

A peer-reviewed version of this preprint was published in PeerJ on 29 March 2016.

[View the peer-reviewed version](http://peerj.com/articles/1839) (peerj.com/articles/1839), which is the preferred citable publication unless you specifically need to cite this preprint.

Delmont TO, Eren AM. 2016. Identifying contamination with advanced visualization and analysis practices: metagenomic approaches for eukaryotic genome assemblies. PeerJ 4:e1839 <https://doi.org/10.7717/peerj.1839>

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

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.

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

Tom O. Delmont¹ and A. Murat Eren^{1,2*}

*Correspondence:

meren@uchicago.edu

¹The Department of Medicine,
The University of Chicago,
Chicago, 60637 IL, USA

²Josephine Bay Paul Center for
Comparative Molecular Biology
and Evolution, Marine Biological
Laboratory, Woods Hole, 02543
MA, USA

Full list of author information is
available at the end of the
article

Abstract

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.

Keywords: genomics; assembly; curation; visualization; contamination; *anvi'o*; HGT

Introduction

Advances in high-throughput sequencing technologies are revolutionizing the field of genomics by allowing researchers to generate large amount of data in a short period of time [1]. These technologies, combined with advances in computational approaches, help us understand the diversity and functioning of life at different scales by facilitating the rapid recovery of bacterial, archaeal, and eukaryotic genomes [2, 3, 4]. Yet, the recovery of genomes is not straightforward, and reconstructing bacterial and archaeal versus eukaryotic genomes present researchers with distinct pitfalls and challenges that result in different molecular and computational workflows.

For instance, difficulties associated with the cultivation of bacterial and archaeal organisms [5] have persuaded microbiologists to reconstruct genomes directly from the environment through assembly-based metagenomics workflows and genome binning. This workflow commonly entails (1) whole sequencing of environmental genetic material, (2) assembly of short reads into contiguous DNA segments (contigs), and (3) identification of draft genomes by binning contigs that originate from the same organism. Due to the extensive diversity of bacteria and archaea in most environmental samples [6, 7], the field of metagenomics has rapidly evolved to accurately delineate genomes in assembly results. Today, microbiologists often exploit two essential properties of bacterial and archaeal genomes to improve the “binning” step: (1) k-mer frequencies that are somewhat preserved throughout a single microbial genome [8], to identify contigs that likely originate from the same genome [9], and (2) a set of genes that occur in the vast majority of bacterial genomes as a single copy, to estimate the level of completion and contamination of genome bins [10, 11, 12]. These properties, along with differential coverage of contigs across multiple samples when such data exist, are routinely used to identify coherent microbial draft genomes in metagenomic assemblies [13, 14, 15, 16].

On the other hand, researchers who study eukaryotic genomes generally focus on the recovery of a single organism, which, in most cases, simplifies the identification of the target genome in assembly results. However, sequences of bacterial origin can contaminate eukaryotic genome assembly results due to their occurrence in samples [17, 18], DNA extraction kits [19], or laboratory environments [20, 21]. One of the major challenges of working with eukaryotic genomes is the extent of repeat regions that complicate the assembly process [22]. To optimize the assembly, researchers often employ multiple library preparations for sequencing [23, 24], which may increase the potential sources of post-DNA extraction contamination. Contaminants in assembly results can eventually contaminate public databases [25], and impair scientific findings [26]. The detection and removal of contaminants poses a major bioinformatics challenge. To identify undesired contigs in a genomic assembly, scientists can simply compare their assembly results to public sequence databases for positive hits to unexpected taxa [23], use k-mer coverage plots to identify distinct genomes [27], or employ scatter plots to partition contigs based on their GC-content and coverage [28]. However, advanced solutions developed for accurate identification of microbial genomes in complex metagenomic assemblies can leverage these approaches further, and offer enhanced curation options for eukaryotic assemblies.

The first release of a tardigrade genome by Boothby et al. [29] demonstrates a striking example of the importance of careful screening for contaminants in eukaryotic genome assemblies. Tardigrades are microscopic animals occurring in a wide range of ecosystems and they exhibit extended capabilities to survive in harsh conditions that would be fatal to most animals [30, 31, 32, 33]. Boothby and his colleagues generated a composite DNA sequencing dataset from a culture of the tardigrade *Hypsibius dujardini* by exploiting some of the best practices of high-throughput sequencing available today [29]. In their assembled tardigrade genome, the authors detected a large number of genes originating from bacteria, making up approximately one-sixth of the gene pool, and suggested that horizontal gene transfers

(HGTs) could explain the unique ability of tardigrades to withstand extreme ranges of temperature, pressure, and radiation. However, Koutsovoulos et al.'s subsequent analysis of Boothby et al.'s assembly suggested that it contained extensive bacterial contamination, casting doubt on the extended HGT hypothesis [34]. By applying two-dimensional scatterplots on their own assembly results (which were also contaminated with bacterial sequences), Koutsovoulos et al. reported a curated draft genome of *H. dujardini*.

Here we re-analyzed the raw sequencing data generated by Boothby et al. [29] and Koutsovoulos et al. [34] using *anvi'o*, an analysis and visualization platform originally designed for the identification and assessment of bacterial genomes in metagenomic assemblies [16]. In our analysis, we relied on bacterial single-copy genes to assess the occurrence of bacterial genomes in assembly results, used k-mer frequencies to organize contigs, combined all sequencing data for each library preparation method from both groups into a single display, and overlaid RNA-Seq data (courtesy of Itai Yanai) over contigs to confirm the origin of contigs.

Material and methods

Genome assemblies, and raw sequencing data for DNA and RNA. Boothby et al. constructed three paired-end Illumina libraries (insert sizes of 0.3, 0.5 and 0.8 kbp) for 2 x 100 paired-end sequencing on a HiSeq2000 and six single-end long-read libraries (five Illumina Moleculo libraries sequenced by the Illumina 'long read' DNA sequencing service, and one PacBio SMRT library sequenced using the P6-C4 chemistry and a 1 X 240 movie), which altogether provided a co-assembly of 252.5 Mbp [29]. The tardigrade genome released by Boothby et al. [29], along with the nine sequencing data used for its assembly, are available at <http://weatherby.genetics.utah.edu/seq-transf>. Independently, Koutsovoulos et al. generated a 0.3 kbp insert library and a 1.1 kbp insert mate-pair library for 2 x 100 paired end sequencing on a HiSeq2000 that provided a co-assembly of 185.8 Mbp [34]. These authors subsequently curated a 135 Mbp draft genome by removing potential bacterial contamination [34]. The tardigrade raw assembly and curated draft genome released by Koutsovoulos et al. [34] are available at <http://badger.bio.ed.ac.uk/H.dujardini>, and their two sequencing datasets are available from the ENA, under study accession PRJEB11910. Itai Yanai (Technion - Israel Institute of Technology, <http://yanailab.technion.ac.il/>) graciously provided RNA-seq data generated from a *H. dujardini* culture, which will be available under the accession ID accession GSE70185 upon their publication. Quality filtering and read mapping. We used *illumina-utils* [35] for quality filtering of short Illumina reads using 'iu-filter-quality-minoche' script with default parameters, which implements the quality filtering described by Minoche et al. [36]. Bowtie2 v2.2.4 [37] with default parameters mapped all reads to assemblies. We used samtools v1.2 [38] to generate BAM files from mapping results.

Processing of contigs, visualization and genome binning. We processed BAM files and raw genome assemblies using *anvi'o* v1.2.2, generated *anvi'o* contig databases, profiled BAM files, and merged resulting profiles using default parameters and following the metagenomic workflow outlined in [16]. In addition,

we mapped and profiled the RNA-seq data to identify scaffolds with transcriptional activity, and exported the table for proportion of each scaffold covered by transcripts using `anvi'o` script `'get-db-table-as-matrix'`. We used the supplementary material published by Boothby et al. [29] (“Dataset S1” in the original publication) to identify scaffolds with proposed HGTs. We included the RNA-seq results and scaffolds with HGTs into our visualization as an additional data file. The URL http://merenlab.org/data/2016_Delmont_et_al_Tardigrade/ reports `anvi'o` files to regenerate Figure 1 and Figure 2, our curation of the tardigrade genome from Boothby et al.’s assembly (which is also available in NCBI via the bioproject ID PRJNA309530), and the FASTA files for bacterial genomes we identified in the Boothby et al. and Koutsovoulos et al. assemblies. To finalize the `anvi'o` generated SVG files for publication, we used Inkscape v0.91 (available from <https://inkscape.org/>).

Predicting number of bacterial genomes. To estimate the number of bacterial genomes in a given collection of scaffolds in a raw assembly or in a curated genome bin, and to visualize the distribution of HMM hits for each bacterial single-copy gene, we used the `anvi'o` script `'gen-stats-for-single-copy-genes'`, which reports the most frequent number in the list of number of hits per single-copy gene as the estimated number of bacterial genomes in a collection of scaffolds. The script uses HMMer v3.1b2 [39] to search for Hidden Markov Profiles (HMMs) of 139 bacterial single-copy genes identified by Campbell et al [11], and the R library `'ggplot'` v1.0.0 [40, 41] to plot results.

Taxonomical and functional annotation of bacterial genomes. After binning, we uploaded bacterial draft genomes recovered from the assembly into the RAST server [42], and used the RAST best taxonomic hits and FigFams to infer the taxonomy of genome bins and functions they harbor.

Results and discussion

Boothby et al. generated sequencing data from a tardigrade culture using three short read (Illumina) and six long read (Moleculo and PacBio) libraries, which altogether provided a co-assembly of 252.5 Mbp [29]. Using this assembly without any curation, authors suggested that 6,663 genes were entered into the tardigrade genome through HGTs. Independently, Koutsovoulos et al. generated sequencing data from another tardigrade culture using two short read Illumina libraries that provided a co-assembly of 185.8 Mbp, from which they could curate a 135 Mbp tardigrade draft genome by removing potential bacterial contamination using two-dimensional scatterplots of scaffolds with respect to their GC-content and coverage [34].

A holistic view of the data

The use of multiple library preparations and sequencing strategies is likely to result in more optimal assembly results [24]. Hence, we focused on the scaffolds generated by Boothby et al. [29] as a foundation to maximize the recovery of the tardigrade

genome. To provide a holistic understanding of the composite sequencing data generated by the two teams, we mapped the raw data from the nine DNA sequencing libraries from Boothby et al., and the two Illumina libraries from Koutsovoulos et al. [34] on this assembly. Anvi'o generated a hierarchical clustering of scaffolds by combining the tetra-nucleotide frequency and coverage of each scaffold across the 11 DNA sequencing libraries [16]. Besides visualizing the coverage of each scaffold in each sample, we highlighted scaffolds with HGTs identified by Boothby et al. on the resulting organization of scaffolds, and visualized RNA-seq mapping results. Figure 1 displays the anvi'o merged profile that represents all this information in a single display.

A larger draft genome for *H. dujardini*

Through the anvi'o interactive interface we selected 14,961 scaffolds from the Boothby et al. assembly that recruited large number of short-reads in a consistent manner (Fig. 1). This 182.2 Mbp selection with consistent coverage (1 in Fig. 1) represents our curation of the tardigrade assembly by Boothby et al. The remaining 7,535 scaffolds, which total about 70 Mbp of the assembly, harbored 96.1% of HGTs identified by Boothby et al. These scaffolds recruited only 0.05% of the reads from the RNA-Seq data, highlighting the extent of contamination in the original assembly. This finding is in agreement with Koutsovoulos et al.'s findings; however, our curated draft genome is 47 Mbp larger than the draft genome released by Koutsovoulos et al. [34]. The portion of scaffolds covered by RNA-Seq data suggests that the additional 47 Mbp still originate from the tardigrade genome. Thus, our selection is likely to be a more complete draft genome for *H. dujardini* than that of Koutsovoulos et al., most probably due to Boothby et al.'s inclusion of longer reads.

The origin of bacterial contamination

Our mapping results indicate the presence of non-target sequences in the assembly that recruit reads only from long-read libraries. One interpretation could be that most of the contamination in Boothby et al.'s assembly originated from Moleculo libraries, post DNA-extraction (Fig. 1). However, a recent study shows that the majority of long reads from Moleculo libraries originated from low-abundance organisms in samples [43], while another study suggests relatively more sequencing bias in Moleculo library preparation results [44]. Therefore another interpretation of the mapping results can be that the bacterial contaminants were present in the sample in low abundances pre-DNA extraction, and individual Moleculo library preparations resulted in long reads originating from different parts of this rare community. Regardless, long reads considerably improved Boothby et al.'s assembly, which resulted in a larger tardigrade genome following the removal of non-target sequences. While these results reiterate that the use of long-read libraries is essential to generate more comprehensive assemblies, they also suggest that extra care should be taken to better mitigate the presence of non-target sequences in assembly results when long-read libraries are used for sequencing.

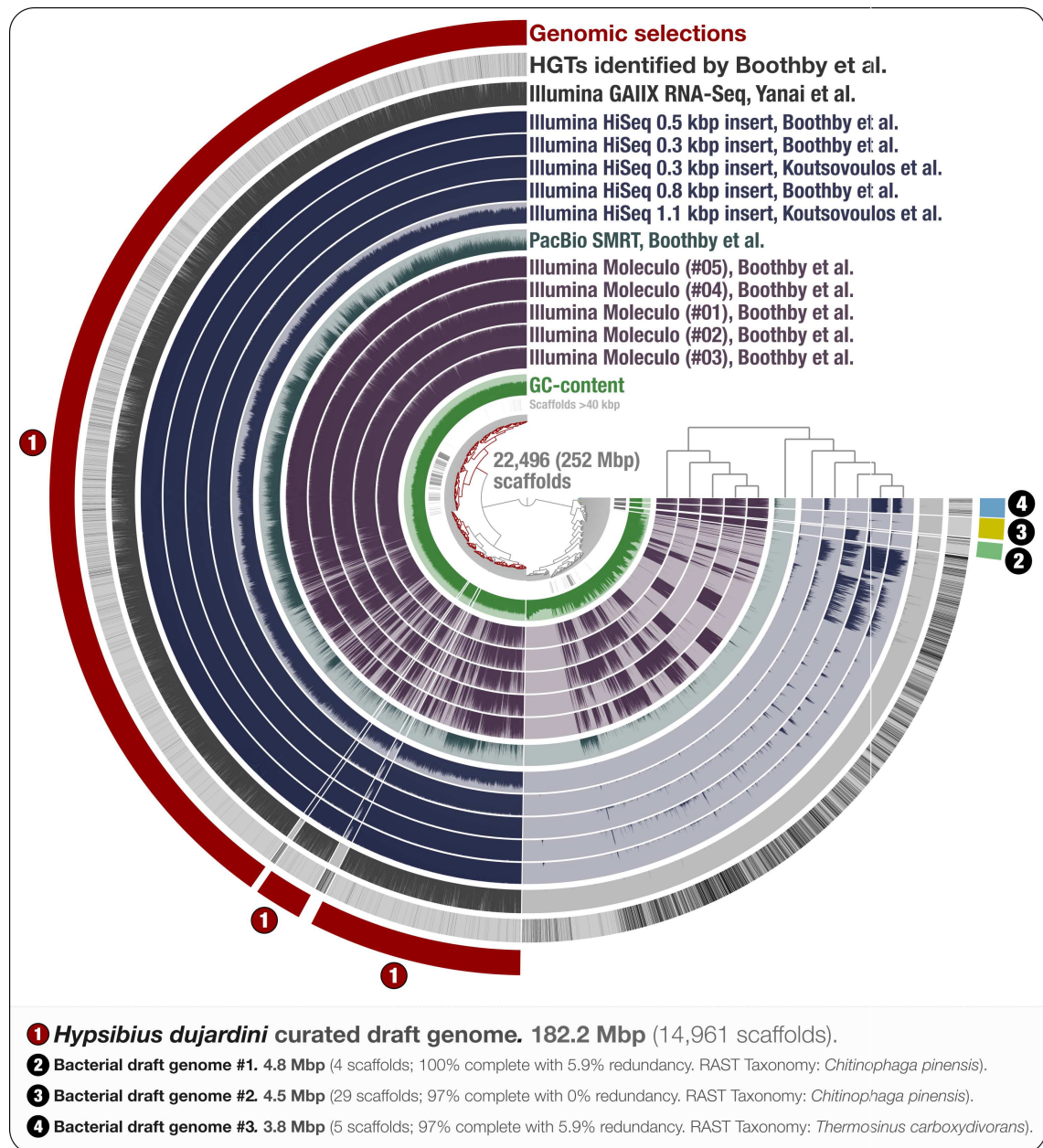


Figure 1 Holistic assessment of the tardigrade genome release from Boothby et al. (2015). Dendrogram in the center organizes scaffolds based on sequence composition and coverage values in data from 11 DNA libraries. Scaffolds larger than 40 kbp were split into sections of 20 kbp for visualization purposes. 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).

We identified three near-complete bacterial genomes affiliated to *Chitinophaga* and *Thermosinus* in Boothby et al.'s assembly (Fig. 1). Surprisingly, Boothby et al. identified only a small portion of these complete bacterial genomes as sources of

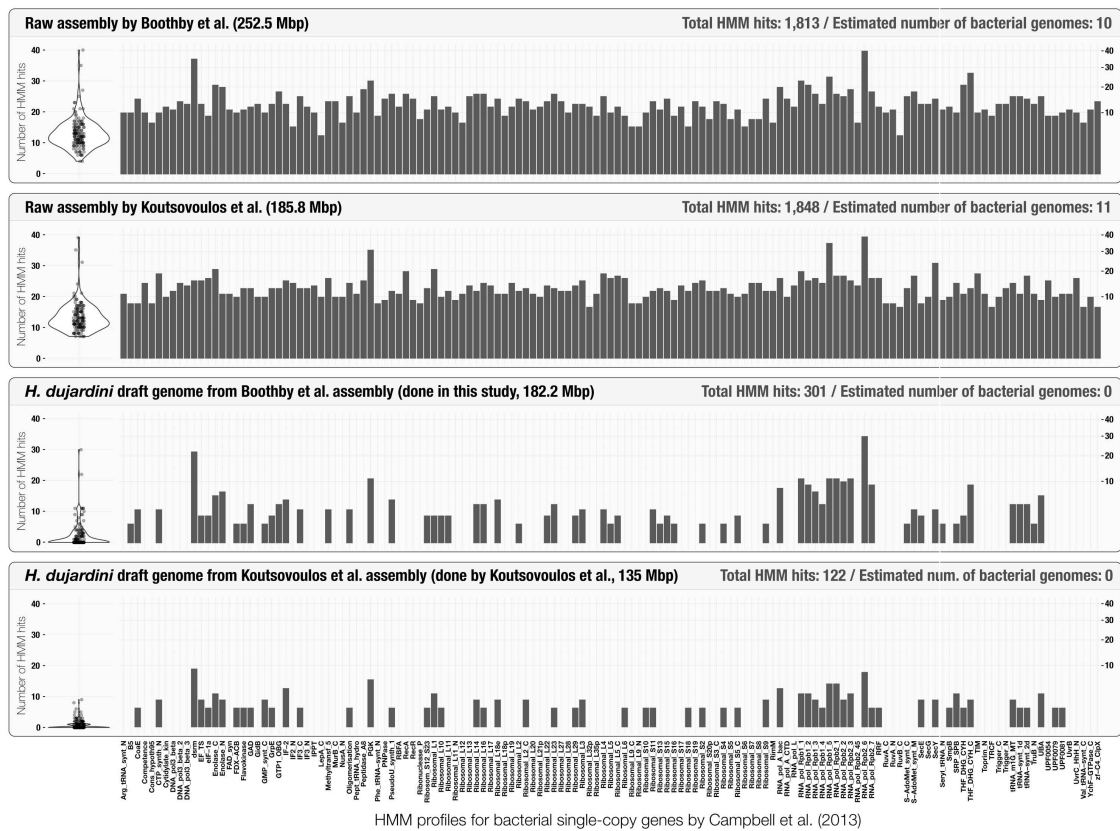
HGTs while applying a metric specifically designed to detect foreign DNA in eukaryotic genomes. For instance, none of the 4,459 genes in bacterial draft genome 2 (selection 3 in Fig. 1) were reported in Boothby et al.'s findings as HGTs. Although this falls outside of the scope of our study, this oddity may indicate a potential flaw in metrics commonly used to quantify foreign DNA in eukaryotic genomes. We also processed and visualized the raw assembly from Koutsovoulos et al. [34] using *anvi'o* (Fig. S1) and recovered eight bacterial genomes, however, we found no taxonomical overlap between high-completion bacterial genomes from the two sequencing projects (Table S1).

Interestingly, one bacterial genome (selection 2 in Fig. 1) was detected in DNA libraries from both groups, as well as in the RNA-seq data, suggesting that the related bacterial population was in all samples prior to the DNA/RNA extraction step. This genome is affiliated to *Chitinophaga*, and harbors genes coding for chitin degradation and utilization (Table S2). Chitin occurs naturally in the feeding apparatus of tardigrades [45], and might be a source of carbon for its microbial inhabitants. The genome also harbors genes coding for the biosynthesis of tryptophan, an essential amino acid for animals [46, 47], proteorhodopsin, host invasion and intracellular resistance, dormancy and sporulation, and oxidative stress. Although this genome may belong to a tardigrade symbiont, the generation of the data does not allow us to rule out the possibility that it may be associated with the food source. Nevertheless, this finding suggests that there may be cases where non-target genomes in an assembly can provide clues about the lifestyle of a given host.

Best practices to assess bacterial contamination

Initial assessment of the occurrence of bacterial single-copy genes in eukaryotic assemblies can provide a quick estimation of the number of bacterial genomes that occur in assembly results. The use of bacterial single-copy genes can give much more accurate representation of potential bacterial contamination than screening for 16S rRNA genes alone, as they are less likely to be found in co-assembly results [48, 49]. Although Boothby et al. reported the lack of 16S rRNA genes in their assembly [29], *anvi'o* estimated that it contained at least 10 complete bacterial genomes (Fig. 2) using a bacterial single-copy gene collection [11]. This simple yet powerful step could identify cases of extensive contamination, and alert researchers to be diligent in identifying scaffolds originating from bacterial organisms. Figure 2 also summarizes the HMM hits in scaffolds found in curated tardigrade genomes from our analysis and Koutsovoulos et al.'s study. We observed that the average significance score for the remaining HMM hits for bacterial single-copy genes in curated genomes was 4.2 times lower in average compared to the HMM hits in assembly results (Table S3). The decrease in the significance scores, and the very similar patterns of occurrence of HMM hits between the two curation efforts suggest that some of the HMM profiles may not be specific enough to be identified only in bacteria.

Two-dimensional scatterplots have a long history of identifying distinct genomes in assembly results [50] and continue to be used for delineating microbial genomes in



HMM profiles for bacterial single-copy genes by Campbell et al. (2013)

Figure 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.

metagenomic assemblies [13, 51], as well as detecting contamination in eukaryotic assembly results [28]. Although scatterplots can describe the organization of contigs in assembly results, they suffer from limited number of dimensions they can display, and their inability to depict complex supporting data that can improve the identification of individual genomes. These limitations are particularly problematic in sequencing projects covering multiple sequencing libraries, where displaying mapping results from each library can help detecting sources of contaminants. Despite their successful applications, two dimensional scatter plots limit researchers to the use of simple characteristics of the data that can be represented on an axis (such as GC-content). In contrast, clustering scaffolds, and overlaying multiple layers of independent information produce more comprehensive 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. The need for

de novo reconstruction of microbial genomes from environmental samples through shotgun metagenomics data has given rise to advanced techniques and software platforms that can make sense of complex assemblies [52, 53, 14, 15, 16]. Our study demonstrates that these approaches can be effectively used in eukaryotic assembly projects for curation purposes.

Acknowledgements

We are grateful to Thomas C. Boothby, Georgios Koutsovoulos, Sujai Kumar, and their colleagues for making their data available and answering our questions. We thank Itai Yanai for providing us with the RNA-Seq data. We also thank Hilary G. Morrison for her invaluable suggestions. This work was supported by the **Frank R. Lillie Research Innovation Award**, and startup funds from the University of Chicago.

Availability of supporting data

The URL http://merenlab.org/data/2016_Delmont_et_al.Tardigrade/ provide access to all essential and supporting data.

Competing interests

The authors declare that they have no competing interests.

Author's contributions

TOD and AME conceived the study, performed the data analyses, and wrote the manuscript.

Acknowledgements

We are grateful to Thomas C. Boothby, Georgios Koutsovoulos, Sujai Kumar, and their colleagues for making their data available and answering our questions. We thank Itai Yanai for providing us with the RNA-Seq data. We also thank Hilary G. Morrison for her invaluable suggestions. This work was supported by the **Frank R. Lillie Research Innovation Award**, and startup funds from the University of Chicago.

Author details

¹The Department of Medicine, The University of Chicago, Chicago, 60637 IL, USA. ²Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, 02543 MA, USA.

References

1. Loman, N.J., Pallen, M.J.: Twenty years of bacterial genome sequencing. *Nature Reviews Microbiology* **13**(12), 787–794 (2015). doi:10.1038/nrmicro3565
2. Venter, J.C., Adams, M.D., Myers, E.W., Li, P.W., Mural, R.J., Sutton, G.G., Smith, H.O., Yandell, M., Evans, C.A., Holt, R.A., Gocayne, J.D., Amanatides, P., Ballew, R.M., Huson, D.H., Wortman, J.R., Zhang, Q., Kodira, C.D., Zheng, X.H., Chen, L., Skupski, M., Subramanian, G., Thomas, P.D., Zhang, J., Gabor Miklos, G.L., Nelson, C., Broder, S., Clark, A.G., Nadeau, J., McKusick, V.A., Zinder, N., Levine, A.J., Roberts, R.J., Simon, M., Slayman, C., Hunkapiller, M., Bolanos, R., Delcher, A., Dew, I., Fasulo, D., Flanigan, M., Florea, L., Halpern, A., Hannenhalli, S., Kravitz, S., Levy, S., Mobarry, C., Reinert, K., Remington, K., Abu-Threideh, J., Beasley, E., Biddick, K., Bonazzi, V., Brandon, R., Cargill, M., Chandramouliswaran, I., Charlab, R., Chaturvedi, K., Deng, Z., Di Francesco, V., Dunn, P., Eilbeck, K., Evangelista, C., Gabrielian, A.E., Gan, W., Ge, W., Gong, F., Gu, Z., Guan, P., Heiman, T.J., Higgins, M.E., Ji, R.R., Ke, Z., Ketchum, K.A., Lai, Z., Lei, Y., Li, Z., Li, J., Liang, Y., Lin, X., Lu, F., Merkulov, G.V., Milshina, N., Moore, H.M., Naik, A.K., Narayan, V.A., Neelam, B., Nuskern, D., Rusch, D.B., Salzberg, S., Shao, W., Shue, B., Sun, J., Wang, Z., Wang, A., Wang, X., Wang, J., Wei, M., Wides, R., Xiao, C., Yan, C., Yao, A., Ye, J., Zhan, M., Zhang, W., Zhang, H., Zhao, Q., Zheng, L., Zhong, F., Zhong, W., Zhu, S., Zhao, S., Gilbert, D., Baumhueter, S., Spier, G., Carter, C., Cravchik, A., Woodage, T., Ali, F., An, H., Awe, A., Baldwin, D., Baden, H., Barnstead, M., Barrow, I., Beeson, K., Busam, D., Carver, A., Center, A., Cheng, M.L., Curry, L., Danaher, S., Davenport, L., Desilets, R., Dietz, S., Dodson, K., Doup, L., Ferriera, S., Garg, N., Gluecksmann, A., Hart, B., Haynes, J., Haynes, C., Heiner, C., Hladun, S., Hosten, D., Houck, J., Howland, T., Ibegwam, C., Johnson, J., Kalush, F., Kline, L.,

- Koduru, S., Love, A., Mann, F., May, D., McCawley, S., McIntosh, T., McMullen, I., Moy, M., Moy, L., Murphy, B., Nelson, K., Pfannkoch, C., Pratts, E., Puri, V., Qureshi, H., Reardon, M., Rodriguez, R., Rogers, Y.H., Romblad, D., Ruhfel, B., Scott, R., Sitter, C., Smallwood, M., Stewart, E., Strong, R., Suh, E., Thomas, R., Tint, N.N., Tse, S., Vech, C., Wang, G., Wetter, J., Williams, S., Williams, M., Windsor, S., Winn-Deen, E., Wolfe, K., Zaveri, J., Zaveri, K., Abril, J.F., Guigó, R., Campbell, M.J., Sjolander, K.V., Karlak, B., Kejarawal, A., Mi, H., Lazareva, B., Hatton, T., Narechania, A., Diemer, K., Muruganujan, A., Guo, N., Sato, S., Bafna, V., Istrail, S., Lippert, R., Schwartz, R., Walenz, B., Yooseph, S., Allen, D., Basu, A., Baxendale, J., Blick, L., Caminha, M., Carnes-Stine, J., Caulk, P., Chiang, Y.H., Coyne, M., Dahlke, C., Mays, A., Dombroski, M., Donnelly, M., Ely, D., Esparham, S., Fosler, C., Gire, H., Glanowski, S., Glasser, K., Glodek, A., Gorokhov, M., Graham, K., Gropman, B., Harris, M., Heil, J., Henderson, S., Hoover, J., Jennings, D., Jordan, C., Jordan, J., Kasha, J., Kagan, L., Kraft, C., Levitsky, A., Lewis, M., Liu, X., Lopez, J., Ma, D., Majoros, W., McDaniel, J., Murphy, S., Newman, M., Nguyen, T., Nguyen, N., Nodell, M., Pan, S., Peck, J., Peterson, M., Rowe, W., Sanders, R., Scott, J., Simpson, M., Smith, T., Sprague, A., Stockwell, T., Turner, R., Venter, E., Wang, M., Wen, M., Wu, D., Wu, M., Xia, A., Zandieh, A., Zhu, X.: The sequence of the human genome. *Science* (New York, N.Y.) **291**(5507), 1304–51 (2001). doi:[10.1126/science.1058040](https://doi.org/10.1126/science.1058040)
3. Brown, C.T., Hug, L.A., Thomas, B.C., Sharon, I., Castelle, C.J., Singh, A., Wilkins, M.J., Wrighton, K.C., Williams, K.H., Banfield, J.F.: Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* **523**(7559), 208–211 (2015). doi:[10.1038/nature14486](https://doi.org/10.1038/nature14486)
 4. Schleper, C., Jurgens, G., Jonuscheit, M.: Genomic studies of uncultivated archaea. *Nature reviews. Microbiology* **3**(6), 479–88 (2005). doi:[10.1038/nrmicro1159](https://doi.org/10.1038/nrmicro1159)
 5. Schloss, P.D., Handelsman, J.: Biotechnological prospects from metagenomics. *Current opinion in biotechnology* **14**(3), 303–10 (2003)
 6. Gans, J., Wolinsky, M., Dunbar, J.: Computational improvements reveal great bacterial diversity and high metal toxicity in soil. *Science* (New York, N.Y.) **309**(5739), 1387–90 (2005). doi:[10.1126/science.1112665](https://doi.org/10.1126/science.1112665)
 7. Rusch, D.B., Halpern, A.L., Sutton, G., Heidelberg, K.B., Williamson, S., Yooseph, S., Wu, D., Eisen, J.A., Hoffman, J.M., Remington, K., Beeson, K., Tran, B., Smith, H., Baden-Tillson, H., Stewart, C., Thorpe, J., Freeman, J., Andrews-Pfannkoch, C., Venter, J.E., Li, K., Kravitz, S., Heidelberg, J.F., Utterback, T., Rogers, Y.-H., Falcón, L.I., Souza, V., Bonilla-Rosso, G., Eguarte, L.E., Karl, D.M., Sathyendranath, S., Platt, T., Bermingham, E., Gallardo, V., Tamayo-Castillo, G., Ferrari, M.R., Strausberg, R.L., Neelson, K., Friedman, R., Frazier, M., Venter, J.C.: The Sorcerer II Global Ocean Sampling expedition: northwest Atlantic through eastern tropical Pacific. *PLoS biology* **5**(3), 77 (2007). doi:[10.1371/journal.pbio.0050077](https://doi.org/10.1371/journal.pbio.0050077)
 8. Pride, D.T., Meinersmann, R.J., Wassenaar, T.M., Blaser, M.J.: Evolutionary implications of microbial genome tetranucleotide frequency biases. *Genome research* **13**(2), 145–58 (2003). doi:[10.1101/gr.335003](https://doi.org/10.1101/gr.335003)
 9. Teeling, H., Meyerdierts, A., Bauer, M., Amann, R., Glöckner, F.O.: Application of tetranucleotide frequencies for the assignment of genomic fragments. *Environmental microbiology* **6**(9), 938–47 (2004). doi:[10.1111/j.1462-2920.2004.00624.x](https://doi.org/10.1111/j.1462-2920.2004.00624.x)
 10. Wu, M., Eisen, J.A.: A simple, fast, and accurate method of phylogenomic inference. *Genome Biology* **9**(10), 151 (2008). doi:[10.1186/gb-2008-9-10-r151](https://doi.org/10.1186/gb-2008-9-10-r151)
 11. Campbell, J.H., O'Donoghue, P., Campbell, A.G., Schwientek, P., Sczyrba, A., Woyke, T., Söll, D., Podar, M.: UGA is an additional glycine codon in uncultured SR1 bacteria from the human microbiota. *Proceedings of the National Academy of Sciences of the United States of America* **110**(14), 5540–5 (2013). doi:[10.1073/pnas.1303090110](https://doi.org/10.1073/pnas.1303090110)
 12. Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., Tyson, G.W.: CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome research* **25**(7), 1043–55 (2015). doi:[10.1101/gr.186072.114](https://doi.org/10.1101/gr.186072.114)
 13. Albertsen, M., Hugenholtz, P., Skarshewski, A., Nielsen, K.L., Tyson, G.W., Nielsen, P.H.: Genome sequences of rare, uncultured bacteria obtained by differential coverage binning of multiple metagenomes. *Nature biotechnology* **31**(6), 533–8 (2013). doi:[10.1038/nbt.2579](https://doi.org/10.1038/nbt.2579)
 14. Alneberg, J., Bjarnason, B.S., de Bruijn, I., Schirmer, M., Quick, J., Ijaz, U.Z., Lahti, L., Loman, N.J., Andersson, A.F., Quince, C.: Binning metagenomic contigs by coverage and composition. *Nature Methods* **11**(11), 1144–1146 (2014). doi:[10.1038/nmeth.3103](https://doi.org/10.1038/nmeth.3103)
 15. Kang, D.D., Froula, J., Egan, R., Wang, Z.: MetaBAT, an efficient tool for accurately reconstructing single genomes from complex microbial communities. *PeerJ* **3**, 1165 (2015). doi:[10.7717/peerj.1165](https://doi.org/10.7717/peerj.1165)
 16. Eren, A.M., Esen, Ö.C., Quince, C., Vineis, J.H., Morrison, H.G., Sogin, M.L., Delmont, T.O.: Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **3**, 1319 (2015). doi:[10.7717/peerj.1319](https://doi.org/10.7717/peerj.1319)
 17. Chapman, J.A., Kirkness, E.F., Simakov, O., Hampson, S.E., Mitros, T., Weinmaier, T., Rattei, T., Balasubramanian, P.G., Borman, J., Busam, D., Disbennett, K., Pfannkoch, C., Sumin, N., Sutton, G.G., Viswanathan, L.D., Walenz, B., Goodstein, D.M., Hellsten, U., Kawashima, T., Prochnik, S.E., Putnam, N.H., Shu, S., Blumberg, B., Dana, C.E., Gee, L., Kibler, D.F., Law, L., Lindgens, D., Martinez, D.E., Peng, J., Wigge, P.A., Bertulat, B., Guder, C., Nakamura, Y., Ozbek, S., Watanabe, H., Khalturin, K., Hemmrich, G., Franke, A., Augustin, R., Fraune, S., Hayakawa, E., Hayakawa, S., Hirose, M., Hwang, J.S., Ikeo, K., Nishimiya-Fujisawa, C., Ogura, A., Takahashi, T., Steinmetz, P.R.H., Zhang, X., Aufschnaiter, R., Eder, M.-K.,

- Gorny, A.-K., Salvenmoser, W., Heimberg, A.M., Wheeler, B.M., Peterson, K.J., Böttger, A., Tischler, P., Wolf, A., Gojobori, T., Remington, K.A., Strausberg, R.L., Venter, J.C., Technau, U., Hobmayer, B., Bosch, T.C.G., Holstein, T.W., Fujisawa, T., Bode, H.R., David, C.N., Rokhsar, D.S., Steele, R.E.: The dynamic genome of Hydra. *Nature* **464**(7288), 592–6 (2010). doi:[10.1038/nature08830](https://doi.org/10.1038/nature08830)
18. Artamonova, I.I., Mushegian, A.R.: Genome sequence analysis indicates that the model eukaryote *Nematostella vectensis* harbors bacterial consorts. *Applied and environmental microbiology* **79**(22), 6868–73 (2013). doi:[10.1128/AEM.01635-13](https://doi.org/10.1128/AEM.01635-13)
 19. Salter, S.J., Cox, M.J., Turek, E.M., Calus, S.T., Cookson, W.O., Moffatt, M.F., Turner, P., Parkhill, J., Loman, N.J., Walker, A.W.: Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biology* **12**(1), 87 (2014). doi:[10.1186/s12915-014-0087-z](https://doi.org/10.1186/s12915-014-0087-z)
 20. Laurence, M., Hatzis, C., Brash, D.E.: Common contaminants in next-generation sequencing that hinder discovery of low-abundance microbes. *PLoS one* **9**(5), 97876 (2014). doi:[10.1371/journal.pone.0097876](https://doi.org/10.1371/journal.pone.0097876)
 21. Strong, M.J., Xu, G., Morici, L., Splinter Bon-Durant, S., Baddoo, M., Lin, Z., Fewell, C., Taylor, C.M., Flemington, E.K.: Microbial contamination in next generation sequencing: implications for sequence-based analysis of clinical samples. *PLoS pathogens* **10**(11), 1004437 (2014). doi:[10.1371/journal.ppat.1004437](https://doi.org/10.1371/journal.ppat.1004437)
 22. Richard, G.-F., Kerrest, A., Dujon, B.: Comparative genomics and molecular dynamics of DNA repeats in eukaryotes. *Microbiology and molecular biology reviews : MMBR* **72**(4), 686–727 (2008). doi:[10.1128/MMBR.00011-08](https://doi.org/10.1128/MMBR.00011-08)
 23. Ekblom, R., Wolf, J.B.W.: A field guide to whole-genome sequencing, assembly and annotation. *Evolutionary Applications* **7**(9), (2014). doi:[10.1111/eva.12178](https://doi.org/10.1111/eva.12178)
 24. Gnerre, S., MacCallum, I., Przybylski, D., Ribeiro, F.J., Burton, J.N., Walker, B.J., Sharpe, T., Hall, G., Shea, T.P., Sykes, S., Berlin, A.M., Aird, D., Costello, M., Daza, R., Williams, L., Nicol, R., Gnirke, A., Nusbaum, C., Lander, E.S., Jaffe, D.B.: High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proceedings of the National Academy of Sciences* **108**(4), 1513–1518 (2010). doi:[10.1073/pnas.1017351108](https://doi.org/10.1073/pnas.1017351108)
 25. Merchant, S., Wood, D.E., Salzberg, S.L.: Unexpected cross-species contamination in genome sequencing projects. *PeerJ* **2**, 675 (2014). doi:[10.7717/peerj.675](https://doi.org/10.7717/peerj.675)
 26. Artamonova, I.I., Lappi, T., Zudina, L., Mushegian, A.R.: Prokaryotic genes in eukaryotic genome sequences: when to infer horizontal gene transfer and when to suspect an actual microbe. *Environmental Microbiology* **17**(7), 2203–2208 (2015). doi:[10.1111/1462-2920.12854](https://doi.org/10.1111/1462-2920.12854)
 27. Percudani, R.: A Microbial Metagenome (*Leucobacter* sp.) in *Caenorhabditis* Whole Genome Sequences. *Bioinformatics and biology insights* **7**, 55–72 (2013). doi:[10.4137/BBI.S11064](https://doi.org/10.4137/BBI.S11064)
 28. Kumar, S., Jones, M., Koutsovoulos, G., Clarke, M., Blaxter, M.: Blobology: exploring raw genome data for contaminants, symbionts and parasites using taxon-annotated GC-coverage plots. *Frontiers in genetics* **4**, 237 (2013). doi:[10.3389/fgene.2013.00237](https://doi.org/10.3389/fgene.2013.00237)
 29. Boothby, T.C., Tenlen, J.R., Smith, F.W., Wang, J.R., Patanella, K.A., Osborne Nishimura, E., Tintori, S.C., Li, Q., Jones, C.D., Yandell, M., Messina, D.N., Glasscock, J., Goldstein, B.: Evidence for extensive horizontal gene transfer from the draft genome of a tardigrade. *Proceedings of the National Academy of Sciences* **112**(52), 201510461 (2015). doi:[10.1073/pnas.1510461112](https://doi.org/10.1073/pnas.1510461112)
 30. Ramløv, H., Westh, P.: Cryptobiosis in the Eutardigrade *Adorybiotus* (*Richtersius*) *coronifer*: Tolerance to Alcohols, Temperature and de novo Protein Synthesis. *Zoologischer Anzeiger - A Journal of Comparative Zoology* **240**(3-4), 517–523 (2001). doi:[10.1078/0044-5231-00062](https://doi.org/10.1078/0044-5231-00062)
 31. Jönsson, K.I., Harms-Ringdahl, M., Torudd, J.: Radiation tolerance in the eutardigrade *Richtersius coronifer*. *International journal of radiation biology* **81**(9), 649–56 (2005). doi:[10.1080/09553000500368453](https://doi.org/10.1080/09553000500368453)
 32. Jönsson, K.I., Rabbow, E., Schill, R.O., Harms-Ringdahl, M., Rettberg, P.: Tardigrades survive exposure to space in low Earth orbit. *Current biology : CB* **18**(17), 729–731 (2008). doi:[10.1016/j.cub.2008.06.048](https://doi.org/10.1016/j.cub.2008.06.048)
 33. Horikawa, D.D., Cumbers, J., Sakakibara, I., Rogoff, D., Leuko, S., Harnoto, R., Arakawa, K., Katayama, T., Kunieda, T., Toyoda, A., Fujiyama, A., Rothschild, L.J.: Analysis of DNA repair and protection in the Tardigrade *Ramazzottius varieornatus* and *Hypsibius dujardini* after exposure to UVC radiation. *PLoS one* **8**(6), 64793 (2013). doi:[10.1371/journal.pone.0064793](https://doi.org/10.1371/journal.pone.0064793)
 34. Koutsovoulos, G., Kumar, S., Laetsch, D.R., Stevens, L., Daub, J., Conlon, C., Maroon, H., Thomas, F., Aboobaker, A., Blaxter, M.: The genome of the tardigrade *Hypsibius dujardini*. Technical report (dec 2015). doi:[10.1101/033464](https://doi.org/10.1101/033464). <http://biorxiv.org/content/early/2015/12/13/033464.abstract>
 35. Eren, A.M., Vineis, J.H., Morrison, H.G., Sogin, M.L.: A Filtering Method to Generate High Quality Short Reads Using Illumina Paired-End Technology. *PLoS ONE* **8**(6), 66643 (2013). doi:[10.1371/journal.pone.0066643](https://doi.org/10.1371/journal.pone.0066643)
 36. Minoche, A.E., Dohm, J.C., Himmelbauer, H.: Evaluation of genomic high-throughput sequencing data generated on Illumina HiSeq and genome analyzer systems. *Genome biology* **12**(11), 112 (2011). doi:[10.1186/gb-2011-12-11-r112](https://doi.org/10.1186/gb-2011-12-11-r112)

37. Langmead, B., Salzberg, S.L.: Fast gapped-read alignment with Bowtie 2. *Nature methods* **9**(4), 357–9 (2012). doi:[10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923)
38. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R.: The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)* **25**(16), 2078–9 (2009). doi:[10.1093/bioinformatics/btp352](https://doi.org/10.1093/bioinformatics/btp352)
39. Eddy, S.R.: Accelerated Profile HMM Searches. *PLoS computational biology* **7**(10), 1002195 (2011). doi:[10.1371/journal.pcbi.1002195](https://doi.org/10.1371/journal.pcbi.1002195)
40. R Development Core Team, R.: R: A Language and Environment for Statistical Computing (2011). doi:[10.1007/978-3-540-74686-7](https://doi.org/10.1007/978-3-540-74686-7). <http://www.r-project.org>
41. Ginestet, C.: ggplot2: Elegant Graphics for Data Analysis. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **174**(1), 245–246 (2011). doi:[10.1111/j.1467-985X.2010.00676](https://doi.org/10.1111/j.1467-985X.2010.00676)
42. Aziz, R.K., Bartels, D., Best, A.A., DeJongh, M., Disz, T., Edwards, R.A., Formsma, K., Gerdes, S., Glass, E.M., Kubal, M., Meyer, F., Olsen, G.J., Olson, R., Osterman, A.L., Overbeek, R.A., McNeil, L.K., Paarmann, D., Paczian, T., Parrello, B., Pusch, G.D., Reich, C., Stevens, R., Vassieva, O., Vonstein, V., Wilke, A., Zagnitko, O.: The RAST Server: rapid annotations using subsystems technology. *BMC genomics* **9**, 75 (2008). doi:[10.1186/1471-2164-9-75](https://doi.org/10.1186/1471-2164-9-75)
43. Sharon, I., Kertesz, M., Hug, L.A., Pushkarev, D., Blauwkamp, T.A., Castelle, C.J., Amirebrahimi, M., Thomas, B.C., Burstein, D., Tringe, S.G., Williams, K.H., Banfield, J.: Accurate, multi-kb reads resolve complex populations and detect rare microorganisms. *Genome Research*, 183012–114 (2015). doi:[10.1101/gr.183012.114](https://doi.org/10.1101/gr.183012.114)
44. Kuleshov, V., Jiang, C., Zhou, W., Jahanbani, F., Batzoglou, S., Snyder, M.: Synthetic long-read sequencing reveals intraspecies diversity in the human microbiome. *Nature Biotechnology* **34**(1), 64–69 (2015). doi:[10.1038/nbt.3416](https://doi.org/10.1038/nbt.3416)
45. Guidetti, R., Bonifacio, A., Altiero, T., Bertolani, R., Rebecchi, L.: Distribution of Calcium and Chitin in the Tardigrade Feeding Apparatus in Relation to its Function and Morphology. *Integrative and comparative biology* **55**(2), 241–52 (2015). doi:[10.1093/icb/icv008](https://doi.org/10.1093/icb/icv008)
46. Crawford, I.P.: Evolution of a biosynthetic pathway: the tryptophan paradigm. *Annual review of microbiology* **43**, 567–600 (1989). doi:[10.1146/annurev.mi.43.100189.003031](https://doi.org/10.1146/annurev.mi.43.100189.003031)
47. Zelante, T., Iannitti, R.G., Cunha, C., De Luca, A., Giovannini, G., Pieraccini, G., Zecchi, R., D'Angelo, C., Massi-Benedetti, C., Fallarino, F., Carvalho, A., Puccetti, P., Romani, L.: Tryptophan catabolites from microbiota engage aryl hydrocarbon receptor and balance mucosal reactivity via interleukin-22. *Immunity* **39**(2), 372–85 (2013). doi:[10.1016/j.immuni.2013.08.003](https://doi.org/10.1016/j.immuni.2013.08.003)
48. Miller, C.S., Baker, B.J., Thomas, B.C., Singer, S.W., Banfield, J.F.: EMIRGE: reconstruction of full-length ribosomal genes from microbial community short read sequencing data. *Genome biology* **12**(5), 44 (2011). doi:[10.1186/gb-2011-12-5-r44](https://doi.org/10.1186/gb-2011-12-5-r44)
49. Delmont, T.O., Eren, A.M., Maccario, L., Prestat, E., Esen, Ö.C., Pelletier, E., Le Paslier, D., Simonet, P., Vogel, T.M.: Reconstructing rare soil microbial genomes using in situ enrichments and metagenomics. *Frontiers in microbiology* **6**, 358 (2015). doi:[10.3389/fmicb.2015.00358](https://doi.org/10.3389/fmicb.2015.00358)
50. Tyson, G.W., Chapman, J., Hugenholtz, P., Allen, E.E., Ram, R.J., Richardson, P.M., Solovyev, V.V., Rubin, E.M., Rokhsar, D.S., Banfield, J.F.: Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**(6978), 37–43 (2004). doi:[10.1038/nature02340](https://doi.org/10.1038/nature02340)
51. Cantor, M., Nordberg, H., Smirnova, T., Hess, M., Tringe, S., Dubchak, I.: Elviz – exploration of metagenome assemblies with an interactive visualization tool. *BMC Bioinformatics* **16**(1), 130 (2015). doi:[10.1186/s12859-015-0566-4](https://doi.org/10.1186/s12859-015-0566-4)
52. Wu, Y.-W., Tang, Y.-H., Tringe, S.G., Simmons, B.A., Singer, S.W.: MaxBin: an automated binning method to recover individual genomes from metagenomes using an expectation-maximization algorithm. *Microbiome* **2**(1), 26 (2014). doi:[10.1186/2049-2618-2-26](https://doi.org/10.1186/2049-2618-2-26)
53. Dick, G.J., Andersson, A.F., Baker, B.J., Simmons, S.L., Thomas, B.C., Yelton, A.P., Banfield, J.F.: Community-wide analysis of microbial genome sequence signatures. *Genome biology* **10**(8), 85 (2009). doi:[10.1186/gb-2009-10-8-r85](https://doi.org/10.1186/gb-2009-10-8-r85)