

AllTheBacteria - all bacterial genomes assembled, available and searchable

Martin Hunt^{1-4,*}, Leandro Lima^{1,*}, Daniel Anderson¹, Jane Hawkey⁵, Wei Shen^{1,6}, John Lees¹, Zamin Iqbal^{1,7,+}

¹European Molecular Biology Laboratory - European Bioinformatics Institute, Hinxton, UK

²Nuffield Department of Medicine, University of Oxford, Oxford, UK

³National Institute of Health Research Oxford Biomedical Research Centre, John Radcliffe Hospital, Headley Way, Oxford, UK

⁴Health Protection Research Unit in Healthcare Associated Infections and Antimicrobial Resistance, University of Oxford, Oxford, UK

⁵Department of Infectious Diseases, School of Translational Medicine, Monash University, Melbourne, Victoria 3004, Australia

⁶Institute for Viral Hepatitis, The Second Affiliated Hospital of Chongqing Medical University, China

⁷Milner Centre for Evolution, University of Bath, UK

*these authors contributed equally

⁺Corresponding author, email: zi245@bath.ac.uk

Abstract

The bacterial sequence data publicly available at the global DNA archives is a vast source of information on the evolution of bacteria and their mobile elements. However, most of it is either unassembled or inconsistently assembled and QC-ed. This makes it unsuitable for large-scale analyses, and inaccessible for most researchers to use. In 2021 Blackwell et al therefore released a uniformly assembled set of 661,405 genomes, consisting of all publicly available whole genome sequenced bacterial isolate data as of November 2018, along with various search indexes. In this study we extend that dataset up to August 2024, more than tripling the number of genomes. We also expand the scope, as we begin a global collaborative project to generate annotations for different species as desired by different research communities.

In this study we describe the project as of release 2024-08, comprising 2,440,377 assemblies (including the 661k dataset). All 2.4 million have been uniformly reprocessed for quality criteria and to give taxonomic abundance estimates with respect to the GTDB phylogeny. We also provide antimicrobial resistance (AMR) gene and mutation annotation via AMRFinderPlus. Using an evolution-informed compression approach, the full set of genomes is just 130Gb in batched xz archives. We also provide multiple search indexes and a method for alignment to the full dataset. Finally, we outline plans for future annotations to be provided in further releases.

Introduction

Bacteria are the dominant cellular organisms on the planet, responsible for the functioning of every biome. As sequencing technology improves and becomes more widely accessible, we are

seeing a rapid expansion in the breadth and depth of sequencing of the bacterial domain. These genomes bear the imprint of millions of years of evolution and constitute a priceless resource for the understanding of their biology, dynamics and the effect on the ecology of our entire planet.

Bacterial genomes evolve both through “vertical” inheritance, as parents fission into pairs of children, and through multiple modes of horizontal gene transfer including those mediated by viruses and mobile genetic elements such as plasmids and transposons. This has profound implications for their plasticity and for the flexibility of their genomes. Members of a single bacterial species can share as little as 50% of their genomes (the core genome), the rest being accessory content, present in only a fraction of the genomes of the species. This “optional extra” content consists of fleetingly present content carried by mobile elements typically purged by selection and therefore rarely observed in the population. It also includes valuable cargo providing vital adaptive traits, observed consistently at intermediate frequencies due to balancing selection. For those who seek to explore the fundamental biology of bacteria, and for those working on clinical microbiology and public health, it is of immense value to be able to study the diversity of bacterial genomes and the dynamics of the functional elements they contain.

Unfortunately, genomes available in the public domain are processed inconsistently or not at all, rendering their use for these purposes inaccessible to most researchers. Even when sequence assemblies are available, specific problems include assembly by a range of different tools and settings; variable quality control (QC); and since many are run together in single projects, there are batch effects caused by blocks of genomes all using the same assembly workflow. As a result, these data are not appropriate for large scale analyses, where uncorrected batch artefacts could masquerade as interesting biology when comparing groups. In order to address this for the community, Blackwell et al[1] set out to uniformly assemble, QC and analyse all bacterial isolate whole genome sequence (WGS) raw data available in the ENA as of November 2018. They released 639,981 high-quality assemblies, along with quality control information and fundamental genome-derived statistics – the most important of which was to check the taxonomic abundance within each putatively single isolate dataset to confirm the species label in the submitted ENA metadata, which is not necessarily sequence-derived. In the process they estimated that 8.1% of the species metadata tags in the ENA were incorrect. They also released multiple search indexes with the assemblies: for whole genome comparison (sourmash[2] and sketchlib[3]), and for k-mer search (COBS[4]).

Reflecting upon this initiative, the more successful aspects of the Blackwell dataset (abbreviated to “661k”) were as follows. It was, to our knowledge, the first uniformly assembled and rigorously QC-ed set of bacterial genomes that set out to encompass all sequenced bacteria. It included assemblies of over 300,000 genomes which had not previously been available (the raw data only had been available). The assemblies and search indexes allowed multiple other studies of plasmids[5, 6], bacterial adaptation[7, 8, 9, 10], and compression/indexing algorithms[11, 12, 13, 14, 15]. However, there were a few limitations. First, the raw data stored at the INSDC has more than doubled since then, and although we realise that keeping up with publicly deposited sequence data is a never-ending task, an update to the dataset would clearly be of great value. Second, the full set of assemblies we produced was almost 1 terabyte in size, even after compression, and the COBS indexes added a further 900Gb - this reduced the accessibility of the data. Third, we wanted to have the taxonomic abundance QC done based on community-supported GTDB[16]. Fourth, we had not provided further useful information on top of the assemblies: gene annotation, species-specific analyses of wide interest (e.g. serotyping, MLST), or built pan-genomes. However to provide all of this was beyond the capacity or expertise of our own research group – to do this properly and best serve the whole community, we realised that we should involve the research communities who focussed

on specific genera/species.

We therefore set up this project, named AllTheBacteria, aiming to update the 661k dataset and improve on the above limitations through a community-centric approach. We advertised the project on Twitter/X and the public microbiology bioinformatics Slack channel and gathered colleagues from across the world keen to work together to produce a valuable public resource.

This paper describes the methodology used for assembly, quality control, taxonomic information and initial search indexes. All software pipelines are open source with permissive licenses, available on GitHub. We also describe the communities which have joined the project and outline plans for future releases. In terms of the data volume reducing accessibility, Brinda et al recently addressed this issue with a general principle called phylogenetic compression[17] – batching data intelligently based on approximate phylogenetic similarity before compressing shrank the 661k assemblies to 20Gb and the indexes to 100Gb. We are able to follow the same approach with this larger dataset.

Methods

Dataset

We downloaded all paired Illumina bacterial isolate whole genome sequence raw sequence meta-data from the ENA, using the query [https://www.ebi.ac.uk/ena/portal/api/search?result=read_run&fields=ALL&query=tax_tree\(2\)&format=tsv](https://www.ebi.ac.uk/ena/portal/api/search?result=read_run&fields=ALL&query=tax_tree(2)&format=tsv). Samples were processed if they were not in the Blackwell 661k dataset, and had metadata “instrument_platform” = “ILLUMINA”, “library_strategy” == “WGS”, “library_source” = “GENOMCIC”, and “library_layout” = “PAIRED”. The samples were processed in two stages: the first (release 0.2) was from metadata downloaded on June 16th 2023, and then a second round of processing from metadata obtained on August 1st 2024 (incremental release 2024-08).

Genome assembly

The genome assembly pipeline used by Blackwell was based around v1.0.4 of Shovill (<https://github.com/tseemann/shovill>) which is a wrapper around Spades[18]. Here we refactored and updated the pipeline (https://github.com/leoisl/bacterial_assembly_pipeline), and used a marginally later version of Shovill (v1.1.0). The difference between v1.0.4 and v1.1.0 was minimal, and therefore there was no need to reassemble the 661k. That pipeline processed all new samples in release 0.2.

All samples in release 2024-08 were processed using a simple Python script (see <https://github.com/AllTheBacteria/AllTheBacteria/tree/main/reproducibility/All-samples/assembly>), again using Shovill v1.1.0, but also including extra stages. The script processes one sample, first downloading the reads, then running Sylph, Shovill, and finally removing contigs matching the human genome (as described later).

Taxonomic abundance estimation

Performing taxonomic analysis on isolate data is considerably simpler than on full metagenomic data - we wanted primarily to establish the major species, its relative abundance, and the nature of contaminants. We therefore ran some simulation experiments with mixtures of different species at different abundances (data not shown here) and determined that sylph[19] was more accurate, faster (~1 minute per sample) and required less RAM (10Gb of RAM for the whole of GTDB) than the tools we used for the 661k (Kraken/Bracken). We therefore used sylph version

0.5.1 with the pre-built GTDB r214 database (<https://storage.googleapis.com/sylph-stuff/v0.3-c200-gtdb-r214.sylpdb>) and default options.

A species call was made from the “Genome_file” column of the sylph output, using a lookup table generated with TaxonKit[20] using GTDB taxonomy data (<https://github.com/shenwei356/gtdb-taxdump>, v0.4.0). The reads from 3,252 samples resulted in no output from sylph, presumably because there were no matches to the reference database.

Human decontamination

After assembly, the contigs output by Shovill were matched to the human genome plus HLA sequences using nucmer from version 4.0.0rc1 of the MUMmer package[21]. We used the T2T CHM13 version 2 assembly (GCA_009914755.4) of the human genome[22, 23]. For HLA sequences, we used the file `hla_gen.fasta` from version 3.55.0 of the IPD-IMGT/HLA database[24, 25, 26]. Any contig that had a single match of at least 99% identity and 90% of its length was removed.

Assembly statistics

The program assembly-stats (<https://github.com/sanger-pathogens/assembly-stats>; git commit 7bdb58b) was run on each assembly to gather basic statistics. Assemblies with a total length of less than 100kbp or greater than 15Mbp were excluded. We found that 21 of the assemblies in the original 661k data set were longer than 15Mbp, and so were removed from our releases, meaning that 661,384 of the samples in AllTheBacteria originate from the 661k data set.

CheckM

CheckM2[27] version 1.0.1 was run on each assembly, using the default downloaded database `uniref100.KO.1.dmnd`. We ran “checkm2 predict” with options `--allmodels --database_path --lowmem uniref100.KO.1.dmnd`. 275 samples did not run to completion, stopping with the error message “No DIAMOND annotation was generated”. This suggests that the assemblies are of low quality, resulting in very few predicted proteins.

MiniPhy

All assembly FASTA files were compressed using MiniPhy[17] commit 7abe08c (this tool (<https://github.com/karel-brinda/MiniPhy>) has been renamed - it was called mof-compress in the original preprint), which uses intelligent batching of genomes to improve compression. The process has 2 steps: Divide the genomes into approximately equal-sized batches, typically done by species. In our case, the highest-abundance species for each sample was previously determined using sylph (see above), and a CSV file was created mapping the filename to species. Batches were auto-created using the `create_batches.py` script from the MiniPhy repository. MiniPhy was then run on each batch; internally it created an approximate phylogenetic tree and reordered the genomes for better compression. The output is one xz compressed archive per batch.

sketchlib.rust

The high-quality assemblies were sketched at k=14 using sketchlib.rust v0.1.0. This database allows sequence similarity search through computing a Jaccard index, either against all the contents, sparse queries returning k-nearest neighbours, below a given distance threshold, or

against a chosen subset of queries. We ran `sketchlib sketch -f 2kk.list.txt -k 14 -s 1000 -o 2kk_sketch --threads 32` to use a sketch size of 1000. The resulting .skd database of sketches is 4.1 GB, and .skm of metadata is 123 MB.

Antimicrobial resistance detection

Antimicrobial resistance determinants were identified using AMRFinderPlus[28] v3.12.8 on all assembly FASTA files with database version v2024-01-31.1. For some specific species (*Acinetobacter baumannii*, *Burkholderia cepacia*, *Burkholderia pseudomallei*, *Campylobacter jejuni*, *Campylobacter coli*, *Citrobacter freundii*, *Clostridioides difficile*, *Enterobacter cloacae*, *Enterobacter asburiae*, *Enterococcus faecalis*, *Enterococcus faecium*, *Enterococcus hirae*, *Escherichia*, *Shigella*, *Klebsiella aerogenes*, *Klebsiella oxytoca*, *Klebsiella pneumoniae*, *Klebsiella pneumoniae species complex*, *Neisseria gonorrhoeae*, *Neisseria meningitidis*, *Pseudomonas aeruginosa*, *Salmonella*, *Serratia marcescens*, *Staphylococcus aureus*, *Staphylococcus pseudintermedius*, *Streptococcus agalactiae*, *Streptococcus mitis*, *Streptococcus pneumoniae*, *Streptococcus pyogenes*, *Vibrio cholerae*, *Vibrio vulnificus*, *Vibrio parahaemolyticus*) we used the GTDB species assigned by sylph for the AMRFinderPlus `--organism` parameter, in accordance with the guidelines at <https://github.com/ncbi/amr/wiki/Running-AMRFinderPlus#--organism-option> (git commit 5f27bbe), thus incorporating known AMR-informative point mutations. This option was omitted for all other species.

Results

This project extends and builds on the 661k dataset, using the same genome assembly pipeline. Thus, we generated 1,778,993 new assemblies, giving a total of 2,440,377 assemblies when combined with the 661k dataset, along with associated taxonomic abundance estimates and quality statistics. We shifted from using the NCBI taxonomy in the 661k project, to using GTDB here, so reprocessed the sequence reads for all samples (including the 661k) in order to get consistent taxonomic estimates. For different use cases, we expect different levels of quality filtering might be needed, but we provide a file list of 2,346,079 “high-quality” assemblies that pass the following criteria: genome size between 100k and 15Mb, no more than 2000 contigs, N50 at least 5000, majority species at above 99% abundance (and the same majority species call for all INSDC sequencing runs from the same sample), CheckM2-completeness at least 90%, and CheckM2-contamination of no more than 5%.

As expected, the data is dominated by the species of high clinical interest - the top 10 species constitute 75% of the high-quality dataset. However, it contains 11,824 different species whereas the 661k contained just 7,003. A comparison of the number of species in the high-quality data set and in the 661k set is shown in Figure 1.

In order to make the data more accessible (i.e. download-able), it was important to compress the assemblies as efficiently as possible, while remaining lossless, and without requiring users to install any special software. Naively gzipping each assembly in its own fasta resulted in a disk usage of 3.9 Terabytes. Applying the MiniPhy tool (renamed, previously mof-compress) to intelligently batch, before compression with xz, reduced the disk use to 130Gb.

Although they are obviously not bacteria, we have also applied the same assembly process to all (illumina) sequenced archaea as of July 31st 2024 (n=815), and make the assemblies also available at OSF.

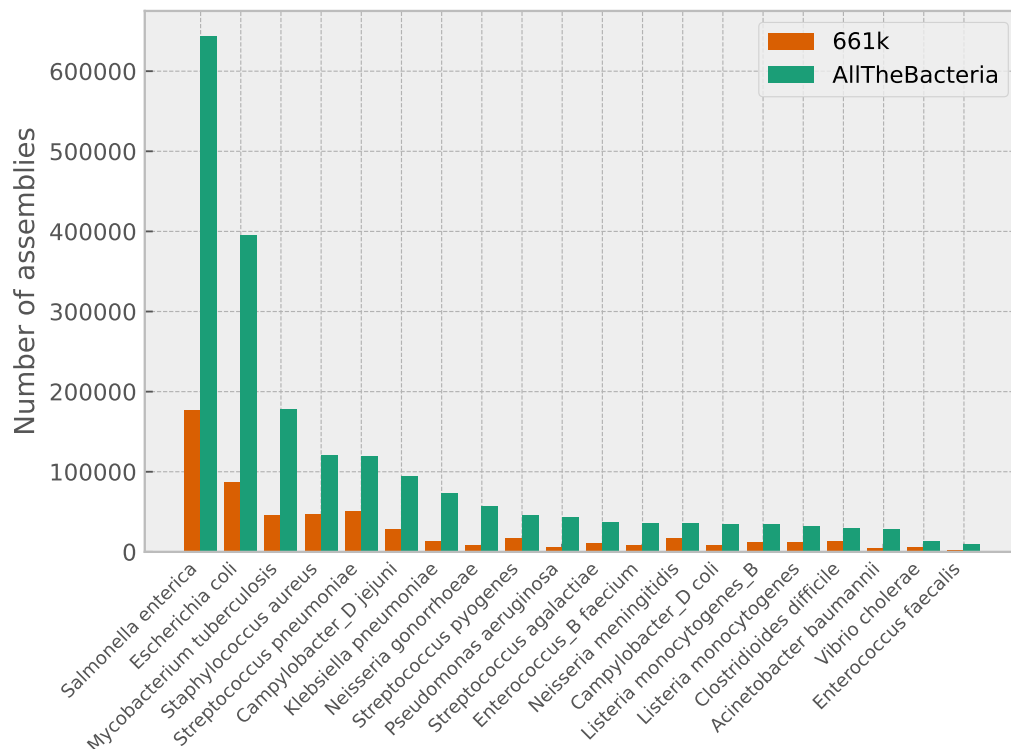


Figure 1: Assembly counts of the 20 most common species in the high quality AllTheBacteria data set, compared with their counts in the 661k set. Species names are from the GTDB.

Discussion

The goal of AllTheBacteria is to generate uniformly QC-ed, high quality genomes and annotations of all sequenced bacteria (and now we also include archaea), and build a community around it who create and share added-value analyses. We have added another 0.5 million assemblies since our first release (current total is 2,440,377), bringing us up to date with all INSDC bacterial and archaeal Illumina data up to August 2024. We now also provide AMR annotation (using AMRFinderPlus) and a sketchlib search index. Multiple other outputs are being generated: the closest to release are gene annotations with Bakta, MLST for multiple species, and defence system annotation using DefenceFinder. In addition, since the last release, we have extensively tested and published LexicMap[29] as a tool for BLAST-like query (min. length 500bp) alignment against the full dataset, with very low RAM requirements (1-2Gb) and extremely quickly (seconds for a rare gene with a few tens of thousand hits, to minutes for a 16s gene which requires alignment to almost every one of the 2 million genomes). The trade-off is that the index is large (around 3Tb); since this would be impractical to download, we do not provide it, but instead recommend downloading the assemblies (130Gb) and indexing them locally.

Future work includes further annotation, including of phage and plasmids, harmonisation of gene annotation to provide consistent id's within a species, and thereby construction of pangenomes (in the standard microbial genomics sense). Also, detection of genes/features of interest to specific research communities.

We want these data to be of use - please use them and publish with them. As our collaborative network continues to grow, we envisage generation of progressively more valuable analytic

outputs for the research community, as well as triggering innovation in search index methods.

Data Availability

Documentation for AllTheBacteria is available at <https://allthebacteria.readthedocs.io/en/latest/>. All data for AllTheBacteria are hosted on the Open Science Framework (OSF) here: <https://osf.io/xv7q9/>. The assembly pipelines for release 0.2 and 2024-08 are at https://github.com/leois1/bacterial_assembly_pipeline and <https://github.com/AllTheBacteria/AllTheBacteria/tree/main/reproducibility/All-samples/assembly>.

Author Contributions

Assembly [LL], taxonomic abundance analysis [SW, MH], compression of assemblies and COBs indexes using miniphy [SW,MH], AMR analysis [DA, JH], sketchlib [JL], all other analyses [MH], planning [LL, MH, SW, JL, ZI], paper writing [ZI, MH, JL, DA].

Acknowledgements

We would like to thank Karel Brinda for help with running MiniPhy. The authorship of this paper is currently very short, as the first phase of this project was originally dependent on the team at EBI/Bath to deliver the assemblies. However many people have volunteered to do future analyses, and their enthusiasm has buoyed us. We would like to thank, for their enthusiasm and probable future contributions: Nabil Fareed-Alikhan, Oliver Schwengers, Laura Carroll, Natacha Couto, Boas van der Putten, Kivumbi Mark Teferi, Sebastian Jaenicke, Conor Meehan, Gultekin Unal, Peter van Heusden, George Bouras, Adrian Cazares, Daniel Cazares, Wendy Figueroa, Michael Hall, Finlay Macguire, Matthew Croxson, Kate Baker, Nick Thomson, Kat Holt, Torsten Seemann and Jo Fothergill.

Funding

This work was supported by the National Institute for Health Research Health Protection Research Unit (NIHR HPRU) in Healthcare Associated Infections and Antimicrobial Resistance at Oxford University in partnership with the UK Health Security Agency (NIHR200915), and the NIHR Biomedical Research Centre, Oxford. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR, the Department of Health or the UK Health Security Agency. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

References

- [1] Grace A. Blackwell, Martin Hunt, Kerri M. Malone, Leandro Lima, Gal Horesh, Blaise T. F. Alako, Nicholas R. Thomson, and Zamin Iqbal. Exploring bacterial diversity via a curated and searchable snapshot of archived DNA sequences. *PLOS Biology*, 19(11):e3001421, November 2021.
- [2] N. Tessa Pierce, Luiz Irber, Taylor Reiter, Phillip Brooks, and C. Titus Brown. Large-scale sequence comparisons with sourmash. *F1000Research*, 8:1006, July 2019.

- [3] John A. Lees, Simon R. Harris, Gerry Tonkin-Hill, Rebecca A. Gladstone, Stephanie W. Lo, Jeffrey N. Weiser, Jukka Corander, Stephen D. Bentley, and Nicholas J. Croucher. Fast and flexible bacterial genomic epidemiology with PopPUNK. *Genome Research*, 29(2):304–316, February 2019.
- [4] Timo Bingmann, Phelim Bradley, Florian Gauger, and Zamin Iqbal. COBS: a Compact Bit-Sliced Signature Index. 2019.
- [5] Florent Lassalle, Salah Al-Shalali, Mukhtar Al-Hakimi, Elisabeth Njamkepo, Ismail Mahat Bashir, Matthew J. Dorman, Jean Rauzier, Grace A. Blackwell, Alyce Taylor-Brown, Mathew A. Beale, Adrián Cazares, Ali Abdullah Al-Somainy, Anas Al-Mahbashi, Khaled Almoayed, Mohammed Aldawla, Abdulalah Al-Harazi, Marie-Laure Quilici, François-Xavier Weill, Ghulam Dhabaan, and Nicholas R. Thomson. Genomic epidemiology reveals multidrug resistant plasmid spread between *Vibrio cholerae* lineages in Yemen. *Nature Microbiology*, 8(10):1787–1798, September 2023.
- [6] Ya Hu, Robert A. Moran, Grace A. Blackwell, Alan McNally, and Zhiyong Zong. Fine-Scale Reconstruction of the Evolution of FII-33 Multidrug Resistance Plasmids Enables High-Resolution Genomic Surveillance. *mSystems*, 7(1):e00831–21, February 2022.
- [7] Kevin O. Tamadonfar, Gisela Di Venzio, Jerome S. Pinkner, Karen W. Dodson, Vasilios Kalas, Maxwell I. Zimmerman, Jesus Bazan Villicana, Gregory R. Bowman, Mario F. Feldman, and Scott J. Hultgren. Structure–function correlates of fibrinogen binding by *Acinetobacter* adhesins critical in catheter-associated urinary tract infections. *Proceedings of the National Academy of Sciences*, 120(4):e2212694120, January 2023.
- [8] Lewis C. E. Mason, David R. Greig, Lauren A. Cowley, Sally R. Partridge, Elena Martinez, Grace A. Blackwell, Charlotte E. Chong, P. Malaka De Silva, Rebecca J. Bengtsson, Jenny L. Draper, Andrew N. Ginn, Indy Sandaradura, Eby M. Sim, Jonathan R. Iredell, Vitali Sintchenko, Danielle J. Ingle, Benjamin P. Howden, Sophie Lefèvre, Elisabeth Njamkepo, François-Xavier Weill, Pieter-Jan Ceyskens, Claire Jenkins, and Kate S. Baker. The evolution and international spread of extensively drug resistant *Shigella sonnei*. *Nature Communications*, 14(1):1983, April 2023.
- [9] Michael Biggel, Nadja Jessberger, Jasna Kovac, and Sophia Johler. Recent paradigm shifts in the perception of the role of *Bacillus thuringiensis* in foodborne disease. *Food Microbiology*, 105:104025, August 2022.
- [10] Tracy M Smith, Madison A Youngblom, John F Kernien, Mohamed A Mohamed, Sydney S Fry, Lindsey L Bohr, Tatum D Mortimer, Mary B O’Neill, and Caitlin S Pepperell. Rapid adaptation of a complex trait during experimental evolution of *Mycobacterium tuberculosis*. *eLife*, 11:e78454, June 2022.
- [11] Barış Ekim, Bonnie Berger, and Rayan Chikhi. Minimizer-space de Bruijn graphs: Whole-genome assembly of long reads in minutes on a personal computer. *Cell Systems*, 12(10):958–968.e6, October 2021.
- [12] Andrea Cracco and Alexandru I. Tomescu. Extremely fast construction and querying of compacted and colored de Bruijn graphs with GGCAT. *Genome Research*, page genome;gr.277615.122v2, May 2023.
- [13] Jamshed Khan, Marek Kokot, Sebastian Deorowicz, and Rob Patro. Scalable, ultra-fast, and low-memory construction of compacted de Bruijn graphs with Cuttlefish 2. *Genome Biology*, 23(1):190, September 2022.

- [14] Sebastian Deorowicz, Agnieszka Danek, and Heng Li. AGC: compact representation of assembled genomes with fast queries and updates. *Bioinformatics*, 39(3):btad097, March 2023.
- [15] Camille Marchet and Antoine Limasset. Scalable sequence database search using Partitioned Aggregated Bloom Comb-Trees. preprint, *Bioinformatics*, February 2022.
- [16] Donovan H Parks, Maria Chuvpina, Christian Rinke, Aaron J Mussig, Pierre-Alain Chaumeil, and Philip Hugenholtz. GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Research*, 50(D1):D785–D794, January 2022.
- [17] Karel Břinda, Leandro Lima, Simone Pignotti, Natalia Quinones-Olvera, Kamil Salikhov, Rayan Chikhi, Gregory Kucherov, Zamin Iqbal, and Michael Baym. Efficient and Robust Search of Microbial Genomes via Phylogenetic Compression. preprint, *Bioinformatics*, April 2023.
- [18] Anton Bankevich, Sergey Nurk, Dmitry Antipov, Alexey A. Gurevich, Mikhail Dvorkin, Alexander S. Kulikov, Valery M. Lesin, Sergey I. Nikolenko, Son Pham, Andrey D. Prjibelski, Alexey V. Pyshkin, Alexander V. Sirotkin, Nikolay Vyahhi, Glenn Tesler, Max A. Alekseyev, and Pavel A. Pevzner. SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. *Journal of Computational Biology*, 19(5):455–477, May 2012.
- [19] Jim Shaw and Yun William Yu. Metagenome profiling and containment estimation through abundance-corrected k-mer sketching with sylph. preprint, *Bioinformatics*, November 2023.
- [20] Wei Shen and Hong Ren. TaxonKit: A practical and efficient NCBI taxonomy toolkit. *Journal of Genetics and Genomics*, 48(9):844–850, September 2021.
- [21] Guillaume Marçais, Arthur L. Delcher, Adam M. Phillippy, Rachel Coston, Steven L. Salzberg, and Aleksey Zimin. MUMmer4: A fast and versatile genome alignment system. *PLOS Computational Biology*, 14(1):e1005944, January 2018.
- [22] Sergey Nurk, Sergey Koren, Arang Rhie, Mikko Rautiainen, Andrey V. Bzikadze, Alla Mikheenko, Mitchell R. Vollger, Nicolas Altemose, Lev Uralsky, Ariel Gershman, Sergey Aganezov, Savannah J. Hoyt, Mark Diekhans, Glennis A. Logsdon, Michael Alonge, Stylianos E. Antonarakis, Matthew Borchers, Gerard G. Bouffard, Shelise Y. Brooks, Gina V. Caldas, Nae-Chyun Chen, Haoyu Cheng, Chen-Shan Chin, William Chow, Leonardo G. De Lima, Philip C. Dishuck, Richard Durbin, Tatiana Dvorkina, Ian T. Fiddes, Giulio Formenti, Robert S. Fulton, Arkarachai Functammasan, Erik Garrison, Patrick G. S. Grady, Tina A. Graves-Lindsay, Ira M. Hall, Nancy F. Hansen, Gabrielle A. Hartley, Marina Haukness, Kerstin Howe, Michael W. Hunkapiller, Chirag Jain, Miten Jain, Erich D. Jarvis, Peter Kerpedjiev, Melanie Kirsche, Mikhail Kolmogorov, Jonas Korlach, Milinn Kremitzki, Heng Li, Valerie V. Maduro, Tobias Marschall, Ann M. McCartney, Jennifer McDaniel, Danny E. Miller, James C. Mullikin, Eugene W. Myers, Nathan D. Olson, Benedict Paten, Paul Peluso, Pavel A. Pevzner, David Porubsky, Tamara Potapova, Evgeny I. Rogaev, Jeffrey A. Rosenfeld, Steven L. Salzberg, Valerie A. Schneider, Fritz J. Sedlazeck, Kishwar Shafin, Colin J. Shew, Alaina Shumate, Ying Sims, Arrian F. A. Smit, Daniela C. Soto, Ivan Sović, Jessica M. Storer, Aaron Streets, Beth A. Sullivan, Françoise Thibaud-Nissen, James Torrance, Justin Wagner, Brian P. Walenz,

- Aaron Wenger, Jonathan M. D. Wood, Chunlin Xiao, Stephanie M. Yan, Alice C. Young, Samantha Zarate, Urvashi Surti, Rajiv C. McCoy, Megan Y. Dennis, Ivan A. Alexandrov, Jennifer L. Gerton, Rachel J. O'Neill, Winston Timp, Justin M. Zook, Michael C. Schatz, Evan E. Eichler, Karen H. Miga, and Adam M. Phillippy. The complete sequence of a human genome. *Science*, 376(6588):44–53, April 2022.
- [23] Arang Rhie, Sergey Nurk, Monika Cechova, Savannah J. Hoyt, Dylan J. Taylor, Nicolas Altemose, Paul W. Hook, Sergey Koren, Mikko Rautiainen, Ivan A. Alexandrov, Jamie Allen, Mobin Asri, Andrey V. Bzikadze, Nae-Chyun Chen, Chen-Shan Chin, Mark Diekhans, Paul Flicek, Giulio Formenti, Arkarachai Functammasan, Carlos Garcia Giron, Erik Garrison, Ariel Gershman, Jennifer L. Gerton, Patrick G. S. Grady, Andrea Guarracino, Leanne Haggerty, Reza Halabian, Nancy F. Hansen, Robert Harris, Gabrielle A. Hartley, William T. Harvey, Marina Haukness, Jakob Heinz, Thibaut Hourlier, Robert M. Hubley, Sarah E. Hunt, Stephen Hwang, Miten Jain, Rupesh K. Kesharwani, Alexandra P. Lewis, Heng Li, Glennis A. Logsdon, Julian K. Lucas, Wojciech Makalowski, Christopher Markovic, Fergal J. Martin, Ann M. Mc Cartney, Rajiv C. McCoy, Jennifer McDaniel, Brandy M. McNulty, Paul Medvedev, Alla Mikheenko, Katherine M. Munson, Terence D. Murphy, Hugh E. Olsen, Nathan D. Olson, Luis F. Paulin, David Porubsky, Tamara Potapova, Fedor Ryabov, Steven L. Salzberg, Michael E. G. Sauria, Fritz J. Sedlazeck, Kishwar Shafin, Valery A. Shepelev, Alaina Shumate, Jessica M. Storer, Likhitha Surapaneni, Angela M. Taravella Oill, Françoise Thibaud-Nissen, Winston Timp, Marta Tomaszewicz, Mitchell R. Vollger, Brian P. Walenz, Allison C. Watwood, Matthias H. Weissensteiner, Aaron M. Wenger, Melissa A. Wilson, Samantha Zarate, Yiming Zhu, Justin M. Zook, Evan E. Eichler, Rachel J. O'Neill, Michael C. Schatz, Karen H. Miga, Kateryna D. Makova, and Adam M. Phillippy. The complete sequence of a human Y chromosome. *Nature*, 621(7978):344–354, September 2023.
- [24] J. Robinson, A. Malik, P. Parham, J.G. Bodmer, and S.G.E. Marsh. IMGT/HLA Database – a sequence database for the human major histocompatibility complex. *Tissue Antigens*, 55(3):280–287, March 2000.
- [25] Dominic J Barker, Giuseppe Maccari, Xenia Georgiou, Michael A Cooper, Paul Flicek, James Robinson, and Steven G E Marsh. The IPD-IMGT/HLA Database. *Nucleic Acids Research*, 51(D1):D1053–D1060, January 2023.
- [26] James Robinson, Dominic J. Barker, and Steven G. E. Marsh. 25 years of the IPD-IMGT/HLA database. *HLA*, 103(6):e15549, June 2024.
- [27] Alex Chklovski, Donovan H. Parks, Ben J. Woodcroft, and Gene W. Tyson. CheckM2: a rapid, scalable and accurate tool for assessing microbial genome quality using machine learning. *Nature Methods*, 20(8):1203–1212, August 2023.
- [28] Michael Feldgarden, Vyacheslav Brover, Narjol Gonzalez-Escalona, Jonathan G. Frye, Julie Haendiges, Daniel H. Haft, Maria Hoffmann, James B. Pettengill, Arjun B. Prasad, Glenn E. Tillman, Gregory H. Tyson, and William Klimke. Amrfinderplus and the reference gene catalog facilitate examination of the genomic links among antimicrobial resistance, stress response, and virulence. *Scientific Reports*, 11(1):12728, June 2021.
- [29] Wei Shen and Zamin Iqbal. LexicMap: efficient sequence alignment against millions of prokaryotic genomes, August 2024.