# Utilizing the Scaffolding Ability of Ultra-Long Oxford Nanopore Reads to Assemble Reference-Grade Genomes Originating from Authenticated ATCC Materials

**ATCC®**
**Credible leads to Incredible™**

Briana Benton,[1] Andrew Frank,[1] Anna McCluskey,[1] Steve King,[1] Nick Greenfield,[2] Juan Lopera[1]

[1]ATCC, Manassas, VA 20110; [2]One Codex, San Francisco, CA 94110

## Background

The advancement and accessibility of next-generation sequencing (NGS) technologies have rapidly transformed microbiological research by providing the ability to analyze and profile microbial communities via metagenomics analyses. These sequencing-based applications have relied on the availability of fully assembled reference genomes for bioinformatics analyses, particularly for variant calling in diagnostic and clinical microbiology. However, despite the availability of existing genome sequences in public databases, the quality, completeness, authenticity, accuracy, and traceability of genomic data are inadequate; the lack of standards for genome quality exacerbates these underlying problems. To address this, ATCC has implemented a robust NGS and genome assembly workflow to advance authentication of bacterial strains in the ATCC collection. Our workflow is accompanied by rigorous quality control methods and criteria to ensure that the data proceeding to the next step are the highest quality. Only data that pass all quality control criteria are published to the ATCC Genome Portal, an online database of reference-grade bacterial genomes.
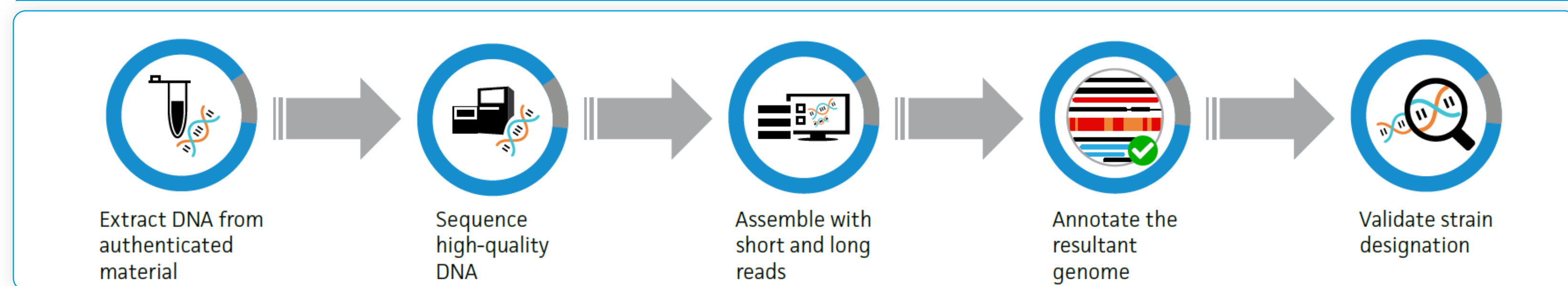
## Whole-Genome Sequencing Workflow



**Figure 1.** Bacterial whole-genome sequencing workflow.

## Extraction of Authenticated Material

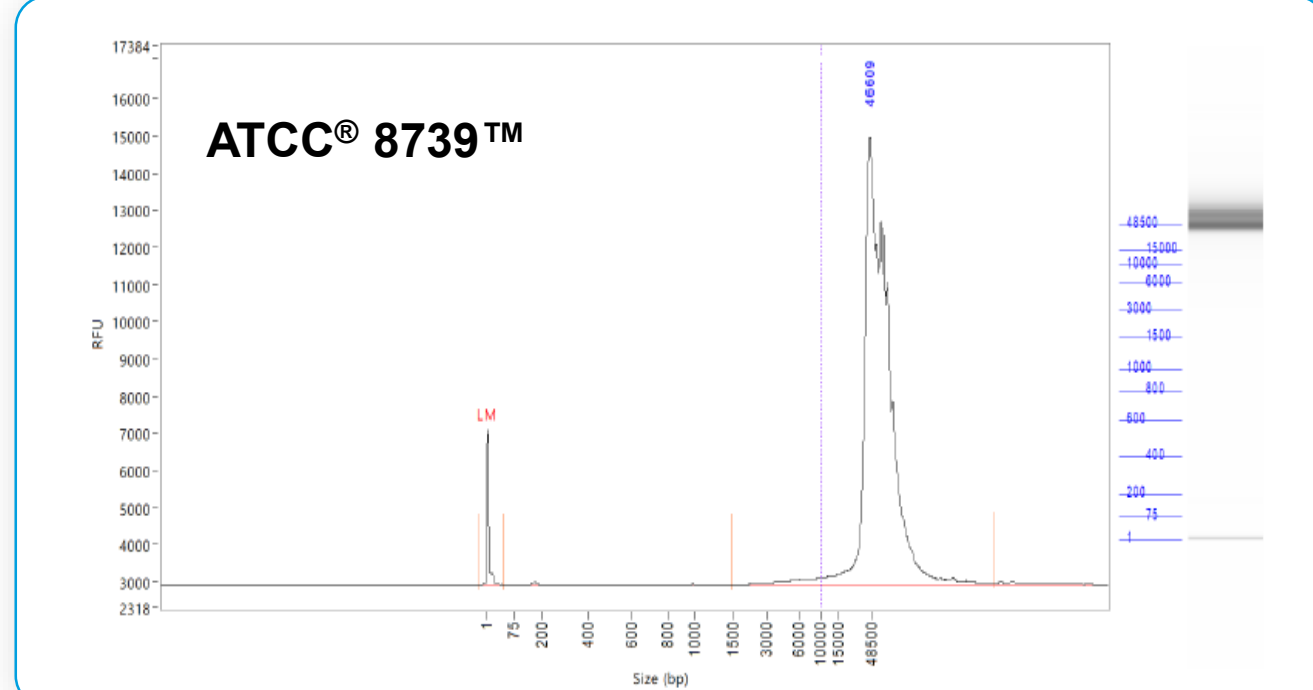| ATCC® No. | Organism | PicoGreen® (ng/µL) | A260/A280 Ratio | Mean Fragment size |
|---|---|---|---|---|
| 8739™ | *Escherichia coli* | 101.9 | 1.92 | >60kb |
| 12228™ | *Staphylococcus epidermidis* | 76.0 | 1.80 | 48kb |
| 13048™ | *Klebsiella aerogenes* | 98.1 | 1.86 | 47kb |
| 14028™ | *Salmonella enterica* subsp. *enterica* | 88 | 1.84 | 59kb |
| 17978™ | *Acinetobacter baumannii* | 133.3 | 1.91 | 46kb |



**Figure 2.** Assessment of quality and quantity of extracted genomic DNA. Fragment size graph obtained from the Agilent Fragment Analyzer platform.
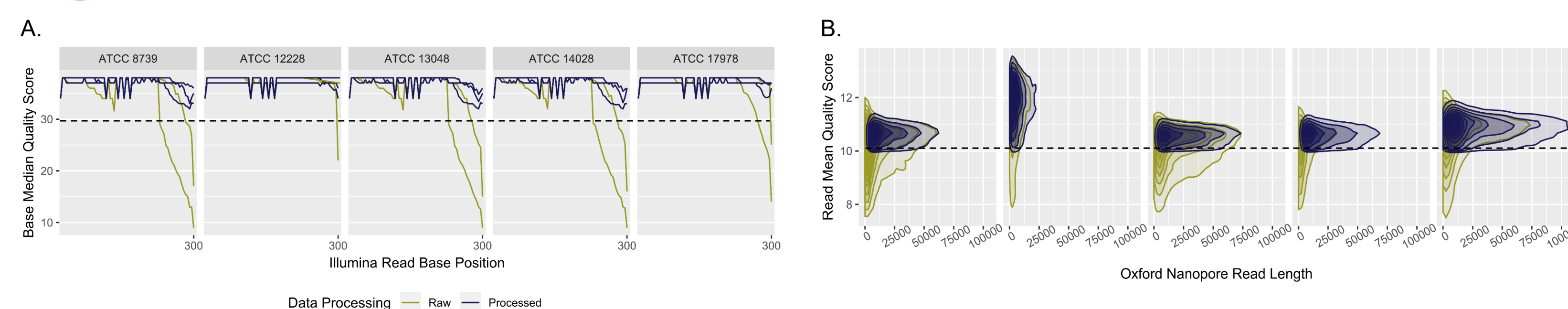
## Sequencing High-Quality DNA



**Figure 3.** ATCC's bacterial genome sequencing quality control (A) substantially improves the quality of Illumina® reads, and (B) improves the length distribution of reads from the Oxford Nanopore Technologies® (ONT) platform. This approach ensures the longest, highest-quality reads are used for assembly. The dashed line indicates the quality score cutoff used for each sequencing technology.
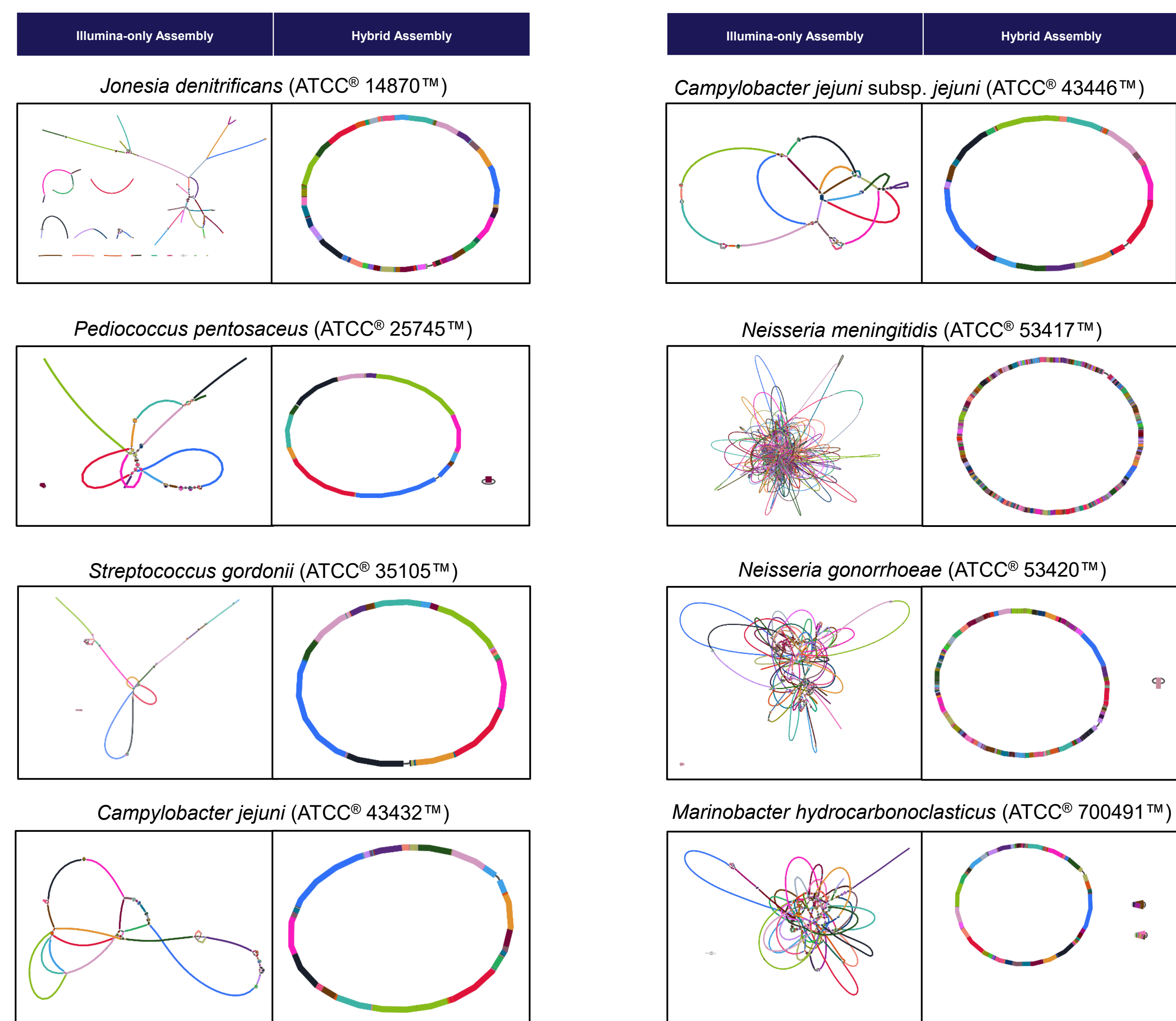
## Assembling with Short and Long Reads



**Figure 4.** Hybrid assembly is a state-of-the-art technique that uses both highly accurate Illumina short reads and ultra-long scaffolding ONT reads. In general, this technique begins with an optimized Illumina assembly. The longest of these resultant contigs are then assembled alongside the ONT reads; this combined assembly then undergoes multiple rounds of both long-read and short-read polishing.

## Analysis and Comparisons

**Table 1.** Reproducibility of ATCC *de-novo* assemblies. Using the ATCC genome assembly workflow, we are able to consistently replicate datasets.

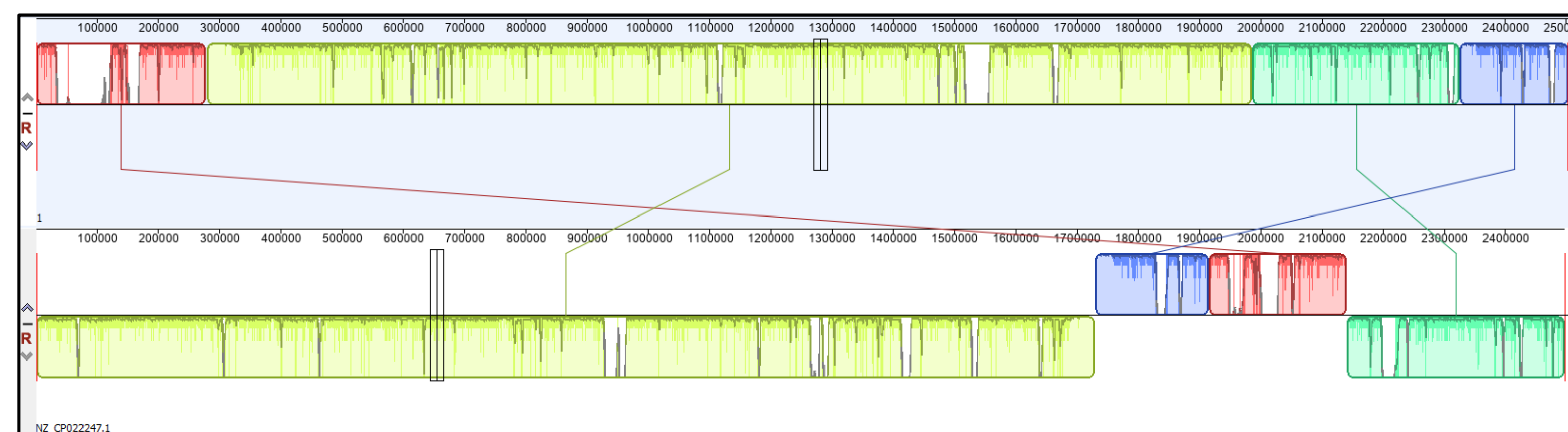| Organism | ATCC® No. | Sequence dataset | total consensus (Mbp) | # of contigs (circular) | N50 (Mbp) | Largest contig (Mbp) | Smallest contig (Mbp) | # of Ns | %GC |
|---|---|---|---|---|---|---|---|---|---|
| *Salmonella enterica* | 14028™ | dataset 1 | 4.96 | 3 | 4.78 | 4.78 | 0.08 | 0 | 52.2 |
| | | dataset 2 | 4.96 | 3 | 4.78 | 4.78 | 0.08 | 0 | 52.2 |
| *Escherichia coli* | 8739™ | dataset 1 | 4.75 | 1 | 4.75 | 4.75 | 4.75 | 0 | 50.9 |
| | | dataset 2 | 4.75 | 1 | 4.75 | 4.75 | 4.75 | 0 | 50.9 |
| *Acinetobacter baumannii* | 17978™ | dataset 1 | 4.08 | 4 | 3.90 | 3.90 | 0.01 | 0 | 38.9 |
| | | dataset 2 | 4.08 | 4 | 3.90 | 3.90 | 0.01 | 0 | 38.9 |



**Figure 5.** Often times, publicly available data is incomplete, lacks in depth traceability, or is incorrect. Here, we show an assembly of ATCC® 12228™ from NCBI compared with our own assembly of ATCC® 12228™. Clear differences are evident in our assembly (top) verses the RefSeq PacBio-only assembly.
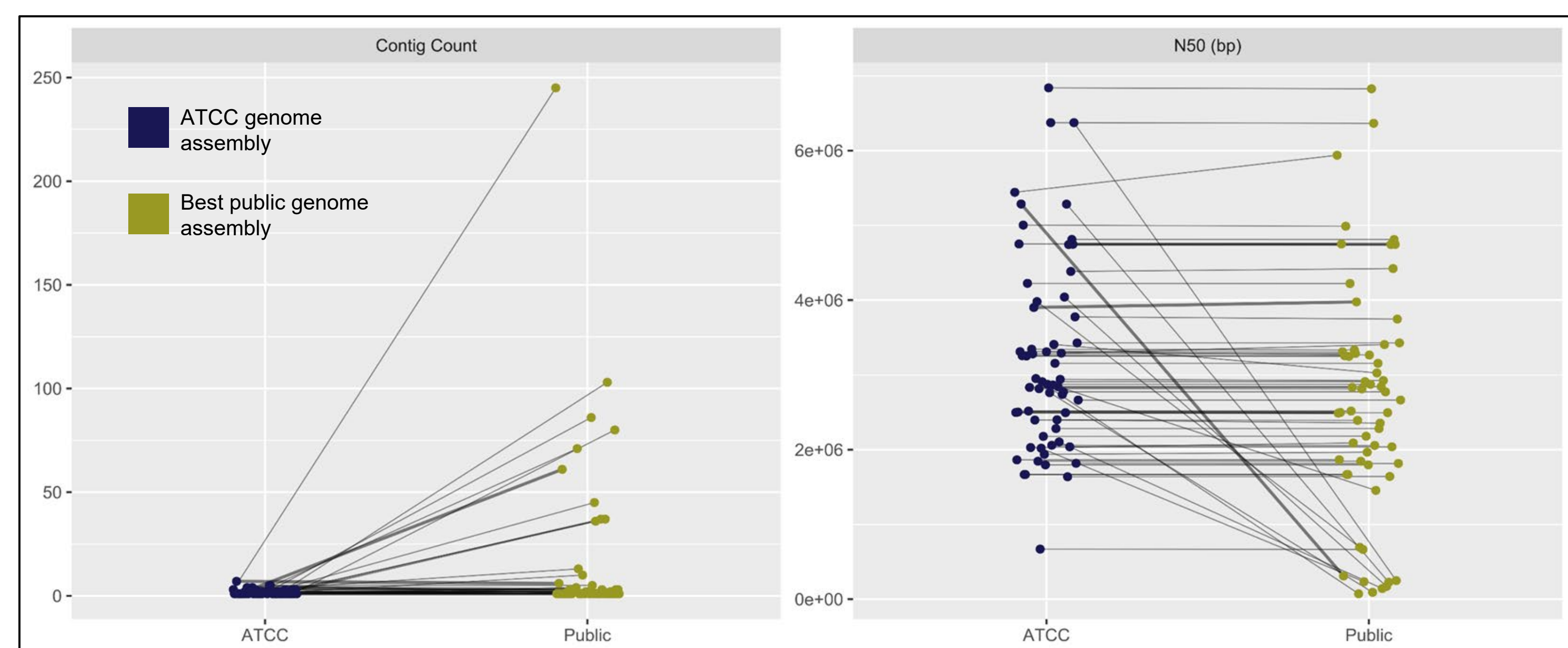


**Figure 6. Downward** trend in contig count and **upward** trend in N50 indicate ATCC is producing more complete genomes than the best publicly available alternative genome.
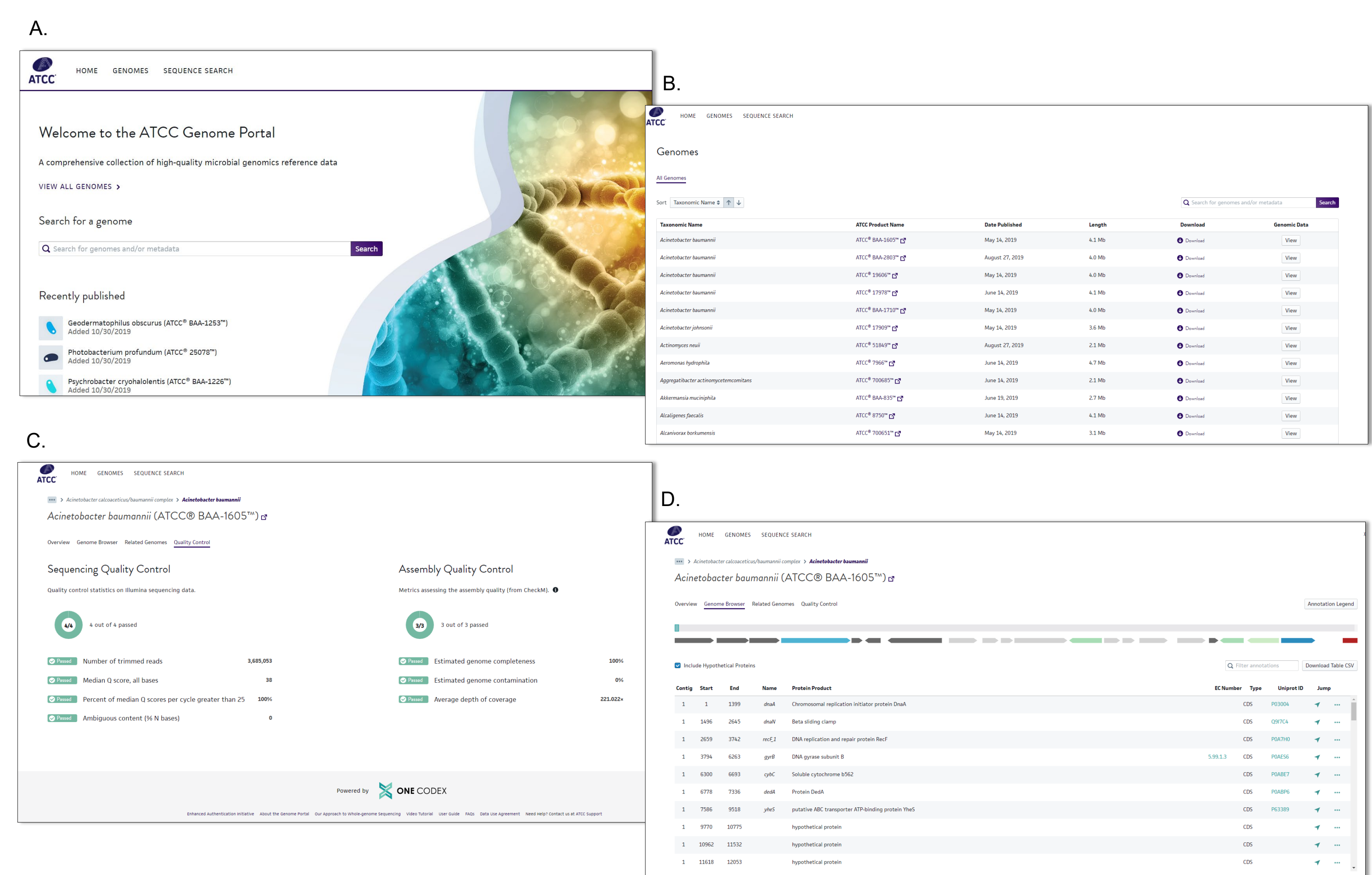


**Figure 7.** Screenshots from the new ATCC Genome Portal. The top left figure illustrates the "Home" page followed by the "Genomes" page, which lists all genomes that pass our rigorous QC criteria. Here, users are able to download the respective fasta and/or gbk files for each organism. (C) The "Quality Control" page, which displays the QC values for each respective genome, and (D) the "Genome Browser" page, which illustrates the annotation of each respective genome. All data is available at genomes.atcc.org.

## Summary

- Our hybrid-assembly method allows us to leverage the power of highly accurate Illumina short-reads with the scaffolding ability of ultra-long ONT reads to generate high-quality reference genomes that are more complete and accurate than what can be generated with each individual technology alone.
- Our standardized and reproducible genome sequencing, assembly, and annotation workflow allows researchers to access higher-quality genomes that are fully authenticated and matched with authenticated ATCC strains.
- ATCC will continue to sequence and assemble the genomes of organisms in our bacteriology collection and make them publicly available to the scientific research community via monthly additions to the ATCC Genome Portal. Data are accessible at genomes.atcc.org.
- We are currently working to expand our applications of the ONT sequencing platform to further enrich the characterization of our virology and mycology collections.