins in which at least one
445 contig was aligned by a *bona fide* *M. vibrans* transcript (RNA-seq *de novo* assembled transcripts
446 whose best hit with euk_prok_db was a *C. owczarzaki* protein [the phylogenetically closest
447 genome to *M. vibrans* available], and whose three best target species are eukaryotes). Thus,
448 CONCOCT most likely put a substantial fraction of *M. vibrans* genome (7104273 bp) in separate
449 bins, some of which include bacterial contigs. Indeed, the bin 56 is a good example of this, as it
450 includes the *M. vibrans* 18S ribosomal gene together with three *bona fide* bacterial contigs and
451 three *bona fide* non-bacterial contigs, this showing the limitations of unsupervised approaches
452 when dealing with complex regions of eukaryotic genomes (e.g., unusual coverage or
453 compositional features).

454

455 Overall, these results suggest that our approach would have been more accurate than CONCOCT
456 with the *M. vibrans* data. On the one hand, although CONCOCT binned most of the *bona fide*
457 bacterial data into two bins (as expected), 7 *bona fide* bacterial contigs were misclassified in bins
458 mostly composed by *bona fide* *M. vibrans* contigs. Moreover, one *bona fide* non-bacterial contig
459 was included in one of the two bins that include most of the bacterial contigs. Despite our binning
460 approach in the first assembly misclassified three contigs as non-bacterial, they were probably
461 eliminated during the decontamination of the second assembly, since we did not found them in
462 *Mvib.gDNA.clean.v1.fasta*. On the other hand, CONCOCT split the *bona fide* *M. vibrans* contigs
463 in 59 bins, whereas we grouped them as a single bin (i.e., *Mvib.gDNA.clean.v1.fasta*). Another
464 relevant feature is that 5.05% of the data was not included in any CONCOCT bin. We thus
465 concluded that our approach was more accurate than the unsupervised approach, and hence we
466 also used it to decontaminate the genomes of *P. atlantis* and the two *Pigoraptor* species.

467

468 **2) *Parvularia atlantis* (Nucleariids, Opisthokonta)**

469

470 **2.1) Cultures, Cell cytometry and DNA sequencing**

471 We started from cultures of *P. atlantis*²¹ (formerly Nuclearia sp. ATCC 50694) growing in ATCC
472 Medium 802 and maintained at 23 °C. As the initial *M. vibrans* culture, *P. atlantis* grow in non-
473 axenic conditions with an undetermined diversity of contaminant bacteria. Because we could not
474 isolate and grow *P. atlantis* in better culture conditions, we directly extracted DNA from a pooling
475 of cells sorted by flow cytometry analyses (see *M. vibrans* section for FACS and DNA extraction
476 protocols). The cell sorting strategy allowed to enrich the *P. atlantis*/Bacteria ratio but not to get
477 rid of bacterial contamination, as shown by the presence of amplified 16S ribosomal PCR-product
478 in the extracted DNA (Supplementary Information 1-Fig. 3B). A total of 264 ng were obtained from

479 the extraction. This yield of DNA was enough to construct a PE but not a MP library (we only did
480 MP library for *M. vibrans*). The PE library was prepared and sequenced in a 50% Illumina HiSeq
481 2500 lane using the sequencing kit HiSeq v4 chemistry (read insert size: 520 bp; read length: 125
482 bp).

483

484 PE reads were preprocessed with *trimmomatic*, using the following parameters: *LEADING:30*
485 *TRAILING:20 SLIDINGWINDOW:2:20 ILLUMINACLIP:2:30:10 MINLEN:80*. As with *M. vibrans*,
486 TruSeq-PE related adapter sequences from *trimmomatic* were used as contaminant database for
487 *ILLUMINACLIP*. The single and paired preprocessed reads were submitted to a read correction
488 step using *SPAdes* [--only-error-correction].

489

490 2.2) Benchmarking of metagenomic assemblers

491 The genomic data of *M. vibrans* was a simple metagenome case, where most of the reads
492 belonged to the eukaryotic species, with also reads belonging to two well defined contaminant
493 bacteria. In the case of *P. atlantis*, despite we enriched the proportion of eukaryotic cells by cell
494 cytometry sorting, we also expected a substantial fraction of reads from an uncertain diversity of
495 contaminant bacterial species (Supplementary Information 1-Fig. 3B). Hence, we decided to
496 benchmark *metaSPAdes* (used for *M. vibrans*) with two other popular metagenome assemblers
497 available at that time: *IDBA-UD*²² and *Ray*²³ *Meta v.2.3.1*. In particular, we evaluated standard
498 contiguity metrics (e.g., N50, L75) and also the tendency to assemble potential chimeric contigs.

499

500 For that, the three metagenome assemblies were ran using the preprocessed and corrected
501 paired and unpaired reads, with default assembly parameters. Because we were interested in the
502 fraction of the assembly corresponding to *P. atlantis*, we ran *BUSCO*²⁴ *v1.22*, using the *all*
503 *Eukaryota* dataset, in order to obtain the *BUSCO orthologs* from our metagenomes. We then
504 aligned these *BUSCO orthologs* with *euk_prok_db*, and we kept those whose best scoring hit was
505 a eukaryotic protein. We expected at least most of the contigs encoding these *bona fide*
506 eukaryotic *BUSCO orthologs* (*BUSCO contigs*) to correspond to *P. atlantis* genome. *BUSCO*
507 *contigs* were annotated using the *BRAKER1* pipeline, and every predicted gene was
508 taxonomically classified into eukaryotic or bacterial following the indirect strategy (see 1.3.2).
509 RNA-seq reads, required by the *BRAKER1* pipeline, were downloaded from NCBI
510 (SRR1617645), preprocessed using *trimmomatic* [*ILLUMINACLIP:2:30:10*
511 *SLIDINGWINDOW:4:5 LEADING:5 TRAILING:5 MINLEN:25*], and corrected using *SEECER*.

512

513 Completeness and contiguity of *BUSCO contigs* were estimated with *QUAST*²⁵ *v4.2* and *BUSCO*.
514 In agreement with the bibliography²⁶, *metaSPAdes* outperformed *IDBA-UD* and *Ray Meta* in most
515 of the metrics (see Supplementary Information 1-Table 1 below). *metaSPAdes* was also the
516 assembler with the highest number of both eukaryotic and bacterial genes in *BUSCO contigs*,
517 and showed the lowest ratio of *BUSCO contigs* without bacterial genes. While these results could

518 be interpreted as *metaSPAdes* being the assembler with more chimeric contigs, they could also
 519 be explained because *BUSCO contigs* were more and longer than in the other two assemblers.
 520 Moreover, because not all genes annotated as bacterial are necessarily *bona fide* contaminants
 521 (e.g., horizontal gene transfer, low scoring *BLAST* hits with bacterial sequences -which are more
 522 represented than eukaryotes in our database-, etc.); not every contig with genes annotated as
 523 bacterial is necessarily a chimera.

524

525 Because each species had its own abundance in the sequenced sample, differences in coverage
 526 between those genomic fragments corresponding to *P. atlantis* and those corresponding to
 527 contaminant bacteria can be expected. If so, and if there are chimeric contigs in the assemblies,
 528 which is uncertain, we would expect them to probably correspond to those contigs showing larger
 529 differences between the coverages of the regions corresponding to eukaryotic genes and those
 530 corresponding to bacterial genes. We hence considered the number of *BUSCO contigs* showing
 531 elevated internal coverage differences as a proxy to estimate the tendency of every assembler to
 532 construct potential chimeric contigs. *BUSCO contigs* with internal coverage differences were
 533 determined in the following manner: for each *BUSCO contig* that contained eukaryotic and
 534 bacterial genes, we computed (i) the coverage of each gene (*cov_gene*), and (ii) the average
 535 coverage of all eukaryotic genes in this specific contig (*cov_all_euk_genes*). Then, we calculated
 536 the absolute distance between the coverage of each gene (eukaryotic or bacterial) and this
 537 average (*dist_gene* = $\text{abs}(\text{cov_gene} - \text{cov_all_euk_genes})$). If the highest distance corresponded
 538 to a bacterial gene, we considered this contig as a potential chimera. The coverage of each gene
 539 was estimated with *samtools*²⁷ 1.3.1 *depth*, previously mapping the clean and preprocessed DNA-
 540 seq reads to *BUSCO contigs* using *bowtie2* [--no-mixed --no-discordant --maxins 750].

541

542 The number of contigs showing potential chimeric features were the same in *metaSPAdes* and
 543 *IDBA-UD*, both lower than in *Ray Meta*. However, *metaSPAdes* showed the proportionally less
 544 potential chimeric contigs (i.e., the best), given that the number of *BUSCO contigs* is higher than
 545 in *IDBA-UD*. Thus, we concluded that *metaSPAdes* was the best assembler for our data because
 546 since it showed better contiguity and completeness measures than the other two assemblers, and
 547 this does not seem to come at the expense of increasing the probability of assembling contigs
 548 showing chimeric-like features. Our results agree with a published benchmarking of metagenome
 549 assemblers, which recommends the use of *metaSPAdes* when the main objective is to retrieve
 550 the genome of a species that is well represented in the sample²⁶.

551

552 **Supplementary Information 1-Table 1.** Benchmarking of three popular metagenomic
 553 assemblers for *P. atlantis* data.

| # All contigs | <i>metaSPAdes</i> | <i>IDBA-UD</i> | <i>Raymeta</i> |
|------------------------------------|-------------------|----------------|----------------|
| Complete BUSCOs (total 429) | 275 | 260 | 265 |

| | | | |
|---|-----------|-----------|-----------|
| Complete and single-copy BUSCOs | 230 | 212 | 218 |
| Complete and duplicated BUSCOs | 45 | 48 | 47 |
| Fragmented BUSCOs | 72 | 85 | 83 |
| Missing BUSCOs | 82 | 84 | 81 |
| Contigs | 61669 | 36486 | 187875 |
| Contigs > 999 bp | 9647 | 14036 | 10078 |
| N50 | 97390 | 40013 | 97458 |
| L75 | 1235 | 2625 | 1132 |
| Genome size (Mb) | 158747712 | 154204313 | 152524304 |
| Genome size, only > 999 bp (Mb) | 142696796 | 143141922 | 117504526 |
| Genome size, only > 49999 bp (Mb) | 88898272 | 69342700 | 73494015 |
| # Contigs with single copy complete BUSCOs euks (<i>BUSCO contigs</i>) | | | |
| Complete BUSCOs (total 429) | 238 | 211 | 217 |
| Complete and single-copy BUSCOs | 237 | 210 | 215 |
| Complete and duplicated BUSCOs | 1 | 1 | 2 |
| Fragmented BUSCOs | 35 | 24 | 34 |
| Missing BUSCOs | 156 | 194 | 178 |
| Contigs | 181 | 179 | 179 |
| Contigs > 24999 bp | 75 | 30 | 45 |
| N50 | 35054 | 20656 | 26013 |
| L75 | 86 | 88 | 85 |
| Genome size (Mb) | 4600459 | 2945403 | 3533535 |
| Genome size, only > 24999 bp (Mb) | 3190783 | 1119512 | 1827248 |
| N° euk genes | 1682 | 1116 | 1316 |
| N° bact genes | 86 | 56 | 73 |
| % euk genes / total genes | 95,14% | 95,22% | 94,74% |
| Contigs with at least 1 bact gene | 118 | 50 | 57 |
| Contigs without bact genes | 63 | 129 | 122 |
| % contigs without bact genes / total contigs | 34,81% | 72,07% | 68,16% |
| Potential chimeric contigs | 16 | 16 | 19 |
| % potential chimeric contigs / total contigs | 8,84% | 8,94% | 10,61% |

554

555

556 2.3) First round of read decontamination

557 The assembled contigs from *metaSPAdes* were evaluated for potential contamination (see 1.3.1).
558 19 contigs were removed because they showed >50% of average identity and >95% average
559 coverage with UniVec sequences. We only found one contig showing signatures of a *bona fide*
560 18S gene, of 10705 bp and high assembly coverage (a metric computed by *SPAdes* for every
561 contig, available in the sequence name). Its best hit with a local 18S database corresponded to
562 the 18S sequence of *P. atlantis*. We also found 5 other contigs aligning to the 18S database, but
563 were all of short length (<302 bp) and had very low assembly coverage values. We also found 71
564 contigs with potential 16S sequences.

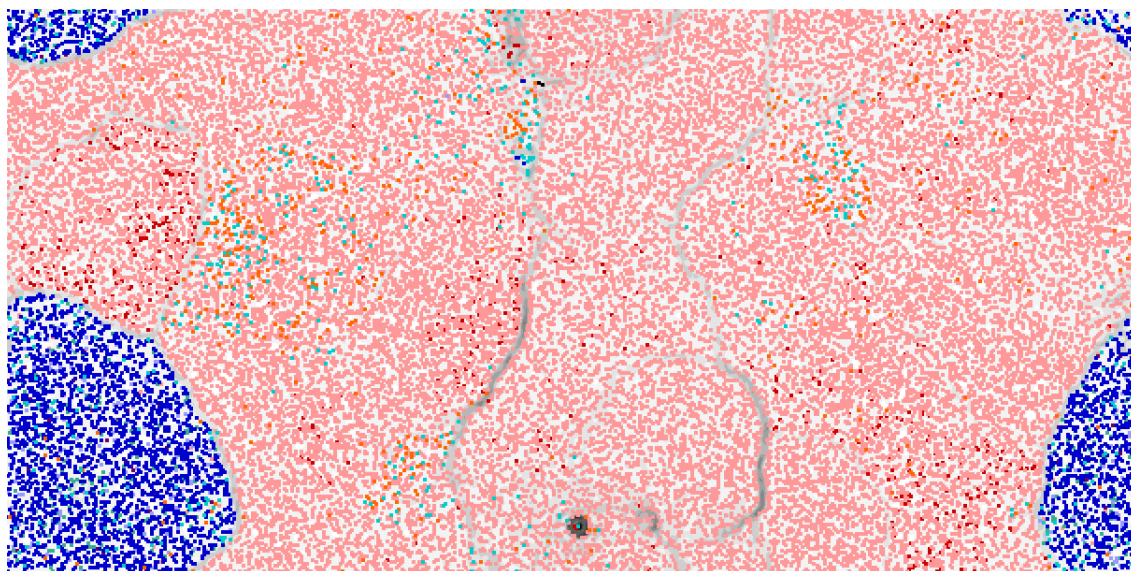
565

566 We also searched for potential mitochondrial sequences in contigs (see 1.3.2). 101 contigs were
567 aligned by the *A. godoyi* sequences, most of which were expected to be bacterial contamination.
568 We used the output from *tBLASTn* to extract the putative amino acid sequences from the
569 alignments, and then aligned the extracted sequences with a database including all sequences
570 from prok_db and the NCBI CDS translations of complete mitochondrial genomes [mito_db],
571 downloaded from <ftp://ftp.ncbi.nlm.nih.gov/blast/db/> [*BLASTp*: -evalue 1e-3]. We considered as
572 potential mitochondrial sequences 5 contigs whose predicted sequences performed the best hit
573 with a mitochondrial protein. We expect two of these contigs to correspond to *P. atlantis*
574 mitochondria. First, because they have a length of 9773 and 13093 bp and a high assembly
575 coverage. Second, because the three sequences predicted from them (COX-I, COX-II, and Cyt-
576 b) did not align with a perfect identity with NCBI nt. The other three contigs showed residual length
577 and assembly coverage values, and none of them aligned with a perfect identity with NCBI nt,
578 suggesting that, as with the spurious 18S contigs, they likely corresponded to miss-assembled
579 mitochondrial regions of *P. atlantis* rather than to eukaryotic contamination.

580

581 Contigs were taxonomically classified using the indirect approach and the automated direct
582 approach (see 1.3.2). All contigs > 2000 bp were submitted to ESOM analyses, splitting contigs
583 larger than 7999 bp into contigs windows of 4000 bp. Before ESOM analyses, contigs were
584 classified into one of the following categories, and colored accordingly in the ESOM map
585 (Supplementary Information 1-Fig. 9): '16S', '18S', 'Bact profile' (excluding contigs in the 16S
586 category), 'BUSCO Automated Bact +' (*BUSCO contigs* classified as bacterial according to the
587 direct approach), 'BUSCO Automated Bact -' (*BUSCO contigs* non-classified as bacterial
588 according to the direct approach), 'Euk profile Automated Bact +' (non-*BUSCO contigs* classified
589 as eukaryotic according to the indirect approach, and as bacterial according to the direct
590 approach), 'Euk profile Automated Bact -' (non-*BUSCO contigs* classified as eukaryotic according
591 to the indirect approach, and non-classified as bacterial according to the direct approach),
592 'Automated Bact +' (other contigs only classified as bacterial according to the direct approach)
593 and 'No data'.

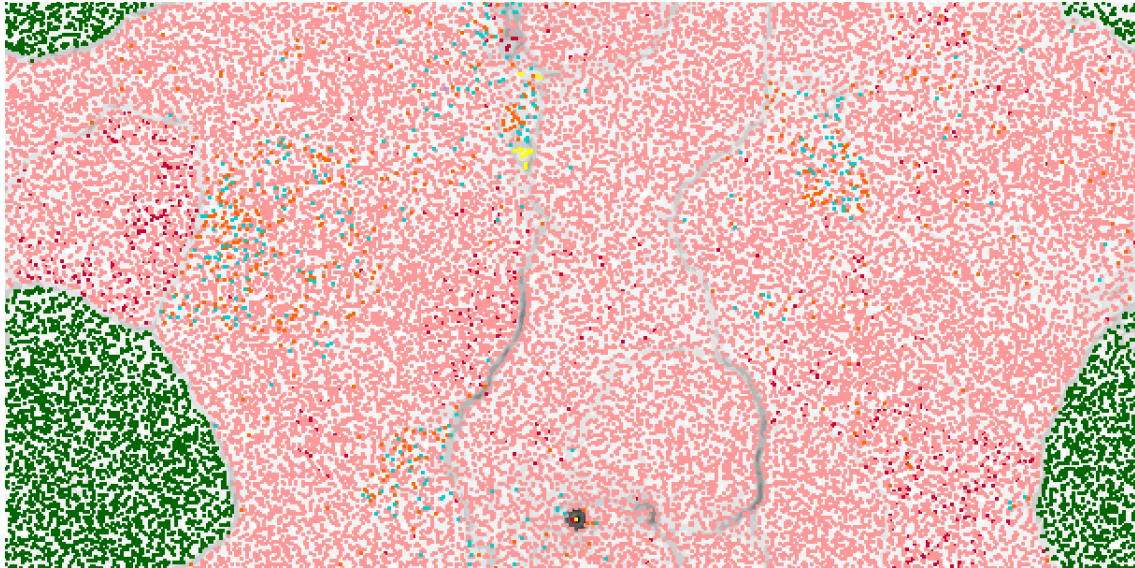
594



595
 596 **Supplementary Information 1-Fig. 9.** ESOM map of the *P. atlantis* contigs from the first assembly. Each
 597 dot in the map correspond to one contig/contig window, which are colored according to the category to which
 598 was classified. Category color code: red for '16S', black for '18S', pink for 'Bact profile', dark purple for
 599 'BUSCO Automated Bact +', turquoise for 'BUSCO Automated Bact -', light purple for 'Euk profile Automated
 600 Bact +', dark blue for 'Euk profile Automated Bact -', orange for 'Automated Bact +', and light blue for 'No
 601 data'. Tetranucleotide frequency distances between neighbour contig/contig windows are represented with a
 602 white/black background gradient for smaller and larger distances, respectively.

603
 604 As expected, *BUSCO contigs* and most of the 'Euk profile' contigs clustered together in the map
 605 (see dark blue dots regions in Supplementary Information 1-Fig. 9), and we hence considered
 606 that region as the *P. atlantis* genome (see dark green dots in Supplementary Information 1-Fig.
 607 10). The 'Non-contaminant' set included all contigs in this region, as well as the two contigs
 608 previously found to contain mitochondrial sequences. A small fraction of 'Euk profile' contigs were
 609 found in a separate region in the north of the map, surrounded by contigs classified as 'Bact
 610 profile', 'Automated Bact+' and 'No data' (Supplementary Information 1-Fig. 9). We considered
 611 these contigs and their neighbours as uncertain (see yellow dots in Supplementary Information
 612 1-Fig. 10). 12 of 53 uncertain contigs were included in the 'Non-contaminant' set because they
 613 corresponded to the mitochondrial or to the 18S contig or because being 'Euk profile'. The other
 614 41 uncertain contigs were included in the 'Contaminant' set together with the other contigs of the
 615 map that were not in the *P. atlantis* dataset.

616



617
 618 **Supplementary Information 1-Fig. 10.** ESOM map of the *P. atlantis* contigs from the first assembly, as in
 619 Supplementary Information 1-Fig. 9, but colored in green those contig/contig windows included in the *P.*
 620 *atlantis* dataset, and colored in yellow those contig/contig windows that further inspected to determine
 621 whether they should be included in the *P. atlantis* dataset, or in the 'Contaminant' set. Other contig/contig
 622 windows were included in the 'Contaminant' set.

623

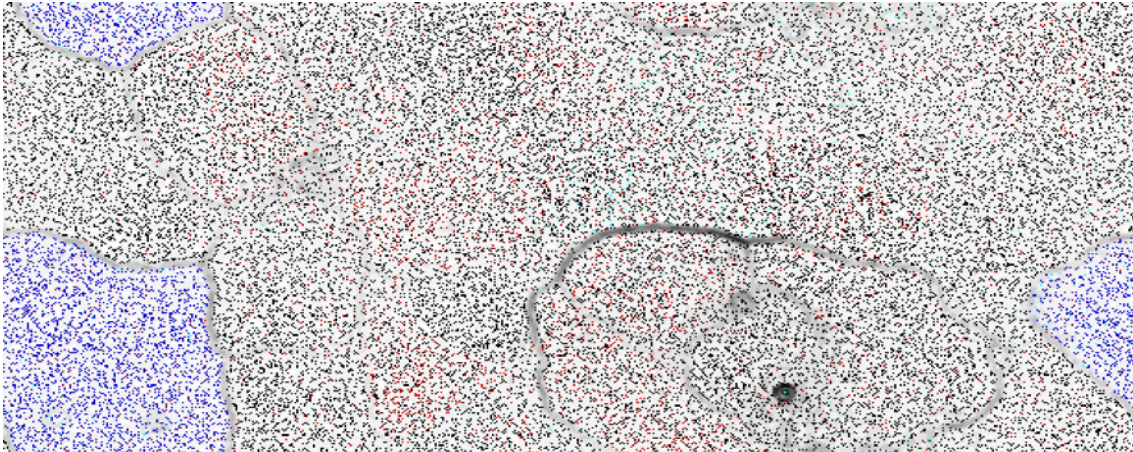
624 2.4) Second round of read decontamination

625 We removed from the PE library the reads that aligned with the 'Contaminant' but not with the
 626 'Non-contaminant' post-ESOM sets (as we did for *M. vibrans*, see 1.4). The surviving reads were
 627 assembled using *SPAdes* with the metagenome disabled, and using the following k-mers: 33, 55,
 628 71, 89 [--careful --cov-cutoff auto]. The average coverage of the 89-mer assembly was 36.28.
 629 Then, as with the first assembly, we screened the assembled scaffolds for potential contamination
 630 of adapter/vector sequences. One scaffold that aligned with Illumina targets with >36.5 *BLASTn*
 631 score (see *M. vibrans* sections) was directly removed because of its short length and low
 632 coverage. We further removed other 7 scaffolds of short length (3470 bp) because they aligned
 633 with vector sequences along most of their sequence. After the adapter/vector decontamination
 634 step, scaffolds were taxonomically classified using the indirect approach and the automated direct
 635 approach (see 1.3.2).

636

637 All scaffolds > 2000 bp were submitted to ESOM analyses, splitting contigs larger than 7999 bp
 638 into contig windows of 4000 bp. Before ESOM analyses, contigs were classified into one of the
 639 following categories, and colored accordingly in the ESOM map (Supplementary Information 1-
 640 Fig. 11): 'Bact profile', 'Euk profile Automated Bact +' (contigs classified as eukaryotic according
 641 to the indirect approach, and as bacterial according to the direct approach), 'Euk profile
 642 Automated Bact -' (contigs classified as eukaryotic according to the indirect approach, and non-
 643 classified as bacterial according to the direct approach), 'Automated Bact +' (contigs only

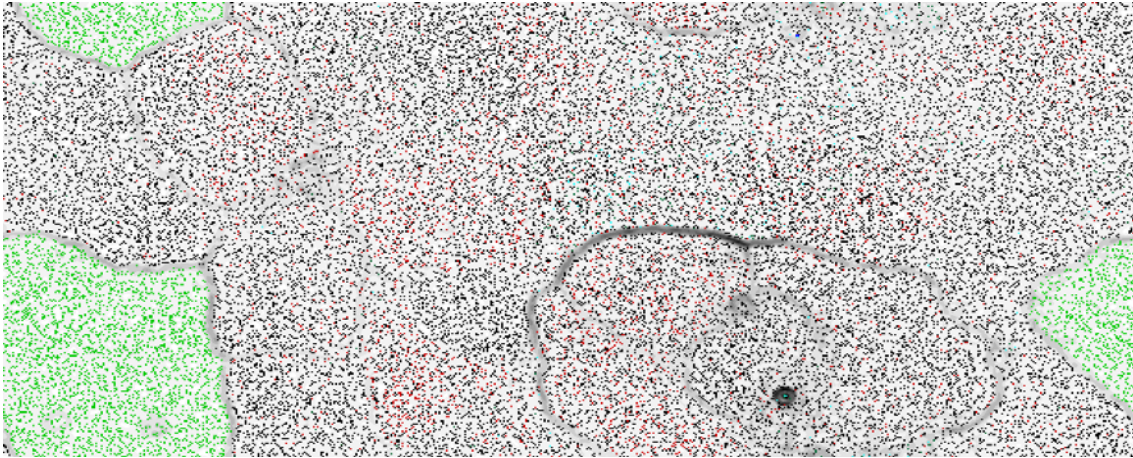
644 classified as bacterial according to the direct approach) and 'No data'. Finally, we also
645 incorporated into ESOM analyses all contaminant contigs from the first assembly.
646



647
648 **Supplementary Information 1-Fig. 11.** ESOM map of the *P. atlantis* scaffolds from the second assembly.
649 Each dot in the map correspond to one scaffold/scaffold window. Tetranucleotide frequency distances
650 between neighbour scaffold/scaffold windows are represented with a white/black background gradient for
651 smaller and larger distances, respectively. Color code: red for 'Bact profile', purple for 'Euk profile Automated
652 Bact +', dark blue for 'Euk profile Automated Bact -', dark green for 'Automated Bact +', light blue 'No data'
653 (see previous sections) purple for 'Euk profile, Automated Bact+' and black for contaminant contigs from the
654 first assembly.

655
656 All except one 'Euk profile' scaffold windows clustered in the same region of the map (dark blue
657 dots in Supplementary Information 1-Fig. 11), and hence all windows within that region were
658 considered to belong to *P. atlantis* genome (light green dots Supplementary Information 1-Fig.
659 12). We inspected the taxonomic classification of the genes from the scaffolds within the *P.*
660 *atlantis* region (colored in light green in Supplementary Information 1-Fig. 12) in order to detect
661 potential misclassified contaminants. Five scaffolds that showed a suspicious pattern (few
662 eukaryotic genes and some bacterial genes predicted) were finally kept because of their ratio of
663 introns per gene (between 3-10.67 introns per gene). We consider the later a good marker of non-
664 bacterial genes as the ratio of introns per gene in the contaminant contigs of the first assembly
665 was 0.13 (10705 introns in 78682 predicted genes).

666



667

668 **Supplementary Information 1-Fig. 12.** ESOM map of the *P. atlantis* scaffolds from the second assembly,
 669 as in Supplementary Information 1-Fig. 11, but colored in light green the scaffold windows within the region
 670 of the map considered to correspond to *P. atlantis* genome.

671

672 To detect potential chimeric scaffolds, we inspected the ESOM class file for scaffolds with
 673 windows within the *P. atlantis* region and also outside this region. This revealed two potential
 674 chimeric scaffolds, a first one (24988 bp) was retained because it contained *bona fide* eukaryotic
 675 genes (verified by aligning them with NCBI nt) and the scaffold region corresponding to the
 676 window outside the *P. atlantis* region did not contain any bacterial gene predicted. The other
 677 potential chimera (8885 bp) was removed because of its short length and also because it did not
 678 include any *bona fide* eukaryotic gene.

679

680 Since scaffolds with less than 2000 of length were not considered in ESOM analyses, we
 681 excluded them from *Patl.gDNA.clean.v1.fasta*, with two exceptions. 1) We kept those aligning
 682 with >95% of identity to the set of *bona fide* *P. atlantis* transcripts (identified with the same
 683 approach as we did for *M. vibrans*), and 2) those containing annotations of eukaryotic genes. In
 684 both cases, scaffolds containing at least one bacterial or potential bacterial gene annotated were
 685 labeled as “_potentialcontaminant”. As a final step to ensure that we did not include bacterial
 686 scaffolds in *Patl.gDNA.clean.v1.fasta*, we aligned them with NCBI nt [*BLASTn*: -task megablast;
 687 -eval 1e-5] and the 4 scaffolds with less than 2000 of length that performed its best hit with a
 688 prokaryotic target (with >25% of average query coverage and > 50% of average identity) were
 689 also labeled as “_potentialcontaminant”. In the end, 13 of 2555 scaffolds (9596 of 19259675 bp)
 690 were labeled as “_potentialcontaminant”. We also labeled the scaffolds with putative
 691 mitochondrial and 18S ribosomal sequences as “_putative-mitochondrial” and “_ribosomal”,
 692 respectively (3 putative mitochondrial scaffolds, 45473 bp in total).

693

694 We finally assessed completeness and contiguity of the decontaminated genome using a set of
 695 *de novo* assembled *bonafide* and non-redundant 5841 *P. atlantis* transcripts. 5461 and 5179 of
 696 these transcripts were complete and contiguous, respectively (5464 and 5147 in the case of the

697 first assembly contigs). We added complete and partial sequences from three contigs of the first
698 assembly (2665 bp, labeled with the suffix “_fromdraftassembly”) found to be incomplete in the
699 second assembly.

700

701 **3) *Pigoraptor* species (Filasterea, Opisthokonta)**

702

703 3.1) Cultures, Cell cytometry, and DNA sequencing

704 As *M. vibrans* and *P. atlantis*, both *Pigoraptor* species grow in non-axenic condition with an
705 undetermined diversity of contaminant bacteria but also with the eukaryotic contaminant
706 *Parabodo caudatus*, used as prey. For both *Pigoraptor* species, we extracted DNA from a pooling
707 of cells sorted by flow cytometry analyses (see *M. vibrans* section for FACS and DNA extraction
708 protocols). Because the FACS protocol is designed to sort eukaryotic cells, we also expected
709 contamination from *P. caudatus*. Hence, to *in silico* decontaminate both *Pigoraptor* libraries, we
710 also obtained DNA from sorted cells from *P. caudatus* cultures. In total, we obtained 14.70 ng,
711 16.14 ng, and 11.73 ng for *P. vietnamica*, *P. chilena* and *P. caudatus* poolings, respectively. The
712 PE libraries were prepared and sequenced each one in a 50% Illumina HiSeq 3000 lane using
713 the sequencing kit HiSeq 3000/4000 chemistry (read insert sizes: 390, 410 and 410 bp for *P.*
714 *vietnamica*, *P. chilena*, and *P. caudatus*, respectively; read lengths: 150 bp).

715

716 PE reads from both *Pigoraptor* species and from *P. caudatus* were preprocessed with
717 *trimmomatic*, using the following parameters: *SLIDINGWINDOW:12:30 LEADING:30*
718 *ILLUMINACLIP:2:30:10 MINLEN:80*. As with *M. vibrans* and *P. atlantis*, TruSeq-PE related
719 adapter sequences from *trimmomatic* were used as contaminant database for *ILLUMINACLIP*.
720 The single and paired preprocessed reads were submitted to a read correction step using *SPAdes*
721 [--only-error-correction].

722

723 3.2) Preliminary decontamination

724 The *P. caudatus* library was sequenced for decontamination purposes. Hence, we first assembled
725 the *P. caudatus* reads using *SPAdes* with the metagenome mode. Then, we aligned the *P.*
726 *vietnamica* and *P. chilena* reads with the *P. caudatus* scaffolds using *bowtie2*; in order to keep
727 only those reads that did not align with *P. caudatus* scaffolds. In particular, we considered as
728 potential *P. caudatus* contaminants all the unpaired reads (UP) that aligned, and all paired-end
729 reads (PE) that aligned concordantly (i.e., both paired reads aligned to the same scaffold with an
730 in-between distance close to the distance expected from the selected insert sizes). In order to
731 discard false positive contaminant cases, we did not remove those reads that also aligned with
732 the *C. owczarzaki* genome (the closest relative genome to *Pigoraptor* available), as they may
733 correspond to highly conserved genomic regions. In total, excluding the reads that aligned with
734 both genomes (0.16% UP and 0.016% PE of the reads that aligned with *P. caudatus* also aligned
735 with *C. owczarzaki*), ~50% of *P. vietnamica* data (49.09% UP, 46.03% PE) aligned with *P.*

736 *caudatus* data (which includes *P. caudatus* and also uncertain Bacteria). In the case of *P.*
737 *chileana*, the percentages of reads that aligned with *P. caudatus* but not with *C. owczarzaki* were
738 61.24% for UP, and 43.79% for PE (0.048% of UP and 0.0078% of PE reads that aligned with *P.*
739 *caudatus* and also with *C. owczarzaki*).

740

741 3.3) First round of read decontamination (*P. vietnamica* and *P. chileana*)

742 For each *Pigoraptor* species, we next assembled the surviving reads using *metaSPAdes*, and the
743 resulting contigs were evaluated for potential contamination as we did for *M. vibrans* and *P.*
744 *atlantis* (i.e., adapters/vectors, 16S, 18S, and mitochondria; see 1.3.1).

745

746 In the case of *P. vietnamica*, we removed a total of 51 contigs (61268 bp) that aligned with > 50%
747 of average query coverage and average identity with VecScreen sequences. We also found 72
748 contigs with potential 16S sequences. The contig with the highest assembly coverage (1980.97
749 cov) among those aligned by 18S sequences corresponded to the *P. vietnamica* 18S sequence.
750 We also found a contig corresponding to the 18S of *P. caudatus*, but it had much less assembly
751 coverage (57.45 cov) than the *P. vietnamica* 18S contig. This suggests that the read mapping-
752 based decontamination approach reduced but did not fully removed all contaminant reads from
753 *P. caudatus*. We also found one contig aligned by 18S sequences (633 bp, 9.30 cov) with 100%
754 query coverage and identity to *H. sapiens* ribosomal sequences. The same contig aligned with
755 lower identity to *C. owczarzaki* genome (77.40%) and with worse alignment metrics with fungal
756 sequences on NCBI nt. Because we also found 1 contig with potential *H. sapiens* mitochondrial
757 sequences (1429 bp, 6.05 cov), these altogether suggested the presence of low coverage *H.*
758 *sapiens* contamination in *P. vietnamica* data.

759

760 In the case of *P. chileana*, we removed a total of 35 vector/adaptor contigs (8140 bp), and we
761 found 44 contigs with potential 16S sequences. As with *P. vietnamica*, we also found contigs
762 corresponding to *P. caudatus* (679 bp, 343.19 cov; 230 bp, 26.92 cov) and possibly to *H. sapiens*
763 (271 bp, 1.74 cov). The screening of mitochondrial sequences revealed two contigs likely
764 containing the *P. chileana* mitochondria (12883 bp and 275.76 cov; 20715 bp and 273.77 cov),
765 but also one potential contaminant from *P. caudatus* (20289 bp and 552.29 cov) and from *H.*
766 *sapiens* (207 bp and 1.07 cov). These suggested as well the presence of potential *H. sapiens* low
767 coverage contamination also in the *P. chileana* data.

768

769 In *M. vibrans* and *P. atlantis* data, the taxonomic classification of contigs into eukaryotic and
770 bacterial allowed to distinguish between the regions of the ESOM map that corresponded to our
771 organism of interest (the only eukaryote), and the regions that corresponded to bacterial
772 contamination. This was not applicable for both *Pigoraptor* species given the presence of *P.*
773 *caudatus* contamination. In both cases, we used an alternative approach that consisted in

774 identifying a set of *bona fide* *P. vietnamica*, *P. chiliana* and *P. caudatus* scaffolds that were used
775 to label the different regions of the map (see below).

776

777 For both *Pigoraptor* species, we aligned the proteins predicted from their transcriptomic data²⁸
778 with euk_prok_db [*BLASTp*: -evalue 1e-3 -task blastp-fast], and we kept those whose best hit was
779 a *C. owczarzaki* protein (the closest relative to *Pigoraptor* in this database). We then aligned these
780 proteins with the *P. caudatus* metagenome scaffolds and also with the contigs of the
781 corresponding *Pigoraptor* species [*tBLASTn*: -evalue 1e-3]. Contigs encoding for those proteins
782 that aligned with the *Pigoraptor* but not with the *P. caudatus* metagenome were considered as
783 *bona fide* *Pigoraptor* contigs. We assumed the contig to which each protein aligned with the
784 highest score as the encoding contig for that protein.

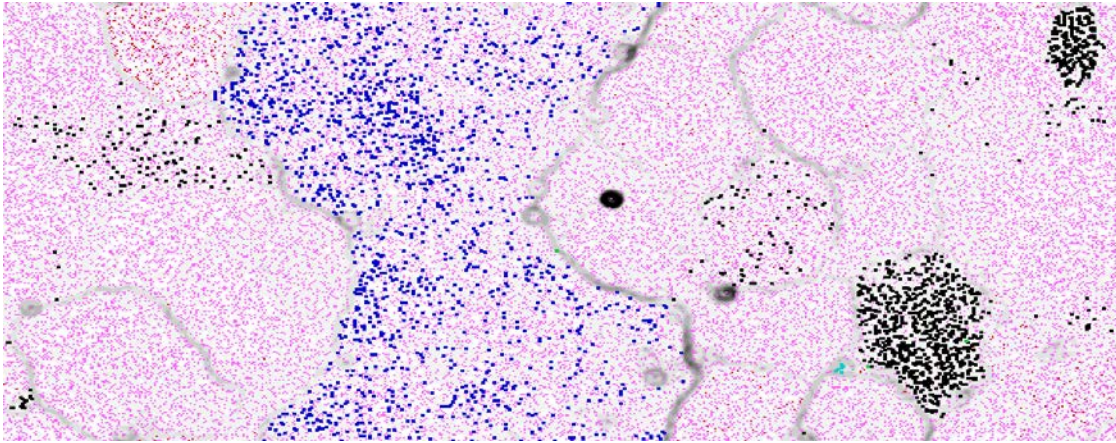
785

786 The set of *bona fide* *P. caudatus* scaffolds was defined in the following manner. We took the
787 proteins from the transcriptomic data of *P. chiliana* (from which *P. caudatus* contamination is
788 expected) that did not align with prokaryotic proteins in the *BLASTp* search with euk_prok_db
789 (see previous paragraph). We then aligned these likely eukaryotic proteins with the *P. caudatus*
790 metagenome [*tBLASTn*: -evalue 1e-49], and we keep as *bona fide* *P. caudatus* scaffolds those in
791 which at least one of these proteins performed aligned with the highest score.

792

793 For each *Pigoraptor* species, contigs > 2000 bp were submitted to ESOM analyses, splitting
794 contigs larger than 7999 bp into contig windows of 4000 bp (SM Figs 13-14 and 15-16 correspond
795 to the ESOM maps of *P. vietnamica* and *P. chiliana*, respectively). Contigs were colored in ESOM
796 maps according to the categories to which they belong: '16S', '18S', 'mitochondria', '*bona fide*
797 *Pigoraptor*', 'Others'. In both ESOM maps we also incorporated the *bona fide* *P. caudatus*
798 scaffolds (see above) in order to detect which regions of both *Pigoraptor* assemblies
799 corresponded to *P. caudatus* contamination. In the case of *P. vietnamica* data, all *bona fide* *P.*
800 *vietnamica* contigs (Supplementary Information 1-Fig. 13) were found within a region that did not
801 include any *bona fide* *P. caudatus* contig (black contigs in Supplementary Information 1-Fig. 13).
802 Accordingly, all contigs with a window within that region were included in the 'Non-contaminant'
803 set (colored in orange in Supplementary Information 1-Fig. 14). In contrast, all contigs in the map
804 without windows within this region were included in the 'Contaminant' set (excluding the putative
805 mitochondrial and 18S contigs of *P. vietnamica*).

806



807

808 **Supplementary Information 1-Fig. 13.** ESOM map of the *P. vietnamica* contigs from the first assembly.

809 Each dot in the map correspond to one contig/contig window, which are colored according to the category

810 to which was classified. Category color code: red for '16S', green for '18S', light blue for 'mitochondria', dark

811 blue for '*bona fide P. vietnamica*', pink for 'Others'. We also incorporated the *bona fide P. caudatus* scaffolds

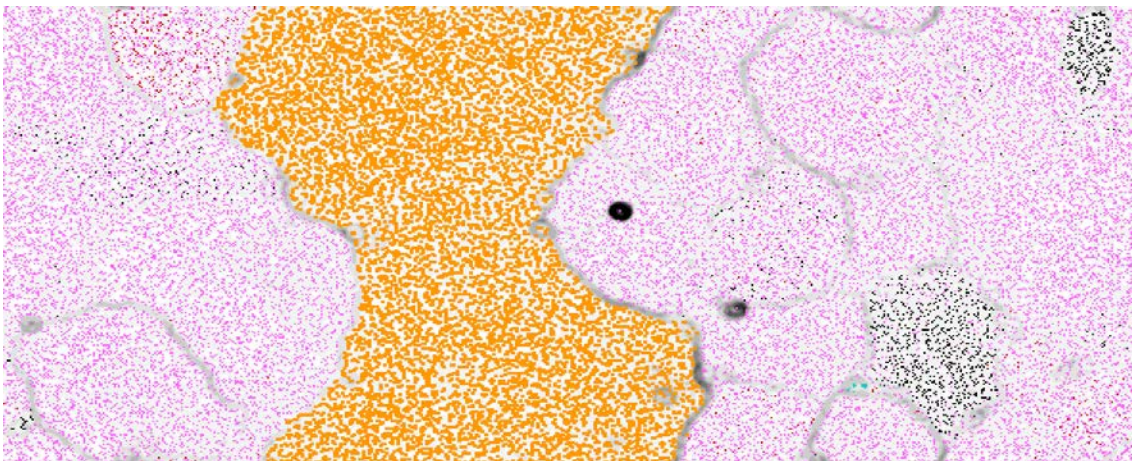
812 (from the *P. caudatus* metagenome, colored in black) in order to detect which regions of the *P. vietnamica*

813 assembly correspond to the *P. caudatus* genome. Tetranucleotide frequency distances between neighbour

814 contig/contig windows are represented with a white/black background gradient for smaller and larger

815 distances, respectively.

816



817

818 **Supplementary Information 1-Fig. 14.** ESOM map of the *P. vietnamica* contigs from the first assembly, as

819 in Supplementary Information 1-Fig. 13, but colored in orange the contigs windows within the region of the

820 map considered to correspond to *P. vietnamica* genome.

821

822 In the case of *P. chiliana*, however, whereas most of the *bona fide P. chiliana* contig windows

823 were proximal in the map (dark blue dots in Supplementary Information 1-Fig. 15), others were in

824 a different region (see top and bottom regions in the middle of the map), which also included *bona*

825 *fide P. caudatus* scaffolds (black dots in Supplementary Information 1-Fig. 15). All these other

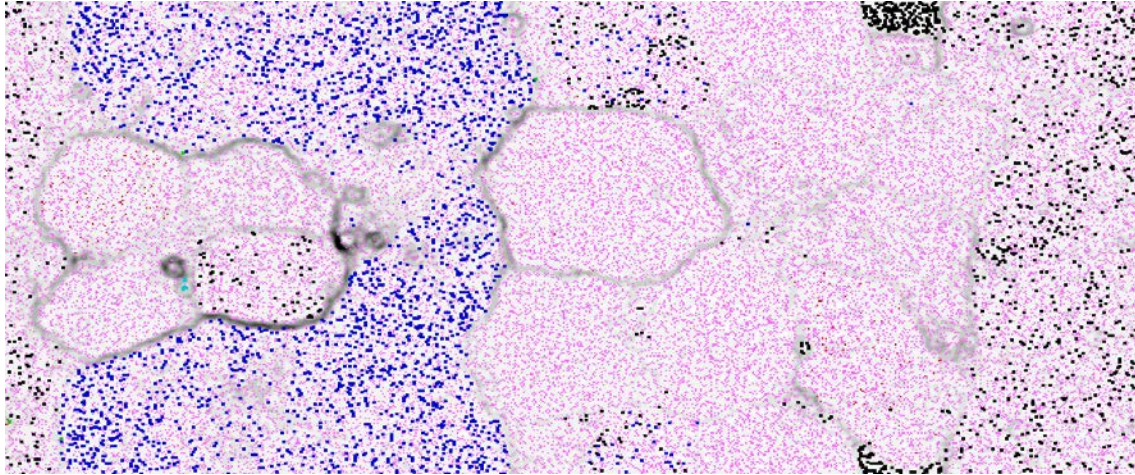
826 windows were from the same contig, suggesting that it was a *P. caudatus* contig misclassified as

827 *P. chiliana*. Consequently, only the region of the map with the majority of dark blue dots was

828 considered as non-contaminant, and hence all contigs with windows within or surrounding this

829 region were considered included in this set (colored in orange in Supplementary Information 1-
830 Fig. 16). All contigs in the map without windows within the selected region were included in the
831 'Contaminant' set.

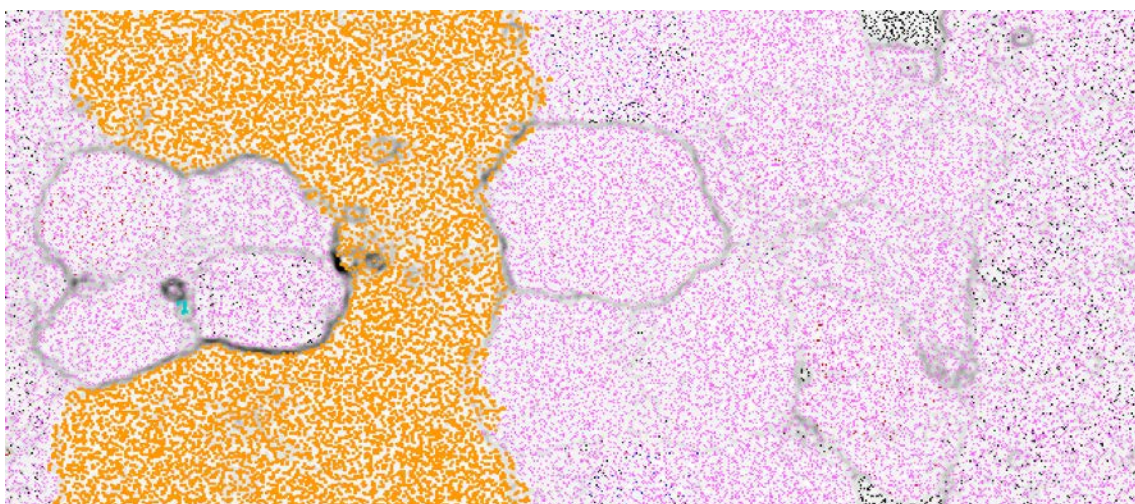
832



833

834 **Supplementary Information 1-Fig. 15.** ESOM map of the *P. chilleana* contigs from the first assembly. Each
835 dot in the map correspond to one contig/contig window, which are colored according to the category to which
836 was classified. Category color code: red for '16S', green for '18S', light blue for 'mitochondria', dark blue for
837 '*bona fide P. chilleana*', pink for 'Others'. We also incorporated the *bona fide P. caudatus* scaffolds (from the
838 *P. caudatus* metagenome, colored in black) in order to detect which regions of the *P. chilleana* assembly
839 correspond to the *P. caudatus* genome. Tetranucleotide frequency distances between neighbour
840 contig/contig windows are represented with a white/black background gradient for smaller and larger
841 distances, respectively.

842



843

844 **Supplementary Information 1-Fig. 16.** ESOM map of the *P. chilleana* contigs from the first assembly, as in
845 Supplementary Information 1-Fig. 15, but colored in orange the contigs windows within the region of the map
846 considered to correspond to *P. chilleana* genome.

847

848 For *M. vibrans* and *P. atlantis*, we directly used the ‘Contaminant’ and ‘Non-contaminant’ contig
849 sets to remove the contaminant reads and then do a second assembly. For both *Pigoraptor*
850 species, before doing this, we first removed potential contaminant reads from *H. sapiens*, as
851 contamination was suggested by the 18S and mitochondrial analyses. For that, we aligned the
852 reads used in the first assembly with the *H. sapiens* genome [bowtie2]. As we did for the
853 alignments with the *P. caudatus* metagenome (see 3.2), we did not consider the reads that also
854 aligned with *C. owczarzaki* genome, as they may correspond to conserved regions. For *P.*
855 *vietnamica*, the 1.80% of UP -unpaired- and 1.48% of PE -paired- reads aligned with *H. sapiens*
856 but not with *C. owczarzaki* (1.71% of UP and 0.06% of PE of the reads that aligned to *H. sapiens*
857 also aligned with *C. owczarzaki*). For *P. chilleana*, the 0.20% of UP and 0.07% of PE reads aligned
858 with *H. sapiens* but not with *C. owczarzaki* (3.03% of UP and 0.74% of PE of the reads that aligned
859 with *H. sapiens* also aligned with *C. owczarzaki*). These results confirmed the hypothesized *H.*
860 *sapiens* contamination, quantitatively lower in *P. chilleana* than in *P. vietnamica*.

861

862 3.4) Second round of read decontamination (*P. vietnamica*)

863 We then also excluded the *P. vietnamica* reads that aligned with the ‘Contaminant’ and ‘Non-
864 contaminant’ post-ESOM sets (see 1.4). 30.08% of the reads aligned with the ‘Contaminant’ but
865 not with the *P. vietnamica* (i.e., non-contaminant) dataset. The second assembly of *P. vietnamica*
866 was done with SPAdes, with the metagenome mode disabled, and using the following k-mers:
867 51, 71, 93, 115 [--careful --cov-cutoff auto]. The average coverage of the 115-mer assembly was
868 42.21.

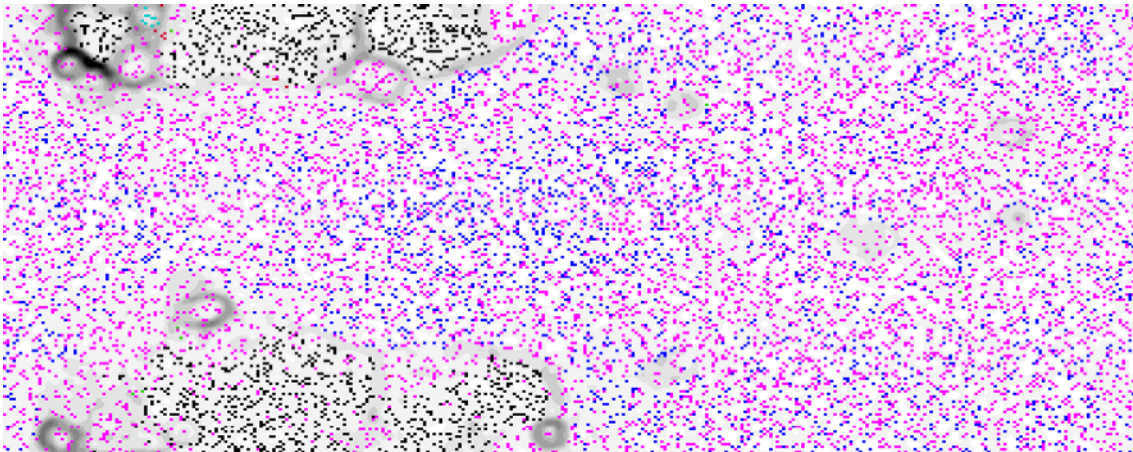
869

870 Then, as with the first assembly, we screened the assembled scaffolds for potential contamination
871 of adapter/vector sequences, and also for 18S, 16S, mitochondrial sequences and *bona fide P.*
872 *vietnamica* scaffolds (see 3.3). We removed 7 scaffolds of short length (4143 bp) because they
873 aligned with vector sequences along most of their sequence. We found 37 scaffolds potentially
874 containing 16S ribosomal genes. We did not find any scaffold with *H. sapiens* 18S or
875 mitochondrial sequences, suggesting that we successfully removed the *H. sapiens*
876 contamination. We also found a scaffold with the putative *P. caudatus* 18S gene (13903 bp, 17.53
877 cov), but it had less assembly coverage than in the first assembly (57.45 cov).

878

879 Scaffolds with > 2000 bp were submitted to ESOM analyses, splitting scaffolds larger than 7999
880 bp into windows of 4000 bp. Scaffolds were colored in ESOM maps according to the categories
881 to which they belong: ‘16S’, ‘18S’, ‘mitochondria’, ‘*bona fide P. vietnamica*’, ‘Others’. We also
882 incorporated the *bona fide P. caudatus* fragments (see 3.3). Contrary to previous results, most of
883 the map corresponds to the *bona fide P. vietnamica* contigs region (see dark blue dots in
884 Supplementary Information 1-Fig. 17), indicating that we successfully removed a substantial
885 fraction of contamination during the decontamination of the first assembly. However, the lack of
886 a clear white/black background gradient separating the *P. vietnamica* and the *P. caudatus* regions

887 (see black dots) left uncertainty surrounding the scaffolds found in-between both regions. From
888 ESOM results, we preliminary classified scaffolds into 'Non-contaminant', 'Contaminant' and
889 'Uncertain' (see orange dots, brown dots and other color dots, respectively, in Supplementary
890 Information 1-Fig. 18). For those scaffolds with windows in more than one category, we used the
891 following criteria: (i) Scaffolds with at least one 'Uncertain' window were classified as 'Uncertain'.
892 (ii) Scaffolds with windows in both 'Non-contaminant' and 'Contaminant' were classified as
893 'Uncertain'. Finally, scaffolds not included in ESOM analyses (i.e., < 2000 bp) were classified as
894 '<2000' (a fourth category).
895



896
897 **Supplementary Information 1-Fig. 17.** ESOM map of the *P. vietnamica* scaffolds from the second
898 assembly. Each dot in the map correspond to one scaffold/scaffold window, which are colored according to
899 the category to which was classified. Category color code: red for '16S', green for '18S', light blue for
900 'mitochondria', dark blue for '*bona fide P. vietnamica*', pink for 'Others'. We also incorporated the *bona fide*
901 *P. caudatus* fragments (from the *P. caudatus* metagenome, colored in black) in order to detect which regions
902 of the *P. vietnamica* assembly correspond to the *P. caudatus* genome. Tetranucleotide frequency distances
903 between neighbour scaffold/scaffold windows are represented with a white/black background gradient for
904 smaller and larger distances, respectively.
905



906
907 **Supplementary Information 1-Fig. 18.** ESOM map of the *P. vietnamica* scaffolds from the second
908 assembly, as in Supplementary Information 1-Fig. 17, but colored in orange the scaffold windows within the

909 region of the map likely corresponding to *P. vietnamica* genome, in brown the region likely corresponding to
910 contaminant fragments. Other scaffolds were considered as 'Uncertain'.

911

912 With *M. vibrans* and *P. atlantis*, we determined which scaffolds corresponded to our genome of
913 interest directly from ESOM results. In this case, given the lack of resolution in-between some *P.*
914 *vietnamica* and *P. caudatus* regions, the preliminary classification from ESOM was redefined by
915 taking into account the following scaffold features: (1) assembly coverage, (2) taxonomic profiling,
916 (3) average number of introns per eukaryotic gene and per bacterial gene, (4) eukaryotic
917 taxonomic profiling, (5) total query coverage in alignments with *P. caudatus* assembly and (6)
918 total query coverage in alignments with NCBI nt.

919

920 (1) The assembly coverage was obtained from the scaffold names (computed by *SPAdes*).

921 (2) The taxonomic profiling was obtained using the indirect taxonomic classification approach
922 (explained in *M. vibrans* section). First scaffolds were preliminary annotated with *BRAKER1*, and
923 then we counted the number of eukaryotic genes ('E'+ 'PE') and bacterial genes ('B'+ 'PB') per
924 scaffold. The RNA-seq reads required for the *BRAKER1* annotation pipeline were kindly provided
925 by Elisabeth Hehenberger²⁸.

926 (3) The average number of introns per eukaryotic gene and per bacterial gene were computed
927 from the *BRAKER1* GFF3 output file.

928 (4) The predicted eukaryotic protein sequences were aligned with a comprehensive local
929 eukaryotic database [*BLASTp*: -evalue 1e-3], and the corresponding genes were later classified
930 into five taxonomic categories: (i) 'F' if the best hit was from Filasterea, (ii) 'E' if the best hit was
931 from Excavata, (iii) 'M' if the best hit was from Metazoa, (iv) 'O' if the best hit was from
932 Opisthokonta but not from Filasterea and Metazoa and (v) 'D' if the best hit was from any other
933 taxonomic group.

934 (5) Scaffolds were aligned with the *P. caudatus* metagenome [*BLASTn*: -evalue 1e-3], and then
935 we computed the total query coverage of every scaffold by dividing the number of positions that
936 aligned with a database sequence by the length of that scaffold.

937 (6) The same approach as (5), but using the NCBI nt database.

938

939 We expected from *P. vietnamica* scaffolds to show (1) higher assembly coverage values than *P.*
940 *caudatus* or other contaminant scaffolds. The median assembly coverage values for 'Non-
941 contaminant' ESOM contigs were 49.42, clearly higher than those for the 'Contaminant' ESOM
942 contigs (3.02); (2) positive E+PE/B+PB ratios; (3) similar number of introns in eukaryotic and in
943 bacterial genes (if any); (4) more F, M, O genes than E or D; (5) spurious query coverage values
944 in the alignment with *P. caudatus* assembly and (6) with the NCBI nt database.

945

946 Compared to *P. vietnamica* scaffolds, we expected from *P. caudatus* scaffolds (1) low assembly
947 coverage values; (4) more E or D genes than F, M, O and (5) high query coverage values in the
948 alignments with *P. caudatus* assembly.

949

950 From bacterial contaminant scaffolds, we expected (1) low assembly coverage values; (2)
951 negative E+PE/B+PB ratios; (3) none or very few introns in genes (false positives); (6) possibly
952 high query coverage values in the alignment with NCBI nt (if the corresponding species is
953 represented in the database).

954

955 From potential *H. sapiens* contaminant scaffolds (most likely excluded in previous
956 decontaminations), we expected (1) low assembly coverage values; (4) more M than other genes,
957 and few or none F genes; (6) possibly high query coverage values in the alignment with NCBI nt
958 database (if the corresponding sequence is represented in the database).

959

960 Based on the above-mentioned expectations, scaffolds were finally classified into '*P. vietnamica*',
961 'Potential contaminant' and 'Contaminant' using a manual decision-tree approach (with the
962 exception of the putative mitochondrial scaffolds, which were directly classified into *P.*
963 *vietnamica*). We explored different combinations of nested conditional *If/Else* statements, which
964 were subsequently improved based on the observed classification outcomes (i.e., after every
965 decision-tree, we evaluated if the parameters of the scaffolds classified within each category
966 disagreed with any of the above-mentioned expectations, and we modified the algorithm
967 accordingly to overcome the observed disagreements). See below the custom-made final
968 decision-tree for *P. vietnamica*.

969

970 **Abbreviations:**

971 cov: assembly coverage value

972 qcPCAU: total query coverage of a scaffold in the alignments with *P. caudatus* scaffolds.

973 qcNCBI: total query coverage of a scaffold in the alignments with NCBI nt.

974 pY: scaffolds preliminary classified as 'Non-contaminant' in ESOM analyses

975 pD: scaffolds preliminary classified as 'Uncertain' in ESOM analyses

976 pN: scaffolds preliminary classified as 'Contaminant' in ESOM analyses

977 less2000: scaffold preliminary classified as '<2000', not included in ESOM analyses

978 Y: scaffolds classified as '*P. vietnamica*'

979 D: scaffolds classified as 'Potential contaminant'

980 N: scaffolds classified as 'Contaminant'

981 Bact: B+PB

982 Euk: E+PE

983 B_introns: average number of introns in Bacterial genes

984 E_introns: average number of introns in Eukaryotic genes

985 Euk_E: eukaryotic genes taxonomically classified as F

986 Euk_M: eukaryotic genes taxonomically classified as M
 987 Euk_O: eukaryotic genes taxonomically classified as O
 988 Euk_E: eukaryotic genes taxonomically classified as E
 989 Euk_D: eukaryotic genes taxonomically classified as D
 990
 991 **Decision-tree for *P. vietnamica*:**
 992 If qcPCAU or If qcNCBI >= 50%: --> N
 993 Else:
 994 If pY: --> by default Y, but:
 995 If cov <= 25:
 996 If Bact > Euk:
 997 If B_introns = 0: --> D
 998 If Euk = 0 --> D
 999 If Euk != 0:
 1000 If Euk_F = 0: --> D
 1001 If Euk_F > 0:
 1002 If Euk_F < Euk_M: --> D
 1003 If Euk_E > 0: --> D
 1004 If (Euk_E+Euk_D) > (Euk_O+Euk_M+Euk_F): --> D
 1005 If cov > 25:
 1006 If Bact > Euk:
 1007 If B_introns = 0: --> D
 1008 If E > 0:
 1009 If Euk_E >= Euk_F:
 1010 If qcPCAU > 5%: --> D
 1011 Else: --> Y
 1012 If pN:
 1013 If Euk > Bact:
 1014 If (Euk_O+Euk_M+Euk_F+Euk_D) > Euk_E:
 1015 If Bact > 0:
 1016 If B_introns > 0: --> D
 1017 If Bact = 0: --> D
 1018 Else: --> N
 1019
 1020 If pD:
 1021 If cov <= 25: --> by default N, but:
 1022 If (Euk_O+Euk_F+Euk_M) >= (Euk_E+Euk_D) or If Euk_F > 0: --> D
 1023 If Euk < Bact or If Euk = 0: --> N
 1024 If cov > 25: --> by default Y, but:
 1025 If Euk < Bact:
 1026 If Euk_F > 0:
 1027 If B_introns > 0: --> Y
 1028 If B_introns < 0: --> D

```

1029             or If B_introns > 0: --> D
1030             Else: --> N
1031         If Euk >= Bact:
1032             If Euk_E > (Euk_O+Euk_F+Euk_M): --> D
1033             Else: --> Y
1034
1035     If less2000: --> by default N, but:
1036         If cov <= 25:
1037             If Euk_F > (Euk_O+Euk_M+Euk_E+Euk_D): --> D
1038         If cov > 25:
1039             If Euk < Bact:
1040                 If (Euk_O+Euk_F+Euk_M) >= (Euk_E+Euk_D):
1041                     If B_introns > 0: --> Y
1042                     Else: --> D
1043             If Bact = 0:
1044                 If (Euk_O+Euk_F+Euk_M) >= (Euk_E+Euk_D): --> Y
1045                 Else: --> D
1046             If Bact > 0 and Bact < Euk:
1047                 If (Euk_O+Euk_F+Euk_M) >= (Euk_E+Euk_D):
1048                     If B_introns > 0: --> Y
1049                     Else: --> D
1050

```

1051 We observed that scaffolds with > 25 assembly coverage tend to be larger, have many eukaryotic
1052 and few bacterial genes, and introns in bacterial genes. Not all genes from the scaffolds classified
1053 as '*P. vietnamica*' performed best hit with *C. owczarzaki* or *M. vibrans* proteins, which could be
1054 expected given the sequence divergence shown from other filastreans respect to the Pigoraptor
1055 clade²⁸. As expected, many scaffolds classified as '*P. vietnamica*' and 'Contaminant' were already
1056 preliminary classified into 'Non-contaminant' and 'Contaminant' after ESOM analyses,
1057 respectively. For the scaffolds preliminary classified as 'Uncertain' and '<2000', the most
1058 determinant parameters were eukaryotic taxonomic profiling and assembly coverage.

1059
1060 Scaffolds classified as '*P. vietnamica*' (3390 scaffolds, 42153970 bp) and as 'Uncertain' (601
1061 scaffolds, 879351) by the decision tree were included in the *Pvie.gDNA.clean.v1.fasta* file, with
1062 the name of the 'Uncertain' scaffolds being labeled with the suffix '_potentialcontaminant'. The
1063 scaffolds with 18S and mitochondrial sequences (1 mitochondrial scaffold, 39208 bp) were
1064 labeled as 'ribosomal' and '_mitochondrial', respectively.

1066 3.5) Second round of read decontamination (*P. chilleana*)

1067 After removing *H. sapiens* contamination, we excluded the reads that aligned with the
1068 'Contaminant' and 'Non-contaminant' post-ESOM sets (see 3.3.3). 66.63% of the reads aligned
1069 with the 'Contaminant' but not with the *P. chilleana* ('Non-contaminant') dataset, more than twice

1070 than in *P. vietnamica* (30.08%). Because the remaining contaminant reads did not allow an
1071 assembly with enough coverage, we sequenced an extra Illumina HiSeq 2500 lane using the
1072 sequencing kit HiSeq v4 chemistry (insert size: 410 bp, read length: 125 bp). Reads from this
1073 extra library were preprocessed, corrected and also decontaminated as with the first library (i.e.,
1074 we removed those reads aligning either with *P. caudatus* scaffolds or with *H. sapiens* but not with
1075 *C. owczarzaki*; as well as those reads aligning to the 'Contaminant' but not to the 'Non-
1076 contaminant' sets).

1077

1078 The second assembly of *P. chilleana* was performed with both libraries, first using *SPAdes* without
1079 the metagenome mode. However, the assembly ended with the warning "Failed to determine
1080 erroneous kmer threshold", which suggested uneven coverage problems occurred because of
1081 substantial amounts of persistent contamination. We found that the coverage problem lied on the
1082 reads corresponding to the first assembly contigs with <2000 bp length, which were not included
1083 in the ESOM decontamination analyses. To avoid potential pitfalls in the assembly related to
1084 uneven coverage, the second assembly was finally performed using the metagenome mode, with
1085 the following k-mers: 21, 33, 45, 57. The average coverage of the 57-mer assembly was 38.40.

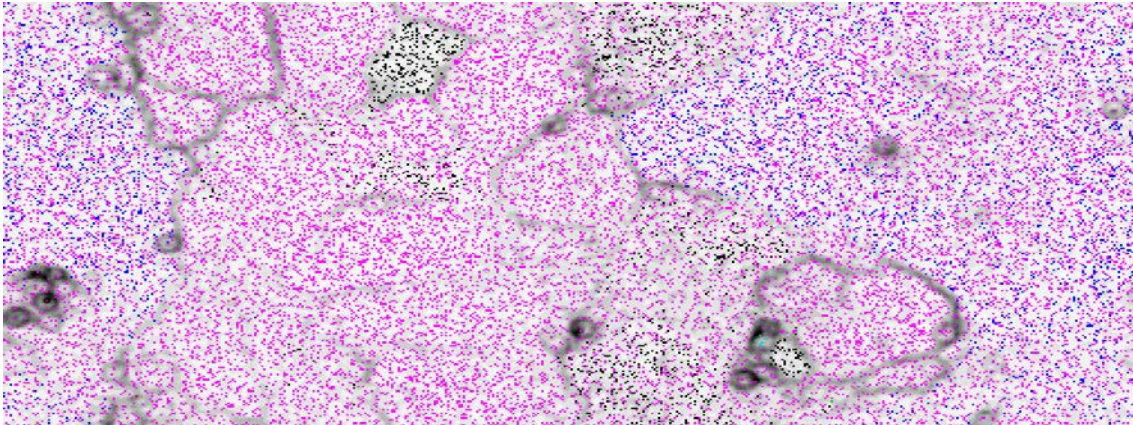
1086

1087 Then, as with the first assembly, we checked for potential remaining adapter/vector sequences,
1088 and also for 18S, 16S, mitochondrial sequences and *bona fide P. chilleana* scaffolds. We removed
1089 other 7 scaffolds of short length (21112 bp) because they aligned against vector sequences along
1090 most of their sequence. We found 64 scaffolds potentially containing 16S ribosomal genes. We
1091 did not find any scaffold with *H. sapiens* 18S or mitochondrial sequences, suggesting that we
1092 successfully removed the *H. sapiens* contamination. Despite we found putative *P. caudatus* 18S
1093 and mitochondrial scaffolds, they had very short lengths (<637 bp) and lower assembly coverage
1094 values than the 18S and mitochondrial *P. chilleana* scaffolds.

1095

1096 Scaffolds with >2000 bp were submitted to ESOM analyses, splitting scaffolds larger than 7999
1097 bp into windows of 4000 bp. Scaffolds were colored in ESOM maps according to the categories
1098 to which they belong: '16S', '18S', 'mitochondria', '*bona fide P. chilleana*', 'Others'. We also
1099 incorporated the *bona fide P. caudatus* scaffolds (from the *P. caudatus* metagenome). The *bona*
1100 *fide P. chilleana* contigs covered an extensive region, approximately half of the map (see dark
1101 blue dots in Supplementary Information 1-Fig. 19). However, as with the ESOM from *P.*
1102 *vietnamica* second assembly, a clear white/black background gradient separating the *P.*
1103 *vietnamica* and the *P. caudatus* regions (see black dots) was missing. Moreover, the bottom left
1104 subregion of the putative *P. vietnamica* region had less blue dots compared to the other parts of
1105 the region, which suggested the presence of contaminant scaffolds within. Because of this, we
1106 did as for *P. chilleana*, and we used again ESOM results to preliminary classify scaffolds into 'Non-
1107 contaminant', 'Uncertain' and 'Contaminant' (see orange dots, brown dots and other color dots,
1108 respectively, in Supplementary Information 1-Fig. 20).

1109

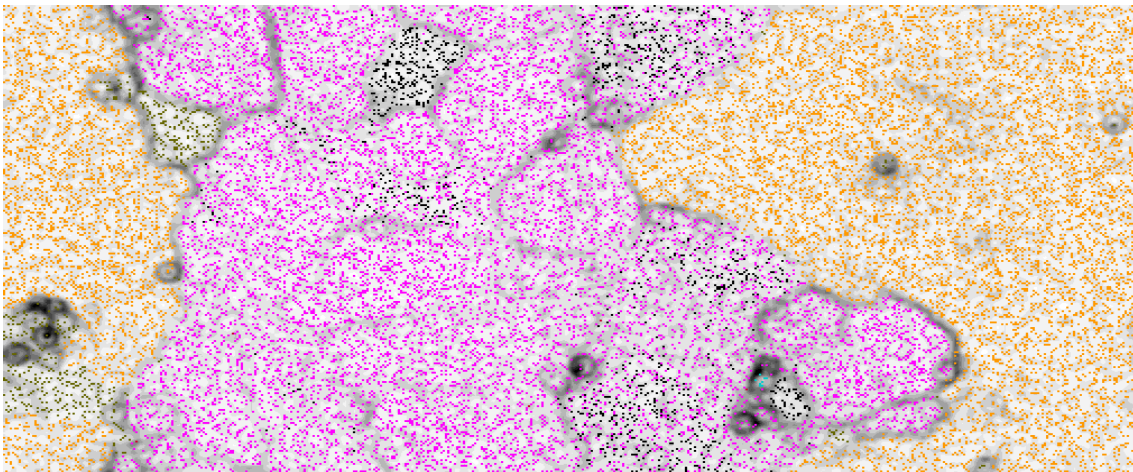


1110

1111

Supplementary Information 1-Fig. 19. ESOM map of the *P. chilleana* scaffolds from the second assembly. Each dot in the map correspond to one scaffold/scaffold window, which are colored according to the category to which was classified. Category color code: red for '16S', green for '18S', light blue for 'mitochondria', dark blue for '*bona fide P. chilleana*', pink for 'Others'. We also incorporated the *bona fide P. caudatus scaffolds* (from the *P. caudatus* metagenome, colored in black) in order to detect which regions of the *P. chilleana* assembly correspond to the *P. caudatus* genome. Tetranucleotide frequency distances between neighbour scaffold/scaffold windows are represented with a white/black background gradient for smaller and larger distances, respectively.

1118
1119



1120

1121

Supplementary Information 1-Fig. 20. ESOM map of the *P. chilleana* scaffolds from the second assembly, as in Supplementary Information 1-Fig. 19, but colored in orange the scaffold windows within the region of the map likely corresponding to *P. chilleana* genome ('Non-contaminant' set), and in brown the region likely corresponding to uncertain fragments.

1124
1125

1126

1127

1128

1129

As with *P. vietnamica*, 'Non-contaminant' scaffolds showed clearly higher assembly coverage values than 'Contaminant' scaffolds (median values were 57.02 and 4.79, respectively). For a final classification of scaffolds into '*P. chilleana*', 'Potential contaminant' and 'Contaminant' we used a similar custom-made decision-tree approach as we did for *P. vietnamica* (see below).

1130

1131 **Abbreviations:**

1132 cov: assembly coverage value

1133 qcPCAU: total query coverage of a scaffold in the alignments with *P. caudatus* scaffolds.

1134 qcNCBI: total query coverage of a scaffold in the alignments with NCBI nt.

1135 pY: scaffolds preliminary classified as 'Non-contaminant' in ESOM analyses

1136 pD: scaffolds preliminary classified as 'Uncertain' in ESOM analyses

1137 pN: scaffolds preliminary classified as 'Contaminant' in ESOM analyses

1138 less2000: scaffold preliminary classified as '<2000', not included in ESOM analyses

1139 Y: scaffolds classified as '*P. chilleana*'

1140 D: scaffolds classified as 'Potential contaminant'

1141 N: scaffolds classified as 'Contaminant'

1142 Bact: B+PB; Euk: E+PE

1143 B_introns: average number of introns in Bacterial genes

1144 E_introns: average number of introns in Eukaryotic genes

1145 Euk_F: eukaryotic genes taxonomically classified as F

1146 Euk_M: eukaryotic genes taxonomically classified as M

1147 Euk_O: eukaryotic genes taxonomically classified as O

1148 Euk_E: eukaryotic genes taxonomically classified as E

1149 Euk_D: eukaryotic genes taxonomically classified as D

1150

1151 **Decision-tree for *P. chilleana*, second round of decontamination:**

1152 If qcPCAU or If qcNCBI \geq 50%: --> N

1153 Else:

1154 If pY: --> by default Y, but:

1155 If cov \leq 25:

1156 If Bact > Euk:

1157 If B_introns = 0: --> D

1158 If B_introns > 0: --> Y

1159 If Euk = 0: --> D

1160 If Euk > 0:

1161 If Euk_F = 0: --> D

1162 If Euk_F < Euk_M: --> D

1163 If Euk_E > 0:

1164 If (Euk_E+Euk_D) \geq (Euk_F+Euk_M+Euk_O): --> D

1165 If cov > 25:

1166 If Bact > Euk:

1167 If B_introns = 0: --> D

1168 If Euk_E > 0:

1169 If Euk_E \geq Euk_F:

1170 If qcPCAU \geq 5%: --> D

1171

1172 If pN: --> by default N, but:

```

1173         If Euk > Bact:
1174             If (Euk_F+Euk_M+Euk_O+Euk_D) > Euk_E:
1175                 If B_introns > 0: --> D
1176
1177     If pD: --> by default N, but:
1178         If cov <= 25:
1179             If Euk > 0:
1180                 If (Euk_M+Euk_F+Euk_O) >= (Euk_D+Euk_E): --> D
1181                 If Euk_F > 1: --> D
1182             If Euk < Bact: --> N
1183             If Euk = 0: --> N
1184         If cov > 25:
1185             If Euk < Bact:
1186                 If Euk_F > 0:
1187                     If B_introns > 0: --> Y
1188                     Else: --> D
1189                 Else: --> N
1190             If Euk >= Bact:
1191                 If Euk > 0:
1192                     If Euk_E > (Euk_F+Euk_M+Euk_O): --> D
1193                     Else: --> Y
1194
1195     If less2000:
1196         If cov <= 25:
1197             If Euk_F > (Euk_M+Euk_O+Euk_D+Euk_E): --> D
1198         If cov > 25:
1199             If Bact = 0:
1200                 If (Euk_M+Euk_O+Euk_F) > (Euk_D+Euk_E): --> Y
1201                 Else: --> D
1202             If Euk < Bact:
1203                 If (Euk_M+Euk_O+Euk_F) > (Euk_D+Euk_E):
1204                     If B_introns > 0: --> Y
1205                     Else: --> D
1206             If Bact > 1 and Euk >= Bact:
1207                 If (Euk_M+Euk_O+Euk_F) > (Euk_D+Euk_E): --> D
1208                 If B_introns > 0: --> Y
1209                 Else: --> D
1210
1211

```

1212 With *P. vietnamica*, the post-decision tree classification of the second assembly was used to
1213 determine which scaffolds were included in the *Pvie.gDNA.clean.v1.fasta* file. However, with *P.*
1214 *chilleana*, the uneven coverage problems limited the quality of the second assembly. We hence

1215 used the classification results from this second assembly to perform a third read decontamination
1216 step to allow a third assembly without uneven coverage problems.

1217

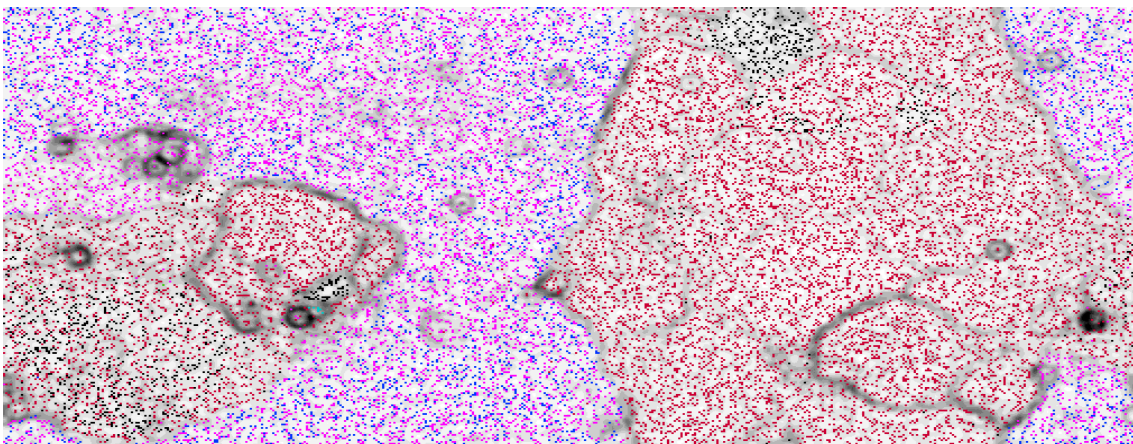
1218 3.6) Third round of read decontamination (*P. chiliana*)

1219 We discarded the reads that aligned with the 'Contaminant' but not with the 'Uncertain' or the
1220 '*bona fide P. chiliana*' scaffolds, this corresponding to the 56.69% of the reads [bowtie2]. For PE
1221 reads, we only removed them if both paired reads satisfied this criteria. We assembled the de-
1222 contaminated reads of *P. chiliana* using SPAdes without the metagenome mode [-k 35,49,63,75
1223 --careful --cov-cutoff auto]. In contrast with the second assembly, this third assembly ended
1224 without warning uneven coverage problems, suggesting that our decision tree classification
1225 approach removed substantial contamination issues. The average coverage of the 75-mer
1226 assembly was 40.55. We found 21 scaffolds with potential 16S sequences, and 3 scaffolds
1227 corresponding to vector/adaptor sequences were removed (957 bp). We did not find any potential
1228 18S or mitochondrial *H. sapiens* sequence. We also identified a set of *bona fide P. chiliana*
1229 scaffolds (see 3.3).

1230

1231 Scaffolds with >2000 bp were submitted to ESOM analyses, splitting scaffolds larger than 7999
1232 bp into windows of 4000 bp. Scaffolds were colored in ESOM maps according to the categories
1233 to which they belong: '16S', '18S', 'mitochondria', '*bona fide P. chiliana*', 'Others', '*bona fide P.*
1234 *caudatus* scaffolds' (from the *P. caudatus* metagenome). We also incorporated the 'Contaminant'
1235 scaffolds from the second assembly. As occurred in the ESOM analyses of the second
1236 decontamination rounds of both *Pigoraptor* species, there were some regions without a clear
1237 white/black background gradient separating the *P. vietnamica* (see blue dots in Supplementary
1238 Information 1-Fig. 21) and the *P. caudatus* and contaminant regions regions (see black and red
1239 dots).

1240



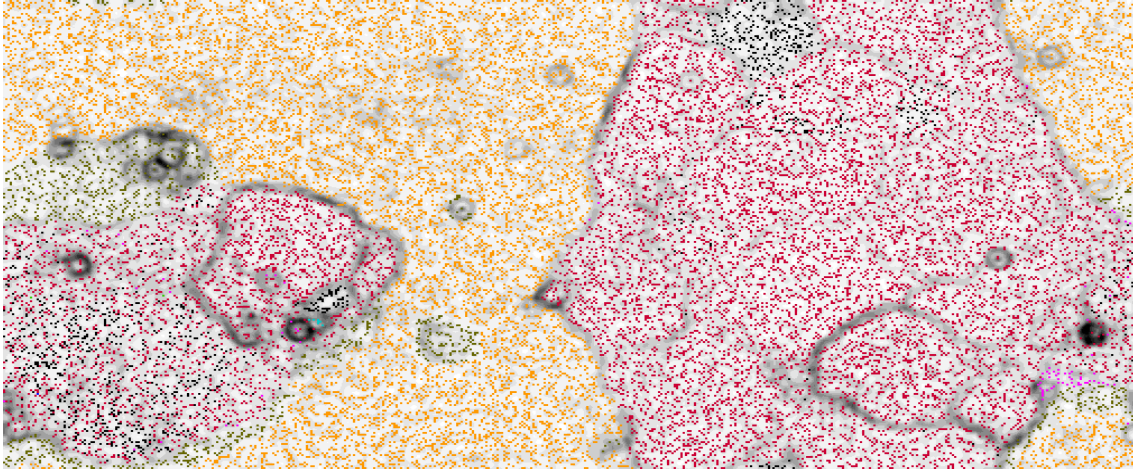
1241

1242 **Supplementary Information 1-Fig. 21.** ESOM map of the *P. chiliana* scaffolds from the second assembly.
1243 Each dot in the map correspond to one scaffold/scaffold window, which are colored according to the category
1244 to which was classified. Category color code: red for '16S', green for '18S', light blue for 'mitochondria', dark

1245 blue for '*bona fide P. chilleana*', pink for 'Others'. We also incorporated the 'Contaminant' scaffolds from the
1246 second assembly and the *bona fide P. caudatus scaffolds* from the *P. caudatus* metagenome (colored in red
1247 in black, respectively). Tetranucleotide frequency distances between neighbour scaffold/scaffold windows
1248 are represented with a white/black background gradient for smaller and larger distances, respectively.

1249

1250 We hence used again ESOM results to preliminary classify scaffolds into 'Non-contaminant',
1251 'Contaminant' and 'Uncertain' (see orange dots, brown dots and other color dots, respectively, in
1252 Supplementary Information 1-Fig. 22).



1253

1254 **Supplementary Information 1-Fig. 22.** ESOM map of the *P. chilleana* scaffolds from the third assembly, as
1255 in Supplementary Information 1-Fig. 21, but colored in orange the scaffold windows within the region of the
1256 map likely corresponding to *P. chilleana* genome, and in brown the region likely corresponding to uncertain
1257 fragments.

1258

1259 For a final classification of scaffolds into '*P. chilleana*', 'Potential contaminant' and 'Contaminant',
1260 we used a custom-made decision-tree approach as in the second assembly.

1261

1262 **Abbreviations:**

1263 cov: assembly coverage value

1264 qcPCAU: total query coverage of a scaffold in the alignments with *P. caudatus* scaffolds.

1265 qcNCBI: total query coverage of a scaffold in the alignments with NCBI nt.

1266 pY: scaffolds preliminary classified as 'Non-contaminant' in ESOM analyses

1267 pD: scaffolds preliminary classified as 'Uncertain' in ESOM analyses

1268 pN: scaffolds preliminary classified as 'Contaminant' in ESOM analyses

1269 less2000: scaffold preliminary classified as '<2000', not included in ESOM analyses

1270 Y: scaffolds classified as '*P. chilleana*'

1271 D: scaffolds classified as 'Potential contaminant'

1272 N: scaffolds classified as 'Contaminant'

1273 Bact: B+PB

1274 Euk: E+PE

1275 B_introns: average number of introns in Bacterial genes

1276 E_introns: average number of introns in Eukaryotic genes
1277 Euk_F: eukaryotic genes taxonomically classified as F
1278 Euk_M: eukaryotic genes taxonomically classified as M
1279 Euk_O: eukaryotic genes taxonomically classified as O
1280 Euk_E: eukaryotic genes taxonomically classified as E
1281 Euk_D: eukaryotic genes taxonomically classified as D
1282
1283 **Decision-tree for *P. chilleana*, third round of decontamination:**
1284 If qcPCAU or If qcNCBI >= 50%: --> N
1285
1286 Else:
1287 If pY: --> by default Y, but:
1288 If cov <= 25:
1289 If (Bact + Euk) > 1:
1290 If Bact > Euk:
1291 If Euk = 0: --> N
1292 If Euk > 1:
1293 If B_introns > 0: --> Y
1294 Else: --> D
1295 or If:
1296 Euk_E > (Euk_F+Euk_M+Euk_O+Euk_D):
1297 Else: --> D
1298 If cov > 25:
1299 If (Bact + Euk) > 1:
1300 If Euk < Bact:
1301 If Euk = 0:
1302 If B_introns > 0:
1303 If Euk_F < Euk_E: --> N
1304 Else: --> D
1305 Else: --> N
1306 Else:
1307 If B_introns > 0: --> Y
1308 Else: --> D
1309 If Euk_F < Euk_E: --> D
1310
1311 If pN: --> by default N, but:
1312 If cov > 25:
1313 If Euk > Bact:
1314 If E_introns > 0:
1315 If Euk_E < (Euk_F+Euk_M+Euk_O+Euk_D):
1316 If Bact > 0:
1317 If B_introns > 0: --> D
1318 Else: --> D

```

1319
1320 If pD: --> by default N, but:
1321     If cov <= 25:
1322         If Euk > 0:
1323             If (Euk_M+Euk_F+Euk_O+Euk_D) >= Euk_E: --> D
1324             If Euk_F > 0: --> D
1325         If Euk + Bact > 0:
1326             If Euk < Bact:
1327                 If B_introns > 0: --> D
1328                 Else: --> N
1329             If Euk = 0: --> N
1330     If cov > 25: --> by default D, but:
1331         If Euk + Bact > 0:
1332             If Euk < Bact: --> by default N, but:
1333                 If Euk_F > 0: --> D
1334                 If B_introns > 0:
1335                     If Euk > 0: --> D
1336                     Else:
1337                         If B_introns > 2: --> D
1338             or If Euk > 0:
1339                 If (Euk_M+Euk_F+Euk_O) > (Euk_D+Euk_E): --> D
1340                 Else: --> Y
1341
1342 If less2000: --> by default N, but:
1343     If cov <= 25:
1344         If Euk > 0:
1345             If Euk_F > (Euk_M+Euk_O+Euk_D+Euk_E): --> D
1346     If cov > 25:
1347         If Euk > 0:
1348             If Euk < Bact:
1349                 If (Euk_F+Euk_M+Euk_O) > (Euk_D+Euk_E):
1350                     If B_introns > 0: --> Y
1351                     Else: --> D
1352             or If Bact = 0: --> by default D, but:
1353                 If (Euk_M+Euk_F+Euk_O) > (Euk_D+Euk_E): --> Y
1354                 Else: --> D
1355             Else:
1356                 If (Euk_F+Euk_M+Euk_O) > (Euk_D+Euk_E):
1357                     If B_introns > 0: --> Y
1358                     Else --> D
1359
1360 Scaffolds classified as 'P. chilleana' (4119 scaffolds, 43905177 bp) and as 'Uncertain' (774
1361 scaffolds, 3157320 bp) by the decision-tree were included in the Pchi.gDNA.clean.v1.fasta, with

```

1362 the suffix ‘_potentialcontaminant’ added to ‘Uncertain’ scaffolds. 18S and mitochondrial scaffolds
1363 (2 putative mitochondrial scaffolds, 38854 bp) were labeled as ‘ribosomal’ and ‘_mitochondrial’,
1364 respectively.

1365

1366 **4) Genome annotation and masking of repetitive regions**

1367

1368 The genomes were annotated using *BRAKER1*, as we did for all preliminary annotations
1369 performed during the decontamination steps (explained in *M. vibrans* section). However, in this
1370 case, we first estimated the maximum intron length of every species for the --max-segment-intron
1371 and --max-intron-length *TopHat* parameters. For that, we aligned the RNA-seq *de novo*
1372 assembled transcripts of each species with its genome [*BLASTn*], and inferred potential intron
1373 positions and lengths from the discontinuities in the alignments between every query transcript
1374 and its best targeting genomic scaffold (only alignments of transcripts showing > 95% of average
1375 query coverage and identity with its best targeting scaffold were considered). The values for both
1376 *TopHat* parameters were set to 17500, 7500, 2500 and 6000 for *M. vibrans*, *P. atlantis*, *P.*
1377 *vietnamica*, and *P. chilleana*, respectively.

1378

1379 We used *PASA*²⁹ v2.0.2 to refine the *BRAKER1* annotations. *PASA* was ran using as input the
1380 transcripts from both *de novo* and genome-guided assemblies, but only those that aligned with its
1381 best targeting genomic scaffold with >90% of average query coverage and identity [*BLASTn*]. *De*
1382 *nov*o transcriptome assemblies were done with *Trinity*, as explained in *M. vibrans* section.
1383 Genome-guided transcriptome assemblies were also done with *Trinity* [--normalize_reads, --
1384 jaccard_clip], using the same *accepted_hits.bam* file as the used for *BRAKER1* annotations, and
1385 also limiting the maximum intron lengths as with *TopHat* alignments. *PASA* annotations were
1386 obtained using both *blat*³⁰ v35x1 and *GMAP*³¹ v2015-12-31 aligners, with the suggested --
1387 stringent_alignment_overlap and --gene_overlap parameters, and after having been
1388 preprocessed the transcripts from adapter/vector sequences using *seqclean*
1389 (<https://sourceforge.net/projects/seqclean/>). *PASA* was run again to add UTR information and
1390 correct some *BRAKER1* annotations by means of two rounds of annotation comparison (as
1391 recommended by the manual). Protein sequences were finally retrieved from corrected *BRAKER1*
1392 annotations, keeping only the longest isoform per gene. Protein sequences corresponding to
1393 genes predicted on scaffolds labeled as “_potentialcontaminant” were equally labeled.

1394

1395 We masked the repetitive regions of the genomes with Ns using *RepeatMasker*³² version open-
1396 4.0.6 (e.g. *Mvib.gDNA.v1.fasta.masked*). For that, we first created for every genome a specific
1397 library of repetitive sequences using *RepeatModeler* v1.0.4 (www.repeatmasker.org).

1398

1399 **5) Correction of false gene fusion/fission events**

1400

1401 Gene fusion and domain rearrangement are important sources of protein innovation in eukaryotic
1402 evolution³³. These events lead to the emergence of composite gene families, which emerge from
1403 the merging of distinct component families³⁴. In a composite gene, the regions corresponding to
1404 the distinct components show similarities at sequence level to distinct sets of proteins (component
1405 families), a pattern that can be detected using sequence-similarity tools (e.g., *BLAST*)³⁵. While
1406 the finding of composite genes may well correspond to true fusion events, they can also occur by
1407 errors during the annotation process (e.g., the software may confuse an intergenic region with an
1408 intron). A preliminary analysis of the *M. vibrans* annotations revealed a *bona fide* example of false
1409 composite (Mvib_g619), a protein of 891 aa with an unexpected Pfam domain architecture (i.e.,
1410 not found in other eukaryotes) that was not supported by the transcriptomic data.

1411

1412 We evaluated the assembled genomes to detect and split those genes showing strong evidences
1413 of being false composites. To do that, we developed a Python script to detect and split the
1414 components (i.e., the true separate genes) within the false composites by analysing results from
1415 *BLAST* alignments of the predicted proteins with two databases: (1) a nucleotide database with
1416 the RNA-seq *de novo* assembled transcripts of the corresponding species; and (2) a
1417 taxonomically-rich database of eukaryotic proteins (euk_db). *BLAST* alignments were performed
1418 separately with each database, using *tBLASTn* and *BLASTP* for (1) and (2), respectively [-evalue
1419 1e-3]. In an alignment of a novel composite (either true or false), we expect to observe distinct
1420 regions of the composite sequence to differentially align with distinct protein sets from euk_db,
1421 each set corresponding to a component family³⁵. However, to be considered as a *bona fide*
1422 composite, the existence of the entire coding sequence must be supported by the transcriptomic
1423 data. In other words, a composite is likely to be false if the discontinuities observed in the
1424 alignments with euk_db proteins are also observed in the alignments with the transcripts.

1425

1426 The criteria to detect false composites and split them into different components (true genes)
1427 consisted in finding those positions (split positions) where all the targets (RNA-seq transcripts and
1428 eukaryotic proteins) that aligned with any upstream position did not align neither with the current
1429 position nor with any downstream position. However, we adapted this criteria to overcome the
1430 limitations of inferring the exact positions that delimit the homologous regions between proteins
1431 from *BLAST* alignments³⁶. In particular, we found cases where a manual inspection of the
1432 alignments strongly suggested the existence of different components, with different regions of the
1433 false composite aligning with differentiated sets of targets. However, a little overlap occurred
1434 between the ending and starting positions of the alignments corresponding to the different
1435 components, most likely because of imprecisions in the alignments. Hence, an algorithm
1436 screening for absolute discontinuities between each pair of consecutive positions will fail to detect
1437 a split position when these overlaps occur. Instead, to skip these misleading contiguity regions,
1438 we evaluated discontinuities between all pairs of positions separated by distances of 15 aa. Thus,
1439 a split position is inferred when the targets that aligned from the 0 to n positions do not align with

1440 any position from n+15 to the end of the sequence; with two exceptions: (1) we discarded split
1441 positions that would lead to the inference of a component shorter than 30 aa; and (2) when the
1442 region delimited by the fifteen upstream and downstream positions to the split position is within
1443 the same predicted Pfam³⁷ domain [*PfamScan*].
1444

1445 Using this script, the false composite Mvib_g619 was separated into two components (true
1446 genes): Mvib_g8200_1-297 and Mvib_g8200_298-891. This division solved the unexpected and
1447 not supported domain architecture shown by Mvib_g619, with the region corresponding to the two
1448 'Uricase' domains and the regions corresponding to the 'ArfGap' domain being now separate in
1449 two distinct genes. Overall, a total of 39, 19, 8 and 2 unexpected domain architectures not
1450 supported by the RNA-seq data were corrected from the genomes of *M. vibrans*, *P. atlantis*, *P.*
1451 *vietnamica* and *P. chileana*; respectively, this representing the 5.63%, 5.48%, 1.50% and 0.45%
1452 of all the unexpected domain architectures found, respectively. The remaining unexpected
1453 domain architectures that were not corrected may correspond either to true or to non-corrected
1454 false composites. The usage of more stringent alignment conditions would have increased the
1455 discontinuities found within transcripts, and hence the number of composites that would have
1456 been split. However, because we used a standard genome-annotation pipeline, we expected an
1457 overall good performance for this automatic annotation approach (a consideration that is
1458 supported by the results from *BUSCO* analyses on protein predictions, see Supplementary
1459 Information 1-Fig. 23). Thus, we decided to use a conservative approach and correct only those
1460 strong candidates of being false composites (the only constraint in the alignments was an E-value
1461 threshold of 1e-3).
1462

1463 In total, we detected 881, 284, 102, and 120 false composites for *M. vibrans*, *P. atlantis*, *P.*
1464 *vietnamica*, and *P. chileana*; respectively, this representing the 7.27%, 3.15%, 0.69%, and 0.83%
1465 of all the genes predicted, respectively. The false composite ratios found for *M. vibrans* and *P.*
1466 *atlantis* are clearly higher than those found for both *Pigoraptor* species. These differences may
1467 be explained by variances in the contiguity of the assemblies, with the genomes of *M. vibrans* and
1468 *P. atlantis*, especially the former, showing better N50 and L75 metrics than the two *Pigoraptor*
1469 genomes (Supplementary Information 1-Fig. 23). The reason is that because only neighbor genes
1470 can be mispredicted as a false composite, the number of potential false composites predictions
1471 should increase proportionally to the contiguity of the genome.
1472

1473 We evaluated the performance of our methodology by counting the differences in the number of
1474 reciprocal best hits (RBH) retrieved between *C. owczarzaki* (Filasterea) and each of the four
1475 species, when using the pre- and the post-corrected protein annotations, based on the following
1476 argument. If the protein A from *C. owczarzaki* (Ca) share orthology with the protein A from *M.*
1477 *vibrans* (Ma), we may expect Ma likely to be the best hit of Ca when aligning all *C. owczarzaki*
1478 proteins with *M. vibrans* proteins; and Ca of Ma when aligning all *M. vibrans* proteins with *C.*

1479 *owczarzaki* proteins. If we erroneously split the protein Ma into two proteins (Ma1 and Ma2), the
1480 protein Ca will still share a RBH relationship with a *M. vibrans* protein, either with Ma1 or Ma2.
1481 However, if we only consider those alignments covering a high percentage of the total query
1482 length (e.g. 75%), it is possible that we did not find any hit between Ca and Ma1 or Ma2 because
1483 both *M. vibrans* proteins could be shorter than Ca. Hence, an erroneous split of a gene may imply
1484 a decrease in the number of RBH recovered. We thus expect the finding and the correction of
1485 false composites to potentially increase the number of RBH recovered. For example, if the
1486 proteins A and B from *M. vibrans* (Ma and Mb) are respective orthologs of the proteins A and B
1487 from *C. owczarzaki* (Ca and Cb), but Ma and Mb were mispredicted as a false composite (Mab),
1488 we would only recover one RBH relation (Ca or Cb to Mab). The correction and split of Mab into
1489 Ma and Mb would increase by one the number of RBH relations (Ca to Ma and Cb to Mb).

1490

1491 We aligned all *C. owczarzaki* proteins with the pre- and post-corrected protein annotations of the
1492 four species and vice versa, and we counted only those RBH relations in which the corresponding
1493 alignments cover at least the 75% of both aligned sequences, with a minimal average identity of
1494 25%. With the post-corrected annotations, we always recovered more RBH relations between *C.*
1495 *owczarzaki* and our target species than with the pre-corrected annotations (19, 9, 2 and 7 for *M.*
1496 *vibrans*, *P. atlantis*, *P. vietnamita* and *P. chiliana*; respectively). Moreover, all the *C. owczarzaki*
1497 proteins involved in a RBH relation with the pre-corrected proteins were also involved in a RBH
1498 relation with the post-corrected proteins. These results altogether suggest that our approach
1499 modestly improved the quality of our annotations by detecting and splitting at least some of the
1500 mispredicted composites, with no evidence of false-positive splits.

1501

1502 Beyond false composites, a miss-prediction in the annotation process can also lead to erroneous
1503 splits of true genes into separate genes (false gene fissions). We also screened the predicted
1504 proteins to find and correct clear cases of false fissions in our annotations. We used *Transdecoder*
1505 (<https://github.com/TransDecoder>) to translate the RNA-seq *de novo* assembled transcripts into
1506 proteins (transcript proteins), and we kept only those transcript proteins that aligned with at least
1507 95% of query coverage and average identity with the same genomic scaffold (*bona fide* transcript
1508 proteins for our organism of interest; to discard potential biases introduced by contaminant
1509 transcripts). *Bona fide* transcript proteins were then aligned with the predicted proteins from the
1510 genomes [*BLASTp*: -evalue 1e-3]; and we screened for transcript proteins with at least two
1511 consecutive regions of its sequence having the best scoring hit with neighbor genes in the
1512 genome. We consider the finding of this alignment pattern as strong evidence of false fission,
1513 given that polycistronic transcripts are rare in eukaryotes³⁸ (the few occurrences of this alignment
1514 pattern corroborate that this assumption is valid for our genomes, see below). Protein sequences
1515 from the false fissioned neighbor genes were joined and aligned with *bona fide* transcript proteins
1516 and with euk_db [*BLASTp*: -evalue 1e-3]. After manual inspection of the alignments, we corrected

1517 2 clear false fissions for *P. chileana* and 1 for *M. vibrans*, *P. atlantis* and *P. vietnamita*. All FASTA
 1518 files produced will be available online as soon as we publish the corresponding manuscript.

1519

1520 **6) Assessment of genome quality**

1521 We used *QUAST v4.2*²⁵ to quantify L75 and N50 metrics for all the genomes of the species
 1522 represented in Supplementary Information 1-Fig. 23A. We also ran *BUSCO v1.22*²⁴ (*all Eukaryota*
 1523 dataset) on the genomes and proteomes of these species; and in the case of *M. vibrans*, *P.*
 1524 *atlantis*, *P. vietnamica* and *P. chileana*, also on proteomes predicted from *de novo* assembled
 1525 transcriptomes using *TransDecoder.LongOrfs v3.0.1* (<https://github.com/TransDecoder>).
 1526 'BUSCO C (%)' metric corresponds to the sum of 'Complete BUSCOs', 'Complete and single-
 1527 copy BUSCOs', 'Complete and duplicated BUSCOs' and 'Fragmented BUSCOs' metrics (in
 1528 percentage); whereas 'BUSCO D (%)' and 'BUSCO F (%)' metrics correspond to the percentages
 1529 of 'Missing BUSCOs' and 'Complete and duplicated BUSCOs', respectively (Supplementary
 1530 Information 1-Fig. 23).

1531

1532

| A | Scaffolds (>499 bp) | Genome size (Mb) | L75 | N50 (kb) | Genes | Proteome (from the genome) | | | Genome | | |
|--------------------------|---------------------|------------------|--------|----------|--------|----------------------------|-------------|-------------|-------------|-------------|-------------|
| | | | | | | BUSCO C (%) | BUSCO D (%) | BUSCO F (%) | BUSCO C (%) | BUSCO D (%) | BUSCO F (%) |
| METAZOA | | | | | | | | | | | |
| <i>M. brevicollis</i> | 219 | 41.71 | 27 | 1,073.6 | 9,233 | 90% | 20% | 8% | 70% | 3% | 9% |
| <i>S. rosetta</i> | 154 | 55.44 | 25 | 1,519.55 | 11,798 | 93% | 25% | 7% | 65% | 2% | 16% |
| <i>C. owczarzakii</i> | 84 | 27.97 | 11 | 1,617.77 | 8,898 | 97% | 28% | 3% | 84% | 8% | 3% |
| <i>M. vibrans</i> | 1,530 | 29.72 | 295 | 64.66 | 12,127 | 96% | 24% | 7% | 78% | 5% | 14% |
| <i>P. vietnamica</i> | 3,823 | 43.03 | 1,332 | 20.69 | 14,822 | 97% | 26% | 9% | 83% | 6% | 13% |
| <i>P. chileana</i> | 4,697 | 47.06 | 1,793 | 16.89 | 14,510 | 96% | 25% | 16% | 79% | 3% | 19% |
| FUNGI | | | | | | | | | | | |
| <i>C. fragrantissima</i> | 83 | 44.82 | 17 | 1,585.96 | 8,644 | 97% | 26% | 0% | 74% | 2% | 13% |
| <i>S. arctica</i> | 15,619 | 121.63 | 1,442 | 64.6 | 18,661 | 90% | 19% | 21% | 68% | 1% | 23% |
| <i>I. hoferi</i> | 1,633 | 88.08 | 515 | 106.77 | 6,351 | 76% | 14% | 16% | 65% | 1% | 16% |
| <i>A. whisleri</i> | 51,561 | 101.91 | 25,133 | 2.355 | 17,283 | 20% | 2% | 12% | 8% | 1% | 4% |
| <i>P. gemmata</i> | 4,697 | 75.27 | 25,440 | 1.853 | 21,835 | 24% | 2% | 19% | 6% | 0 | 4% |
| <i>C. perkinsii</i> | 3,994 | 36.91 | 187 | 120.165 | 12,463 | 97% | 26% | 3% | 89% | 8% | 6% |
| <i>C. limacisporum</i> | 287 | 24.11 | 86 | 180.47 | 7,535 | 96% | 25% | 2% | 87% | 12% | 2% |
| <i>F. alba</i> | 214 | 31.3 | 8 | 2,529.56 | 6,465 | 90% | 19% | 9% | 70% | 2% | 11% |
| <i>P. atlantis</i> | 1,308 | 19.24 | 452 | 27.01 | 9,028 | 98% | 25% | 2% | 85% | 1% | 18% |

BUSCO C(%): Complete + Fragmented + Multiple
 BUSCO D(%): Multiple
 BUSCO F(%): Fragmented

1533

| B | Proteome (from the genome) | | | Proteome (from the transcriptome) | | | Difference (Prot. genome. - Prot. transcript.) | | | C Genome (transcriptome-based metrics) | |
|----------------------|----------------------------|-------------|-------------|-----------------------------------|-------------|-------------|--|-------------|-------------|--|------------|
| | BUSCO C (%) | BUSCO D (%) | BUSCO F (%) | BUSCO C (%) | BUSCO D (%) | BUSCO F (%) | BUSCO C (%) | BUSCO D (%) | BUSCO F (%) | Completeness | Contiguity |
| <i>M. vibrans</i> | 96% | 24% | 7% | 92% | 26% | 21% | +4% | -2% | -14% | 94% (+16%) | 88% (+2%) |
| <i>P. vietnamica</i> | 97% | 26% | 9% | 56% | 16% | 20% | +41% | +10% | -11% | 94% (+11%) | 84% (-3%) |
| <i>P. chileana</i> | 96% | 25% | 16% | 80% | 29% | 13% | +16% | -4% | +3% | 91% (+12%) | 77% (-3%) |
| <i>P. atlantis</i> | 98% | 25% | 2% | 97% | 22% | 25% | +1% | +3% | -23% | 94% (+9%) | 89% (+7%) |

1534 **Supplementary Information 1-Fig. 23.** Completeness and contiguity metrics of the genomes from the four
 1535 species sequenced (highlighted in bold font) and also from other unicellular relatives of Metazoa and Fungi
 1536 with publicly available genome data. BUSCO software metrics are expressed in % (429 BUSCO markers),
 1537 and were computed for the genomic protein predictions and for the genomic scaffolds, and (B) also for the
 1538 protein predictions obtained from the transcriptomic data of *M. vibrans*, *P. vietnamica*, *P. chileana* and *P.*
 1539 *atlantis* using Transdecoder software (<https://github.com/TransDecoder>). The white/blue gradients in (A) are

1540 column-specific and represent differences between metric values (from the lowest to the highest). The
1541 blue/red code in (B) illustrates when the metrics are better or worse for the protein predictions coming from
1542 the genome than for those predictions coming from the transcriptome, respectively. (C) Genome
1543 completeness and contiguity metrics estimated from alignments of *bona fide* transcripts of each species to
1544 the corresponding genome. In particular, we estimated completeness by counting how many *bona fide*
1545 transcripts of the species aligned with the genome with an average identity of >95% and with a total query
1546 coverage of >95%. Contiguity was estimated as completeness, but only hits with the genomic fragment to
1547 which each transcript aligned with the highest score were considered (e.g., a transcript which sequence is
1548 complete but split into distinct genomic fragments will sum for completeness but not for contiguity).
1549 Differences between this transcriptome-based metrics and those found by BUSCO are indicated within the
1550 parenthesis, with the square being colored in blue or red according to whether the transcriptome-based
1551 metrics indicated a greater or a worse quality for the genome than the BUSCO metrics. Note that
1552 transcriptome-based completeness and contiguity metrics should be compared to 'BUSCO C(%)' and [100
1553 - 'BUSCO F(%)'], respectively. Number of bona fide transcript sequences (i.e., markers) per species: *M.*
1554 *vibrans* 10,056; *P. vietnamica*: 746; *P. chileana*: 2,200; *P. atlantis*: 5,841
1555

1556 An alternative transcriptome-based approach to assess completeness and contiguity
1557 (Supplementary Information 1-Fig. 23C) was applied to *M. vibrans*, *P. atlantis*, *P. vietnamica* and
1558 *P. chileana* genomes. It consists of aligning [*BLASTn*⁸: -evaluate 1e-3] a set of *bona fide* transcripts
1559 from this species to the corresponding genome, and computing the completeness and contiguity
1560 according to alignment results. In particular, we estimated completeness by counting how many
1561 transcripts aligned with the genome with an average identity of >95% and with a total query
1562 coverage of >95%. Contiguity was estimated as completeness, but only hits with the best scoring
1563 target genomic fragment were considered (e.g., a transcript in which the sequence is complete
1564 but split into distinct genomic fragments will sum for completeness but not for contiguity). For *M.*
1565 *vibrans* and *P. atlantis*, we used the set of *bona fide* transcripts defined during the
1566 decontamination process (see Supplementary Material 1). For both *Pigoraptor* species, the set
1567 was constructed in the following manner: we first performed a *de novo* transcriptome assembly
1568 of RNA-seq reads of these species using *Trinity* v2.2.0³⁹ [--jaccard_clip, --normalize_reads],
1569 previously preprocessed using *trimmomatic* v0.36⁴ [TruSeq3-PE-2.fa:2:30:10,
1570 SLIDINGWINDOW:4:5, LEADING:5, TRAILING:5, MINLEN:25]. The raw RNA-seq reads
1571 previously used to produce the transcriptomic data used in²⁸ were kindly provided by Elisabeth
1572 Hehenberger. We then used *TransDecoder.LongOrfs* to keep only transcripts corresponding to
1573 complete coding sequences, which were aligned to the metagenomes of *Parabodo caudatus* and
1574 of the corresponding *Pigoraptor* species [*BLASTn*: -evaluate 1e-3]. Transcripts that aligned to
1575 *Pigoraptor* but not to the *P. caudatus* metagenome were kept and aligned to euk_prok_db
1576 [*BLASTx*: -task blastx-fast, -evaluate 1e-3] (see Supplementary Material 1), and only those whose
1577 best scoring hit was a protein from *Capsaspora owczarzakii* (the only filasterean in the dataset)
1578 were considered as the *bona fide* transcripts, which were lastly processed for redundancy removal
1579 using *CD-HIT* v4.6¹⁸.

1580

1581 **References**

- 1582 1. Tong, S. M. Heterotrophic flagellates and other protists from Southampton Water, U.K. *Ophelia* **47**,
1583 71–131 (1997).
- 1584 2. Marron, A. O., Akam, M. & Walker, G. A Duplex PCR-Based Assay for Measuring the Amount of
1585 Bacterial Contamination in a Nucleic Acid Extract from a Culture of Free-Living Protists. *PLoS One*
1586 **8**, e61732 (2013).
- 1587 3. Torruella, G., de Mendoza, A., Grau-Bové, X. & Ruiz-Trillo, I. Phylogenomics Reveals Convergent
1588 Evolution of Lifestyles in Close Relatives of Animals and Fungi. *Curr. Biol.* **25**, 2404–2410 (2015).
- 1589 4. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data.
1590 *Bioinformatics* **30**, 2114–2120 (2014).
- 1591 5. O’Connell, J. *et al.* NxTrim: Optimized trimming of Illumina mate pair reads. *Bioinformatics* **31**,
1592 2035–2037 (2015).
- 1593 6. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-
1594 Cell Sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
- 1595 7. Grau-Bové, X. *et al.* Dynamics of genomic innovation in the unicellular ancestry of animals. *Elife* **6**,
1596 e26036 (2017).
- 1597 8. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search
1598 tool. *J. Mol. Biol.* **215**, 403–410 (1990).
- 1599 9. Burger, G., Gray, M. W., Forget, L. & Lang, B. F. Strikingly Bacteria-Like and Gene-Rich
1600 Mitochondrial Genomes throughout Jakobid Protists. *Genome Biol. Evol.* **5**, 418–438 (2013).
- 1601 10. Hoff, K. J., Lange, S., Lomsadze, A., Borodovsky, M. & Stanke, M. BRAKER1: Unsupervised RNA-
1602 Seq-based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics* **32**, 767–769
1603 (2015).
- 1604 11. Besemer, J. & Borodovsky, M. GeneMark: web software for gene finding in prokaryotes,
1605 eukaryotes and viruses. *Nucleic Acids Res.* **33**, W451–W454 (2005).
- 1606 12. Stanke, M. *et al.* AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res.* **34**,
1607 W435–439 (2006).
- 1608 13. Le, H. S., Schulz, M. H., Mccauley, B. M., Hinman, V. F. & Bar-Joseph, Z. Probabilistic error
1609 correction for RNA sequencing. *Nucleic Acids Res.* **41**, e109 (2013).
- 1610 14. Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence of insertions,
1611 deletions and gene fusions. *Genome Biol.* **14**, R36 (2013).
- 1612 15. Adl, S. M. *et al.* The Revised Classification of Eukaryotes. *J. Eukaryot. Microbiol.* **59**, 429–514
1613 (2012).
- 1614 16. Dick, G. J. *et al.* Community-wide analysis of microbial genome sequence signatures. *Genome*
1615 *Biol.* **10**, R85 (2009).
- 1616 17. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–
1617 359 (2012).
- 1618 18. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation
1619 sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
- 1620 19. Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat. Methods* **11**,
1621 1144–1146 (2014).
- 1622 20. Zaremba-Niedzwiedzka, K. *et al.* Asgard archaea illuminate the origin of eukaryotic cellular

- 1623 complexity. *Nature* **541**, 353–358 (2017).
- 1624 21. López-Escardó, D., López-García, P., Moreira, D., Ruiz-Trillo, I. & Torruella, G. Parvularia atlantis
1625 gen. et sp. nov., a Nucleariid Filose Amoeba (Holomycota, Opisthokonta). *J. Eukaryot. Microbiol.*
1626 **65**, 170–179 (2018).
- 1627 22. Peng, Y., Leung, H. C. M., Yiu, S. M. & Chin, F. Y. L. IDBA-UD: A de novo assembler for single-
1628 cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics* **28**, 1420–1428
1629 (2012).
- 1630 23. Boisvert, S., Laviolette, F. & Corbeil, J. Ray: Simultaneous Assembly of Reads from a Mix of High-
1631 Throughput Sequencing Technologies. *J. Comput. Biol.* **17**, 1519 (2010).
- 1632 24. Simao, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO:
1633 Assessing genome assembly and annotation completeness with single-copy orthologs.
1634 *Bioinformatics* **31**, 3210–3212 (2015).
- 1635 25. Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUASt: Quality assessment tool for genome
1636 assemblies. *Bioinformatics* **29**, 1072–1075 (2013).
- 1637 26. Vollmers, J., Wiegand, S. & Kaster, A. K. *Comparing and evaluating metagenome assembly tools*
1638 *from a microbiologist's perspective - Not only size matters! PLoS ONE* **12**, (2017).
- 1639 27. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–9 (2009).
- 1640 28. Hehenberger, E. *et al.* Novel Predators Reshape Holozoan Phylogeny and Reveal the Presence of
1641 a Two-Component Signaling System in the Ancestor of Animals. *Curr. Biol.* **27**, 2043–2050 (2017).
- 1642 29. Haas, B. J. *et al.* Improving the Arabidopsis genome annotation using maximal transcript alignment
1643 assemblies. *Nucleic Acids Res.* **31**, 5654–5666 (2003).
- 1644 30. Kent, W. J. BLAT--the BLAST-like alignment tool. *Genome Res.* **12**, 656–64 (2002).
- 1645 31. Wu, T. D. & Watanabe, C. K. GMAP: a genomic mapping and alignment program for mRNA and
1646 EST sequences. *Bioinformatics* **21**, 1859–1875 (2005).
- 1647 32. Chen, N. Using RepeatMasker to Identify Repetitive Elements in Genomic Sequences. *Curr.*
1648 *Protoc. Bioinforma.* **5**, 4.10.1-4.10.14 (2004).
- 1649 33. Basu, M. K., Carmel, L., Rogozin, I. B. & Koonin, E. V. Evolution of protein domain promiscuity in
1650 eukaryotes. *Genome Res.* **18**, 449–461 (2008).
- 1651 34. Haggerty, L. S. *et al.* A pluralistic account of homology: Adapting the models to the data. *Mol. Biol.*
1652 *Evol.* **31**, 501–516 (2014).
- 1653 35. Pathmanathan, J. S., Lopez, P., Lapointe, F.-J. & Baptiste, E. CompositeSearch: A Generalized
1654 Network Approach for Composite Gene Families Detection. *Mol. Biol. Evol.* **35**, 252–255 (2017).
- 1655 36. Slater, G. & Birney, E. Automated generation of heuristics for biological sequence comparison.
1656 *BMC Bioinformatics* **6**, 31 (2005).
- 1657 37. Finn, R. D. *et al.* Pfam: The protein families database. *Nucleic Acids Res.* **42**, D222–D230 (2014).
- 1658 38. Blumenthal, T. Gene clusters and polycistronic transcription in eukaryotes. *BioEssays* **20**, 480–487
1659 (1998).
- 1660 39. Haas, B. J. *et al.* De novo transcript sequence reconstruction from RNA-seq using the Trinity
1661 platform for reference generation and analysis. *Nat. Protoc.* **8**, 1494–1512 (2013).
- 1662
- 1663
- 1664
- 1665