



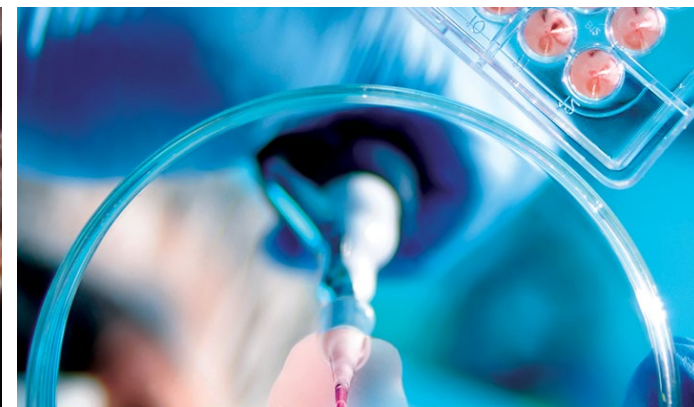
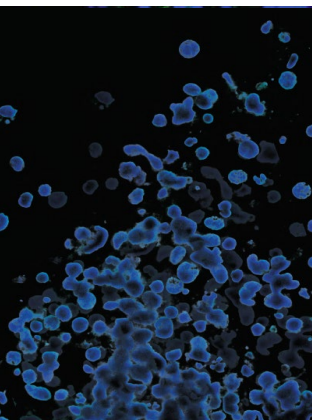
# The ATCC Genome Portal

Authenticity and Traceability for Microbial Genomes

World Microbe Forum 2021  
Industry & Science Symposium

Jonathan Jacobs, PhD  
Senior Director, Bioinformatics  
Principal Scientist  
Sequencing and Bioinformatics Center, ATCC

Credible Leads to Incredible™



# Overview

- **The ATCC Genome Portal**
- Traceability and authentication of microbial genomes
- Standards for authenticated reference genomes
- Development roadmap preview



# The ATCC Genome Portal

The ATCC Genome Portal is a cloud-based platform that enables users to easily browse genomic data and metadata by simply logging into the portal



Download whole-genome sequences and annotations of ATCC materials

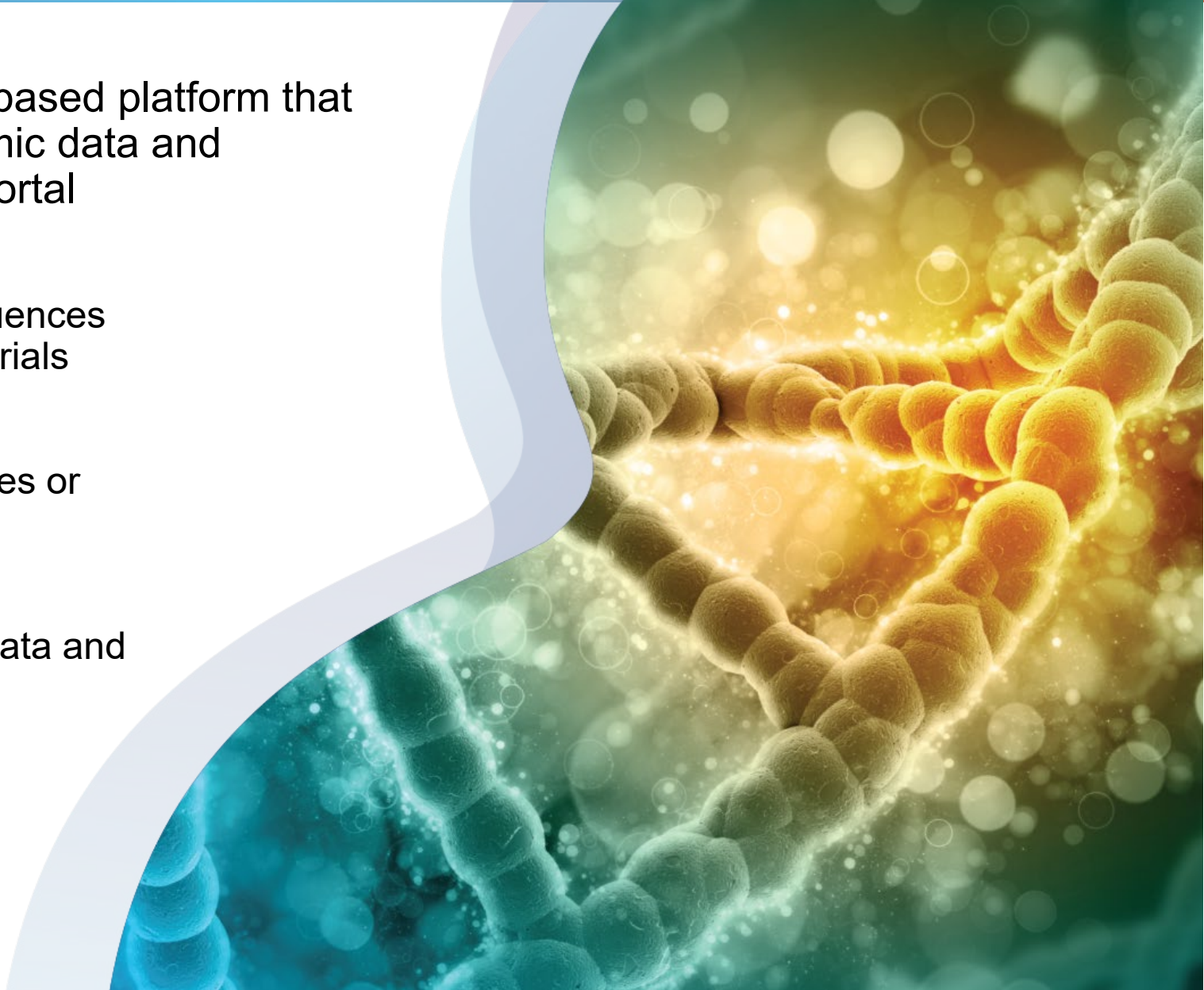


Search for nucleotide sequences or genes within genomes



View genome assembly metadata and quality metrics

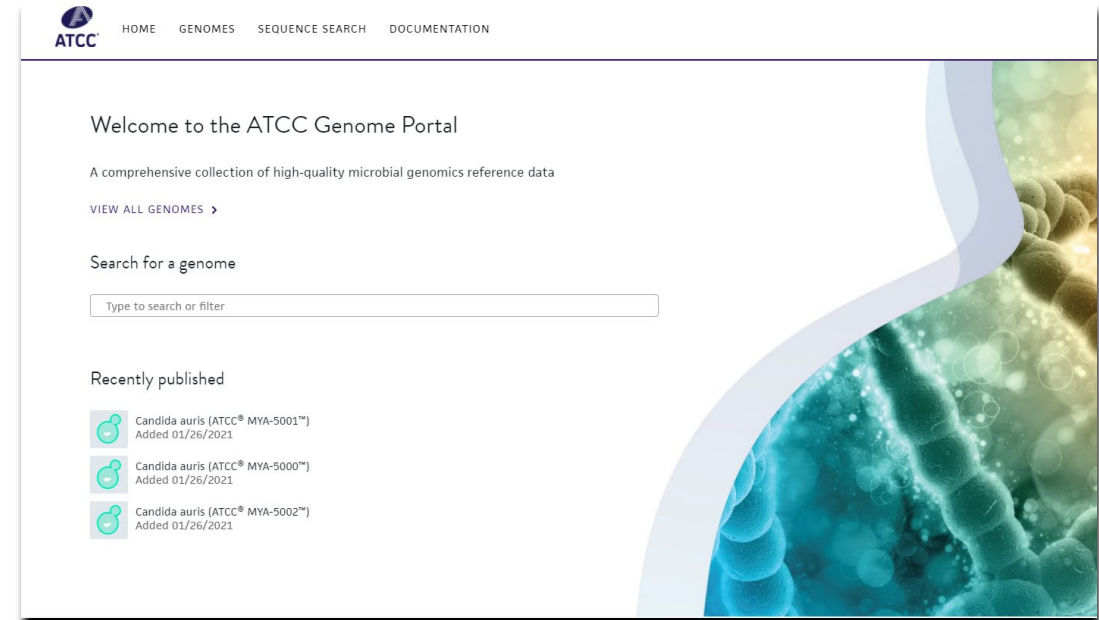
**[genomes.atcc.org](https://genomes.atcc.org)**



# The ATCC Genome Portal

*ATCC's authenticated reference genomes*

- **2017 – 2018** – Planning and proof-of-concept experiments
- **2018** – ATCC Commitment
  - Laboratory, staffing, resources, instrumentation, and bioinformatics development
- **2019** – ATCC Enhanced Authentication Initiative
  - June 2019 – beta launch at ASM Microbe 2019
  - Sept 2019 – formal launch of the ATCC Genome Portal
- **2020+** – Expansion
  - 1200+ authenticated, reference-grade whole-genome assemblies
  - Inclusion of viral, bacterial, and fungal genomes



# Providing reference-quality genomes

## Why - Challenge # 1



- Public databases routinely host genomic data that is cited as “ATCC,” but...
  - Often no traceability back to genuine ATCC cultures
  - ATCC cannot authenticate 3rd party genomes in public databases
- So, how do researchers \*know\* which data set to use?
  - Which is the “correct” one?
  - Close enough?
- How do researchers have confidence in their selection?

# Providing reference-quality genomes

## Why - Challenge # 2



- How do we bring authentication into the genomics era while maintaining our commitment to our customers that we've fully and accurately authenticated our material?
- Typically, authentication\* may refer to:
  - Morphology
  - Purity
  - Viability
  - Phenotypic testing
  - Genotypic testing
    - 16S ribosomal gene
    - ITS and D1D2

**ATCC** **CERTIFICATE OF ANALYSIS**

ATCC\* Number: 12228D-5™  
Lot Number: 70006950  
Designation: Staphylococcus epidermidis genomic DNA  
Fill Volume Prior to Drying: 80 µL  
Product Format: Dried microbial DNA  
Expiration Date: Not applicable  
Storage Conditions: 2°C to 8°C

Test / Method	Specification	Result
OD <sub>600</sub> /OD <sub>660</sub> ratio (Spectrophotometer method)	1.6 to 2.1	1.8
Total amount of DNA (PicoGreen® measurement)	≥ 5 µg per vial	6 µg/vial
Agarose gel electrophoresis	High molecular weight chromosomal DNA; No visible RNA	High molecular weight chromosomal DNA; No visible RNA See photograph below
PCR Functionality	Successful PCR amplification of selected gene(s)	Successful PCR amplification of selected gene(s)
Sequencing of selected gene(s)	Consistent with source organism	Consistent with source organism

1 2  
Lane 1: Invitrogen™ TrackIt™ 1 Kb Plus DNA Ladder  
Lane 2: 12228D-5™

Quality Assurance Specialist: Quality Assurance  
ATCC hereby represents and warrants that the material provided under this certificate is pure and has been subjected to the tests and procedures specified and that the results described, along with any other data provided in this certificate, are true and correct to the best of our knowledge.

ATCC 12228D-5™ or 703-365-2700  
12228D-5™ or 703-365-2700  
Manassas, VA 20115-2209 USA Fax: 703-365-2700  
www.atcc.org Email: tech@atcc.org  
or contact your local distributor

- Page 1 of 2 -  
Template Revision: 4  
Template Effective Date: 12/15/2014

\*not an inclusive list

# Providing reference-quality genomes

Why - Challenge # 3

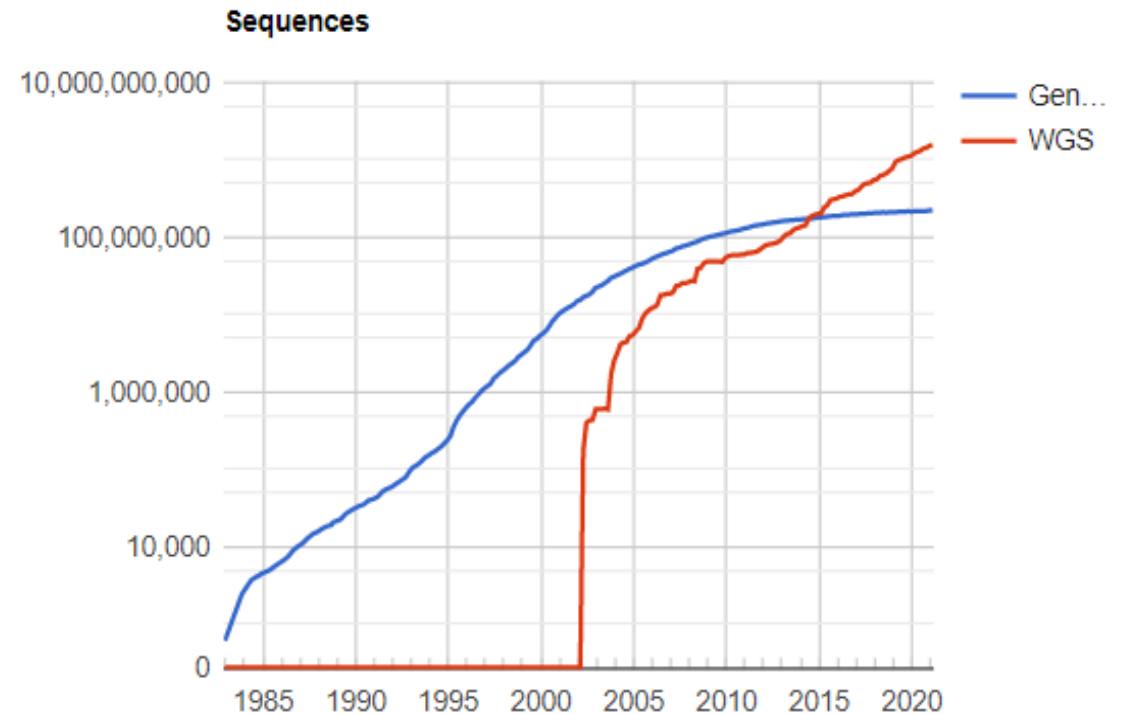


- Acknowledge there is a problem with reference genomes ✓
- Work through a plan to address the problem ✓
- **How do we effectively and easily provide customers with genomic data while not diluting it or burying it in a public database?**

# Reference genomes

Where can researchers turn to for “reference” genomes?

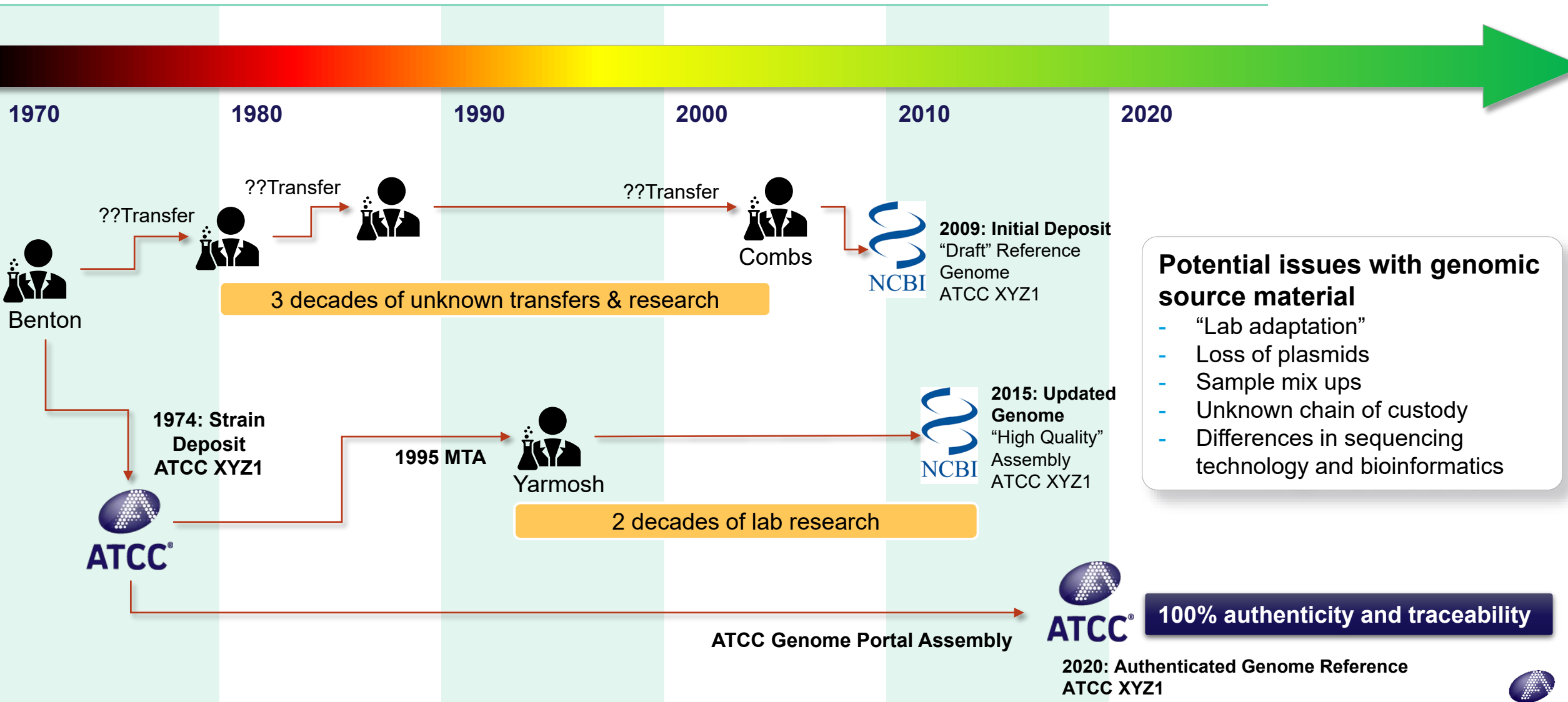
- De facto standard
  - The sequence database for the entire public scientific community
  - Contains numerous genomes
  - Genomes submitted by a variety of labs
- Relatively little curation
- Highly variable quality
- **NEVER** authenticated by ATCC



<https://www.ncbi.nlm.nih.gov/genbank/statistics/>

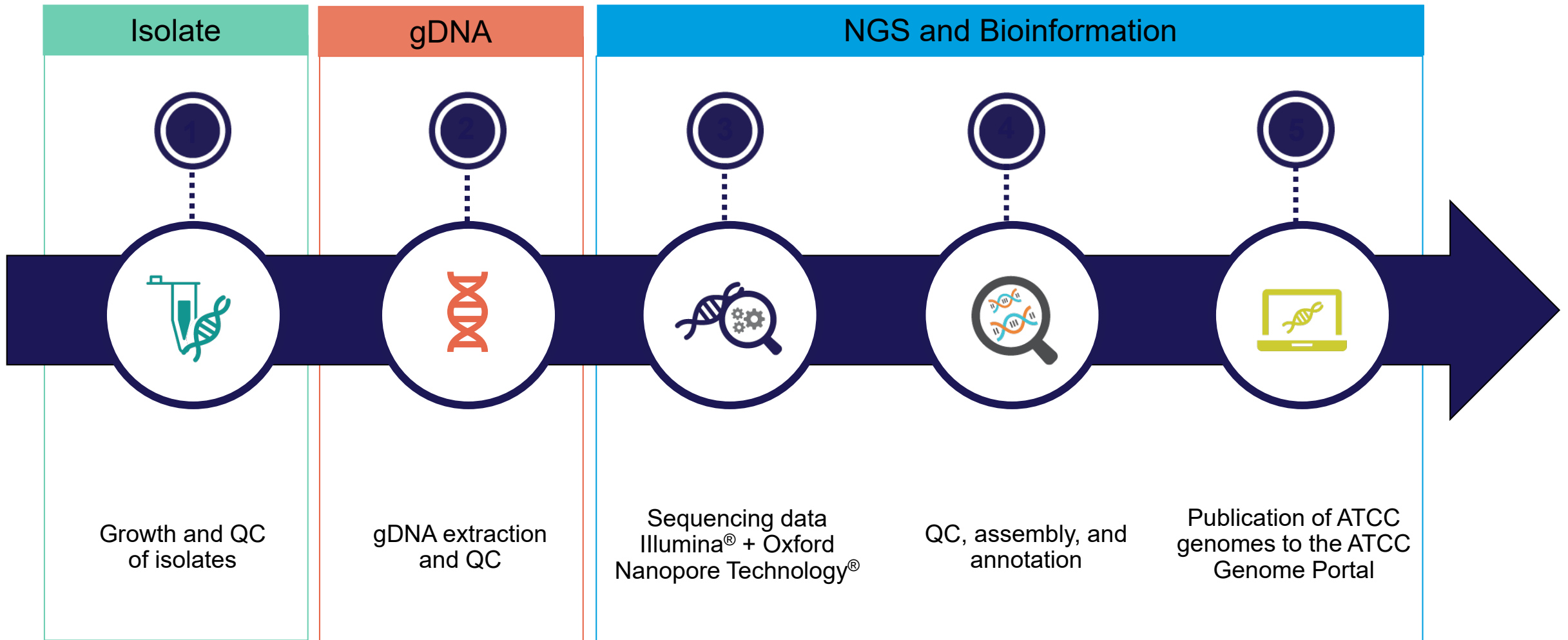


# The elephant in the room: Authenticated reference genomes

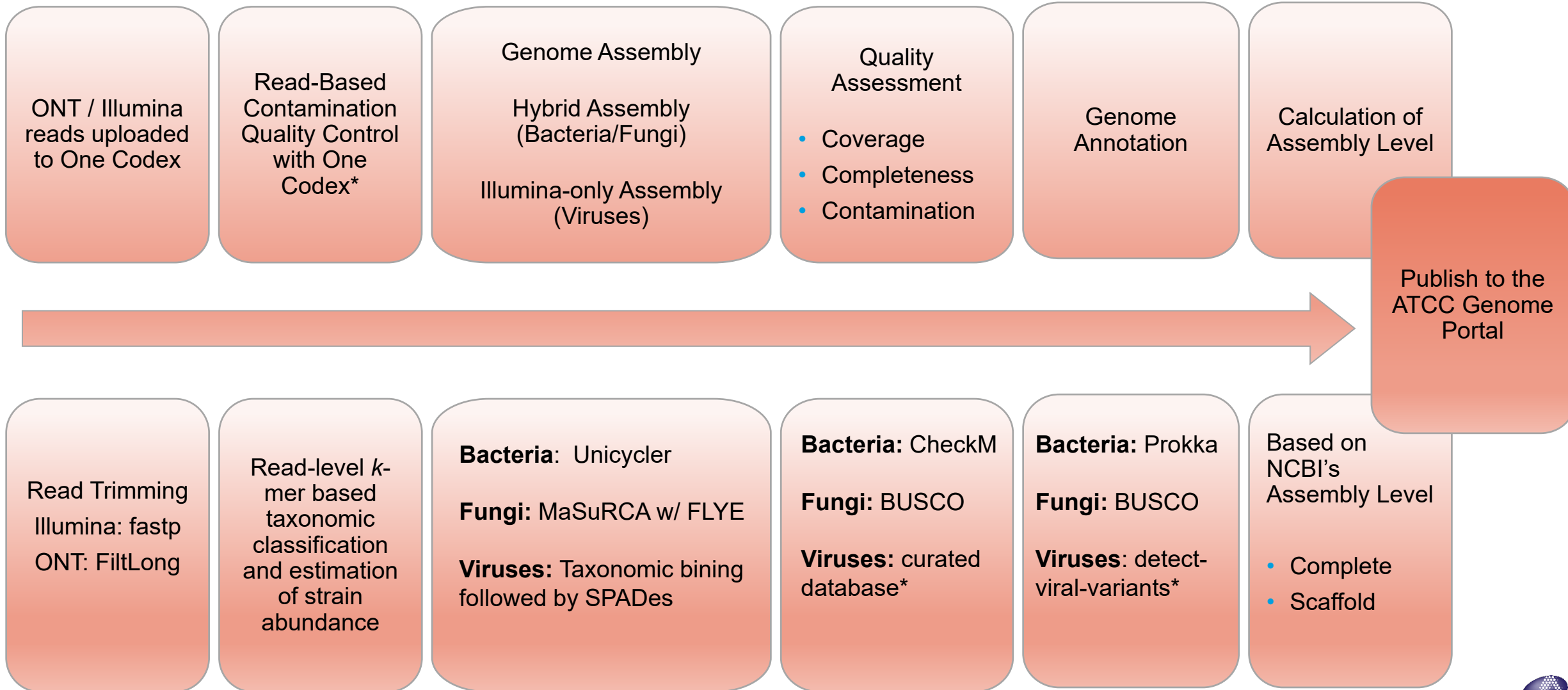


- Potential issues with genomic source material**
- "Lab adaptation"
  - Loss of plasmids
  - Sample mix ups
  - Unknown chain of custody
  - Differences in sequencing technology and bioinformatics

# Our process: Authenticated physical material coupled with reference-quality genome sequences



# ATCC genome assembly process

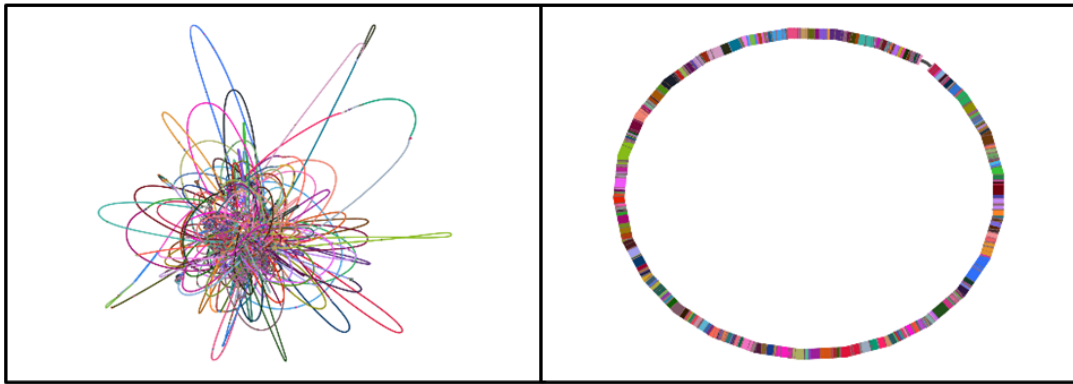


\* One Codex proprietary software

# Advantage of hybrid assemblies

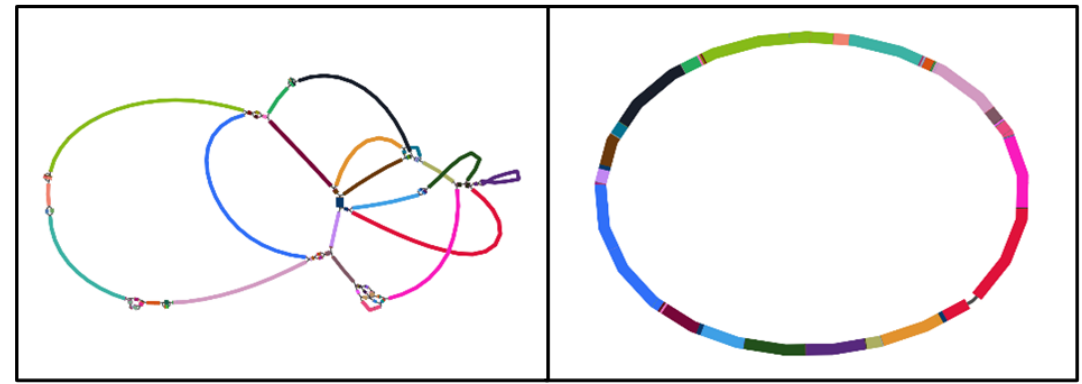
Illumina-only assembly      Hybrid assembly

*Neisseria meningitidis* (ATCC® 53417™)

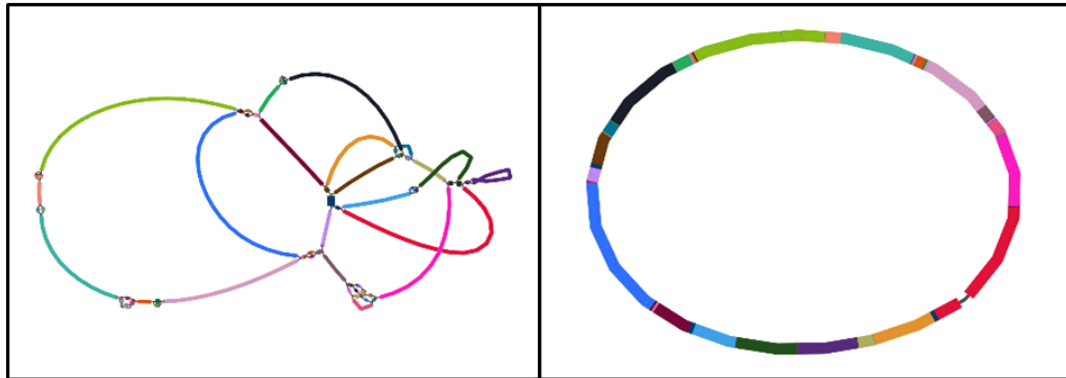


Illumina-only assembly      Hybrid assembly

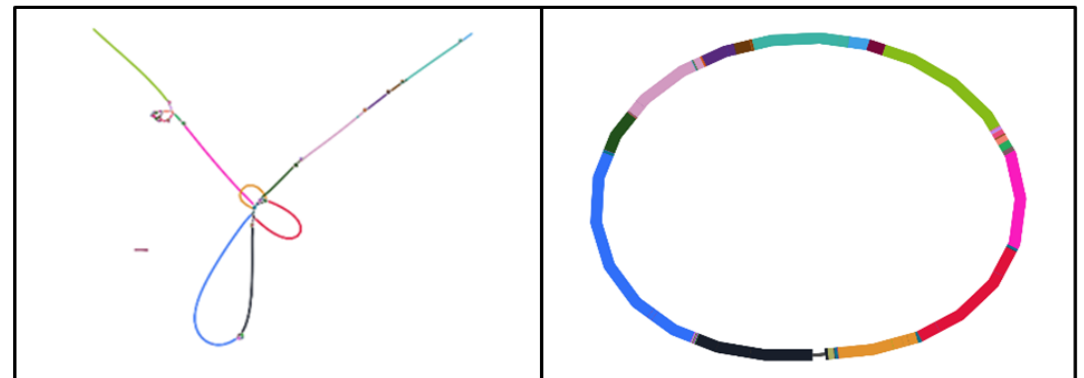
*Campylobacter jejuni* subsp. *jejuni* (ATCC® 43446™)



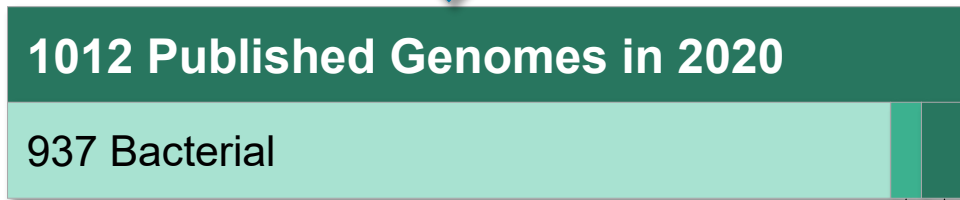
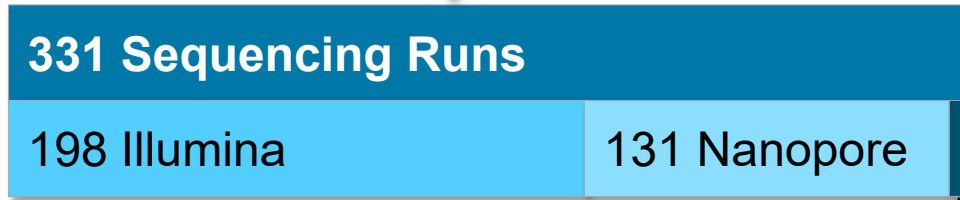
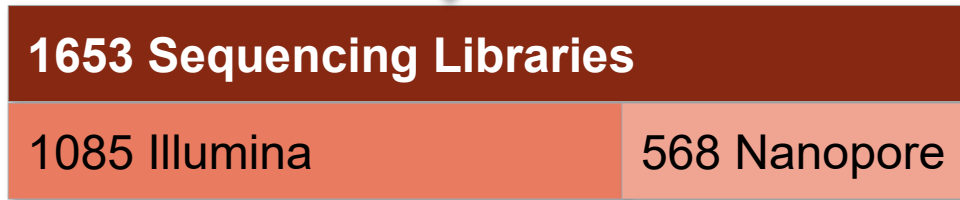
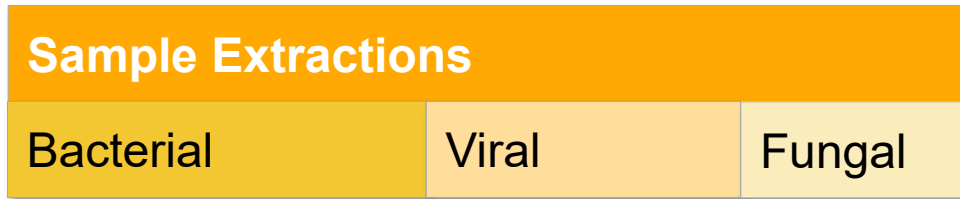
*Campylobacter jejuni* subsp. *jejuni* (ATCC® 43446™)



*Streptococcus gordonii* (ATCC® 35105™)



# The ATCC Genome Portal



**1,251**  
Authenticated  
Genomes

← Today

- 1,118 bacterial genomes (739 complete circularized, 391 type strains)
- 59 viral genomes
- 74 mycology genomes

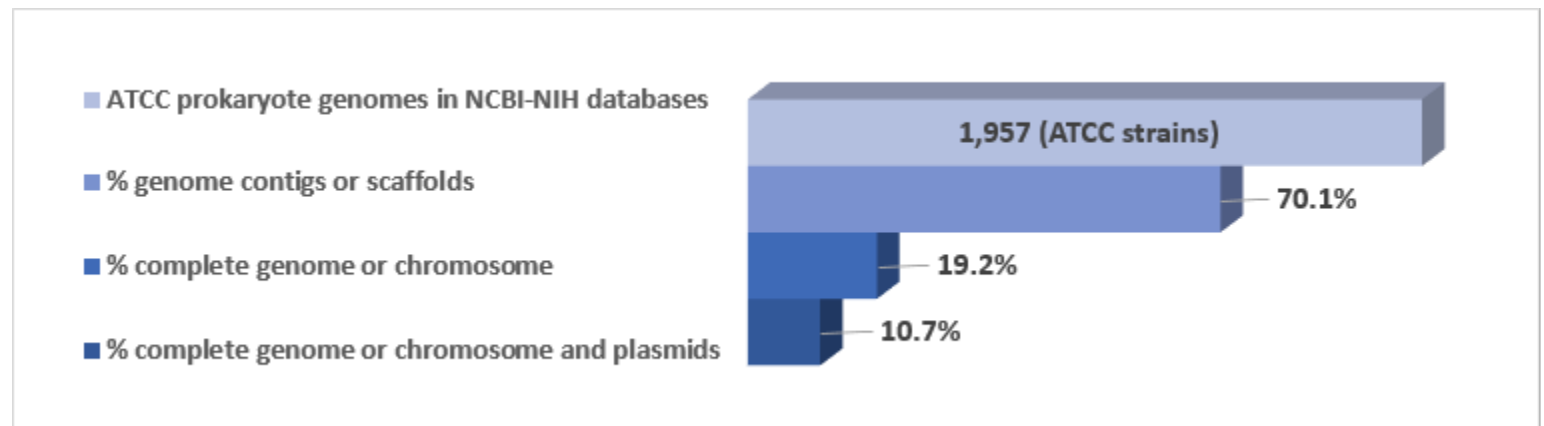
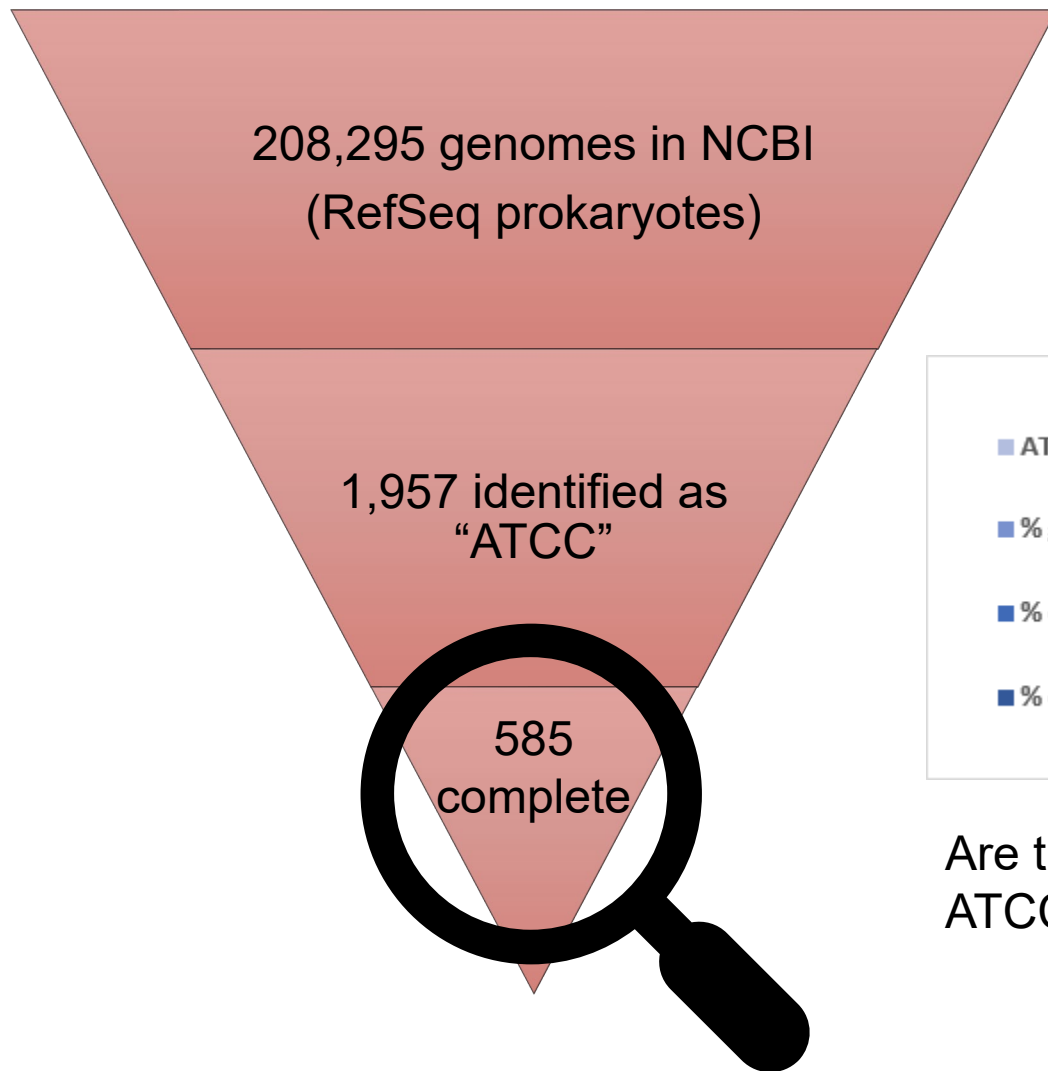
- Monthly updates
- All genomes are traceable to ATCC's biomaterials
- Hybrid assemblies for all bacterial & fungal genomes
- All genomes annotated
- Additional improvements to fungal and viral genome annotations coming

# Overview

- The ATCC Genome Portal
- **Traceability and authentication of reference genomes**
- Standards for authenticated reference genomes
- Development roadmap preview



# Reference genomes

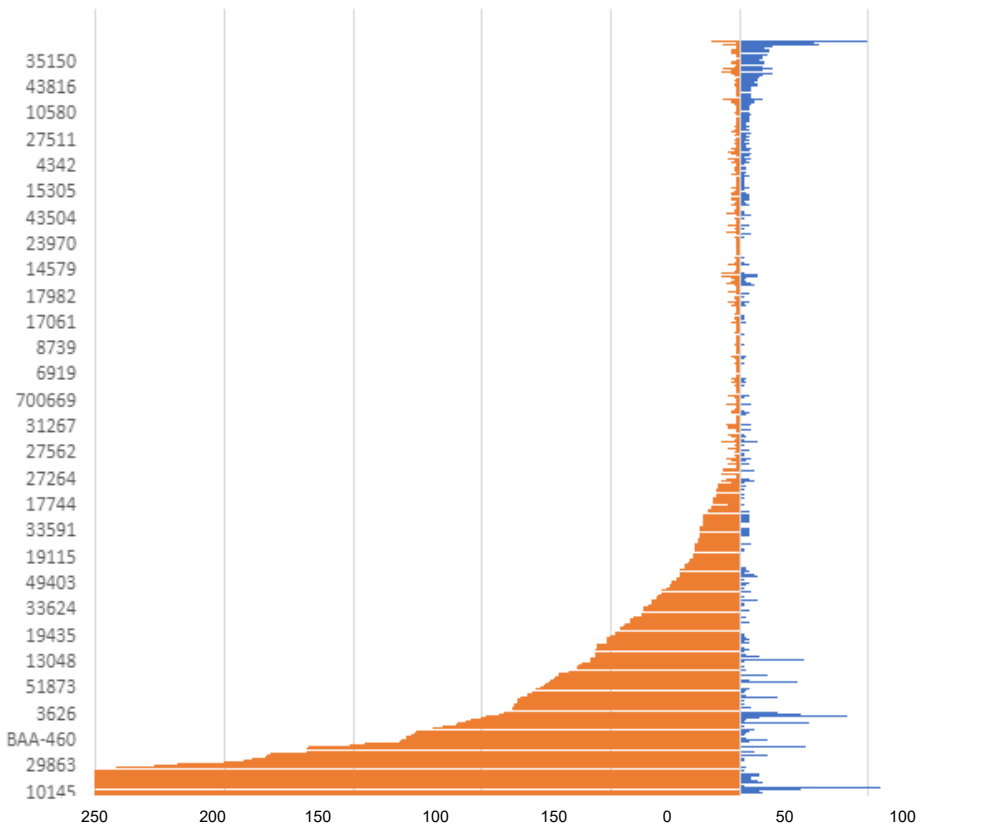


Are these 585 RefSeq genomes traceable back to authenticated ATCC cultures with well-documented growth and storage conditions?

# Genome assembly quality

*Equivalency analysis of ATCC Genome Portal assemblies vs. those from public databases*

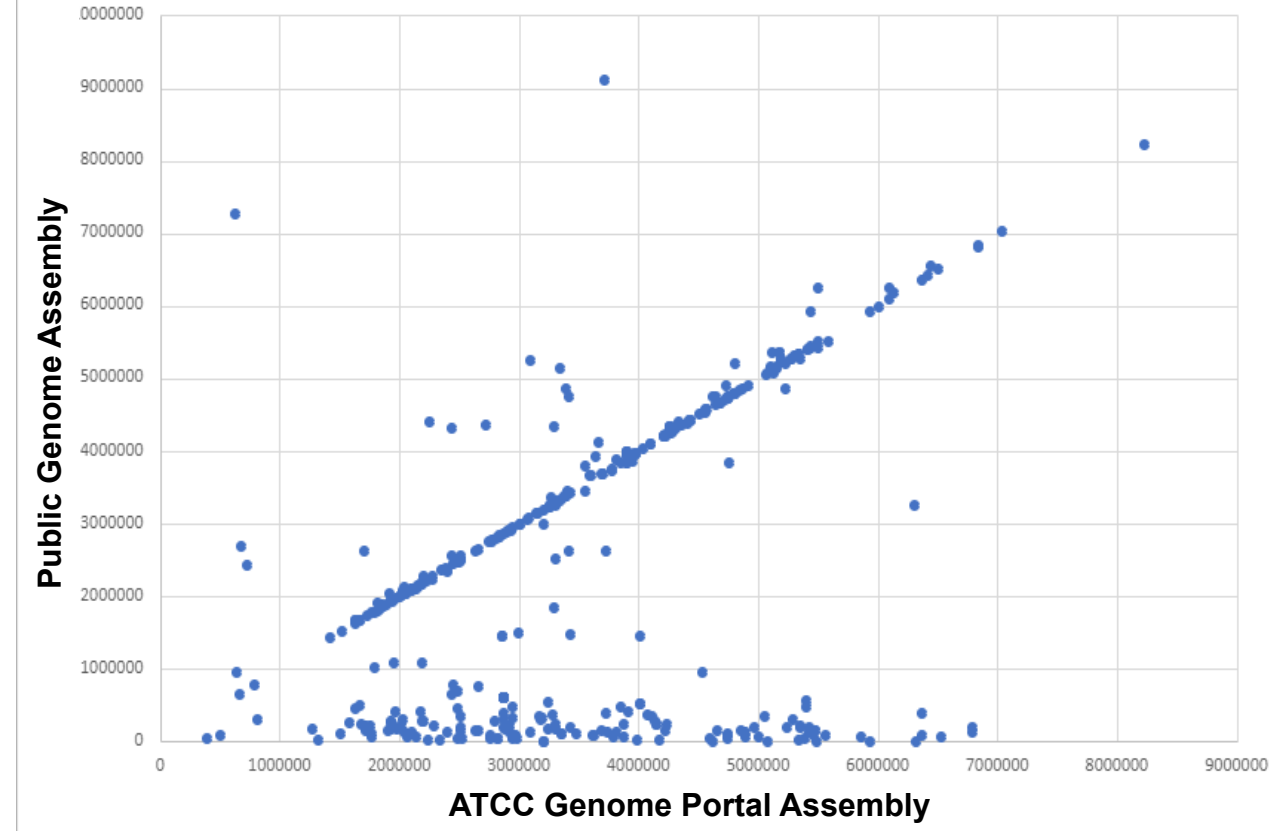
Bacteriology Products Contig Counts



Number of Contigs per assembly

■ #PUBLIC DATA      ■ ATCC GENOME PORTAL

Bacteriology Assembly N50 Comparisons



The downward trend in contig count and the upward trend in N50 indicate the ATCC produced genomes are of higher quality



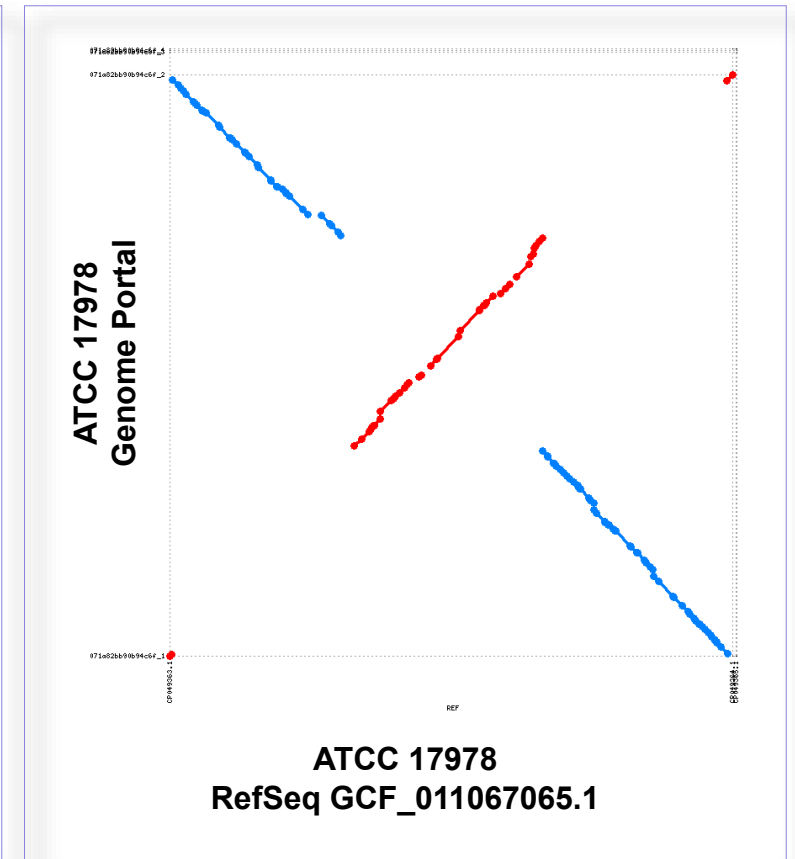
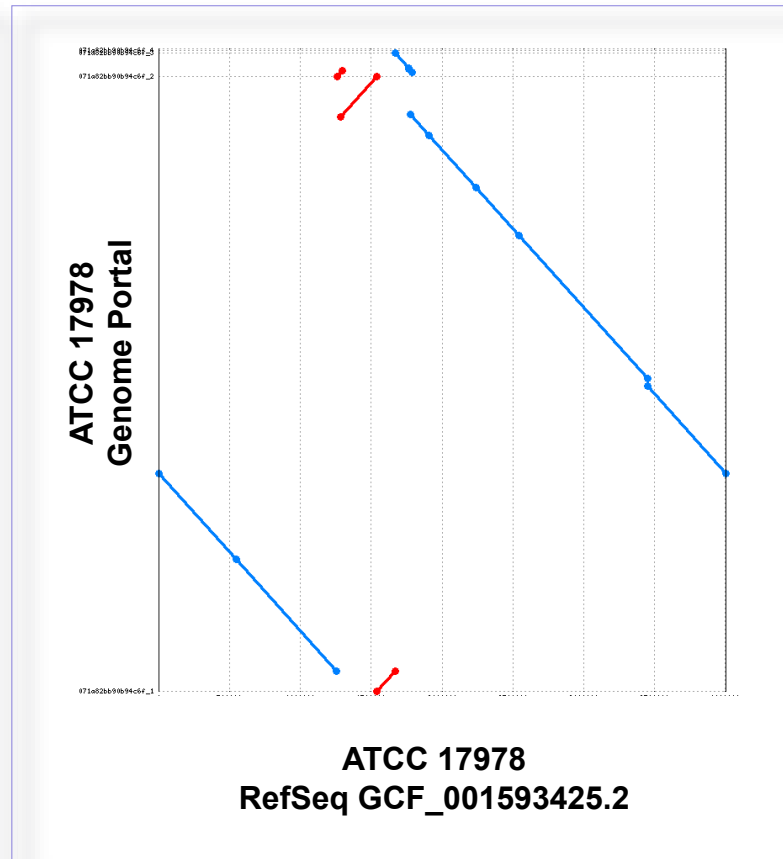
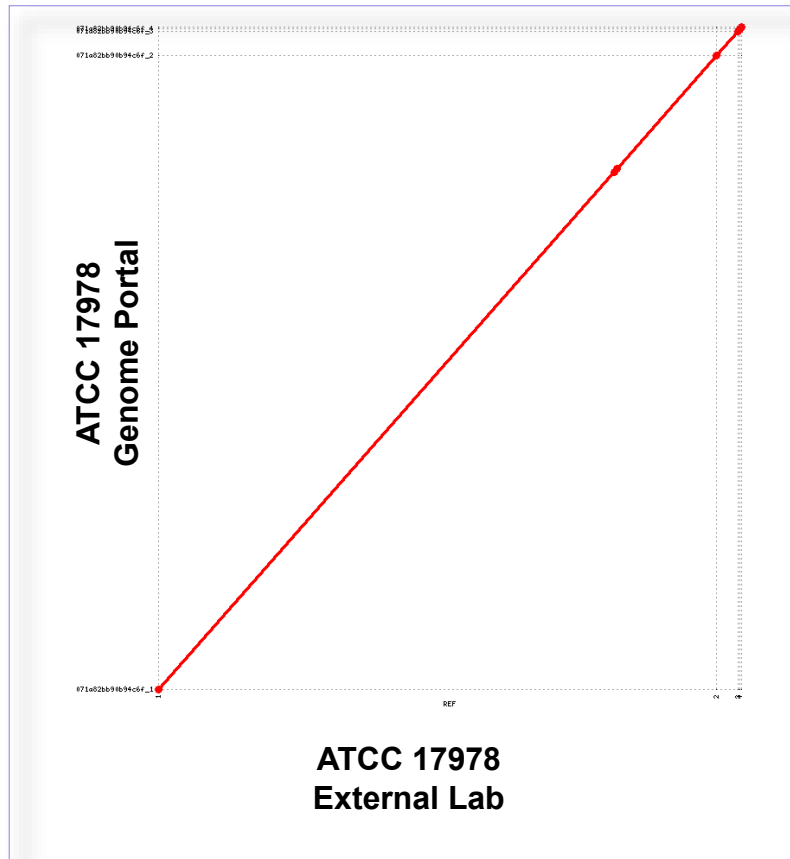
# Evaluation of genome sequences from public databases

Product	NCBI existing reference genomes	NCBI assembly level (plasmids)	Sequencing technology and coverage	# of SNPs	# of indels	Average coverage (variants)	
<i>Acinetobacter baumannii</i> (ATCC® 17978™)	GCA_001593425.2	Complete Genome	Illumina (300.0x)	14	5	210.1	1 strain
	GCA_000015425.1*	Complete Genome (2)	Not available	118	656	152.7	
	GCA_014672775.1	Complete Genome (1)	PacBio (399.24x)	15	87	170.4	7 assemblies Unknown origin of materials
	GCA_013372085.1	Complete Genome (2)	Illumina, Nanopore (80x)	14	2	210.2	
	GCA_004797155.2	Complete Genome (2)	PacBio (247.19x)	28	62	162.1	
	GCA_001077675.1	Complete Genome (1)	Illumina, PacBio (153x)	15	6	135.9	
	GCA_011067065.1	Complete Genome (2)	PacBio (231.08x)	60227	2486	165.6	
<i>Candida albicans</i> (ATCC® 10231™)	GCA_015227795.1	3, 081 Contigs	NovaSeq (16x)	10174	1573	265.6	
	GCA_002276455.1	2,219 Scaffolds	HiSeq (95x)	13408	2390	274.6	
<i>Meyerozyma guilliermondii</i> (ATCC® 6260™)	GCF_000149425.1	9 RefSeq Scaffolds	Not available	505	1973	278.2	
	GCA_006942155.1	9 Contigs	ONT+MiSeq (240x)	74	386	223.3	
<i>Clavispora lusitaniae</i> (ATCC® 42720™)	GCF_000003835.1	9 RefSeq Scaffolds	Not available	587	2336	265.6	
	GCA_003675505.1	109 Scaffolds	NextSeq (182x)	102	5142	236.9	

# Evaluation of public sequences for ATCC 17978

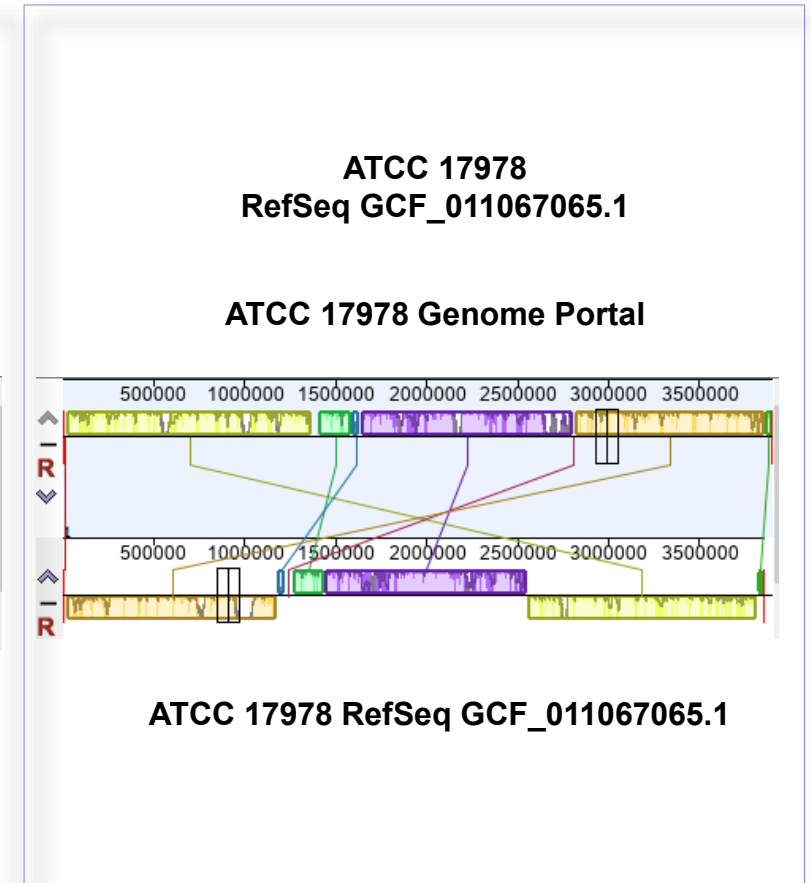
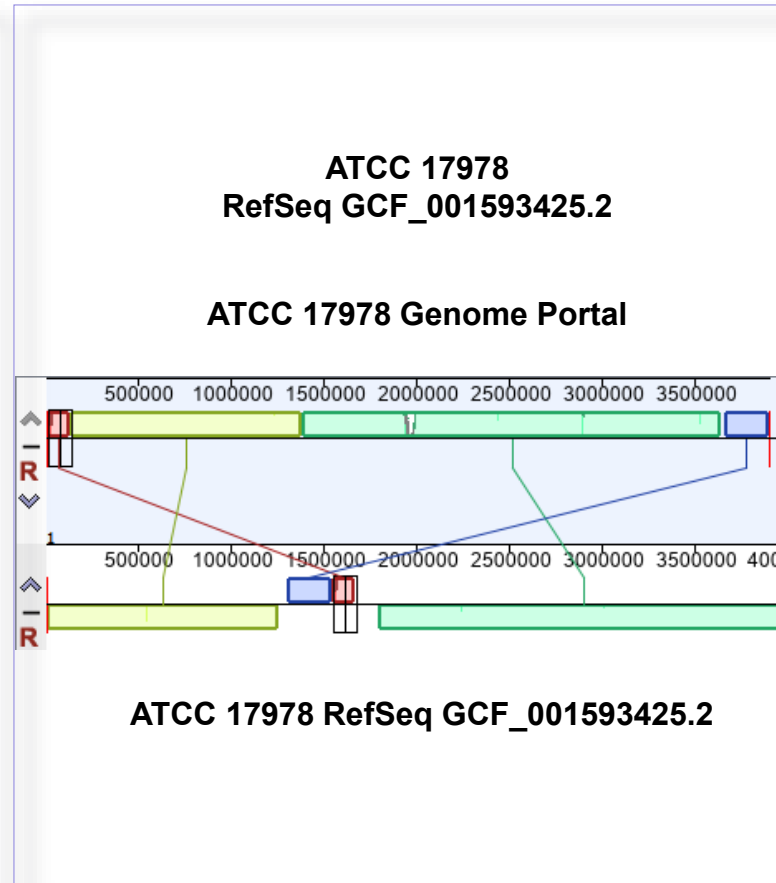
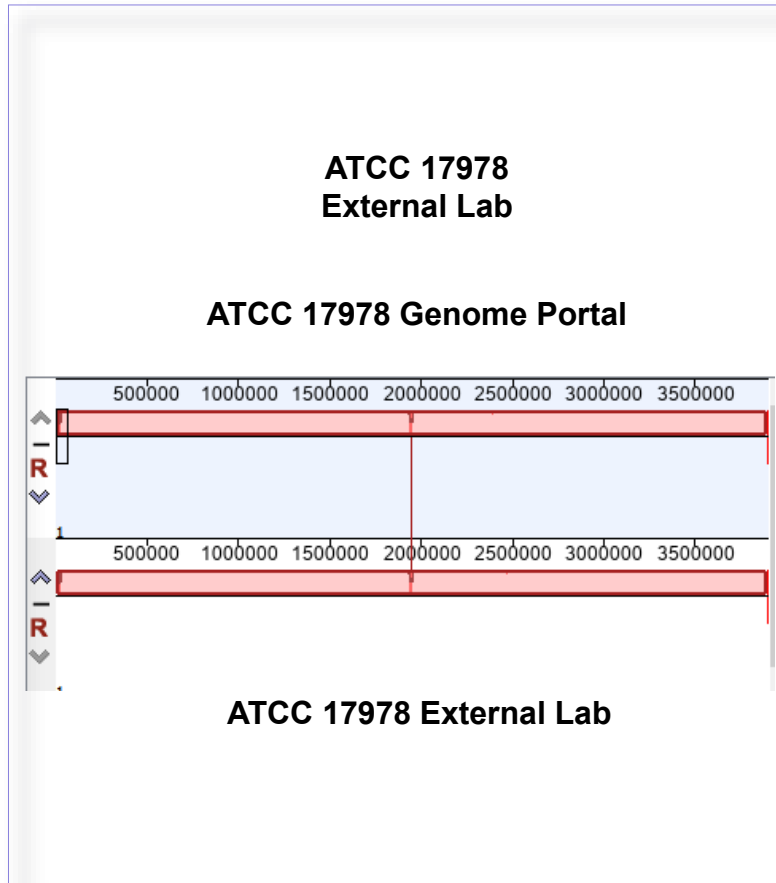
MUMmer alignment with the de novo ATCC 17978 versus GenBank RefSeq genome assemblies GCF\_001593425.2 and GCF\_011067065.1

*Acinetobacter baumannii* strain 5377 (ATCC 17978)



# Evaluation of public sequences for ATCC 17978

*Acinetobacter baumannii* strain 5377 (ATCC 17978)

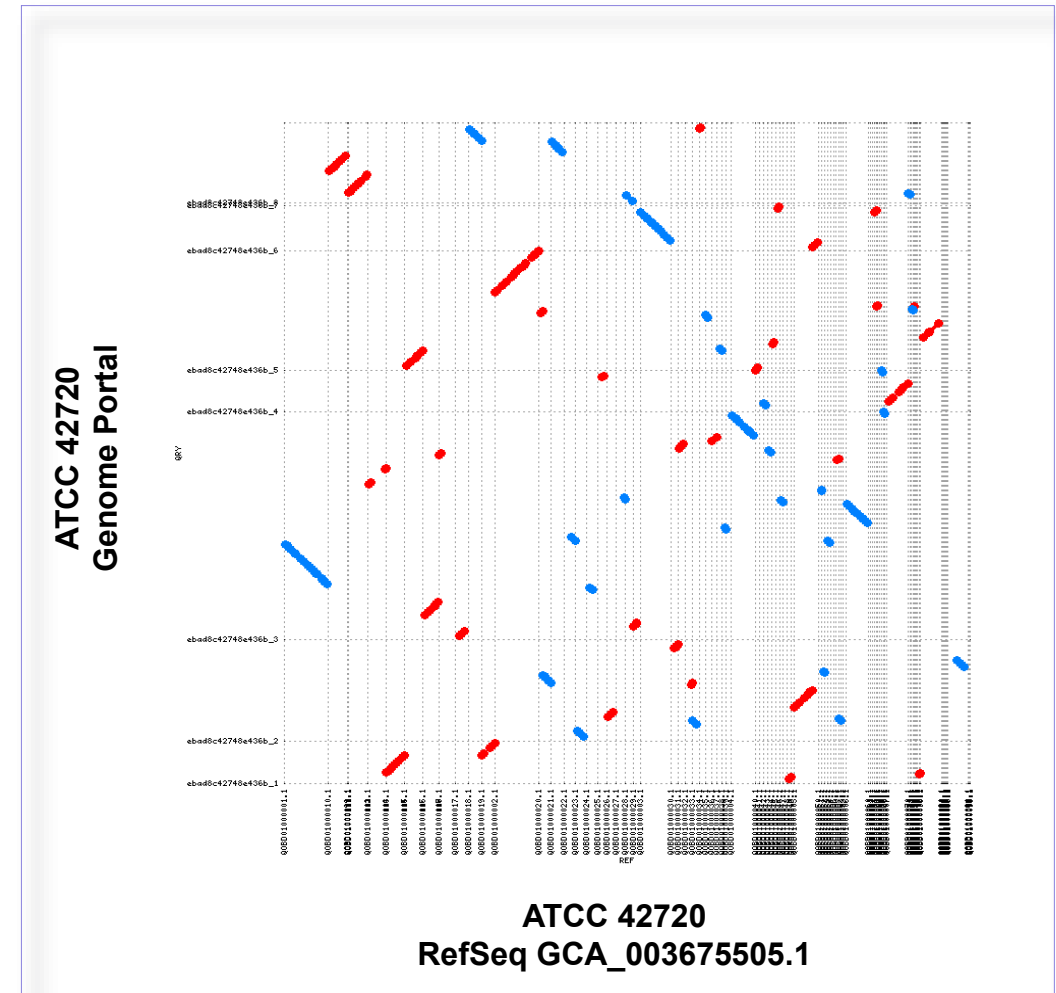
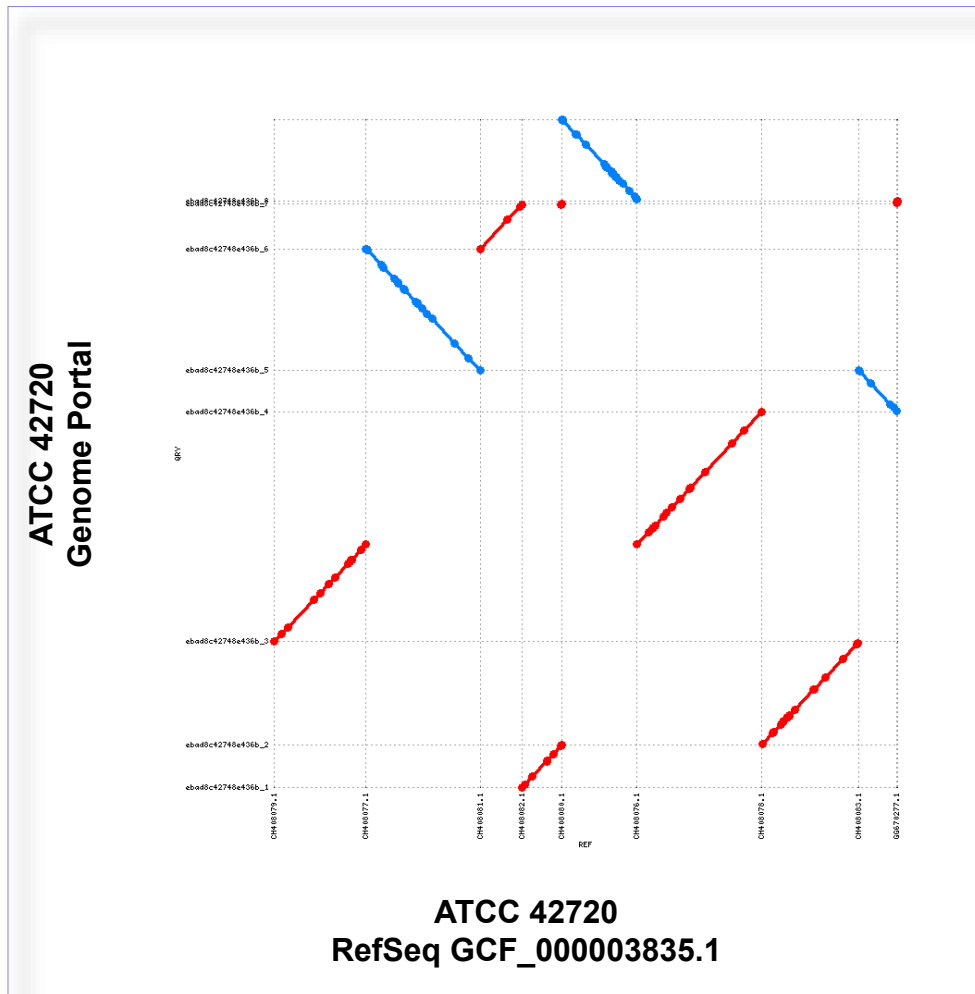


# Evaluation of genome sequences from public databases

Product	NCBI existing reference genomes	NCBI assembly level (plasmids)	Sequencing technology and coverage	# of SNPs	# of indels	Average coverage (variants)
<i>Acinetobacter baumannii</i> (ATCC 17978)	GCA_001593425.2	Complete Genome	Illumina (300.0x)	14	5	210.1
	GCA_000015425.1*	Complete Genome (2)	Not available	118	656	152.7
	GCA_014672775.1	Complete Genome (1)	PacBio (399.24x)	15	87	170.4
	GCA_013372085.1	Complete Genome (2)	Illumina, Nanopore (80x)	14	2	210.2
	GCA_004797155.2	Complete Genome (2)	PacBio (247.19x)	28	62	162.1
	GCA_001077675.1	Complete Genome (1)	Illumina, PacBio (153x)	15	6	135.9
	GCA_011067065.1	Complete Genome (2)	PacBio (231.08x)	60227	2486	165.6
<i>Candida albicans</i> (ATCC 10231)	GCA_015227795.1	3,081 Contigs	NovaSeq (16x)	10174	1573	265.6
	GCA_002276455.1	2,219 Scaffolds	HiSeq (95x)	13408	2390	274.6
<i>Meyerozyma guilliermondii</i> (ATCC 6260)	GCF_000149425.1	9 RefSeq Scaffolds	Not available	505	1973	278.2
	GCA_000942155.1	9 Contigs	ONT+MiSeq (240x)	74	386	223.3
<i>Clavispora lusitaniae</i> (ATCC 42720)	GCF_000003835.1	9 RefSeq Scaffolds	Not available	587	2336	265.6
	GCA_003675505.1	109 Scaffolds	NextSeq (182x)	102	5142	236.9

# Evaluation of public sequences for ATCC 42720

*MUMmer* whole genome alignments of ATCC de-novo genome assembly of ATCC 42720 versus GenBank RefSeq genome assemblies GCF\_000003835.1 and GCA\_003675505.1




# Overview

- The ATCC Genome Portal
- Traceability and authentication of reference genomes
- **Standards for authenticated reference genomes**
- Development roadmap preview



# Selected timeline for (microbial) genomics standards



1970s	1980s	1990s	2000s	2010s	2020s
	<b>1982</b> – GenBank and ENA created		<b>2005</b> – Genomic Standards Consortium established  <b>2008</b> – Minimal Information on Genome Sequence (MIGS) specification  <b>2009</b> – Genome Project Standards published by GSC	<b>2012</b> – CDC NGS Standards for Clinical Testing (Nex-StoCT)  <b>2014</b> – Viral Genome Reference Standards  <b>2016</b> – FDA Draft Guidance on NGS for Pathogen Identification	<b>2020</b> – ATCC Enhanced Authentication Initiative  <b>2020</b> – ATCC Genome Portal Launch
<b>1974</b> – Complete RNA genome bacteriophage MS2	<b>1984</b> – Complete Epstein Barr virus genome	<b>1995</b> – Complete genome of <i>Haemophilus influenzae</i>	<b>2001</b> – Draft Human Genome <b>2007</b> - Genomic Encyclopedia of Bacteria and Archaea (GEBA) and Human Microbiome Project (HMP) launch.	<b>2011</b> – GEBA II Launched	<b>2020</b> – First “end to end” gapless genome for Human Chr. X

# Recognition of the importance of traceability to biomaterials

nature  
biotechnology

## PERSPECTIVE

### The minimum information about a genome sequence (MIGS) specification

Dawn Field\*<sup>1</sup>, George Garrity<sup>2</sup>, Tanya Gray<sup>1</sup>, Norman Morrison<sup>3,4</sup>, Jeremy Selengut<sup>5</sup>, Peter Sterk<sup>6</sup>, Tatiana Tatusova<sup>7</sup>, Nicholas Thomson<sup>8</sup>, Michael J Allen<sup>9</sup>, Samuel V Angiuoli<sup>5,10</sup>, Michael Ashburner<sup>11</sup>, Nelson Axelrod<sup>5</sup>, Sandra Baldauf<sup>12</sup>, Stuart Ballard<sup>13</sup>, Jeffrey Boore<sup>14</sup>, Guy Cochrane<sup>6</sup>, James Cole<sup>2</sup>, Peter Dawyndt<sup>15</sup>, Paul De Vos<sup>16,17</sup>, Claude dePamphilis<sup>18</sup>, Robert Edwards<sup>19,20</sup>, Nadeem Faruque<sup>6</sup>, Robert Feldman<sup>21</sup>, Jack Gilbert<sup>9</sup>, Paul Gilna<sup>22</sup>, Frank Oliver Glöckner<sup>23</sup>, Philip Goldstein<sup>24</sup>, Robert Guralnick<sup>24</sup>, Dan Haft<sup>5</sup>, David Hancock<sup>3,4</sup>, Henning Hermjakob<sup>6</sup>, Christiane Hertz-Fowler<sup>8</sup>, Phil Hugenholtz<sup>25</sup>, Ian Joint<sup>9</sup>, Leonid Kagan<sup>5</sup>, Matthew Kane<sup>26</sup>, Jessie Kennedy<sup>27</sup>, George Kowalchuk<sup>28</sup>, Renzo Kottmann<sup>23</sup>, Eugene Kolker<sup>29–31</sup>, Saul Kravitz<sup>5</sup>, Nikos Kyrpides<sup>32</sup>, Jim Leebens-Mack<sup>33</sup>, Suzanna E Lewis<sup>34</sup>, Kelvin Li<sup>5</sup>, Allyson L Lister<sup>35,36</sup>, Phillip Lord<sup>35</sup>, Natalia Maltsev<sup>20</sup>, Victor Markowitz<sup>37</sup>, Jennifer Martiny<sup>38</sup>, Barbara Methe<sup>5</sup>, Ilene Mizrahi<sup>7</sup>, Richard Moxon<sup>39</sup>, Karen Nelson<sup>5,40</sup>, Julian Parkhill<sup>8</sup>, Lita Proctor<sup>26</sup>, Owen White<sup>10</sup>, Susanna-Assunta Sansone<sup>6</sup>, Andrew Spiers<sup>42</sup>, Robert Stevens<sup>3</sup>, Paul Swift<sup>1</sup>, Chris Taylor<sup>6</sup>, Yoshio Tateno<sup>43</sup>, Adrian Tett<sup>1</sup>, Sarah Turner<sup>1</sup>, David Ussery<sup>44</sup>, Bob Vaughan<sup>6</sup>, Naomi Ward<sup>45</sup>, Trish Whetzel<sup>46</sup>, Ingio San Gil<sup>41</sup>, Gareth Wilson<sup>1</sup> & Anil Wipat<sup>35,36</sup>

With the quantity of genomic data increasing at an exponential rate, it is imperative that these data be captured electronically, in a standard format. Standardization activities must proceed within the auspices of open-access and international working bodies. To tackle the issues surrounding the development of better descriptions of genomic investigations, we have formed the Genomic Standards Consortium (GSC). Here, we introduce the minimum information about a genome sequence (MIGS) specification with the intent of promoting participation in its development and discussing the resources that will be required to develop improved mechanisms of metadata capture and exchange. As part of its wider goals, the GSC also supports

can manipulate it to provide new solutions to critical problems. Such solutions include therapies and cures for disease, industrial products, approaches for biodegradation of xenobiotic compounds and renewable energy sources. With improvements in sequencing technologies, the growing interest in metagenomic approaches and the proven power of comparative analysis of groups of related genomes, we can envision the day when it will be commonplace to sequence tens to hundreds of genomes or more as part of a single study. At current rates of genome sequencing, it has been estimated that >4,000 bacterial genomes will be available soon after 2010 (ref. 1).

Given the importance of the growing genome collection, the capital investment in its creation and the benefits of leveraging its value

**“Source material identifier” is an exception; the GSC recommends this be a core descriptor, but as of yet, physical archives are not yet routinely created for all cases or types of biological material subjected to genome sequencing ...**

**This was in 2008.**

We agree.

But, 12 years later “physical archives are [still] not yet routinely created” by groups doing whole genome sequencing.

Chain of custody of biomaterials is rarely or poorly documented.

Field, D. *et al.* (2008) ‘The minimum information about a genome sequence (MIGS) specification’, *Nature Biotechnology*, 26(5), pp. 541–547. doi: [10.1038/nbt1360](https://doi.org/10.1038/nbt1360).

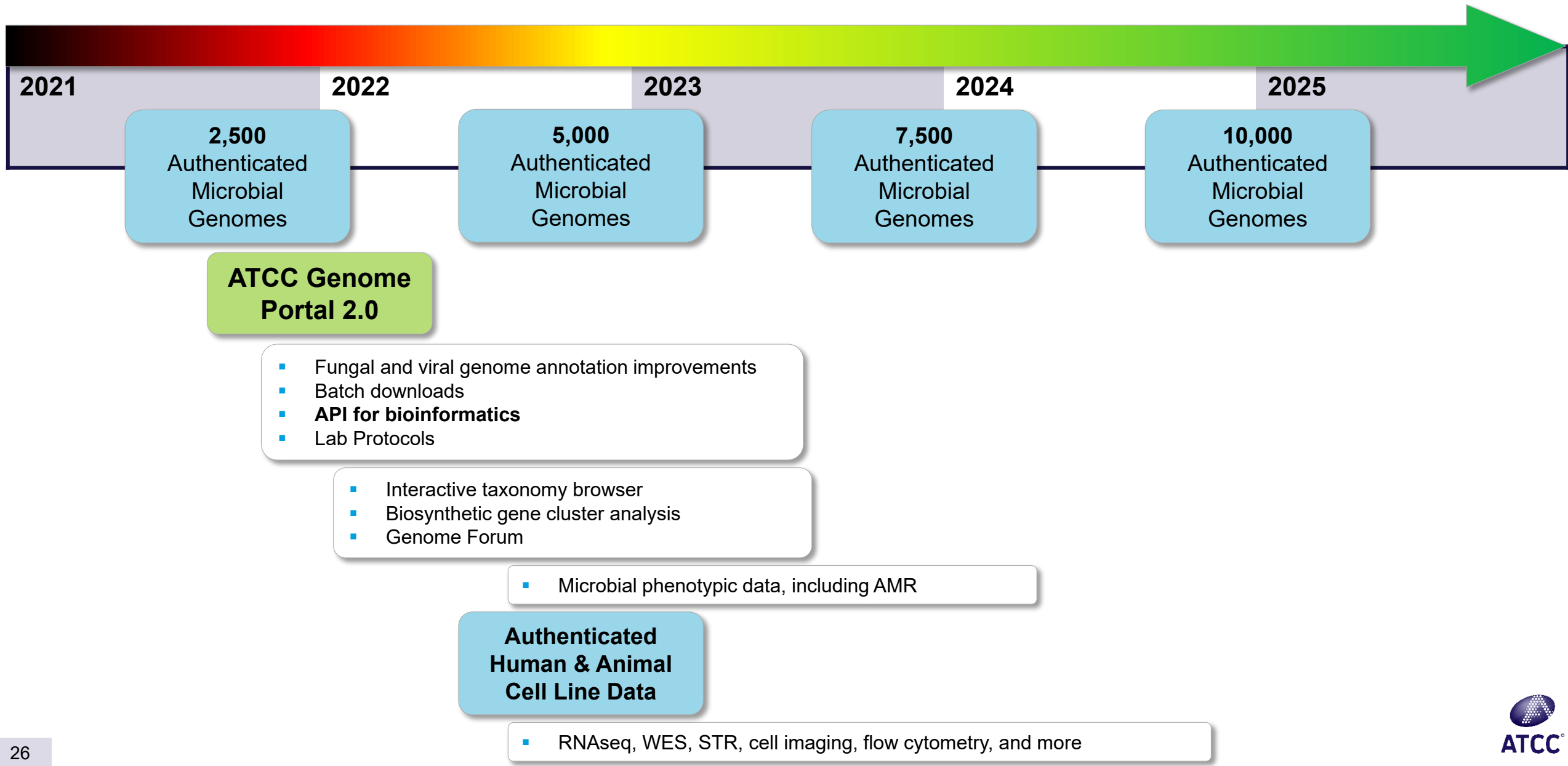


# Overview

- The ATCC Genome Portal
- Traceability and authentication of reference genomes
- Importance of authenticated reference genomes
- **Development roadmap preview**



# ATCC Genome Portal development goals



# The ATCC Genome Portal Team

**Jonathan Jacobs, PhD**

jjacobs@atcc.org

## Genomics Lab

**Briana Benton**

Stephen King, MSc

James Duncan, MSc

Robert Marlow

Samuel Greenfield

Corina Tabron

Fabio Martinez

Amanda Pierola

## Bioinformatics Lab

**John Bagnoli**

David Yarmosh, MSc

Nikhita Puthaveetil,  
MSc

P. Ford Combs, MSc

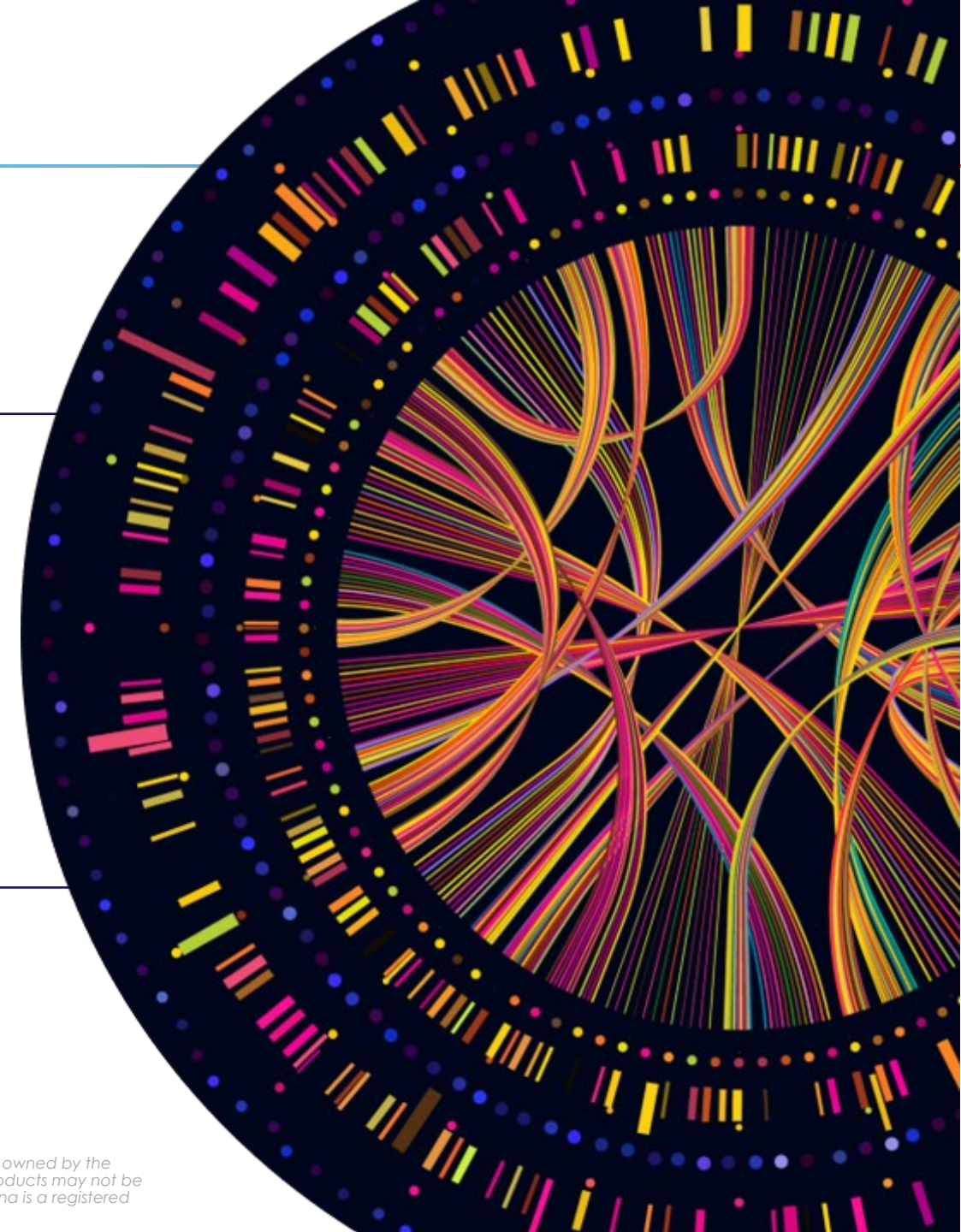
## Partners

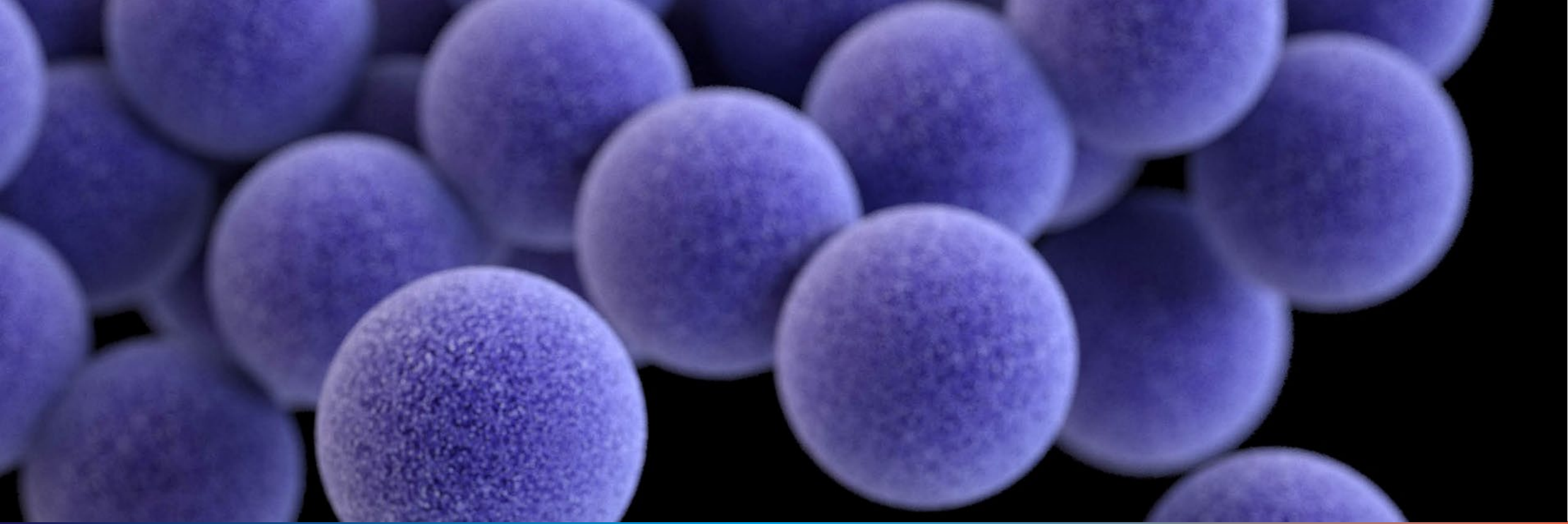
Juan Lopera, PhD

Marco Riojas, PhD

*... and One Codex!*

**JOIN OUR TEAM! We're hiring!**





Thank you