# Identifying contamination with advanced visualization and analysis practices: metagenomic approaches for eukaryotic genome assemblies

Tom O Delmont, A. Murat Eren

High-throughput sequencing provides a fast and cost effective mean to recover genomes of organisms from all domains of life. However, adequate curation of the assembly results against potential contamination of non-target organisms requires advanced bioinformatics approaches and practices. Here, we re-analyzed the sequencing data generated for the tardigrade *Hypsibius dujardini* using approaches routinely employed by microbial ecologists who reconstruct bacterial and archaeal genomes from metagenomic data. We created a holistic display of the eukaryotic genome assembly using DNA data originating from two groups and eleven sequencing libraries. By using bacterial single-copy genes, k-mer frequencies, and coverage values of scaffolds we could identify and characterize multiple near-complete bacterial genomes, and curate a 182 Mbp draft genome for *H. dujardini* supported by RNA-Seq data. Our results indicate that most contaminant scaffolds were assembled from Moleculo long-read libraries, and most of these contaminants have differed between library preparations. Our re-analysis shows that visualization and curation of eukaryotic genome assemblies can benefit from tools designed to address the needs of today's microbiologists, who are constantly challenged by the difficulties associated with the identification of distinct microbial genomes in complex environmental metagenomes.

1 # Identifying contamination with advanced
2 # visualization and analysis practices:
3 # metagenomic approaches for eukaryotic
4 # genome assemblies

5 Tom O. Delmont[1] and A. Murat Eren[1,2]

6 [1] Department of Medicine, The University of Chicago, Chicago, IL , United States.

7 [2] Josephine Bay Paul Center, Marine Biological Laboratory, Woods Hole, MA , United States.

8 Corresponding Author:

9     A. Murat Eren[1]
10     *Knapp Center for Biomedical Discovery,*
11     *900 E. 57th St., MB 9, RM 9118, Chicago, IL 60637 USA*
12
13     Email address: meren@uchicago.edu
14
15

16  ## Abstract

17  High-throughput sequencing provides a fast and cost effective mean to recover genomes of
18  organisms from all domains of life. However, adequate curation of the assembly results
19  against potential contamination of non-target organisms requires advanced bioinformatics
20  approaches and practices. Here, we re-analyzed the sequencing data generated for the
21  tardigrade *Hypsibius dujardini,* and created a holistic display of the eukaryotic genome
22  assembly using DNA data originating from two groups and eleven sequencing libraries. By
23  using bacterial single-copy genes, k-mer frequencies, and coverage values of scaffolds we
24  could identify and characterize multiple near-complete bacterial genomes from the raw
25  assembly, and curate a 182 Mbp draft genome for *H. dujardini* supported by RNA-Seq data.
26  Our results indicate that most contaminant scaffolds were assembled from Moleculo long-
27  read libraries, and most of these contaminants have differed between library preparations.
28  Our re-analysis shows that visualization and curation of eukaryotic genome assemblies can
29  benefit from tools designed to address the needs of today's microbiologists, who are
30  constantly challenged by the difficulties associated with the identification of distinct
31  microbial genomes in complex environmental metagenomes.

32  ## Introduction

33  Advances in high-throughput sequencing technologies are revolutionizing the field of
34  genomics by allowing researchers to generate large amount of data in a short period of
35  time (Loman & Pallen, 2015). These technologies, combined with advances in
36  computational approaches, help us understand the diversity and functioning of life at
37  different scales by facilitating the rapid recovery of bacterial, archaeal, and eukaryotic
38  genomes (Venter et al., 2001; Schleper, Jurgens & Jonuscheit, 2005; Brown et al., 2015). Yet,
39  the recovery of genomes is not straightforward, and reconstructing bacterial and archaeal
40  versus eukaryotic genomes present researchers with distinct pitfalls and challenges that
41  result in different molecular and computational workflows.

42  For instance, difficulties associated with the cultivation of bacterial and archaeal organisms
43  (Schloss & Handelsman, 2003) have persuaded microbiologists to reconstruct genomes
44  directly from the environment through assembly-based metagenomics workflows and
45  genome binning. This workflow commonly entails (1) whole sequencing of environmental
46  genetic material, (2) assembly of short reads into contiguous DNA segments (contigs), and
47  (3) identification of draft genomes by binning contigs that originate from the same
48  organism. Due to the extensive diversity of bacteria and archaea in most environmental
49  samples (Gans, Wolinsky & Dunbar, 2005; Rusch et al., 2007), the field of metagenomics
50  has rapidly evolved to accurately delineate genomes in assembly results. Today,
51  microbiologists often exploit two essential properties of bacterial and archaeal genomes to
52  improve the "binning" step: (1) k-mer frequencies that are somewhat preserved
53  throughout a single microbial genome (Pride et al., 2003), to identify contigs that likely
54  originate from the same genome (Teeling et al., 2004), and (2) a set of genes that occur in
55  the vast majority of bacterial genomes as a single copy, to estimate the level of completion
56  and contamination of genome bins (Wu & Eisen, 2008; Campbell et al., 2013; Parks et al.,

57 2015). These properties, along with differential coverage of contigs across multiple
58 samples when such data exist, are routinely used to identify coherent microbial draft
59 genomes in metagenomic assemblies (Dick et al., 2009; Albertsen et al., 2013; Wu et al.,
60 2014; Alneberg et al., 2014; Kang et al., 2015; Eren et al., 2015).
61
62 On the other hand, researchers who study eukaryotic genomes generally focus on the
63 recovery of a single organism, which, in most cases, simplifies the identification of the
64 target genome in assembly results. However, sequences of bacterial origin can contaminate
65 eukaryotic genome assembly results due to their occurrence in samples (Chapman et al.,
66 2010; Artamonova & Mushegian, 2013), DNA extraction kits (Salter et al., 2014), or
67 laboratory environments (Laurence, Hatzis & Brash, 2014; Strong et al., 2014). One of the
68 major challenges of working with eukaryotic genomes is the extent of repeat regions that
69 complicate the assembly process (Richard, Kerrest & Dujon, 2008). To optimize the
70 assembly, researchers often employ multiple library preparations for sequencing (Gnerre
71 et al., 2010; Ekblom & Wolf, 2014), which may increase the potential sources of post-DNA
72 extraction contamination. Contaminants in assembly results can eventually contaminate
73 public databases (Merchant, Wood & Salzberg, 2014), and impair scientific findings
74 (Artamonova et al., 2015). The detection and removal of contaminants poses a major
75 bioinformatics challenge. To identify undesired contigs in a genomic assembly, scientists
76 can simply compare their assembly results to public sequence databases for positive hits to
77 unexpected taxa (Ekblom & Wolf, 2014), use k-mer coverage plots to identify distinct
78 genomes (Percudani, 2013), or employ scatter plots to partition contigs based on their GC-
79 content and coverage (Kumar et al., 2013). However, advanced solutions developed for
80 accurate identification of microbial genomes in complex metagenomic assemblies can
81 leverage these approaches further, and offer enhanced curation options for eukaryotic
82 assemblies.
83
84 The first release of a tardigrade genome by Boothby et al. (2015) demonstrates a striking
85 example of the importance of careful screening for contaminants in eukaryotic genome
86 assemblies. Tardigrades are microscopic animals occurring in a wide range of ecosystems
87 and they exhibit extended capabilities to survive in harsh conditions that would be fatal to
88 most animals (Ramløv & Westh, 2001; Jönsson, Harms-Ringdahl & Torudd, 2005; Jönsson
89 et al., 2008; Horikawa et al., 2013). Boothby and his colleagues generated a composite DNA
90 sequencing dataset from a culture of the tardigrade *Hypsibius dujardini* by exploiting some
91 of the best practices of high-throughput sequencing available today (Boothby et al., 2015).
92 In their assembled tardigrade genome, the authors detected a large number of genes
93 originating from bacteria, making up approximately one-sixth of the gene pool, and
94 suggested that horizontal gene transfers (HGTs) could explain the unique ability of
95 tardigrades to withstand extreme ranges of temperature, pressure, and radiation. However,
96 Koutsovoulos et al.'s subsequent analysis of Boothby et al.'s assembly suggested that it
97 contained extensive bacterial contamination, casting doubt on the extended HGT
98 hypothesis (Koutsovoulos et al., 2015). By applying two-dimensional scatterplots on their
99 own raw assembly results, Koutsovoulos et al. also reported a curated draft genome of *H.*
100 *dujardini*.

101  Here we re-analyzed the raw sequencing data generated by Boothby et al. (2015) and
102  Koutsovoulos et al. (2015), in combination with an independent RNA-Seq dataset
103  generated by Levin et al. (2016) for *H. dujardini*. Using anvi'o, an analysis and visualization
104  platform originally designed for the identification of bacterial genomes in metagenomic
105  assemblies (Eren et al., 2015), we employed bacterial single-copy genes to assess the
106  occurrence of bacterial genomes in the raw and curated assembly results, utilized k-mer
107  frequencies and coverage values across multiple sequencing libraries to organize scaffolds,
108  and visualized our findings in in a single display.


## Material and methods

110  **Genome assemblies, and raw sequencing data for DNA and RNA.** Boothby et al.
111  constructed three paired-end Illumina libraries (insert sizes of 0.3, 0.5 and 0.8 kbp) for 2 x
112  100 paired-end sequencing on a HiSeq2000 and six single-end long-read libraries (five
113  Illumina Moleculo libraries sequenced by the Illumina "long read" DNA sequencing service,
114  and one PacBio SMRT library sequenced using the P6-C4 chemistry and a 1 X 240 movie),
115  which altogether provided a co-assembly of 252.5 Mbp (Boothby et al., 2015). The
116  tardigrade genome released by Boothby et al. (2015), along with the nine sequencing data
117  used for its assembly, are available at http://weatherby.genetics.utah.edu/seq_transf.
118  Independently, Koutsovoulos et al. generated a 0.3 kbp insert library and a 1.1 kbp insert
119  mate-pair library for 2 x 100 paired end sequencing on a HiSeq2000 that provided a co-
120  assembly of 185.8 Mbp (nHd.1.0) (Koutsovoulos et al. 2015). These authors subsequently
121  curated a 135 Mbp draft genome (nHd.2.3) by removing potential contamination and re-
122  assembling filtered short reads (Koutsovoulos et al., 2015). The tardigrade raw assembly
123  and curated draft genome released by Koutsovoulos et al. (2015) are available at
124  http://badger.bio.ed.ac.uk/H_dujardini, and their two sequencing datasets are available
125  from the ENA, under study accession PRJEB11910.

126  **RNA-seq data**. We obtained the RNA-seq data using the NCBI accession id PRJNA272543
127  (Levin et al., 2016). Briefly, Levin et al. isolated RNA from *H. dujardini* using the Trizol
128  reagent (Invotrogen), constructed paired-end Illumina libraries according to the TruSeq
129  RNA-seq protocol, and sequenced their cDNA libraries with a read length of 100 bp.

130  **Quality filtering and read mapping**. We used illumina-utils (Eren et al., 2013) (available
131  from http://github.com/meren/illumina-utils) for quality filtering of short Illumina reads
132  using 'iu-filter-quality-minoche' script with default parameters, which implements the
133  quality filtering described by Minoche et al. (Minoche, Dohm & Himmelbauer, 2011).
134  Bowtie2 v2.2.4 (Langmead & Salzberg, 2012) with default parameters mapped all reads to
135  the scaffolds, and we used samtools v1.2 (Li et al., 2009) to convert reported SAM files to
136  BAM files.

137  **Overview of the anvi'o workflow.** Our workflow with anvi'o to identify and remove
138  contamination from a given collection of scaffolds consists of four main steps. The first step
139  is the processing of the FASTA file of scaffolds to create an anvi'o contigs database (CDB).
140  The resulting database holds basic information about each scaffold in the assembly (such as

141    the k-mer frequency, or GC-content). The second step is the profiling of each BAM file with
142    respect to the CDB we generated in the previous step. Each anvi'o profile describes
143    essential statistics for each scaffold in a given BAM file, including their average coverage,
144    and the portion of each scaffold covered by at least one read. The third step is to merge all
145    anvi'o profiles. Merging step combines all statistics from individual profiles, and uses them
146    to compute hierarchical clusterings of scaffolds. The default organization of scaffolds is
147    determined by the average coverage information from individual profiles, and the
148    sequence composition information from the CDB. This organization makes it possible to
149    identify scaffolds that distribute similarly across different library preparations. The final
150    step is to visualize the merged data on the anvi'o interactive interface. The anvi'o
151    interactive interface provides a holistic perspective of the combined data, and allows the
152    identification of draft genome bins, and removal of contaminants.

153    **Processing of scaffolds, and mapping results.** We used anvi'o v1.2.2 (available from
154    http://github.com/meren/anvio) to process scaffolds and mapping results, visualize the
155    distribution of scaffolds, and identify draft genomes following the workflow outlined in the
156    previous section, and detailed in Eren et al (2015). We created an anvi'o contigs database
157    CDB for each scaffold collection using the 'anvi-gen-contigs-database' program with default
158    parameters (where k equals 4 for k-mer frequency analysis). We then annotated scaffolds
159    with myRAST (available from http://theseed.org/) and imported these results into the CDB
160    using the program 'anvi-populate-genes-table' to store the information about the locations
161    of open reading frames (ORFs) in scaffolds, and their taxonomical and functional inference.
162    We profiled individual BAM files using the program 'anvi-profile' with a minimum contig
163    length of 1 kbp, and the program 'anvi-merge' combined resulting profiles with default
164    parameters. For the analysis of Boothby et al. (2015) assembly, we also profiled the RNA-
165    Seq data published by Levin et al. (2016) to identify scaffolds with transcriptomic activity,
166    and exported the table for proportion of each scaffold covered by transcripts using the
167    script 'get-db-table-as-matrix'. We used the supplementary material published by Boothby
168    et al. (2015) ("Dataset S1" in the original publication) to identify scaffolds with proposed
169    HGTs. Finally, we used the program 'anvi-interactive' visualize the merged data, and to
170    identify genome bins. We included RNA-Seq results and scaffolds with HGTs into our
171    visualization using the ` --additional-layers` flag. To finalize the anvi'o generated SVG files
172    for publication, we used Inkscape v0.91 (available from https://inkscape.org/).

173    **Predicting the number of bacterial genomes in an assembly**. We used the occurrence of
174    bacterial single-copy genes as a proxy to the expected number of bacterial genomes in a
175    raw assembly or in a curated genome bin. First, we ran on each CDB generated in this study
176    the anvi'o program 'anvi-populate-search-tables' to search using HMMer v3.1b2 (Eddy,
177    2011) for bacterial single-copy genes Campbell et al. (2013) published. Then, we used the
178    anvi'o script 'gen-stats-for-single-copy-genes' to report the number of hits per single-copy
179    gene as an array of integers from each CDB. We finally used mode (i.e., the most frequently
180    occurring number) of this array as the expected number of complete bacterial genomes in a
181    given collection of scaffolds. The script 'gen-stats-for-single-copy-genes' also used the R
182    library 'ggplot' v1.0.0 (R Development Core Team, 2011; Ginestet, 2011) to plot the
183    occurrence of single-copy genes.

184   **Taxonomical and functional annotation of bacterial genomes.** We uploaded bacterial
185   draft genomes identified from the raw tardigrade genomic assembly results into the RAST
186   server (Aziz et al., 2008), and used the RAST best taxonomic hits and FigFams to infer the
187   taxonomy of genome bins and functions they harbor.

188   **Data availability:** The URL <u>http://merenlab.org/data/</u> reports (1) anvi'o files to
189   regenerate Figure 1 and Figure 2, (2) our curation of the tardigrade genome from Boothby
190   et al.'s assembly (which is also available through the NCBI under the bioproject ID
191   PRJNA309530), and (3) the FASTA files for bacterial genomes we identified in the raw
192   assemblies from Boothby et al. and Koutsovoulos et al..

## Results and Discussion

194   Boothby et al. generated sequencing data from a tardigrade culture using three short read
195   (Illumina) and six long read (Moleculo and PacBio) libraries, which altogether provided a
196   co-assembly of 252.5 Mbp (Boothby et al., 2015). Using this assembly, the authors
197   suggested that 6,663 genes were entered into the tardigrade genome through HGTs.
198   Independently, Koutsovoulos et al. generated sequencing data from another tardigrade
199   culture using two short read Illumina libraries that provided a co-assembly of 185.8 Mbp,
200   from which they could curate a 135 Mbp tardigrade draft genome by removing potential
201   bacterial contamination using two-dimensional scatterplots of scaffolds with respect to
202   their GC-content and coverage (Koutsovoulos et al., 2015).

### A holistic view of the data

204   The use of multiple library preparations and sequencing strategies is likely to result in
205   more optimal assembly results (Gnerre et al., 2010). Hence, we focused on the scaffolds
206   generated by Boothby et al. (2015) as a foundation to maximize the recovery of the
207   tardigrade genome. To provide a holistic understanding of the composite sequencing data
208   generated by the two teams, we mapped the raw data from the nine DNA sequencing
209   libraries from Boothby et al., and the two Illumina libraries from Koutsovoulos et al. (2015)
210   on this assembly. Anvi'o generated a hierarchical clustering of scaffolds by combining the
211   tetra-nucleotide frequency and coverage of each scaffold across the 11 DNA sequencing
212   libraries (Eren et al., 2015). Besides visualizing the coverage of each scaffold in each
213   sample, we highlighted scaffolds with HGTs identified by Boothby et al. on the resulting
214   organization of scaffolds, and visualized RNA-seq mapping results. Figure 1 displays the
215   anvi'o merged profile that represents all this information in a single display.

### A draft genome for *H. dujardini*

217   Through the anvi'o interactive interface we selected 14,961 scaffolds from the Boothby et
218   al. assembly that recruited large number of short-reads in a consistent manner (Fig. 1).
219   This 182.2 Mbp selection with consistent coverage (#1 in Fig. 1) represents our curation of
220   the tardigrade draft genome from Boothby et al.'s assembly. The remaining 7,535 scaffolds,
221   which total about 70 Mbp of the assembly, harbored 96.1% of HGTs identified by Boothby

222  et al. These scaffolds recruited only 0.05% of the reads from the RNA-Seq data, highlighting
223  the extent of contamination in the original assembly. This finding is in agreement with
224  Koutsovoulos et al.'s findings; however, our curated draft genome from the Boothby et al.'s
225  assembly is 47 Mbp larger than the draft genome released by Koutsovoulos et al. (2015),
226  most probably due to Boothby et al.'s inclusion of longer reads from Moleculo libraries.
227  While the portion of scaffolds covered by RNA-Seq data suggests that this additional 47
228  Mbp still originate from the tardigrade genome, the biological relevance of this information
229  (or lack thereof) for the characterization of the tardigrade genome falls outside of the scope
230  of our study.

231  ## The origin of bacterial contamination

232  Our mapping results indicate the presence of non-target sequences in the assembly that
233  recruit reads only from long-read libraries. One interpretation could be that most of the
234  contamination in Boothby et al.'s assembly originated from Moleculo libraries, post DNA-
235  extraction (Fig. 1). However, while a recent study shows that the majority of long reads
236  from Moleculo libraries originated from low-abundance organisms in the analyzed samples
237  (Sharon et al., 2015), another study suggests relatively more sequencing bias in Moleculo
238  library preparation results (Kuleshov et al., 2015). Therefore, an alternative interpretation
239  of the mapping results can be that the bacterial contaminants were present in the sample
240  pre-DNA extraction at very low abundances, and each Moleculo library preparation
241  included long reads originating from different parts of this rare community. Regardless,
242  long reads considerably improved Boothby et al.'s assembly, which resulted in a larger
243  tardigrade genome following the removal of non-target sequences. While these results
244  reiterate that the use of long-read libraries is essential to generate more comprehensive
245  assemblies, they also suggest that extra care should be taken to better mitigate the
246  presence of non-target sequences in assembly results when long-read libraries are used for
247  sequencing.

248  We identified three near-complete bacterial genomes affiliated to *Chitinophaga* and
249  *Thermosinus* in Boothby et al.'s assembly (Fig. 1). Surprisingly, Boothby et al. identified only
250  a small portion of these complete bacterial genomes as sources of HGTs while applying a
251  metric specifically designed to detect foreign DNA in eukaryotic genomes. For instance,
252  none of the 4,459 genes in bacterial draft genome #2 (selection #3 in Fig. 1) were reported
253  in Boothby et al.'s findings as HGTs. We also processed and visualized the raw assembly
254  (nHd.1.0) from Koutsovoulos et al. (2015) using anvi'o (Figure S1), and recovered eight
255  bacterial genomes. However, we found no taxonomical overlap between high-completion
256  bacterial genomes from the two sequencing projects (Table S1).
257
258  Interestingly, one bacterial genome (selection #2 in Fig. 1) was detected in DNA libraries
259  from both groups, as well as in the RNA-seq data, suggesting that the related bacterial
260  population was in all samples prior to the DNA/RNA extraction step. This genome is
261  affiliated to *Chitinophaga*, and harbors genes coding for chitin degradation and utilization
262  (Table S2). Chitin occurs naturally in the feeding apparatus of tardigrades (Guidetti et al.,
263  2015), and might be a source of carbon for its microbial inhabitants. The genome also
264  harbors genes coding for the biosynthesis of proteorhodopsin, host invasion and

265  intracellular resistance, dormancy and sporulation, oxidative stress, and tryptophan, which
266  is an essential amino acid for animals (Crawford, 1989; Zelante et al., 2013). Although this
267  genome may belong to a tardigrade symbiont, the generation of the data does not allow us
268  to rule out the possibility that it may be associated with the food source. Nevertheless, this
269  finding suggests that there may be cases where non-target genomes in an assembly can
270  provide clues about the lifestyle of a given host.

271  **Best practices to assess bacterial contamination**

272  Initial assessment of the occurrence of bacterial single-copy genes in eukaryotic assemblies
273  can provide a quick estimation of the number of bacterial genomes that occur in assembly
274  results. The use of bacterial single-copy genes can give much more accurate representation
275  of potential bacterial contamination than screening for 16S rRNA genes alone, as they are
276  less likely to be found in co-assembly results (Miller et al., 2011; Delmont et al., 2015).
277  Although Boothby et al. reported the lack of 16S rRNA genes in their assembly (Boothby et
278  al., 2015), anvi'o estimated that it contained at least 10 complete bacterial genomes (Fig. 2)
279  using a bacterial single-copy gene collection (Campbell et al., 2013). This simple yet
280  powerful step could identify cases of extensive contamination, and alert researchers to be
281  diligent in identifying scaffolds originating from bacterial organisms. Figure 2 also
282  summarizes the HMM hits in scaffolds found in curated tardigrade genomes from our
283  analysis and Koutsovoulos et al.'s study. We observed that the average significance score
284  for the remaining HMM hits for bacterial single-copy genes in curated genomes was 4.2
285  times lower in average compared to the HMM hits in assembly results (Table S3). The
286  decrease in the significance scores, and the very similar patterns of occurrence of HMM hits
287  between the two curation efforts suggest that some of the HMM profiles may not be specific
288  enough to be identified only in bacteria.

289  Two-dimensional scatterplots have a long history of identifying distinct genomes in
290  assembly results (Tyson et al., 2004) and continue to be used for delineating microbial
291  genomes in metagenomic assemblies (Albertsen et al., 2013; Cantor et al., 2015), as well as
292  detecting contamination in eukaryotic assembly results (Kumar et al., 2013). Although
293  scatterplots can describe the organization of contigs in assembly results, they suffer from
294  limited number of dimensions they can display, and their inability to depict complex
295  supporting data that can improve the identification of individual genomes. These
296  limitations are particularly problematic in sequencing projects covering multiple
297  sequencing libraries, where displaying mapping results from each library can help
298  detecting sources of contaminants. Despite their successful applications, two dimensional
299  scatter plots limit researchers to the use of simple characteristics of the data that can be
300  represented on an axis (such as GC-content). In contrast, clustering scaffolds, and
301  overlaying multiple layers of independent information produce more comprehensive
302  visualizations that display multiple aspects of the data.

## Conclusions

The field of genomics requires advanced computational approaches to take best advantage of constantly evolving ways to generate sequencing data, and to identify and remove contamination from genome assemblies. Our study indicates that some of these advanced approaches may emerge from the field of metagenomics, where the need for *de novo* reconstruction of microbial genomes from environmental samples has given raise to techniques and software platforms that can make sense of complex assemblies. Here we used k-mer frequencies to organize scaffolds, the occurrence of bacterial single-copy genes to estimate the extent of contamination, and advanced visualization strategies to detect and remove contamination in a eukaryotic assembly project while simultaneously characterizing the sources of contamination. Our results also suggest that metagenomic binning strategies can be used to recover near-complete bacterial genomes from raw eukaryotic assemblies, which can provide insights into the potential host-microbe interactions during the curation step.

## Acknowledgments

## References

Albertsen M., Hugenholtz P., Skarshewski A., Nielsen KL., Tyson GW., Nielsen PH. 2013. Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. *Nature biotechnology* 31:533–8.

Alneberg J., Bjarnason BS., de Bruijn I., Schirmer M., Quick J., Ijaz UZ., Lahti L., Loman NJ., Andersson AF., Quince C. 2014. Binning metagenomic contigs by coverage and composition. *Nature Methods* 11:1144–1146.

Artamonova II., Lappi T., Zudina L., Mushegian AR. 2015. Prokaryotic genes in eukaryotic genome sequences: when to infer horizontal gene transfer and when to suspect an actual microbe. *Environmental Microbiology* 17:2203–2208.

Artamonova II., Mushegian AR. 2013. Genome sequence analysis indicates that the model eukaryote Nematostella vectensis harbors bacterial consorts. *Applied and environmental microbiology* 79:6868–73.

Aziz RK., Bartels D., Best AA., DeJongh M., Disz T., Edwards RA., Formsma K., Gerdes S., Glass EM., Kubal M., Meyer F., Olsen GJ., Olson R., Osterman AL., Overbeek RA., McNeil LK.,

338         Paarmann D., Paczian T., Parrello B., Pusch GD., Reich C., Stevens R., Vassieva O.,
339         Vonstein V., Wilke A., Zagnitko O. 2008. The RAST Server: rapid annotations using
340         subsystems technology. *BMC genomics* 9:75.

341   Boothby TC., Tenlen JR., Smith FW., Wang JR., Patanella KA., Osborne Nishimura E., Tintori
342         SC., Li Q., Jones CD., Yandell M., Messina DN., Glasscock J., Goldstein B. 2015. Evidence
343         for extensive horizontal gene transfer from the draft genome of a tardigrade.
344         *Proceedings of the National Academy of Sciences* 112:201510461.

345   Brown CT., Hug LA., Thomas BC., Sharon I., Castelle CJ., Singh A., Wilkins MJ., Wrighton KC.,
346         Williams KH., Banfield JF. 2015. Unusual biology across a group comprising more than
347         15% of domain Bacteria. *Nature* 523:208–211.

348   Campbell JH., O'Donoghue P., Campbell AG., Schwientek P., Sczyrba A., Woyke T., Söll D.,
349         Podar M. 2013. UGA is an additional glycine codon in uncultured SR1 bacteria from the
350         human microbiota. *Proceedings of the National Academy of Sciences of the United States
351         of America* 110:5540–5.

352   Cantor M., Nordberg H., Smirnova T., Hess M., Tringe S., Dubchak I. 2015. Elviz – exploration
353         of metagenome assemblies with an interactive visualization tool. *BMC Bioinformatics*
354         16:130.

355   Chapman JA., Kirkness EF., Simakov O., Hampson SE., Mitros T., Weinmaier T., Rattei T.,
356         Balasubramanian PG., Borman J., Busam D., Disbennett K., Pfannkoch C., Sumin N.,
357         Sutton GG., Viswanathan LD., Walenz B., Goodstein DM., Hellsten U., Kawashima T.,
358         Prochnik SE., Putnam NH., Shu S., Blumberg B., Dana CE., Gee L., Kibler DF., Law L.,
359         Lindgens D., Martinez DE., Peng J., Wigge PA., Bertulat B., Guder C., Nakamura Y., Ozbek
360         S., Watanabe H., Khalturin K., Hemmrich G., Franke A., Augustin R., Fraune S.,
361         Hayakawa E., Hayakawa S., Hirose M., Hwang JS., Ikeo K., Nishimiya-Fujisawa C., Ogura
362         A., Takahashi T., Steinmetz PRH., Zhang X., Aufschnaiter R., Eder M-K., Gorny A-K.,
363         Salvenmoser W., Heimberg AM., Wheeler BM., Peterson KJ., Böttger A., Tischler P., Wolf
364         A., Gojobori T., Remington KA., Strausberg RL., Venter JC., Technau U., Hobmayer B.,
365         Bosch TCG., Holstein TW., Fujisawa T., Bode HR., David CN., Rokhsar DS., Steele RE.
366         2010. The dynamic genome of Hydra. *Nature* 464:592–6.

367   Crawford IP. 1989. Evolution of a biosynthetic pathway: the tryptophan paradigm. *Annual
368         review of microbiology* 43:567–600.

369   Delmont TO., Eren AM., Maccario L., Prestat E., Esen ÖC., Pelletier E., Le Paslier D., Simonet
370         P., Vogel TM. 2015. Reconstructing rare soil microbial genomes using in situ
371         enrichments and metagenomics. *Frontiers in microbiology* 6:358.

372   Dick GJ., Andersson AF., Baker BJ., Simmons SL., Thomas BC., Yelton AP., Banfield JF. 2009.
373         Community-wide analysis of microbial genome sequence signatures. *Genome biology*
374         10:R85.

375   Eddy SR. 2011. Accelerated Profile HMM Searches. *PLoS computational biology* 7:e1002195.

376    Ekblom R., Wolf JBW. 2014. A field guide to whole-genome sequencing, assembly and
377        annotation. *Evolutionary Applications* 7:n/a–n/a.

378    Eren AM., Vineis JH., Morrison HG., Sogin ML. 2013. A Filtering Method to Generate High
379        Quality Short Reads Using Illumina Paired-End Technology. *PLoS ONE* 8:e66643.

380    Eren AM., Esen ÖC., Quince C., Vineis JH., Morrison HG., Sogin ML., Delmont TO. 2015.
381        Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* 3:e1319.

382    Gans J., Wolinsky M., Dunbar J. 2005. Computational improvements reveal great bacterial
383        diversity and high metal toxicity in soil. *Science (New York, N.Y.)* 309:1387–90.

384    Ginestet C. 2011. ggplot2: Elegant Graphics for Data Analysis. *Journal of the Royal Statistical
385        Society: Series A (Statistics in Society)* 174:245–246.

386    Gnerre S., MacCallum I., Przybylski D., Ribeiro FJ., Burton JN., Walker BJ., Sharpe T., Hall G.,
387        Shea TP., Sykes S., Berlin AM., Aird D., Costello M., Daza R., Williams L., Nicol R., Gnirke
388        A., Nusbaum C., Lander ES., Jaffe DB. 2010. High-quality draft assemblies of
389        mammalian genomes from massively parallel sequence data. *Proceedings of the
390        National Academy of Sciences* 108:1513–1518.

391    Guidetti R., Bonifacio A., Altiero T., Bertolani R., Rebecchi L. 2015. Distribution of Calcium
392        and Chitin in the Tardigrade Feeding Apparatus in Relation to its Function and
393        Morphology. *Integrative and comparative biology* 55:241–52.

394    Horikawa DD., Cumbers J., Sakakibara I., Rogoff D., Leuko S., Harnoto R., Arakawa K.,
395        Katayama T., Kunieda T., Toyoda A., Fujiyama A., Rothschild LJ. 2013. Analysis of DNA
396        repair and protection in the Tardigrade Ramazzottius varieornatus and Hypsibius
397        dujardini after exposure to UVC radiation. *PloS one* 8:e64793.

398    Jönsson KI., Rabbow E., Schill RO., Harms-Ringdahl M., Rettberg P. 2008. Tardigrades
399        survive exposure to space in low Earth orbit. *Current biology : CB* 18:R729–R731.

400    Jönsson KI., Harms-Ringdahl M., Torudd J. 2005. Radiation tolerance in the eutardigrade
401        Richtersius coronifer. *International journal of radiation biology* 81:649–56.

402    Kang DD., Froula J., Egan R., Wang Z. 2015. MetaBAT, an efficient tool for accurately
403        reconstructing single genomes from complex microbial communities. *PeerJ* 3:e1165.

404    Koutsovoulos G., Kumar S., Laetsch DR., Stevens L., Daub J., Conlon C., Maroon H., Thomas F.,
405        Aboobaker A., Blaxter M. 2015. *The genome of the tardigrade Hypsibius dujardini*. Cold
406        Spring Harbor Labs Journals.

407    Kuleshov V., Jiang C., Zhou W., Jahanbani F., Batzoglou S., Snyder M. 2015. Synthetic long-
408        read sequencing reveals intraspecies diversity in the human microbiome. *Nature
409        Biotechnology* 34:64–69.

410    Kumar S., Jones M., Koutsovoulos G., Clarke M., Blaxter M. 2013. Blobology: exploring raw
411         genome data for contaminants, symbionts and parasites using taxon-annotated GC-
412         coverage plots. *Frontiers in genetics* 4:237.

413    Langmead B., Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature*
414         *methods* 9:357–9.

415    Laurence M., Hatzis C., Brash DE. 2014. Common contaminants in next-generation
416         sequencing that hinder discovery of low-abundance microbes. *PloS one* 9:e97876.

417    Levin M., Anavy L., Cole AG., Winter E., Mostov N., Khair S., Senderovich N., Kovalev E., Silver
418         DH., Feder M., Fernandez-Valverde SL., Nakanishi N., Simmons D., Simakov O., Larsson
419         T., Liu S-Y., Jerafi-Vider A., Yaniv K., Ryan JF., Martindale MQ., Rink JC., Arendt D.,
420         Degnan SM., Degnan BM., Hashimshony T., Yanai I. 2016. The mid-developmental
421         transition and the evolution of animal body plans. *Nature* advance on.

422    Li H., Handsaker B., Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin
423         R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford,*
424         *England)* 25:2078–9.

425    Loman NJ., Pallen MJ. 2015. Twenty years of bacterial genome sequencing. *Nature Reviews*
426         *Microbiology* 13:787–794.

427    Merchant S., Wood DE., Salzberg SL. 2014. Unexpected cross-species contamination in
428         genome sequencing projects. *PeerJ* 2:e675.

429    Miller CS., Baker BJ., Thomas BC., Singer SW., Banfield JF. 2011. EMIRGE: reconstruction of
430         full-length ribosomal genes from microbial community short read sequencing data.
431         *Genome biology* 12:R44.

432    Minoche AE., Dohm JC., Himmelbauer H. 2011. Evaluation of genomic high-throughput
433         sequencing data generated on Illumina HiSeq and genome analyzer systems. *Genome*
434         *biology* 12:R112.

435    Parks DH., Imelfort M., Skennerton CT., Hugenholtz P., Tyson GW. 2015. CheckM: assessing
436         the quality of microbial genomes recovered from isolates, single cells, and
437         metagenomes. *Genome research* 25:1043–55.

438    Percudani R. 2013. A Microbial Metagenome (Leucobacter sp.) in Caenorhabditis Whole
439         Genome Sequences. *Bioinformatics and biology insights* 7:55–72.

440    Pride DT., Meinersmann RJ., Wassenaar TM., Blaser MJ. 2003. Evolutionary implications of
441         microbial genome tetranucleotide frequency biases. *Genome research* 13:145–58.

442    R Development Core Team R. 2011. R: A Language and Environment for Statistical
443         Computing. *R Foundation for Statistical Computing* 1:409.

PeerJ

444   Ramløv H., Westh P. 2001. Cryptobiosis in the Eutardigrade Adorybiotus (Richtersius)
445       coronifer: Tolerance to Alcohols, Temperature and de novo Protein Synthesis.
446       *Zoologischer Anzeiger - A Journal of Comparative Zoology* 240:517–523.

447   Richard G-F., Kerrest A., Dujon B. 2008. Comparative genomics and molecular dynamics of
448       DNA repeats in eukaryotes. *Microbiology and molecular biology reviews : MMBR*
449       72:686–727.

450   Rusch DB., Halpern AL., Sutton G., Heidelberg KB., Williamson S., Yooseph S., Wu D., Eisen
451       JA., Hoffman JM., Remington K., Beeson K., Tran B., Smith H., Baden-Tillson H., Stewart
452       C., Thorpe J., Freeman J., Andrews-Pfannkoch C., Venter JE., Li K., Kravitz S., Heidelberg
453       JF., Utterback T., Rogers Y-H., Falcón LI., Souza V., Bonilla-Rosso G., Eguiarte LE., Karl
454       DM., Sathyendranath S., Platt T., Bermingham E., Gallardo V., Tamayo-Castillo G.,
455       Ferrari MR., Strausberg RL., Nealson K., Friedman R., Frazier M., Venter JC. 2007. The
456       Sorcerer II Global Ocean Sampling expedition: northwest Atlantic through eastern
457       tropical Pacific. *PLoS biology* 5:e77.

458   Salter SJ., Cox MJ., Turek EM., Calus ST., Cookson WO., Moffatt MF., Turner P., Parkhill J.,
459       Loman NJ., Walker AW. 2014. Reagent and laboratory contamination can critically
460       impact sequence-based microbiome analyses. *BMC Biology* 12:87.

461   Schleper C., Jurgens G., Jonuscheit M. 2005. Genomic studies of uncultivated archaea.
462       *Nature reviews. Microbiology* 3:479–88.

463   Schloss PD., Handelsman J. 2003. Biotechnological prospects from metagenomics. *Current*
464       *opinion in biotechnology* 14:303–10.

465   Sharon I., Kertesz M., Hug LA., Pushkarev D., Blauwkamp TA., Castelle CJ., Amirebrahimi M.,
466       Thomas BC., Burstein D., Tringe SG., Williams KH., Banfield J. 2015. Accurate, multi-kb
467       reads resolve complex populations and detect rare microorganisms. *Genome*
468       *Research*:gr.183012.114.

469   Strong MJ., Xu G., Morici L., Splinter Bon-Durant S., Baddoo M., Lin Z., Fewell C., Taylor CM.,
470       Flemington EK. 2014. Microbial contamination in next generation sequencing:
471       implications for sequence-based analysis of clinical samples. *PLoS pathogens*
472       10:e1004437.

473   Teeling H., Meyerdierks A., Bauer M., Amann R., Glöckner FO. 2004. Application of
474       tetranucleotide frequencies for the assignment of genomic fragments. *Environmental*
475       *microbiology* 6:938–47.

476   Tyson GW., Chapman J., Hugenholtz P., Allen EE., Ram RJ., Richardson PM., Solovyev V V.,
477       Rubin EM., Rokhsar DS., Banfield JF. 2004. Community structure and metabolism
478       through reconstruction of microbial genomes from the environment. *Nature* 428:37–
479       43.

480   Venter JC., Adams MD., Myers EW., Li PW., Mural RJ., Sutton GG., Smith HO., Yandell M.,

481  Evans CA., Holt RA., Gocayne JD., Amanatides P., Ballew RM., Huson DH., Wortman JR.,
482  Zhang Q., Kodira CD., Zheng XH., Chen L., Skupski M., Subramanian G., Thomas PD.,
483  Zhang J., Gabor Miklos GL., Nelson C., Broder S., Clark AG., Nadeau J., McKusick VA.,
484  Zinder N., Levine AJ., Roberts RJ., Simon M., Slayman C., Hunkapiller M., Bolanos R.,
485  Delcher A., Dew I., Fasulo D., Flanigan M., Florea L., Halpern A., Hannenhalli S., Kravitz
486  S., Levy S., Mobarry C., Reinert K., Remington K., Abu-Threideh J., Beasley E., Biddick K.,
487  Bonazzi V., Brandon R., Cargill M., Chandramouliswaran I., Charlab R., Chaturvedi K.,
488  Deng Z., Di Francesco V., Dunn P., Eilbeck K., Evangelista C., Gabrielian AE., Gan W., Ge
489  W., Gong F., Gu Z., Guan P., Heiman TJ., Higgins ME., Ji RR., Ke Z., Ketchum KA., Lai Z., Lei
490  Y., Li Z., Li J., Liang Y., Lin X., Lu F., Merkulov G V., Milshina N., Moore HM., Naik AK.,
491  Narayan VA., Neelam B., Nusskern D., Rusch DB., Salzberg S., Shao W., Shue B., Sun J.,
492  Wang Z., Wang A., Wang X., Wang J., Wei M., Wides R., Xiao C., Yan C., Yao A., Ye J., Zhan
493  M., Zhang W., Zhang H., Zhao Q., Zheng L., Zhong F., Zhong W., Zhu S., Zhao S., Gilbert D.,
494  Baumhueter S., Spier G., Carter C., Cravchik A., Woodage T., Ali F., An H., Awe A.,
495  Baldwin D., Baden H., Barnstead M., Barrow I., Beeson K., Busam D., Carver A., Center
496  A., Cheng ML., Curry L., Danaher S., Davenport L., Desilets R., Dietz S., Dodson K., Doup
497  L., Ferriera S., Garg N., Glucksmann A., Hart B., Haynes J., Haynes C., Heiner C., Hladun
498  S., Hostin D., Houck J., Howland T., Ibegwam C., Johnson J., Kalush F., Kline L., Koduru S.,
499  Love A., Mann F., May D., McCawley S., McIntosh T., McMullen I., Moy M., Moy L.,
500  Murphy B., Nelson K., Pfannkoch C., Pratts E., Puri V., Qureshi H., Reardon M.,
501  Rodriguez R., Rogers YH., Romblad D., Ruhfel B., Scott R., Sitter C., Smallwood M.,
502  Stewart E., Strong R., Suh E., Thomas R., Tint NN., Tse S., Vech C., Wang G., Wetter J.,
503  Williams S., Williams M., Windsor S., Winn-Deen E., Wolfe K., Zaveri J., Zaveri K., Abril
504  JF., Guigó R., Campbell MJ., Sjolander K V., Karlak B., Kejariwal A., Mi H., Lazareva B.,
505  Hatton T., Narechania A., Diemer K., Muruganujan A., Guo N., Sato S., Bafna V., Istrail S.,
506  Lippert R., Schwartz R., Walenz B., Yooseph S., Allen D., Basu A., Baxendale J., Blick L.,
507  Caminha M., Carnes-Stine J., Caulk P., Chiang YH., Coyne M., Dahlke C., Mays A.,
508  Dombroski M., Donnelly M., Ely D., Esparham S., Fosler C., Gire H., Glanowski S., Glasser
509  K., Glodek A., Gorokhov M., Graham K., Gropman B., Harris M., Heil J., Henderson S.,
510  Hoover J., Jennings D., Jordan C., Jordan J., Kasha J., Kagan L., Kraft C., Levitsky A., Lewis
511  M., Liu X., Lopez J., Ma D., Majoros W., McDaniel J., Murphy S., Newman M., Nguyen T.,
512  Nguyen N., Nodell M., Pan S., Peck J., Peterson M., Rowe W., Sanders R., Scott J., Simpson
513  M., Smith T., Sprague A., Stockwell T., Turner R., Venter E., Wang M., Wen M., Wu D., Wu
514  M., Xia A., Zandieh A., Zhu X. 2001. The sequence of the human genome. *Science (New
515  York, N.Y.)* 291:1304–51.

516  Wu Y-W., Tang Y-H., Tringe SG., Simmons BA., Singer SW. 2014. MaxBin: an automated
517  binning method to recover individual genomes from metagenomes using an
518  expectation-maximization algorithm. *Microbiome* 2:26.

519  Wu M., Eisen JA. 2008. A simple, fast, and accurate method of phylogenomic inference.
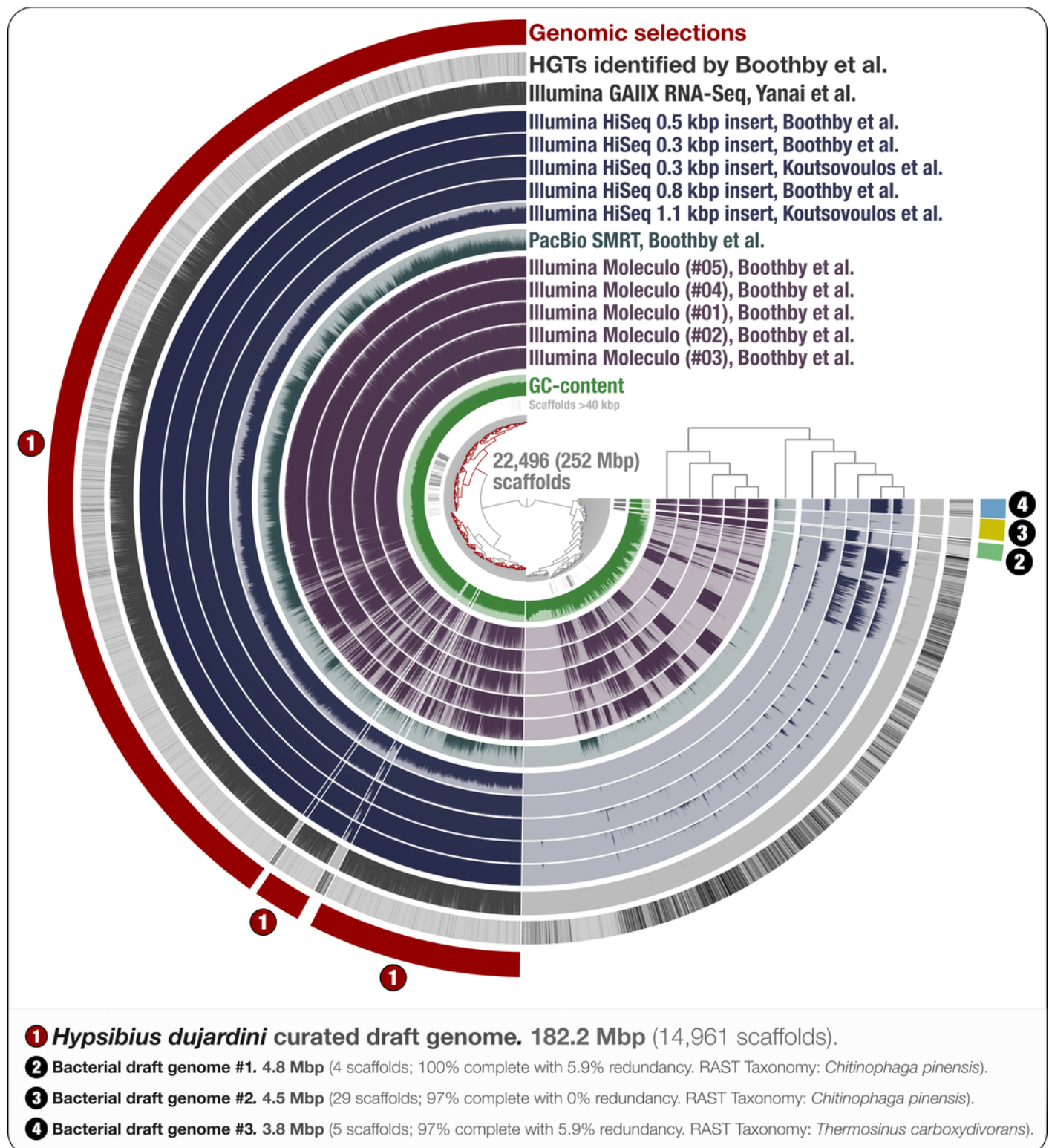520  *Genome Biology* 9:R151.

521  Zelante T., Iannitti RG., Cunha C., De Luca A., Giovannini G., Pieraccini G., Zecchi R., D'Angelo
522  C., Massi-Benedetti C., Fallarino F., Carvalho A., Puccetti P., Romani L. 2013.
523  Tryptophan catabolites from microbiota engage aryl hydrocarbon receptor and

524      balance mucosal reactivity via interleukin-22. *Immunity* 39:372–85.

# 1

Holistic assessment of the tardigrade genome release from Boothby et al. (2015).

Dendrogram in the center organizes scaffolds based on sequence composition and coverage values in data from 11 DNA libraries. Scaffolds larger than 40 kbp were split into sections of 20 kbp for visualization purposes. Splits are displayed in the first inner circle and GC-content (0-71%) in the second circle. In the following 11 layers, each bar represents the portion of scaffolds covered by short reads in a given sample. The next layer shows the same information for RNA-Seq data. Scaffolds harboring genes used by Boothby et al. to support the expended HGT hypothesis is shown in the next layer. Finally, the outermost layer shows our selections of scaffolds as draft genome bins: the curated tardigrade genome (selection #1), as well as three near-complete bacterial genomes originating from various contamination sources (selection #2, #3, and #4).

**Genomic selections**

**HGTs identified by Boothby et al.**

**Illumina GAIIX RNA-Seq, Yanai et al.**

**Illumina HiSeq 0.5 kbp insert, Boothby et al.**
**Illumina HiSeq 0.3 kbp insert, Boothby et al.**
**Illumina HiSeq 0.3 kbp insert, Koutsovoulos et al.**
**Illumina HiSeq 0.8 kbp insert, Boothby et al.**
**Illumina HiSeq 1.1 kbp insert, Koutsovoulos et al.**

**PacBio SMRT, Boothby et al.**

**Illumina Moleculo (#05), Boothby et al.**
**Illumina Moleculo (#04), Boothby et al.**
**Illumina Moleculo (#01), Boothby et al.**
**Illumina Moleculo (#02), Boothby et al.**
**Illumina Moleculo (#03), Boothby et al.**

**GC-content**
Scaffolds >40 kbp

22,496 (252 Mbp)
scaffolds

**❶** *Hypsibius dujardini* curated draft genome. 182.2 Mbp (14,961 scaffolds).

**❷** **Bacterial draft genome #1.** 4.8 Mbp (4 scaffolds; 100% complete with 5.9% redundancy. RAST Taxonomy: *Chitinophaga pinensis*).

**❸** **Bacterial draft genome #2.** 4.5 Mbp (29 scaffolds; 97% complete with 0% redundancy. RAST Taxonomy: *Chitinophaga pinensis*).

**❹** **Bacterial draft genome #3.** 3.8 Mbp (5 scaffolds; 97% complete with 5.9% redundancy. RAST Taxonomy: *Thermosinus carboxydivorans*).

# 2

Occurrence of the 139 bacterial single-copy genes reported by Campbell et al. (2013) across scaffold collections.

The top two plots display the frequency and distribution of single-copy genes in the raw tardigrade genomic assembly generated by Boothby et al. (2015), and Koutsovoulos et al. (2015), respectively. The bottom two plots display the same information for each of the curated tardigrade genomes. Each bar represents the squared-root normalized number of significant hits per single-copy gene. The same information is visualized as box-plots on the left side of each plot.