

# The ATCC<sup>®</sup> Genome Portal: Authenticated Microbial Reference Genomes

Joseph R. Petrone, PhD; Nikhita P. Puthuveetil, MS; David A. Yarmosh, MS; Amy L. Reese, MS; Corina Tabron, MS; Jade Kirkland, BS; Kaitlyn Gaffney, MS; Noah Wax, MS; James Duncan, MS; Robert Marlow, BS; Stephen King, MS; Scott Nguyen, PhD; John Bagnoli, BS; Briana Benton, BS; Jonathan L. Jacobs, PhD | ATCC, Manassas, VA 20110

## Background

The ATCC<sup>®</sup> Genome Portal was developed as part of an initiative to produce high-quality reference genomes for the ATCC<sup>®</sup> microbial collection. Currently, the portal includes fully authenticated and annotated genome assemblies for over 4,500 microbes, encompassing 3,828 bacterial, 362 viral, 306 fungal, and 4 protist genomes. The ATCC<sup>®</sup> Genome Portal (AGP), which warehouses these data, is meticulously curated and links all the metadata and sequencing data to physical materials in ATCC's collection. As the sole ISO 9001-compliant reference database for ATCC<sup>®</sup> microbial strains, the portal ensures full end-to-end data provenance and authentication. Accessible for research purposes, the ATCC<sup>®</sup> Genome Portal and its data can be explored online (<https://genomes.atcc.org>) or through an authenticated REST-API.

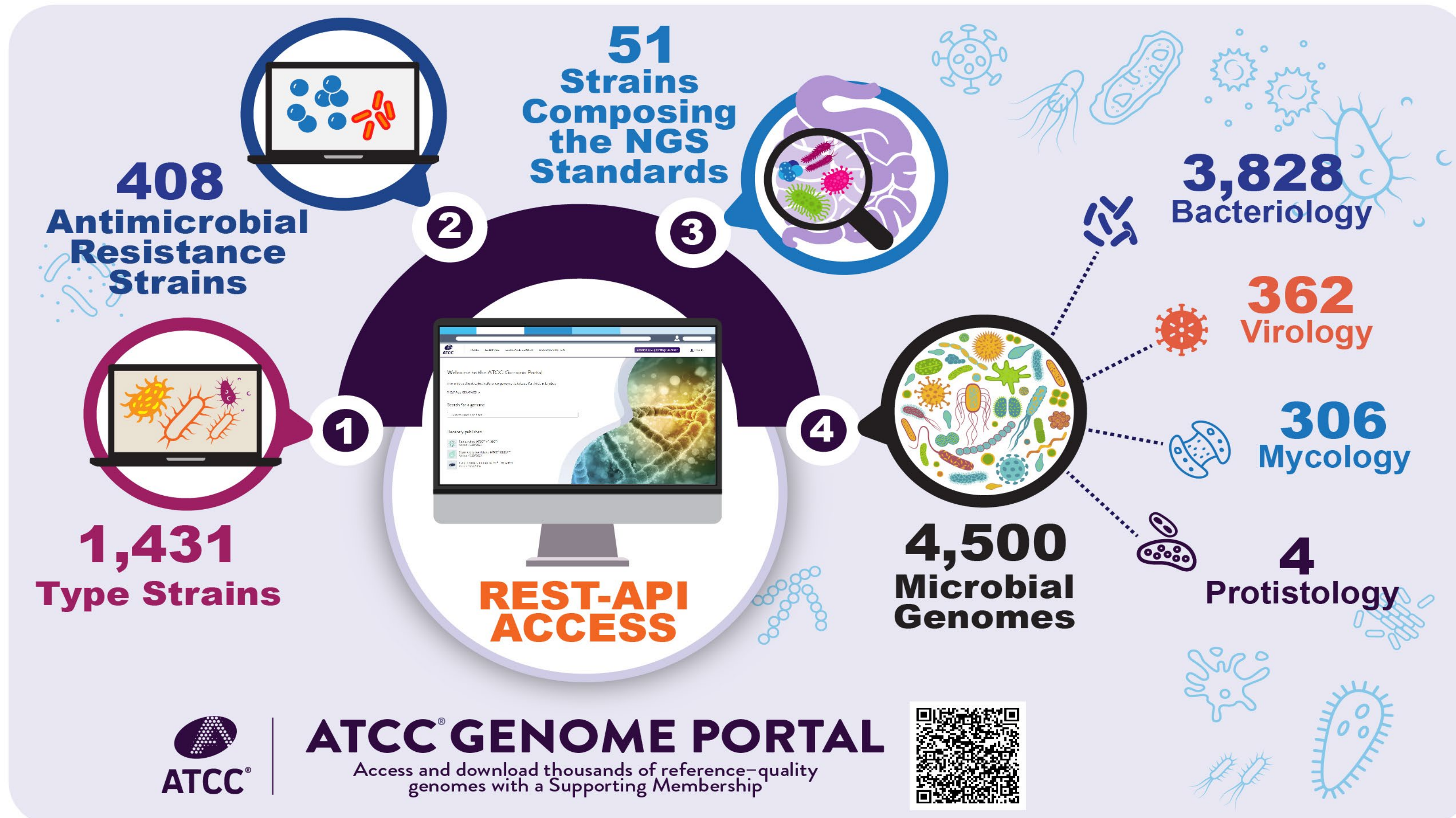


Figure 1: Reference genomes available through the ATCC<sup>®</sup> Genome Portal. The statistics are current as of May 1, 2024.

## AGP Supporting Membership

Table 1: ATCC<sup>®</sup> Genome Portal Supporting Membership plans

|  | Free Plan                   | Individual Member        | Research Member          | Institutional Member |
|--|-----------------------------|--------------------------|--------------------------|----------------------|
| View organism and genome metadata, assemblies, and annotations                     | ✓                           | ✓                        | ✓                        | ✓                    |
| Search for genes of interest   | ✓                           | ✓                        | ✓                        | ✓                    |
| Download genome assemblies and annotations   | Only for purchased products | ✓                        | ✓                        | ✓                    |
| Access the REST API  | Not available               | ✓                        | ✓                        | ✓                    |
| Compare your sequencing data to ATCC <sup>®</sup> genomes with Discrepancy Reports | Fee for each report         | 12 free reports per year | 60 free reports per year | Inquire for details  |
| Members with full access   | 0                           | 1                        | 5                        | Unlimited            |

## Methodology

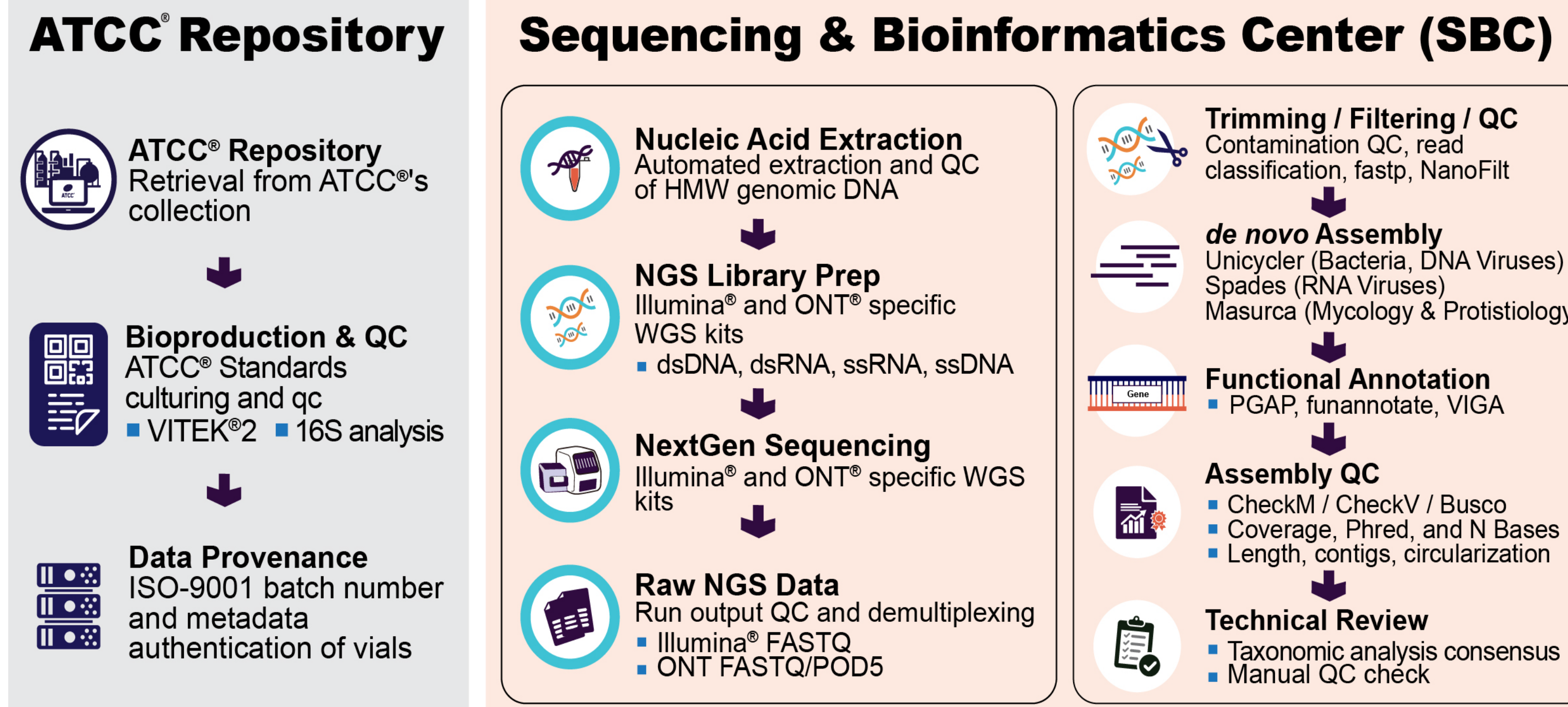


Figure 2: Standard pipeline for end-to-end genome portal publication.

## Authenticated ATCC<sup>®</sup> Genomes

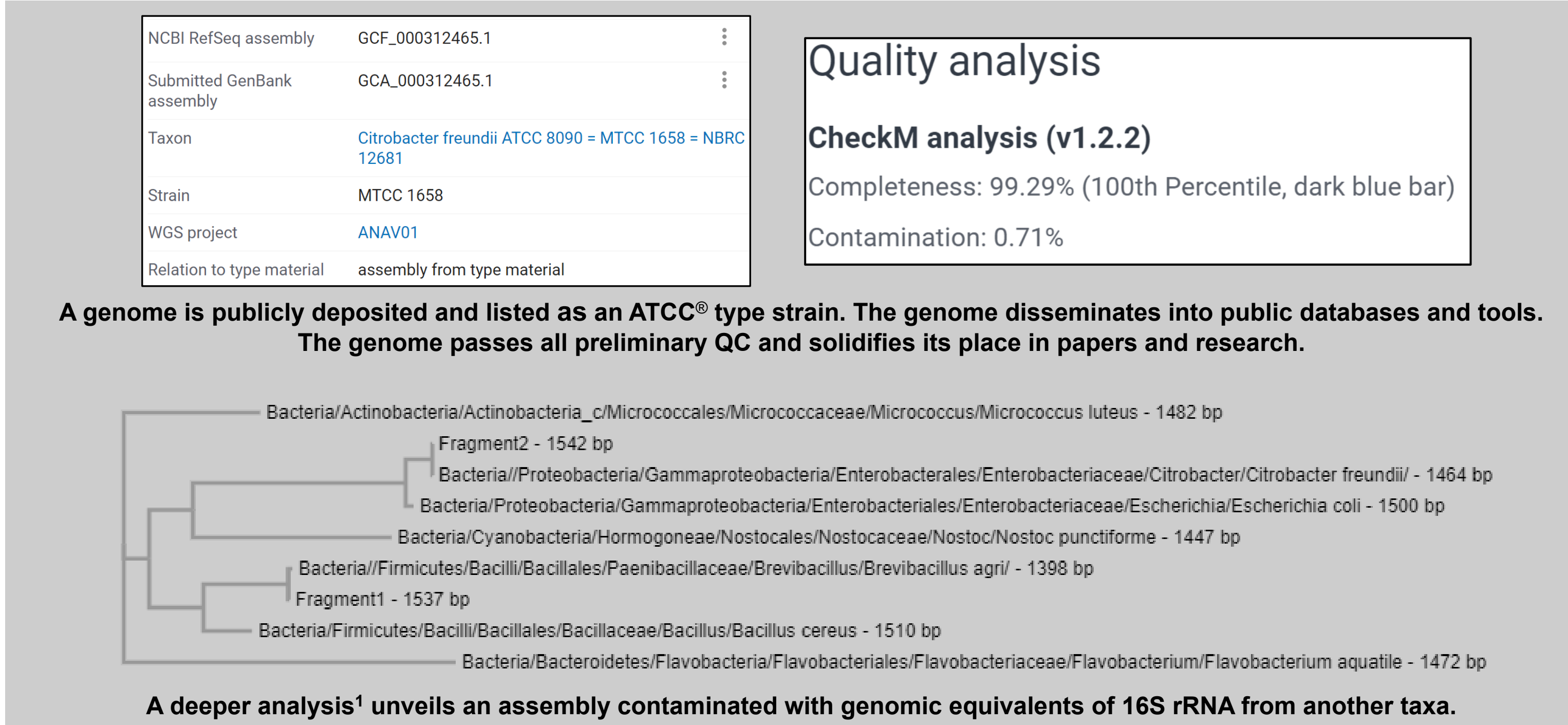


Figure 3: The genome for an ATCC<sup>®</sup> item deposited into public databases that is not authenticated and contains critical errors as compared to an ATCC<sup>®</sup>-curated assembly for that same item.

```

genome_stats (3)
  filtered_contig_count : 3
  filtered_contig_length : 4818689
  number_of_n_bases : 0
  catalog_details (2)
    ATCC_catalog_number : ATCC 780561
    ATCC_lot_number : 947825
  illumina_metadata (5)
    sequencer : Illumina NextSeq
    barcoding_kit : IDT1mDNARNAID1SetATagmentation
    library_kit : IlluminaDNAprep
    basecaller_version : 1.4.1.39716
    basecaller_model : NextSeq 1800/2000 Control Software
  ont_metadata (6)
    sequencer : Oxford Nanopore GridION X5
    f1bcell_type : FLB-HM106
    barcoding_kit : EXP-HBB106
    library_kit : SQK-LSK109
    basecaller_version : 6.3.4-d9e0f64
    basecaller_model : HAC
    
```

Table 2: Additional JSON fields available to access.

| Additional JSON Fields | Values |
|------------------------|--------|
| Input Reads Summary    | 6      |
| Annotations Summary    | 8      |
| AMR Summary            | 4      |
| Genome Provider        | 6      |
| NGS Read QC (ILM/ONT)  | 4/4    |
| Assembly QC            | 4      |
| Product Page Info      | >18    |

Figure 4: Snapshot of embedded genome JSON metadata. All collected information about the cataloged item and QC are stored and accessible.

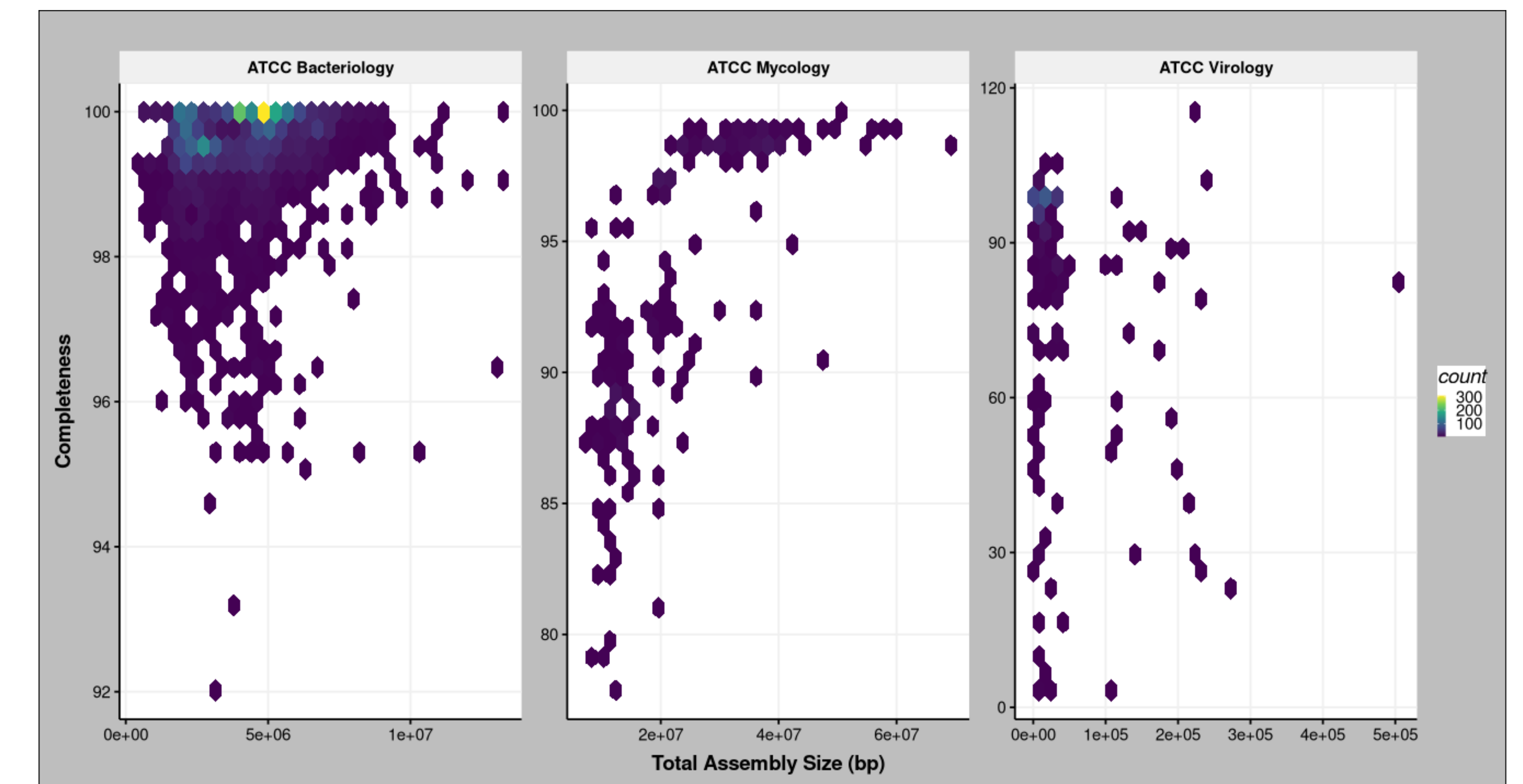


Figure 5: Completeness metrics of the ATCC<sup>®</sup> Genome Portal. CheckM, CheckV, and Busco were used for Bacteriology, Virology, and Mycology products, respectively.<sup>2,3</sup>

## Conclusions

- Public genome databases lack the authentication and traceability that is provided by the ATCC<sup>®</sup> Genome Portal. This gap in data provenance potentially complicates downstream bioinformatics applications and research objectives.<sup>4</sup>
- ATCC<sup>®</sup> is producing ultra-high quality reference genomes for all microbial species in our collection and providing research-use only access to these data via the ATCC<sup>®</sup> Genome Portal.
- All genomes from the collection follow a standardized protocol and rigorous QC from deposition and extraction to sequencing and assembly and are released quarterly to the portal.<sup>5</sup>

### References

- Lee I, et al. ContEst16S: an algorithm that identifies contaminated prokaryotic genomes using 16S RNA gene sequences. In J Syst Evol Microbiol 67(6): 2053-2057, 2017.
- Parks DH, et al. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. Genome Res 25(7): 1043-1055, 2015.
- Simão FA, et al. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. Bioinformatics 31(19): 3210-3212, 2015.
- Yarmosh DA, et al. Comparative analysis and data provenance for 1,113 bacterial genome assemblies. mSphere 7(3): e00077-22, 2022.
- Nguyen SV, et al. The ATCC Genome Portal: 3,938 authenticated microbial reference genomes. Microbiol Resour Anounc 13(2): e0104523, 2024.

