

1 Comparative Analysis and Data Provenance for 1,113 Bacterial Genome 2 Assemblies

3 David A. Yarmosh¹, Juan G. Lopera¹†, Nikhita P. Puthuveetil¹, Patrick Ford Combs¹, Amy L.
4 Reese¹, Corina Tabron¹, Amanda E. Pierola¹, James Duncan¹, Samuel R. Greenfield¹, Robert
5 Marlow¹, Stephen King¹, Marco A. Riojas^{1,2}, John Bagnoli¹, Briana Benton¹, Jonathan L.
6 Jacobs^{1*}

7 ¹American Type Culture Collection (ATCC); Manassas, VA (USA)

8 ²BEI Resources; Manassas, VA (USA)

9 *Corresponding author. Email: jjacobs@atcc.org

10 †Present Address: EDAN Diagnostics; Madison, WI (USA)

11
12 **The quality and traceability of microbial genomics data in public databases is deteriorating**
13 **as they rapidly expand and struggle to cope with data curation challenges. While the**
14 **availability of public genomic data has become essential for modern life sciences research,**
15 **the curation of the data is a growing area of concern that has significant real-world impacts**
16 **on public health epidemiology, drug discovery, and environmental biosurveillance**
17 **research¹⁻⁶. While public microbial genome databases such as NCBI's RefSeq database**
18 **leverage the scalability of crowd sourcing for growth, they do not require data provenance**
19 **to the original biological source materials or accurate descriptions of how the data was**
20 **produced⁷. Here, we describe the *de novo* assembly of 1,113 bacterial genome references**
21 **produced from authenticated materials sourced from the American Type Culture**
22 **Collection (ATCC), each with full data provenance. Over 98% of these ATCC Standard**
23 **Reference Genomes (ASRGs) are superior to assemblies for comparable strains found in**
24 **NCBI's RefSeq database. Comparative genomics analysis revealed significant issues in**
25 **RefSeq bacterial genome assemblies related to genome completeness, mutations, structural**
26 **differences, metadata errors, and gaps in traceability to the original biological source**
27 **materials. For example, nearly half of RefSeq assemblies lack details on sample source**
28 **information, sequencing technology, or bioinformatics methods. We suggest there is an**
29 **intrinsic connection between the quality of genomic metadata, the traceability of the data,**

30 **and the methods used to produce them with the quality of the resulting genome assemblies**
31 **themselves. Our results highlight common problems with “reference genomes” and**
32 **underscore the importance of data provenance for precision science and reproducibility.**
33 **These gaps in metadata accuracy and data provenance represent an “elephant in the**
34 **room” for microbial genomics research, but addressing these issues would require raising**
35 **the level of accountability for data depositors and our own expectations of data quality.**

36

37 The National Center for Biotechnology Information’s (NCBI) RefSeq database has become an
38 essential cornerstone of the global genomics research community, but declining data quality and
39 the increasing cost of manual data curation by end-users are growing areas of concern ^{8,1,3,2}. As
40 RefSeq continues to expand, so too does the risk for data errors, omission, obfuscation, or
41 falsification to go undetected and to potentially damage trust in this enormously important public
42 resource ^{9,10}. RefSeq contains over 357,657 prokaryotic genome assemblies. It is the largest
43 collection of non-redundant, annotated genome assemblies available, and it is built exclusively
44 from crowd-sourced data. However, despite extensive efforts to create automated curation
45 pipelines and tools to improve RefSeq data, significant quality issues remain in genome
46 assemblies found within RefSeq ^{11–13}. For example, while all newly deposited prokaryote
47 genome assemblies are automatically annotated, the associated metadata records (i.e.,
48 BioSample, BioProject, SRA, Assembly data) are submitted by depositors who are not required
49 to provide attribution for the biological materials behind each genome ^{14,15}. In fact, the
50 International Nucleotide Sequence Database Collaboration (INSDC) policy states “*the quality*
51 *and accuracy of the record are the responsibility of the submitting author, not of the database,*”
52 which is to say that metadata, which are often crucial for comparative genomics research, are not
53 curated or verified for accuracy ¹⁶. This is further complicated by data omissions, lack of
54 controlled vocabulary for terms, variable taxonomic naming conventions, and competing
55 metadata package formats. In many cases, tracing the provenance of an individual assembly to its
56 source material in order to verify its authenticity becomes challenging, and manual curation is
57 frequently required to detect and correct RefSeq metadata errors ¹⁷.

58 In this study, we present the results of an ongoing whole-genome sequencing initiative at the
59 American Type Culture Collection (ATCC) to provide end-to-end data provenance from source
60 materials to reference-grade microbial genomes, hereafter referred to as “ATCC Standard

61 Reference Genomes” (ASRGs). Presented here are 1,113 bacterial ASRGs from authenticated
62 materials that were produced via a hybrid *de novo* assembly approach. We compared them to
63 assemblies in RefSeq where metadata indicated they were produced by 3rd party labs from
64 materials sourced from ATCC For 366 ASRGs (~33%), we were able use metadata to compare
65 them to one or more assemblies in RefSeq. The remaining 747 ASRGs (~66%) represented
66 potentially novel assemblies. All ASRGs described here are available for research use via the
67 ATCC Genome Portal (<https://genomes.atcc.org>)¹⁸.

68 **Whole-genome Sequencing of 1,113 ATCC Bacterial Strains**

69 High-molecular-weight genomic DNA (HMW-gDNA) was extracted from 1,113 bacterial strains
70 obtained from ATCC’s biorepository. Each strain was cultured using strain-specific protocols
71 and subjected to quality control (QC) for contamination, viability, purity, phenotype, and
72 taxonomic identity (Figure 1). For whole-genome sequencing (WGS), HMW-gDNA was split
73 and subjected to sequencing using both Illumina and Oxford Nanopore Technologies (ONT)
74 next-generation sequencing (NGS) platforms (Figure 1). Next, reads were taxonomically
75 classified using One Codex’s metagenomics platform to assess the purity of each NGS library
76 prior to *de novo* assembly¹⁹. Read sets were then down-sampled to predetermined coverage
77 depths (Illumina, 100x; ONT, 60x) expected to be optimal for bacterial genome assemblies^{20–22}.
78 Lastly, a hybrid-assembly pipeline incorporating reads from both platforms produced *de novo*
79 assemblies for each strain²³. High-level summary metrics for each ASRG are shown in Table S1
80 and Fig. 2. All 1,113 ASRG assemblies were estimated to be over 95% complete by *CheckM*;
81 1,015 were found to be over 99% complete and 329 are 100% complete²⁴. A total of 617 are
82 considered high-quality, closed genome references.

83 **Survey of Bacterial Genome Assemblies in RefSeq**

84 We compared the ASRG assemblies to those in NCBI’s RefSeq Bacteria Database
85 (<https://ftp.ncbi.nlm.nih.gov/genomes/refseq/bacteria/>) labeled as representing ATCC bacterial
86 strains, i.e., assemblies where the ATCC strain name (or a synonymous name) was indicated in
87 the title, description, or other metadata field in the GenBank assembly record. We intentionally
88 did not search RefSeq using a traditional comparative genomics approach (i.e., by sequence
89 homology, BLAST, etc.) since this would require arbitrary thresholds for determining strain
90 identity, and metadata descriptors are intended to be useful for these types of queries. Using this
91 approach, we found 2,701 genome assemblies in RefSeq, which collectively comprised 1,960

92 different ATCC strains (Table S2, Fig. 3A). Interestingly, RefSeq had numerous examples of
93 bacterial strains represented by multiple assemblies or submitted by different groups, and it often
94 included “strains” resulting from intentional genetic modification (i.e., there are 33 different
95 RefSeq assemblies for *Serratia marcescens* subsp. *marcescens* ATCC® 13880™). This is despite
96 it representing a “non-redundant” database. Overall, we found one or more duplicate assemblies
97 in RefSeq for 158 strains for which we also produced an ASRG, including instances of
98 assemblies for genetically modified strains mislabeled as representing “type strains” (See Table
99 S2). These errors and strain duplications create risks for researchers who may unwittingly use
100 these data in their own research yet remain unaware of these issues.

101 Further examination of the metadata for the 2,701 RefSeq assemblies labeled as ATCC strains
102 also revealed numerous records with incomplete, missing, or obscured descriptor fields (Figure
103 S1). For example, “Assembly type” is present in every assembly record but the value is “na” for
104 all. “Sequencing technology” is not included or has a value of “Unknown” for 1,088 assemblies
105 (~40%, Table S2), and spelling and nonstandard abbreviations further complicate the rest.
106 “Assembly method” is not included for 1,082 assemblies, contains the value “Unknown” for 88
107 assemblies or “other” for 4 assemblies, and has numerous misspellings for various
108 bioinformatics tools (i.e., “Velevt” or “Velveth” for the Velvet assembler). One particularly poor
109 example includes an assembly for *Streptomyces clavuligerus* ATCC® 27064™
110 (GCF_015708605.1) that indicates the “Assembly method” as “Several assembly pipelines,
111 manual curation v. 2018-09-27.” Underutilized fields included “Description,” “Isolate,” and
112 “Relation to type material,” which had no values in 99%, 98%, and 38% of the assembly records,
113 respectively. The damaging impact that inconsistent depositor metadata has on scientific research
114 and reproducibility has been extensively covered elsewhere^{1,3,25}.

115 Of the 2,701 RefSeq assemblies for ATCC bacterial strains, 708 had a counterpart ASRG (Table
116 S2, Figure 3A). Of these, 303 (43%) are labeled “complete genome” or “chromosome” level
117 assemblies. Despite this, N50 values were largely inferior when compared to their ASRG
118 counterparts (Figure 3B). While 241 RefSeq assemblies had the same number of scaffolds as
119 their corresponding ASRGs, 341 were more fragmented. Altogether, 662 ASRGs had equivalent
120 or superior N50 values to their RefSeq counterparts (ATCC N50 / RefSeq N50 \geq 0.95), while 46
121 ASRG assemblies were more fragmented (Figure 3D). The greatest difference was observed for

122 a RefSeq assembly for *Pseudomonas aeruginosa* ATCC[®] 700888[™] (GCF_000297315.1), which
123 comprised 600 contigs while the ASRG equivalent is closed, containing only one contig.

124 **Comparative Genomics of 303 RefSeq Assemblies**

125 Next, we compared the 303 complete RefSeq assemblies to their corresponding ASRGs for the
126 same strains (represented by 212 ASRGs). First, we found that the pairwise average nucleotide
127 identity (ANI) ranged from 97% to 100% for identical strains, which at first glance suggested a
128 high level of similarity²⁶. Although large differences in the high-level assembly metrics were
129 previously observed (e.g., N50, GC content), after conducting pairwise whole-genome
130 alignments with *Mummer4* for all 303 RefSeq assemblies against ASRGs for the same strain, we
131 found 292 had over 95% of their sequence aligned. Next, we examined pairwise structural
132 variations and found significant differences in sequence repeats, inversions, indels, and
133 translocations between RefSeq assemblies and ASRGs for the same strains (Tables S3, S4)²⁷.
134 Analysis with *dnaDiff* of all 303 RefSeq assemblies revealed an average 6.73 structural
135 rearrangements in comparison to ASRGs, the worst of which was GCF_000160895.1 for
136 *Bacillus cereus* ATCC[®] 10876[™] with 232 structural differences (despite both assemblies having
137 over 99% reciprocally aligned bases). Structural relocations were the most common, with 256
138 RefSeq assemblies having at least one per assembly (average 4.3 per assembly). Structural
139 inversions were found in 74 RefSeq assemblies (average 2.2). Translocations were relatively
140 rare, with only 9 RefSeq assemblies having structural translocations relative to the ASRG
141 assembly for the same strain (Table S4). We also found that RefSeq assemblies with the greatest
142 number of structural differences from the ATCC assemblies corresponded to those submitted to
143 NCBI prior to 2010, and for which sequencing technology or assembly method were not
144 indicated in the RefSeq metadata. The distribution of structural variations in the 303 complete
145 RefSeq assemblies compared to their corresponding ASRGs is shown in Figure S2.

146 **Variants in 303 RefSeq Assemblies**

147 Next, we sought to investigate the prevalence of single-nucleotide polymorphisms (SNPs) and
148 insertions/deletions (InDels) that would arise by using RefSeq assemblies as a reference genome
149 against which Illumina sequencing data would be mapped—a common approach used by labs
150 without the resources or expertise for *de novo* assembly and annotation. For each of the 303
151 complete RefSeq assemblies described above, we mapped the same Illumina reads used in
152 creating the corresponding ASRGs for the same strain. Variant calling from the resulting

153 consensus genomes was carried out on all 303 references to detect SNPs and InDels in each (see
154 Materials & Methods). Overall, the number of SNPs and InDels per assembly ranged from zero
155 (none detected) to as many as 60,064 SNPs (*Acinetobacter baumannii* ATCC[®] 17978[™],
156 GCF_011067065.1) and 2,699 InDels for a given assembly (*Parabacteroides distasonis* ATCC[®]
157 8503[™], GCF_900683725.1) (Table S5). The median level of SNPs and InDels was 7 SNPs and
158 8 InDels per assembly, with 7 of the 303 mappings having no detectable SNPs and InDels. These
159 results were promising overall, yet significant outliers were detected, and 26 strains had SNPs
160 and InDels beyond an extreme-outlier boundary, i.e., greater than 3-times interquartile range
161 (IQL) above the median with 9 of them having over 1,000 SNPs and InDels each (Figs. S3, S4a,
162 S4b).

163 A total of 111 assemblies had fewer than 10 variants, while 15 assemblies had more than 500
164 variants (SNPs, Indels). Not surprisingly, as the number of SNPs increased, so too did the
165 number of InDels (Figure S3). Of these, 52 of the 303 assemblies had no expected non-
166 synonymous mutations, but 87 had at least 10 non-synonymous variants per genome (Figure
167 S4b). Importantly, 52 RefSeq assemblies identified as “assembled from type material” were
168 found to have at least 10 non-synonymous variants, and seven assemblies had over 100; this
169 could have potentially deleterious impacts on future comparative genomics studies utilizing
170 those reference assemblies (Table S5).

171 We found that complete RefSeq assemblies without the label “reference genome” or
172 “representative genome” (250 genomes) were enriched for SNPs (7.6-fold) and InDels (9.6-fold)
173 compared to reference RefSeq genomes (53 assemblies). Furthermore, type strain assemblies in
174 RefSeq (i.e., labeled as “assembly designated as neotype,” “assembly from synonym type
175 material,” or “assembly from type material”) had marginally fewer SNPs and InDels than other
176 assemblies overall, but some significant exceptions to this were also observed (see above). No
177 statistically significant enrichment for SNPs or InDels was detectable by taxonomic clade or G:C
178 content. Collectively, these results underscore the importance of data provenance of the
179 originating materials (e.g. “type-strains”) and assembly quality (e.g. “reference genome” or
180 “representative genome”), and that they are both important drivers in reducing variability and
181 improving genome assembly quality.

182 **Discussion**

183 Over the last 20 years, several non-commercial and government initiatives have specifically tried
184 to address issues relating to the quality and standardization of metadata for microbial genomics,
185 which has had some benefit for end-users, but substantial work remains to be done^{15,28,29}. As the
186 unmet need for curated, high-quality microbial genomics data continues to grow, we will no
187 doubt continue to see a variety of commercial initiatives be successful in developing solutions
188 designed to address gaps in quality, content, and reliability, such as QIAGEN's CLC Microbial
189 Reference Database, ARES Genetics' ARESdb, and the One Codex platform. While these public
190 and private efforts have been largely successful, by some measures the overall quality of public
191 microbial genomics data has been declining over the last decade, carrying a potentially great cost
192 to the broader research community^{2,3,5,13,30}. We propose that widespread gaps in the traceability
193 of genome assemblies to their originating biological materials, lab protocols, and bioinformatics
194 methods represent fundamental weaknesses in these data that will hinder research and increase
195 costs unless it is addressed.

196 At the outset of the work described here, we sought to develop methods to systematically
197 sequence ATCC's bacterial collection and share that data with the research community alongside
198 the physical strain materials. However, during the course of our work we found that bacterial
199 genome assemblies in RefSeq labeled as representing ATCC strains compared poorly against
200 ASRGs. More broadly, our analysis uncovered disparities in the quality, accuracy, and
201 completeness of metadata associated with assemblies in RefSeq, suggesting that gaps in data
202 provenance may be playing a role in the decline of data quality. As an example, over 33%
203 (1,087) of the RefSeq assemblies included in our study completely lacked any description for
204 how they were sequenced or assembled.

205 There remain significant gaps in the quality of "typical" genome assemblies available from
206 crowd-sourced databases such as RefSeq. Researchers should be cautious about the data they use
207 and avoid blindly ingesting reference genome data without first being curious about the origins
208 of the data and the methods used to produce them. Further studies are needed to better
209 understand the importance of establishing data provenance in genomics data and the impact its
210 absence has on the research of those who use it. It is our hope that initiatives focused on genomic
211 data provenance, such as the ATCC Genome Portal (<https://genomes.atcc.org>), will serve to
212 highlight the value of establishing higher standards of traceability and accountability for
213 genomics data in the public domain.

214

215 **References**

- 216 1. Gonçalves, R. S. & Musen, M. A. The variable quality of metadata about biological samples
217 used in biomedical experiments. *Sci Data* **6**, 190021 (2019).
- 218 2. Pettengill, J. B. *et al.* Interpretative labor and the bane of non-standardized metadata in
219 public health surveillance and food safety. *Clinical Infectious Diseases* ciab615 (2021)
220 doi:10.1093/cid/ciab615.
- 221 3. Rajesh, A. *et al.* Improving the completeness of public metadata accompanying omics
222 studies. *Genome Biol* **22**, 106, s13059-021-02332-z (2021).
- 223 4. Vangay, P. *et al.* Microbiome Metadata Standards: Report of the National Microbiome Data
224 Collaborative’s Workshop and Follow-On Activities. *mSystems* **6**, e01194-20,
225 /msystems/6/1/mSys.01194-20.atom (2021).
- 226 5. Toczydlowski, R. H. *et al.* Poor data stewardship will hinder global genetic diversity
227 surveillance. *Proc Natl Acad Sci USA* **118**, e2107934118 (2021).
- 228 6. Leipzig, J. *et al.* The role of metadata in reproducible computational research.
229 *arXiv:2006.08589 [cs]* (2021).
- 230 7. Li, W. *et al.* RefSeq: expanding the Prokaryotic Genome Annotation Pipeline reach with
231 protein family model curation. *Nucleic Acids Research* **49**, D1020–D1028 (2021).
- 232 8. Pruitt, K. D., Tatusova, T. & Maglott, D. R. NCBI Reference Sequence (RefSeq): a curated
233 non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res* **33**,
234 D501-504 (2005).

- 235 9. Gopalakrishna, G. *et al.* *Prevalence of responsible research practices and their potential*
236 *explanatory factors: a survey among academic researchers in The Netherlands.*
237 <https://osf.io/xsn94> (2021) doi:10.31222/osf.io/xsn94.
- 238 10. Caswell, J. *et al.* Defending Our Public Biological Databases as a Global Critical
239 Infrastructure. *Front. Bioeng. Biotechnol.* **7**, 58 (2019).
- 240 11. Steinegger, M. & Salzberg, S. L. Terminating contamination: large-scale search identifies
241 more than 2,000,000 contaminated entries in GenBank. *Genome Biol* **21**, 115 (2020).
- 242 12. Segerman, B. The Most Frequently Used Sequencing Technologies and Assembly Methods
243 in Different Time Segments of the Bacterial Surveillance and RefSeq Genome Databases.
244 *Front. Cell. Infect. Microbiol.* **10**, 527102 (2020).
- 245 13. Smits, T. H. M. The importance of genome sequence quality to microbial comparative
246 genomics. *BMC Genomics* **20**, 662, s12864-019-6014-5 (2019).
- 247 14. Tatusova, T. *et al.* NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res* **44**,
248 6614–6624 (2016).
- 249 15. Field, D. *et al.* The minimum information about a genome sequence (MIGS) specification.
250 *Nature Biotechnology* **26**, 541–547 (2008).
- 251 16. Karsch-Mizrachi, I., Nakamura, Y., Cochrane, G., & on behalf of the International
252 Nucleotide Sequence Database Collaboration. The International Nucleotide Sequence
253 Database Collaboration. *Nucleic Acids Research* **40**, D33–D37 (2012).

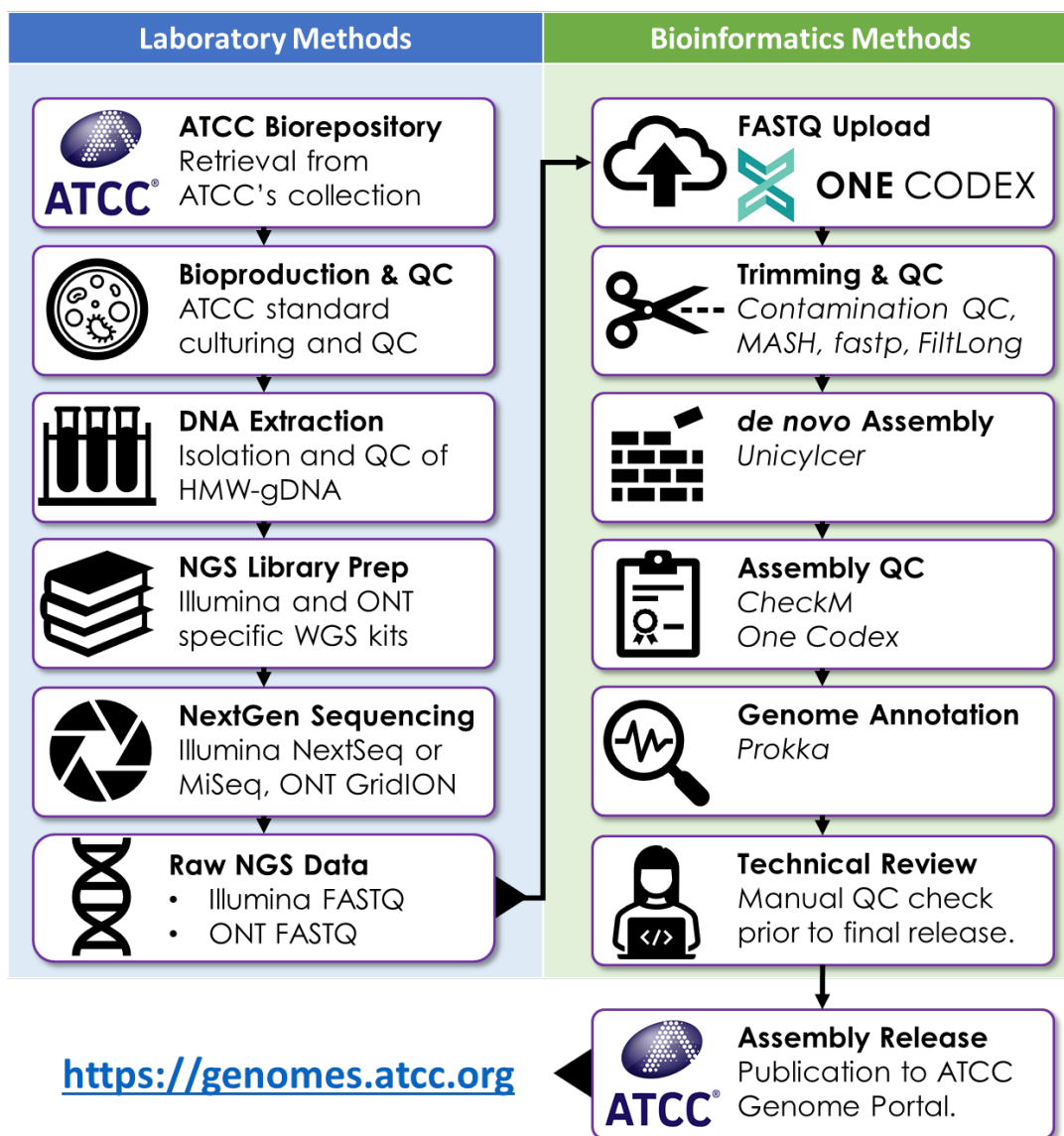
- 254 17. Schmedes, S. E., King, J. L. & Budowle, B. Correcting Inconsistencies and Errors in
255 Bacterial Genome Metadata Using an Automated Curation Tool in Excel (AutoCurE). *Front.*
256 *Bioeng. Biotechnol.* **3**, (2015).
- 257 18. Benton, B. *et al.* The ATCC Genome Portal: Microbial Genome Reference Standards with
258 Data Provenance. *Microbiology Resource Announcements* (2022).
- 259 19. Minot, S. S., Krumm, N. & Greenfield, N. B. *One Codex: A Sensitive and Accurate Data*
260 *Platform for Genomic Microbial Identification.*
261 <http://biorxiv.org/lookup/doi/10.1101/027607> (2015) doi:10.1101/027607.
- 262 20. Desai, A. *et al.* Identification of Optimum Sequencing Depth Especially for De Novo
263 Genome Assembly of Small Genomes Using Next Generation Sequencing Data. *PLoS ONE*
264 **8**, e60204 (2013).
- 265 21. Chen, Z., Erickson, D. L. & Meng, J. Benchmarking hybrid assembly approaches for
266 genomic analyses of bacterial pathogens using Illumina and Oxford Nanopore sequencing.
267 *BMC Genomics* **21**, 631 (2020).
- 268 22. Sims, D., Sudbery, I., Ilott, N. E., Heger, A. & Ponting, C. P. Sequencing depth and
269 coverage: key considerations in genomic analyses. *Nat Rev Genet* **15**, 121–132 (2014).
- 270 23. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome
271 assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595 (2017).

- 272 24. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
273 assessing the quality of microbial genomes recovered from isolates, single cells, and
274 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 275 25. Marc, D. T., Beattie, J., Herasevich, V., Gatewood, L. & Zhang, R. Assessing Metadata
276 Quality of a Federally Sponsored Health Data Repository. *AMIA Annu Symp Proc* **2016**,
277 864–873 (2016).
- 278 26. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
279 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat*
280 *Commun* **9**, 5114 (2018).
- 281 27. Marçais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLOS*
282 *Computational Biology* **14**, e1005944 (2018).
- 283 28. Sichtig, H. *et al.* FDA-ARGOS is a database with public quality-controlled reference
284 genomes for diagnostic use and regulatory science. *Nat Commun* **10**, 3313 (2019).
- 285 29. Yilmaz, P. *et al.* The genomic standards consortium: bringing standards to life for microbial
286 ecology. *ISME J* **5**, 1565–1567 (2011).
- 287 30. Breitwieser, F. P., Pertea, M., Zimin, A. V. & Salzberg, S. L. Human contamination in
288 bacterial genomes has created thousands of spurious proteins. *Genome Res.* **29**, 954–960
289 (2019).
- 290 31. Bacteriology Culture Guide. *American Type Culture Collection*
291 <https://www.atcc.org/resources/culture-guides/bacteriology-culture-guide> (2021).

- 292 32. ATCC Ready-to-Use Nucleic Acids. *American Type Culture Collection*
293 <https://www.atcc.org/microbe-products/bacteriology-and-archaea/nucleic-acids>.
- 294 33. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor.
295 *Bioinformatics* **34**, i884–i890 (2018).
- 296 34. Wingett, S. W. & Andrews, S. FastQ Screen: A tool for multi-genome mapping and quality
297 control. *Fl000Res* **7**, 1338 (2018).
- 298 35. Wick, R. & Menzel, P. *Filtlong*. (2020).
- 299 36. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using MinHash.
300 *Genome Biology* **17**, (2016).
- 301 37. Salzberg, S. L. *et al.* GAGE: A critical evaluation of genome assemblies and assembly
302 algorithms. *Genome Research* **22**, 557–567 (2012).
- 303 38. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069
304 (2014).
- 305 39. ATCC Genome Portal. *American Type Culture Collection* <https://genomes.atcc.org> (2021).
- 306 40. Ole Tange. GNU Parallel - The Command-Line Power Tool. *The USENIX Magazine* 42–47
307 (2011).
- 308 41. Khelik, K., Lagesen, K., Sandve, G. K., Rognes, T. & Nederbragt, A. J. NucDiff: in-depth
309 characterization and annotation of differences between two sets of DNA sequences. *BMC*
310 *Bioinformatics* **18**, 338 (2017).

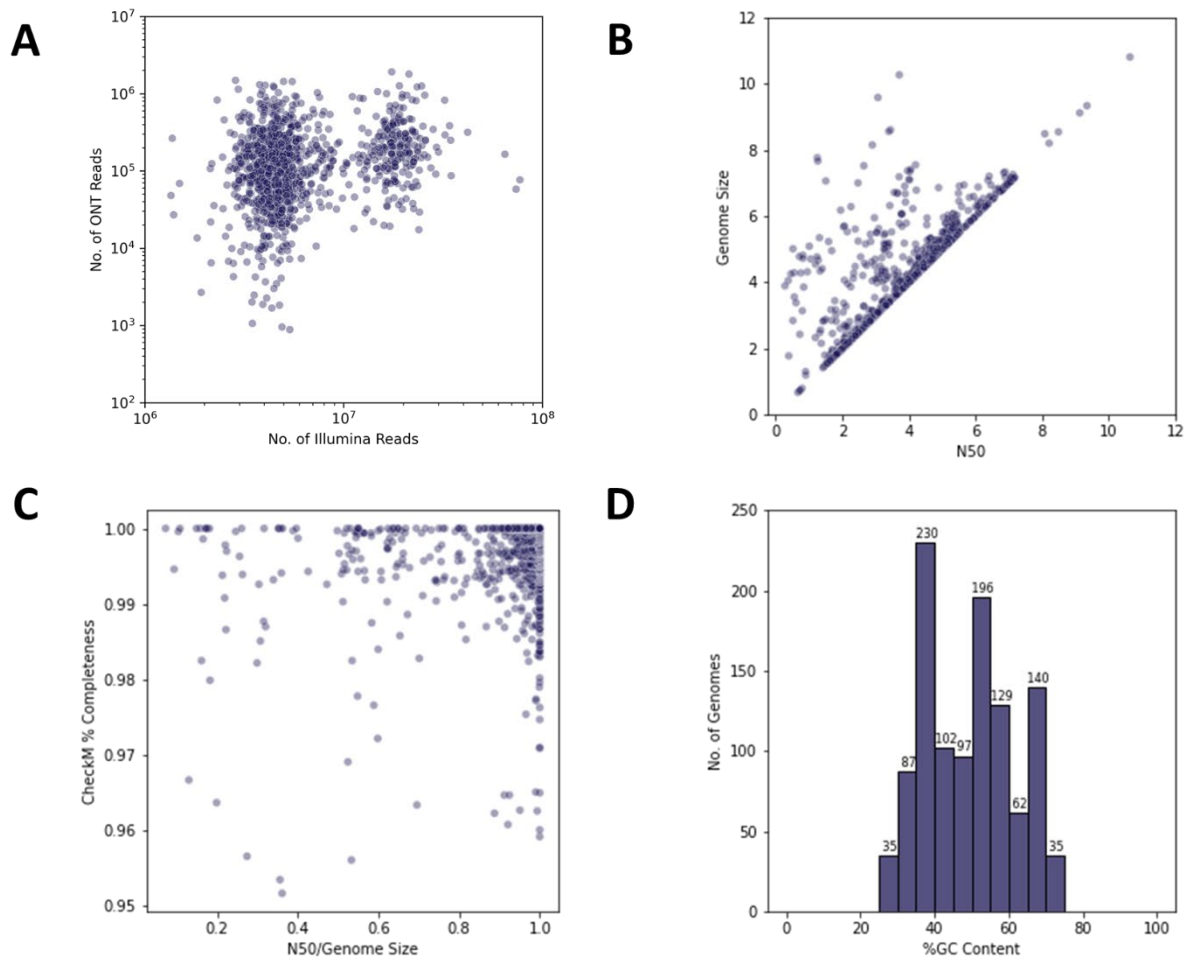
311

FIGURES



313 **Fig 1. A Pipeline for End-to-End Genomic Data Provenance.** Source materials are obtained
314 directly from the ATCC biorepository and tracked through to the final assembly and genome
315 annotation. Upfront culture conditions varied depending on the species cultured, but downstream
316 process steps were performed using standardized protocols for DNA extraction, library prep,
317 sequencing, and bioinformatics. Each pipeline is hosted on One Codex's cloud infrastructure.

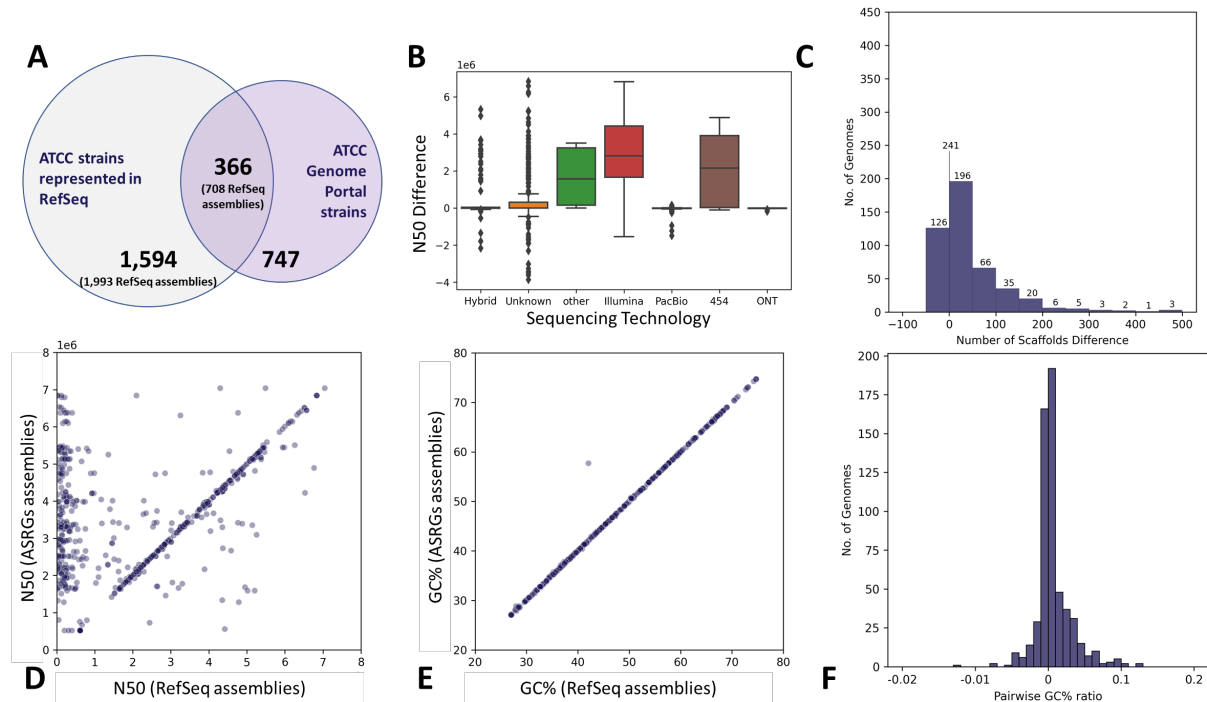
318



320 **Fig 2. Sequencing & Quality Metrics for 1,113 Bacterial Genome Assemblies.** (A) Illumina
321 vs. ONT reads for ASRGs before down-sampling. (B) N50 metrics vs. genome size. (C) N50
322 normalized by genome size vs *CheckM* genome completion estimates. (D) Diversity of G:C%
323 content for all 1,113 ASRG assemblies.

324

325



327 **Fig 3. Comparative Metrics for 1,113 ASRGs vs. RefSeq Assemblies.** (A) Intersection of
 328 ASRGs vs. RefSeq for strains labeled as being from ATCC. In parentheses are the total number
 329 of RefSeq assemblies, allowing for strain-redundancy. (B) N50 variability of RefSeq vs. ASRGs
 330 by sequencing technology. Note the scale is 1E6. (C) Differences in contig counts for ASRG vs.
 331 RefSeq assemblies. Positive values indicating the RefSeq assembly had more contigs. (D) Ratios
 332 of ASRG N50 values (y-axis) to RefSeq N50 values (“public,” x-axis). Density along the
 333 diagonal indicates many assemblies are similar, while the density along the y-axis indicates
 334 ASRGs with higher N50 value. (E) GC%-content for ASRGs (y-axis) to RefSeq (x-axis). Nearly
 335 all assemblies have less than 0.1% difference in GC-content. (F) Pairwise GC% differences
 336 between ASRGs and comparable RefSeq assemblies for the same strain.
 337

338 **Materials and Methods**

339 *Sample Acquisition and Culture Conditions*

340 All the bacterial cell cultures and genomic DNA used in this study met or exceeded ATCC's
341 quality standards (<https://www.atcc.org/about-us/quality-commitment>), underwent extensive
342 phenotypic and genotypic characterization to ensure accurate strain identification, and were
343 extensively tested for contamination before being accepted for use in this study. ATCC is
344 certified by the ANSI National Accreditation Board (ANAB) to meet both ISO 17034:2016
345 standards as a reference material producer and ISO/IEC 17025:2017 as a testing and calibration
346 reference laboratory. Each bacterial strain included in this study is available from ATCC's
347 biorepository and was authenticated according to protocols executed in accordance with ATCC's
348 quality management system (see above). The specific protocols for each strain varied depending
349 on the specific species in question. In general, strain identification and authentication included
350 assessment of colony morphology, gram staining, culture purity, metabolic profiling, antibiotic
351 susceptibility testing (AST), broad-spectrum biochemical reactivity testing, 16S rRNA gene
352 sequencing, ribotyping, matrix-assisted laser desorption/ionization time-of-flight mass
353 spectrometry (e.g., BioMérieux VITEK MS™ system), and whole-genome next-generation
354 sequencing (NGS). Additional details used for culturing, growth conditions, and authentication
355 of each bacterial strain are available online in each bacterial strain's catalog page at ATCC.org,
356 and by visiting ATCC's Bacterial Cell Culture portal³¹.

357 *DNA Templates and Quality Control*

358 To facilitate the successful NGS library preparation for multiple sequencing platforms (long- and
359 short-read sequences), both high-quality and high-quantity input DNA was obtained from
360 authenticated genomic DNA (gDNA) available in ATCC Bacterial Nucleic Acids repository³².
361 ATCC uses several commercially available extraction kits and in-house validated protocols to
362 obtain pure high-molecular-weight DNA depending on the biological characteristics of the
363 organism undergoing extraction. For strains with no preexisting genomic DNA in ATCC's
364 repository, total high molecular weight genomic DNA (HMW gDNA) was extracted from
365 thawed or resuspended frozen cultures with 10^7 - 10^9 cells/mL using the QIAGEN Genomic-
366 Tip™ 20/g or 100/g kit and analyzed for purity, concentration and fragment size. HMW-gDNA
367 samples meeting or exceeding the following criteria were subjected to sequencing; median

368 fragment size larger than 20 kb, optical density A260/280 between 1.75 – 2.00, and a final
369 elution concentration over 20ng/μL per extraction.

370 ***Short-Read Next Generation Sequencing***

371 High-quality gDNA from each strain was subjected to whole-genome sequencing using a short-
372 read next generation sequencing (NGS) workflow. Briefly, sequencing libraries from each
373 extraction were prepared using the DNA Prep kit and indexed using DNA/RNA UD indexes
374 (Illumina), and subsequently subjected to paired-end sequencing on either an Illumina MiSeq®
375 or NextSeq 2000® instrument. Sample multiplexing was based on achieving a minimum 100X
376 average depth of coverage for each genome. Base-calling and adapter trimming was initially
377 done using onboard Illumina instrument software and followed by an additional round of
378 trimming and quality-score filtering using *fastp* and *FastQC*^{33,34}. Illumina reads accepted for
379 further use passed the following quality control thresholds: median Q score, all bases > 30,
380 median Q score, per base > 25, ambiguous content (% N bases) < 5%.

381 ***Long-Read Next Generation Sequencing***

382 Long-read sequencing was carried out using the Oxford Nanopore Technologies (ONT) GridION
383 platform. ONT Ligation Sequencing Kit (Oxford Nanopore, UK, SQK-LSK109) sequencing
384 libraries were prepared from the same physical samples of HMW gDNA used for Illumina
385 sequencing above, multiplexed using the ONT Native Barcoding Expansion kit (Oxford
386 Nanopore, UK, EXP-NBD104 or EXP-NBD114), and sequenced using GridION flow cells
387 (Oxford Nanopore, UK, R9.4.1). As with Illumina sequencing, the number of samples
388 multiplexing was based on the estimated genome size of a given organism and sequencing was
389 performed for a minimum of 48 hours per flow cell. Using the most up to date version of
390 *MinKNOW*, reads were base-called, using the high accuracy settings, demultiplexed, and barcode
391 trimmed. Furthermore, ONT sequencing reads were quality trimmed and filtered using *Filtlong* to
392 meet the following minimum acceptance criteria: minimum mean Q score per read > 10,
393 minimum read length > 5000 bp³⁵.

394 ***Assembly of ATCC Standard Reference Genomes***

395 For genome references deposited to the ATCC Genome Portal, genome assembly size was first
396 estimated from raw reads using *MASH*, and this estimate was used to down-sample the Illumina
397 and ONT raw sequencing libraries to a maximum 100x and 40x coverage respectively³⁶. These

398 coverage requirements were selected to maximize accuracy for individual consensus base-calls
399 in the final assemblies^{20,22}. After down-sampling each sequencing library, a hybrid *de novo*
400 assembly approach was taken using *Unicycler*²³. Briefly, Illumina libraries were first assembled
401 individually into contigs. The longest contigs in the initial set were then scaffolded with reads
402 from the ONT library. The combined hybrid-assembly was then iteratively polished using both
403 long and short reads from both input libraries, resulting in highly contiguous or closed reference
404 genomes. Sequencing and assembly artifacts of less than 1000 bp that also had significantly
405 different coverage depth (e.g., “chaff” contigs) were removed from the final draft reference³⁷.
406 These draft assemblies were subsequently checked using One Codex to confirm the species¹⁹.
407 Finally, each draft assembly was assessed for completeness and potential contamination with
408 *CheckM* v1.12, which is based on orthologous gene copy numbers present in an assembly²⁴.
409 Assemblies which were determined to have a *CheckM* “completeness” score above 95% and a
410 contamination value below 5% were deemed final assemblies. Each final assembly was
411 subsequently annotated using *Prokka* v1.14 for CDS, rRNA, tRNA, signal leader peptides, and
412 non-coding RNA identification³⁸. Finally, each complete and annotated genome was deposited
413 into the ATCC Genome Portal and is referred to herein as an ATCC Standard Reference Genome
414 (ASRG)³⁹.

415 ***Characterization of Public Genome Assemblies***

416 To gather the public assemblies of ATCC bacterial strains, the “assembly_summary_refseq.txt”
417 file was downloaded from the NCBI Bacterial RefSeq ftp site
418 (<https://ftp.ncbi.nlm.nih.gov/genomes/refseq/bacteria/>). This file contains accession numbers and
419 metadata, such as “Isolate”, “Assembly Level”, and “Tax ID,” for every assembly in NCBI
420 Bacterial RefSeq. First, this file was filtered to keep all records that contained either the “ATCC”
421 or “NCTC” keyword. This was done because many strains have synonymous ATCC and NCTC
422 IDs, though often only one of the two is present in a record. Of the records containing “ATCC”
423 or “NCTC,” all that included the “ATCC” were kept, but records containing “NCTC” were
424 filtered to keep only those with a synonymous ATCC ID. This final set of records contained the
425 2,701 public assemblies of ATCC strains. While “assembly_summary_refseq.txt” does contain
426 metadata, the complete set of metadata was collected by downloading the “assembly_report.txt”
427 for each assembly from the NCBI ftp site. Metadata comparisons were performed using the
428 *compare.all.levels.py* script after appending the RefSeq assembly data with a GC content

429 column, calculated by *bbnorm_stats.sh*, all of which was paralleled with GNU Parallels⁴⁰.
430 ATCC's Genome Portal does not distinguish between contigs and scaffolds, which RefSeq
431 defines as contigs that are connected across gaps. For this, all data comparing ASRGs in terms of
432 contiguity uses RefSeq scaffold information.

433 ***Comparisons of NCBI and ATCC Genome Assembly Metrics***

434 For each of the bacterial strains included in the ATCC Genome Portal, we identified and
435 downloaded all 2,701 genome assemblies that had the same name or similar names from NCBI's
436 RefSeq and Genome Assembly databases. For the 303 NCBI assemblies with a finished assembly
437 status of "Complete" or "Chromosome" and representation in ATCC's Genome Portal, we
438 carried out pairwise whole genome alignments for each NCBI and ASRG using *MUMmer4* and
439 its associated suite of tools for comparative genomics²⁷. In some cases, due to duplications in
440 RefSeq and NCBI's Genome Assembly database, multiple NCBI assemblies were compared
441 against the same ASRG assembly. Following the creation of the alignments, we identified
442 genome-wide variants for each NCBI assembly as compared to the ASRG assembly, including
443 single nucleotide polymorphisms (SNPs), insertions and deletions (InDels), and structural
444 variants (SV). Genome-wide comparisons using *dnaDiff* v1.3 included assembly length, number
445 of contigs, pairwise percent aligned, and N50 values⁴¹ (*SVs_and_ANI.sh*). Furthermore,
446 *MUMmer4's dnadiff* tool was run with default settings using the ASRG assemblies against each
447 NCBI RefSeq assembly, and relocations, translocations, and inversions are reported alongside
448 total and aligned bases²⁷. Prior to running *MUMmer4's dnadiff* tool on these assemblies, each
449 was filtered to remove contigs <1kb in length to prevent short sequences from exaggerating SVs
450 between assemblies. Structural variants included breakpoints, relocations, translocations, and
451 inversions, and summarized as rearrangements.

452
453 **Data Availability:** ATCC Standard Reference Genomes (ASRGs), metadata, and raw (FASTQ)
454 data are subject to controlled access, but may be used for any non-commercial research-use only
455 purposes by meeting the requirements outlined below. Data can be obtained directly from the
456 ATCC Genome Portal (<https://genomes.atcc.org>), via our REST-API (access and details
457 available upon request), or via URLs found in [https://github.com/ATCC-Bioinformatics/AGP-
458 Raw-Data/blob/main/AGP_Raw-Data-Access.txt](https://github.com/ATCC-Bioinformatics/AGP-Raw-Data/blob/main/AGP_Raw-Data-Access.txt). Downloading these data requires a ATCC
459 Web User Profile (<https://www.atcc.org/web-profile/create-a-web-profile>) and acceptance of

460 ATCC's Data Use Agreement (<https://www.atcc.org/policies/product-use-policies/data-use->
461 [agreement](#)). Any commercial use of ATCC genomics data requires express permission of ATCC
462 (please contact licensing@atcc.org for details). MIT Licensed, open-source code for scripts used
463 in this manuscript are available at <https://github.com/ATCC->
464 [Bioinformatics/Equivalency Analysis](#) .

465 **Acknowledgements:** We thank One Codex for contributing to the development of the ATCC
466 Genome Portal. We also thank Drs. Raymond Cypess, Mindy Goldsborough, and Rebecca
467 Bradford for critical comments and review prior to submission.

468 **Author Contributions:**

469 Conceptualization: JGL, BB, JB, JLJ

470 Data curation: DAY, JGL, NPP, PFC, ALR, MAR

471 Formal Analysis: DAY, NPP, PFC, ALR

472 Investigation: CT, AEP, JD, SRG, SK, RM, BB

473 Project administration: BB, JB

474 Software: DAY, NPP, PFC, ALR, JB

475 Supervision: BB, JB, JLJ

476 Visualization: PFC, BB, JLJ

477 Writing – original draft: DAY, JGL, JLJ

478 Writing – review & editing: DAY, NPP, PFC, BB, JB, MAR, JLJ

479 **Funding:** The work described herein was financially supported entirely by the American Type
480 Culture Collection.

481 **Competing Interests:** All authors are employees of the American Type Culture Collection,
482 which solely funded the work presented here and provided all the bacterial strain materials
483 needed for the research. No other competing interests are claimed.

484 **Supplemental Data:**

485 All Supplemental Tables are found in a single MS Excel document, with each worksheet labeled
486 accordingly for each table.

487 **Supplementary Materials**

488 Supplementary Information is available for this paper. Figs. S1 to S4, and Tables S1 to S8 (as
489 separate Excel document).

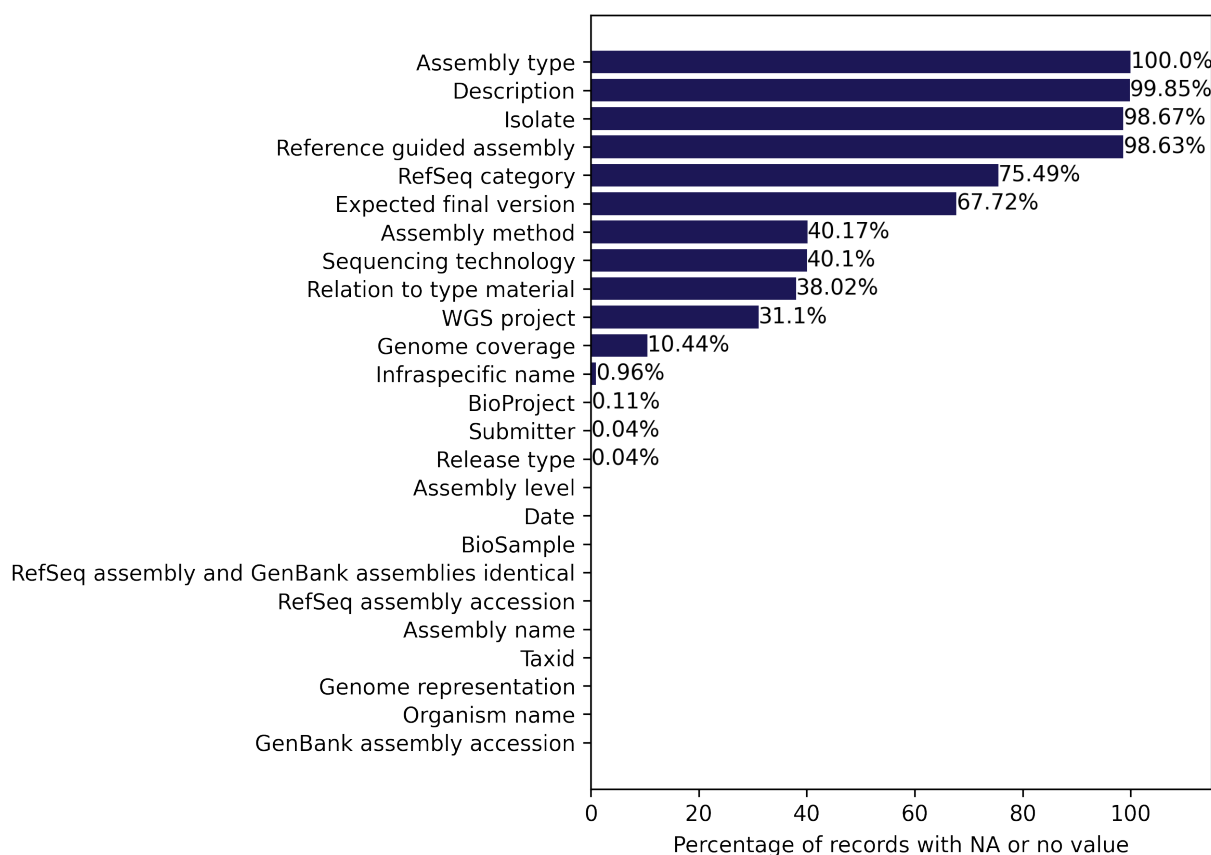
490

491 Correspondence and requests for materials should be addressed to Jonathan L Jacobs

492 (jjacobs@atcc.org).

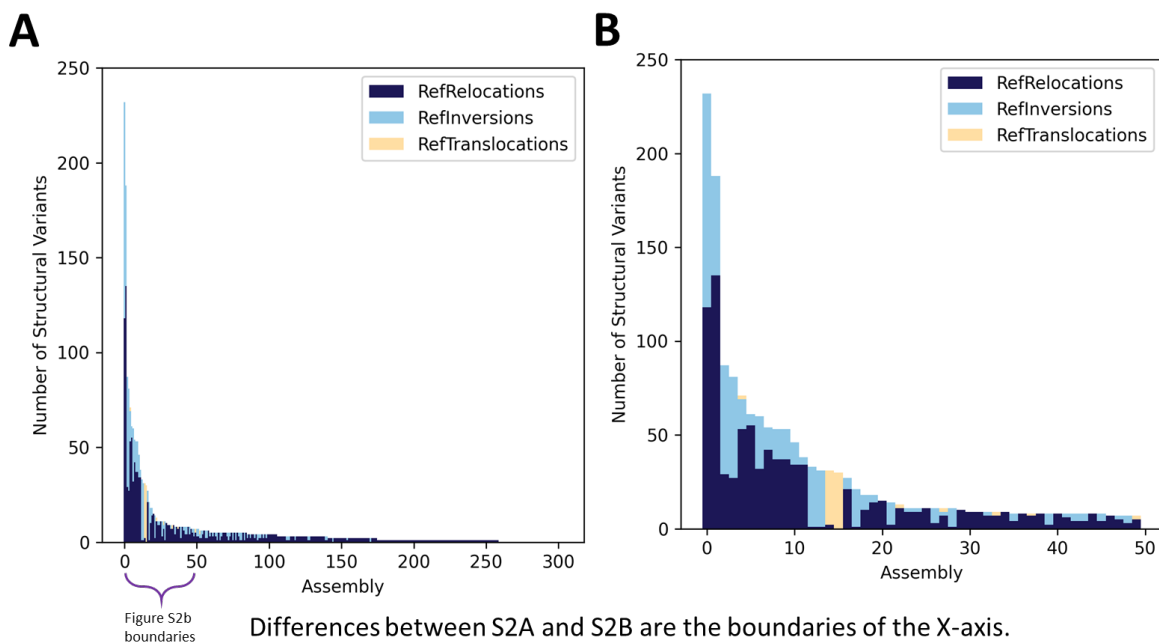
493

Percentage of unused or empty fields in RefSeq assembly reports



495 **Fig. S1.** Bar chart demonstrating the percentage of RefSeq assembly report fields that are left
496 empty or contain “na” as a value. While some of these, such as RefSeq category, have implicit
497 definitions for empty fields, others, such as Relation to type material, are potentially crucial
498 pieces of information.

499

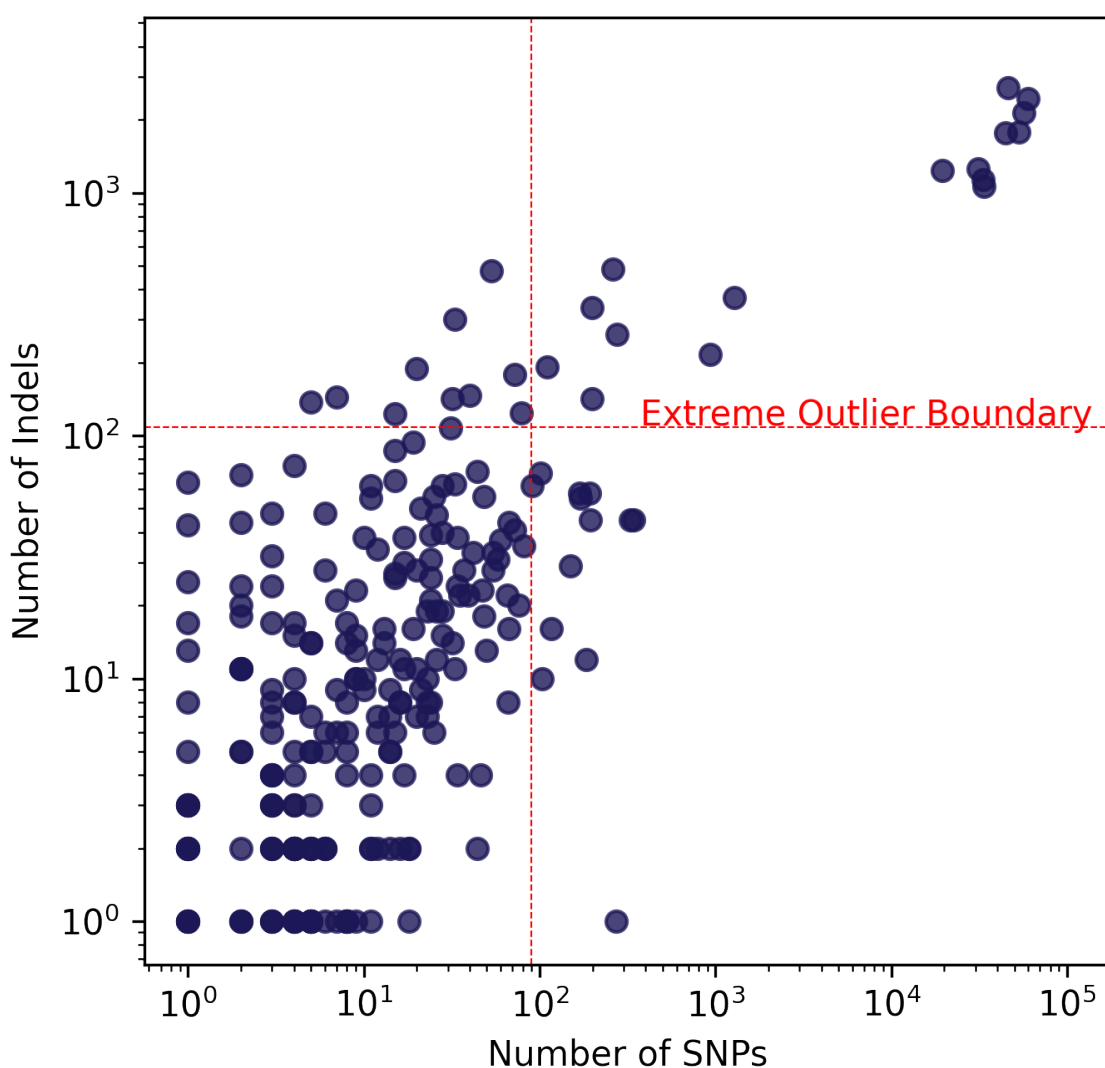


501

502 **Fig. S2. (A)** Stacked bar chart showing relocation, inversion, and translocation structural variants
503 between all ATCC assemblies and assemblies generated from mapping ATCC's read data of
504 specific strains to assemblies of those strains. **(B)** Stacked bar chart showing relocation,
505 inversion, and translocation structural variants between ATCC assemblies and assemblies
506 generated from mapping ATCC's read data of specific strains to assemblies of those strains, for
507 the 50 ATCC products with the greatest total of structural variants.

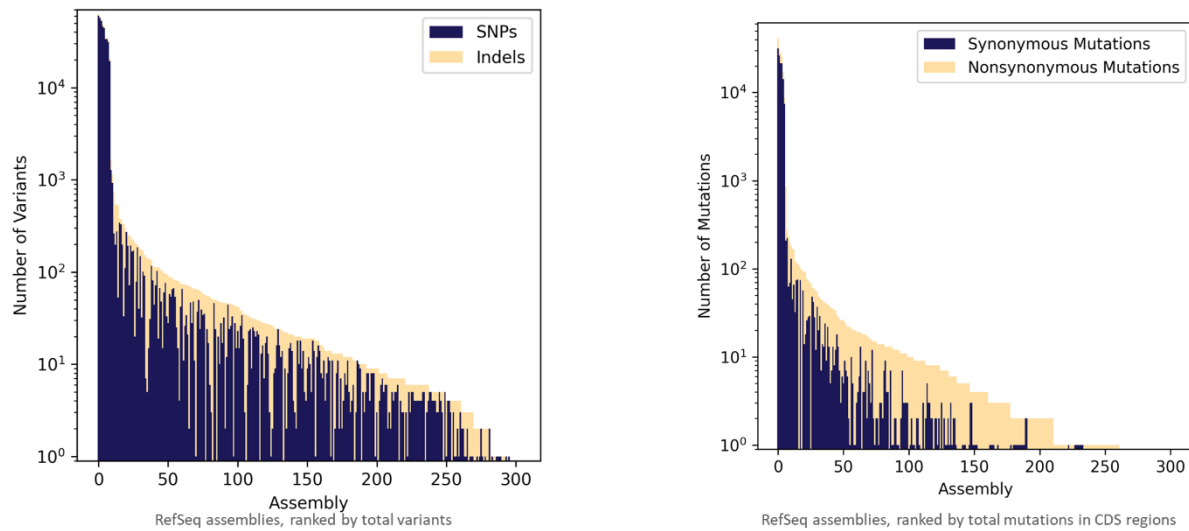
508

509



511 **Fig. S3. Single-nucleotide polymorphisms (SNPs) and insertions/deletions (InDels) of ASRG**
512 **raw data mapped to RefSeq references.** Each data point represents a read-mapping of ASRG
513 raw data (Illumina only) to a RefSeq genome assembly for the same bacterial strain. In cases
514 where multiple RefSeq assemblies exist for the same bacterial strain, ASRG reads were mapped
515 to each and is represented above by multiple data points. The extreme outlier boundary (red) is
516 determined is 3x the interquartile range above median for both SNPs and Indels (See Materials &
517 Methods).
518

519



521 **Fig. S4. A visualization of the number and types of variants found when mapping the**
522 **trimmed Illumina reads for an ATCC product to its corresponding RefSeq**
523 **assembly/assemblies. (A) The total number of variants, the number of SNPs, and the number of**
524 **indels found across this mapping. (B) The total number of variants and the characterization of**
525 **those variants into either Synonymous or Nonsynonymous, as determined by VEP. Synonymous**
526 **variants represent alterations to a coding sequence that does not change the amino acid upon**
527 **translation. Nonsynonymous variants represent alterations to a coding sequence that does change**
528 **the amino acid upon translation. Variants outside of coding regions were calculated as well, but**
529 **are not shown here.**

530