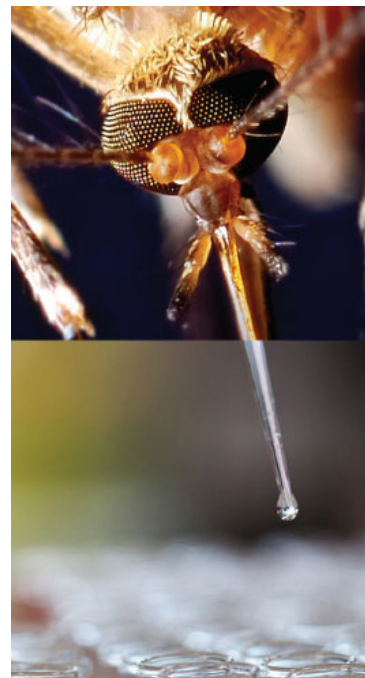
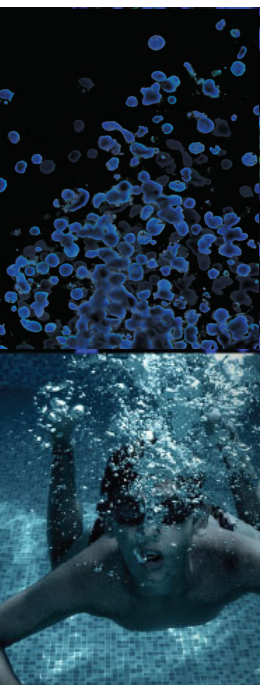




Addressing the Need for Accuracy and Traceability in Microbial Genomic Data: The ATCC Genome Portal

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

Credible Leads to Incredible™

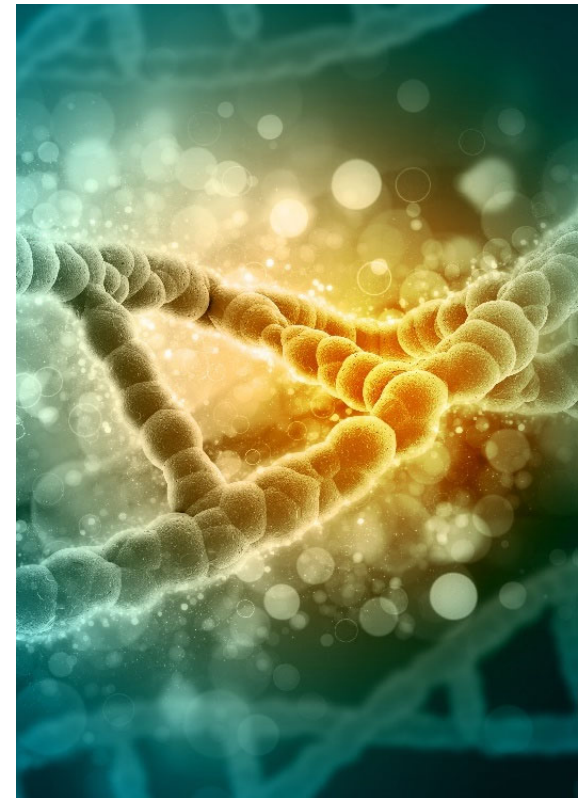


About ATCC

- Founded in 1925, ATCC is a not-for-profit organization with HQ in Manassas, VA, and an R&D and Services center in Gaithersburg, MD
- World's largest, most diverse biological materials and information resource for microbes – the “*gold standard*”
- Innovative R&D company featuring gene editing, microbiome, NGS, and advanced models
- cGMP biorepository
- Partner with government, industry, and academia
- Leading global supplier of authenticated cell lines, microorganisms, and molecular standards
- Sales and distribution in 150 countries, 19 international distributors
- Talented team of 450+ employees, over one-third with advanced degrees

Overview

- The ATCC Genome Portal
- Traceability and authentication of microbial genomes
- Importance of 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

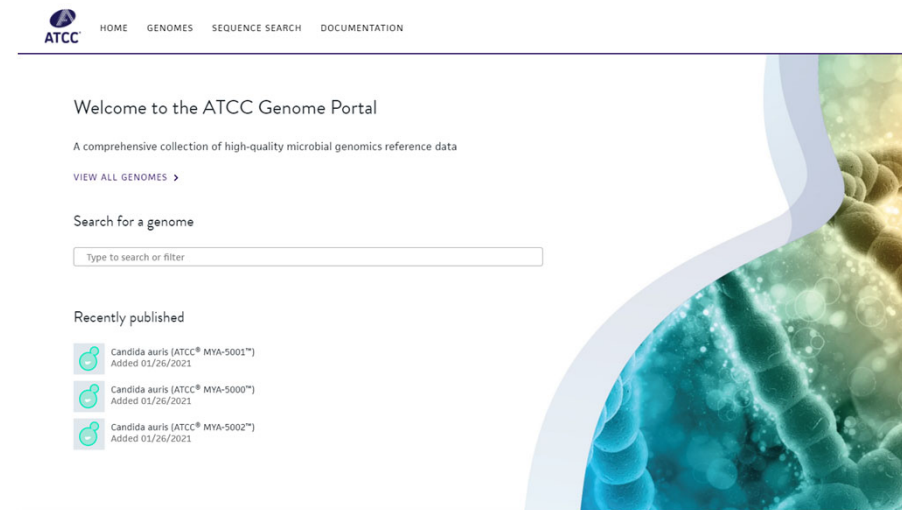
genomes.atcc.org



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
- 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: 12290-5™
Lot Number: 7000050

Designation: *Staphylococcus epidermidis* genomic DNA
FBI 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 ₆₀₀ /OD ₇₅₀ ratio (Spectrophotometer method)	1.8 to 2.1	1.8
Total amount of DNA (PicoGreen® measurement)	≥ 5 µg per vial	7 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: 12290-5™

Quality Assurance: Specificity, 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
10801 University Boulevard
Manassas, VA 20108-2209 USA
www.atcc.org

800-635-0047 or 703-365-2700
Fax: 703-365-2700
E-mail: help@atcc.org
or contact your local distributor

Page 1 of 2
Template Number: 1
Template (Rev. 04-2016)

*not an inclusive list

Providing reference-quality genomes

Why - Challenge # 3

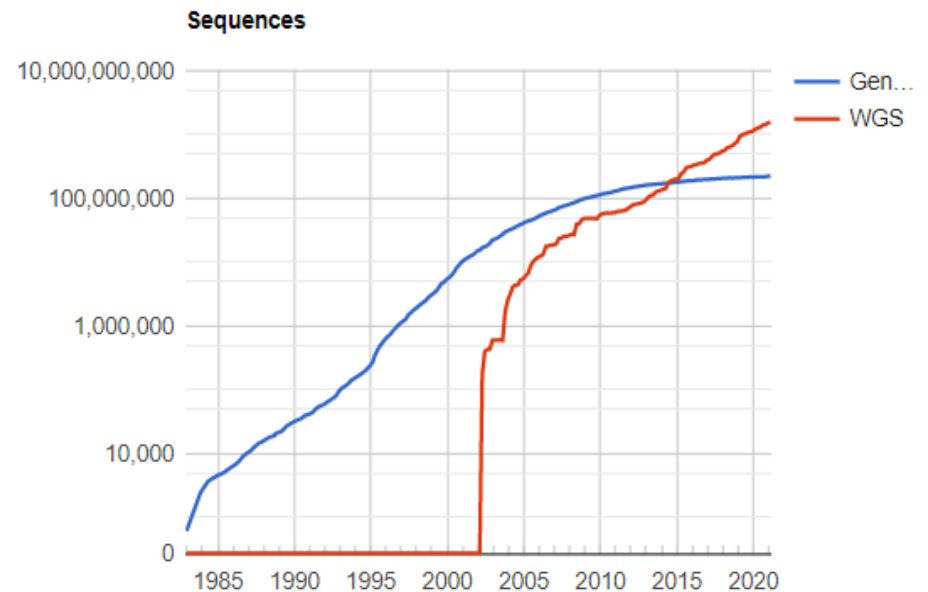


- Acknowledge there is a problem with reference genomes ✓
- Work through a plan to address the problem ✓
- **How do we effectively and easily provide customers with genomic data while not diluting it or burying it in a public database?**

Reference genomes

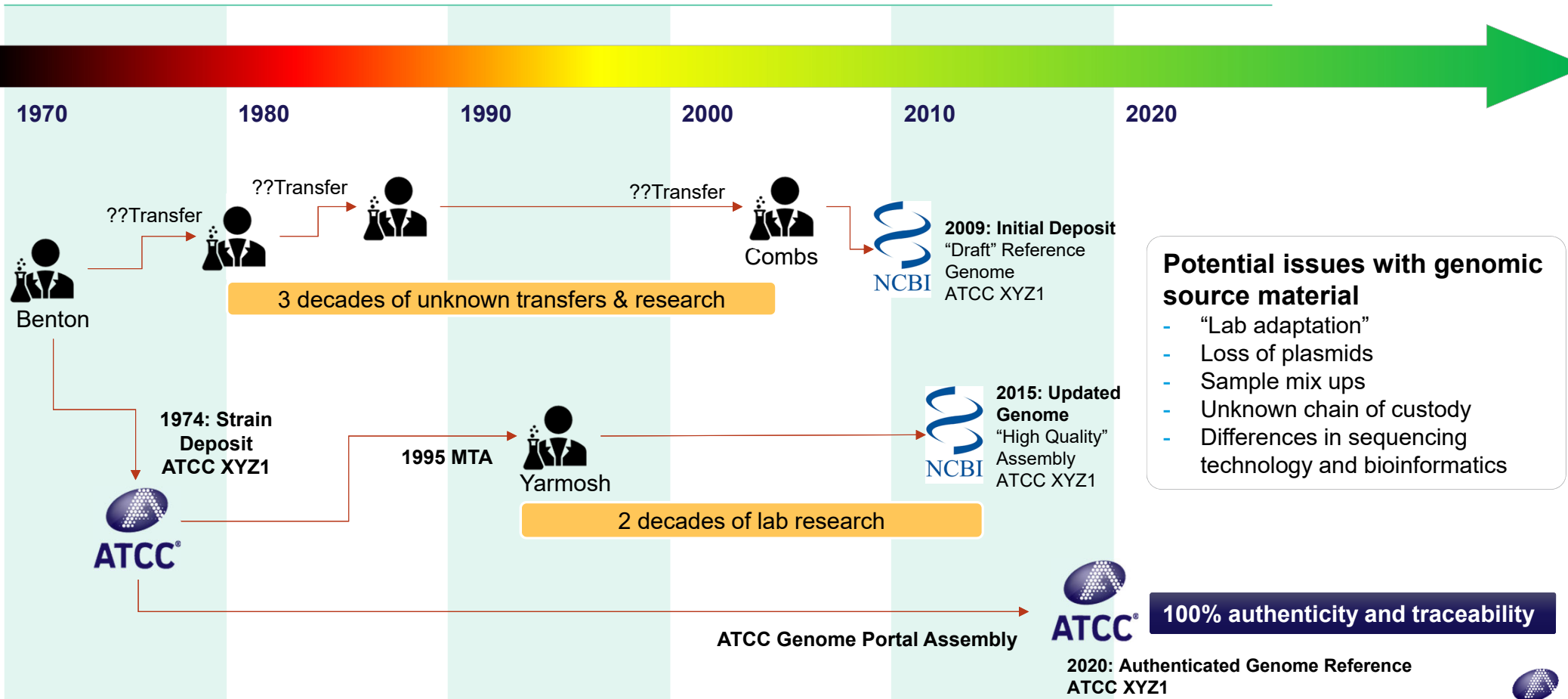
Where can researchers turn to for “reference” genomes?

- 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

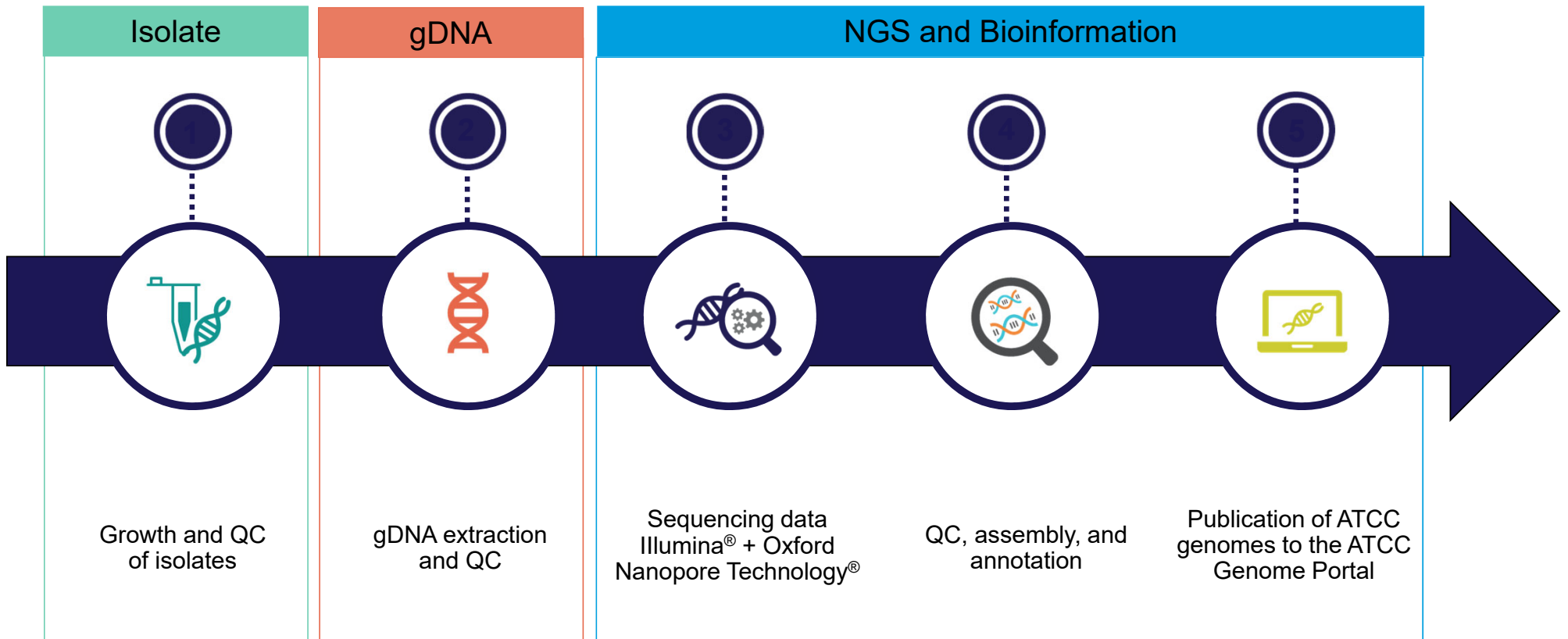


<https://www.ncbi.nlm.nih.gov/genbank/statistics/>

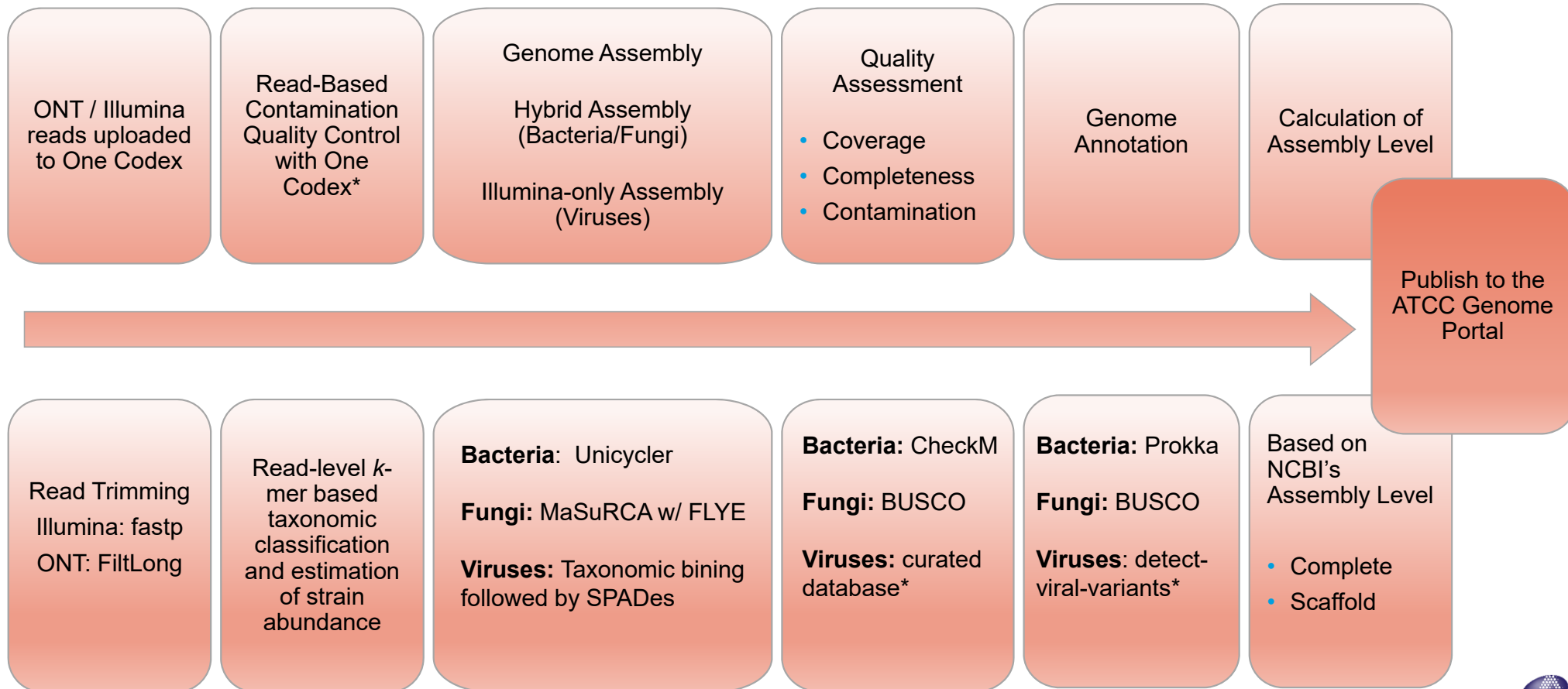
The elephant in the room: Authenticated reference genomes



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



ATCC genome assembly process

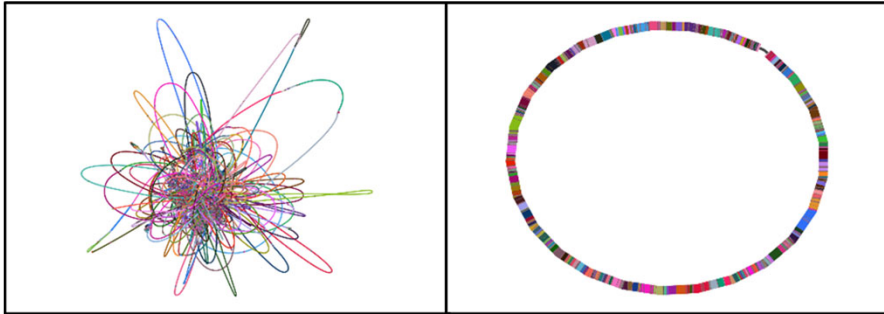


* One Codex proprietary software

Advantage of hybrid assemblies

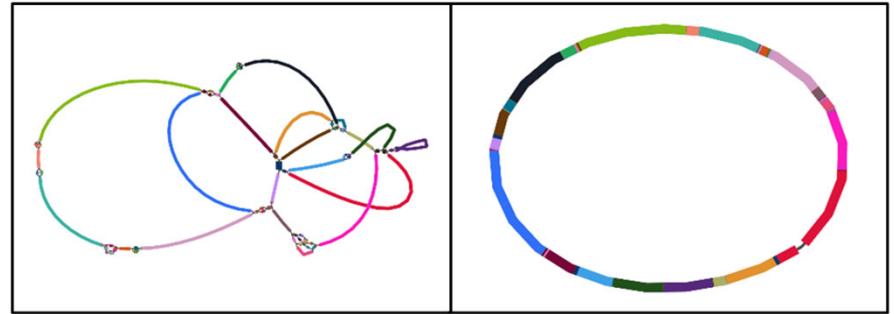
Illumina-only assembly Hybrid assembly

Neisseria meningitidis (ATCC® 53417™)

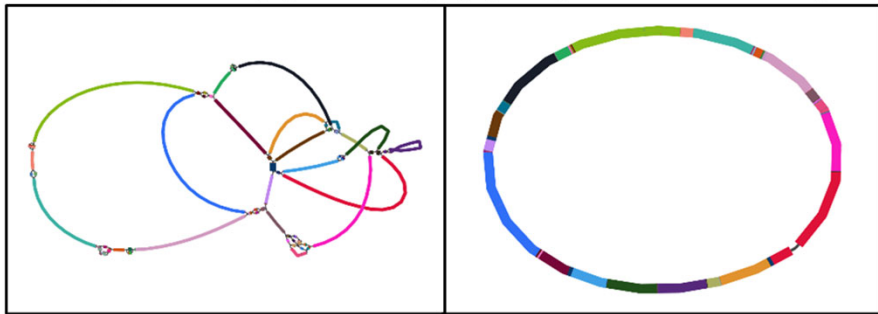


Illumina-only assembly Hybrid assembly

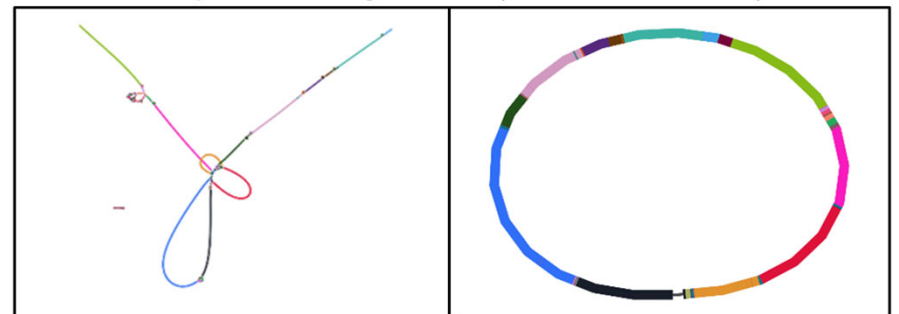
Campylobacter jejuni subsp. *jejuni* (ATCC® 43446™)



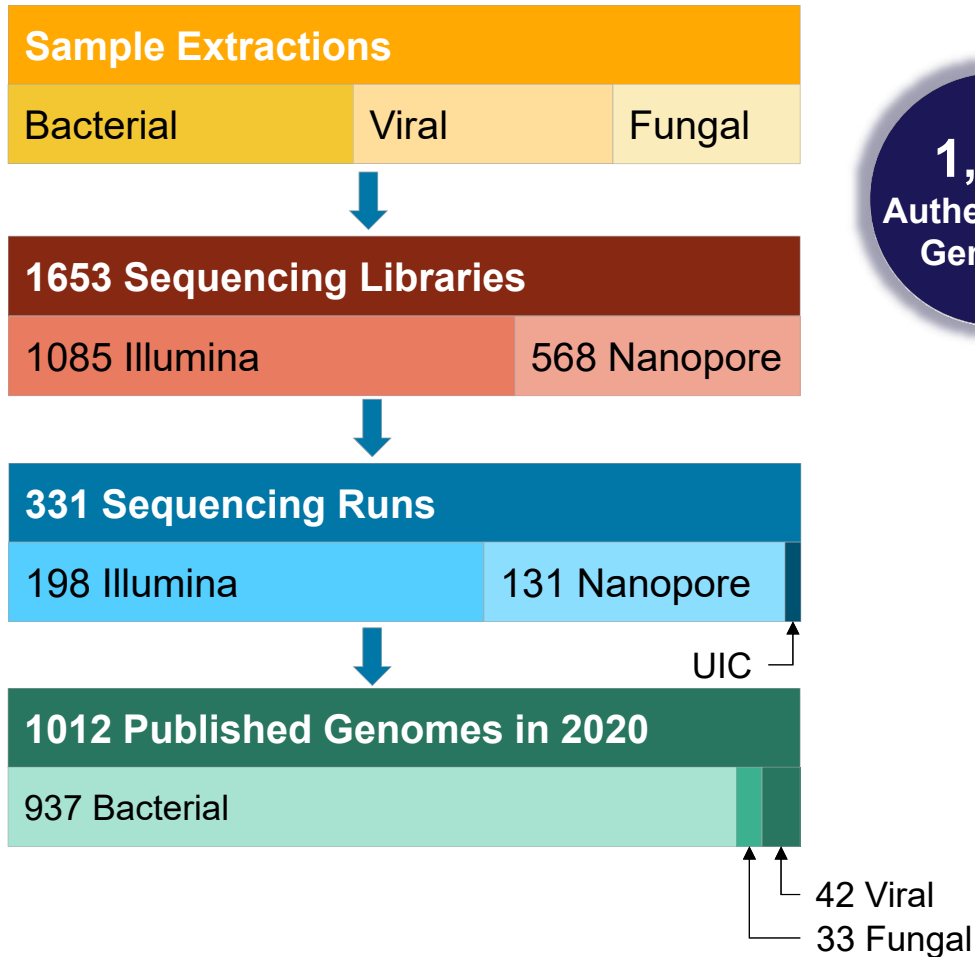
Campylobacter jejuni subsp. *jejuni* (ATCC® 43446™)



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



Reference genomes

208,295 genomes in NCBI
(RefSeq prokaryotes)

1,957 identified as
"ATCC"

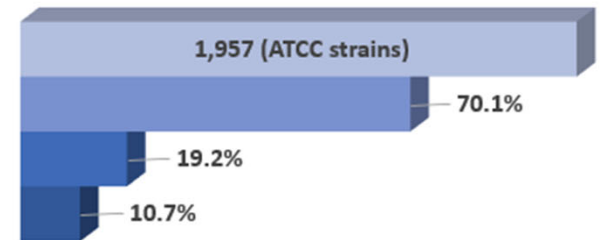
585
complete

■ ATCC prokaryote genomes in NCBI-NIH databases

■ % genome contigs or scaffolds

■ % complete genome or chromosome

■ % complete genome or chromosome and plasmids

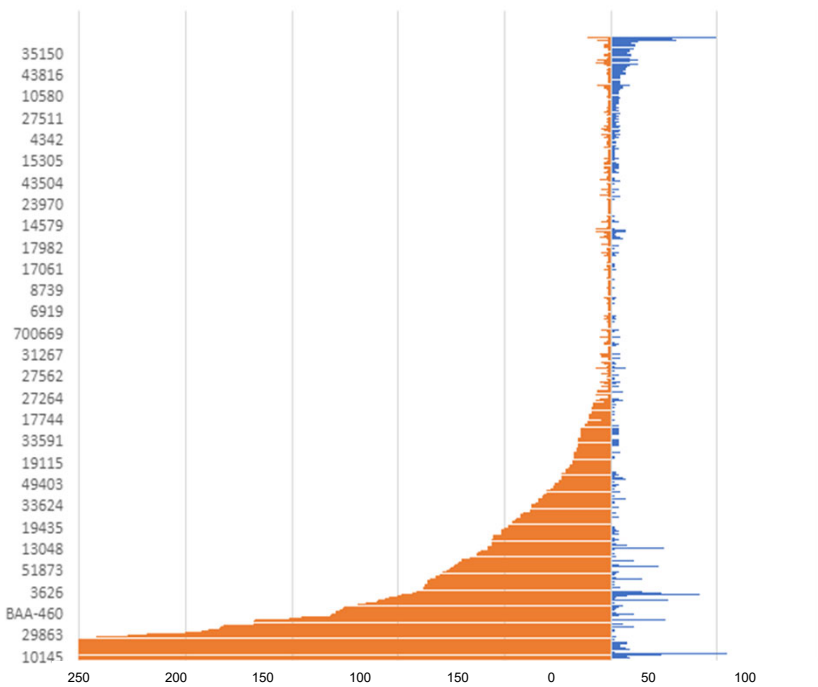


Are these 585 RefSeq genomes traceable back to authenticated ATCC cultures with well-documented growth and storage conditions?

Genome assembly quality

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

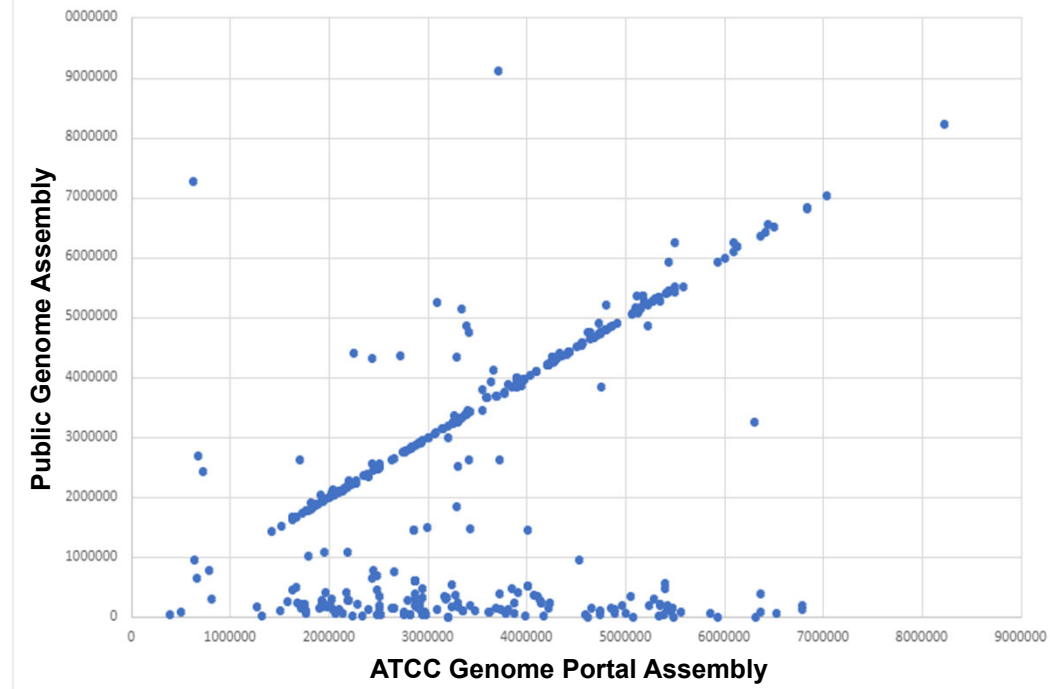
Bacteriology Products Contig Counts



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



ATCC

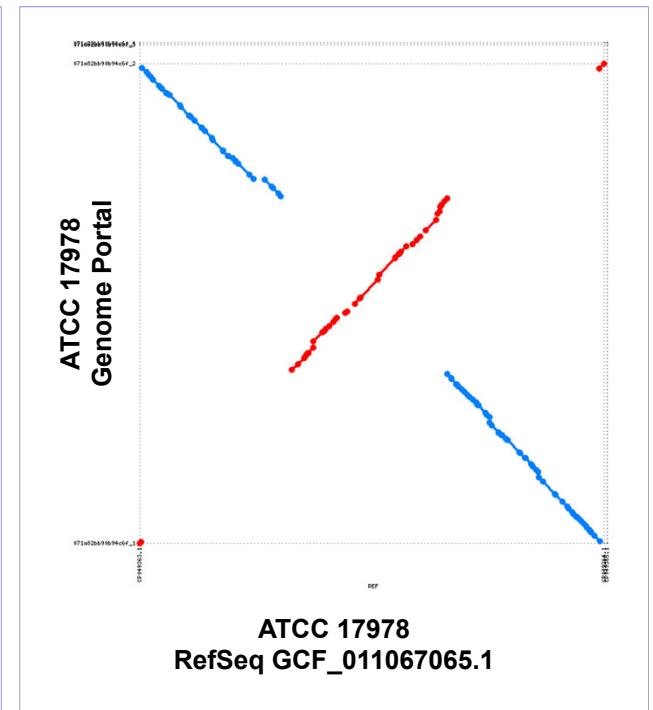
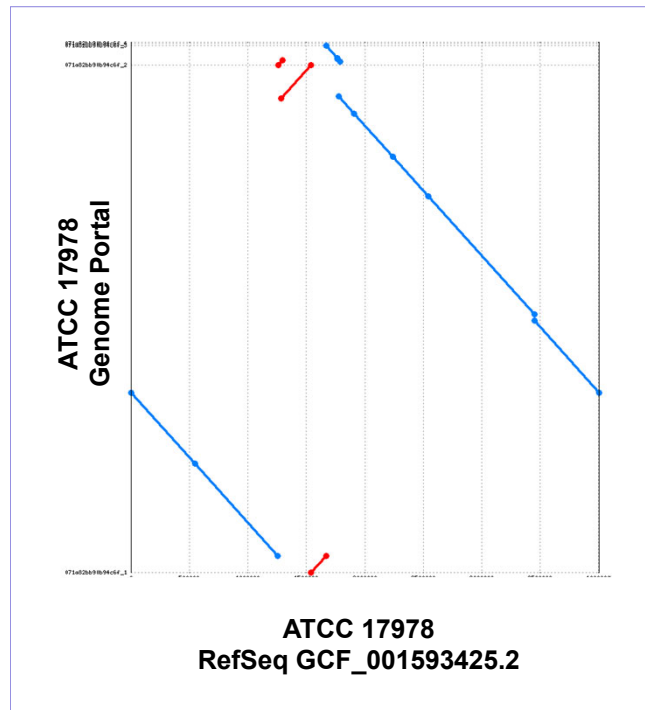
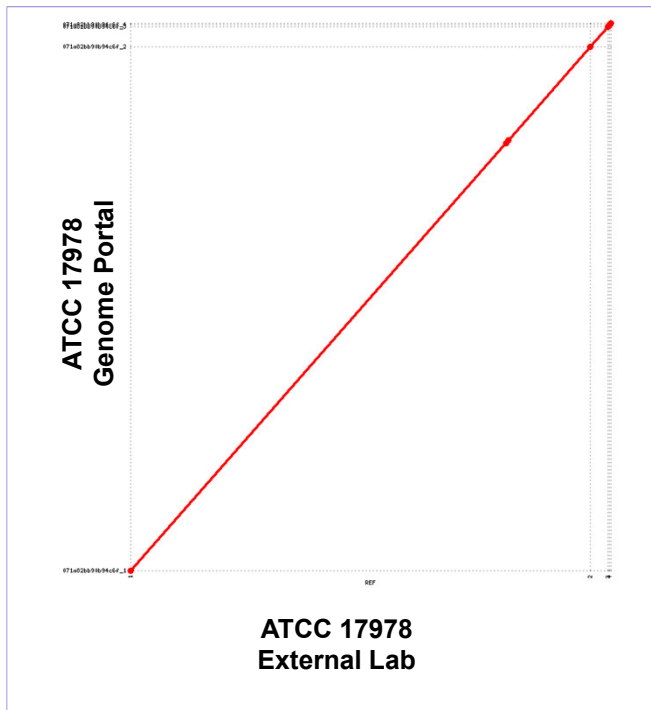
Evaluation of genome sequences from public databases

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 7 assemblies Unknown origin of materials
	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_006942155.1	9 Contigs	ONT+MiSeq (240x)	74	386	223.3	
<i>Clavispora lusitanae</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

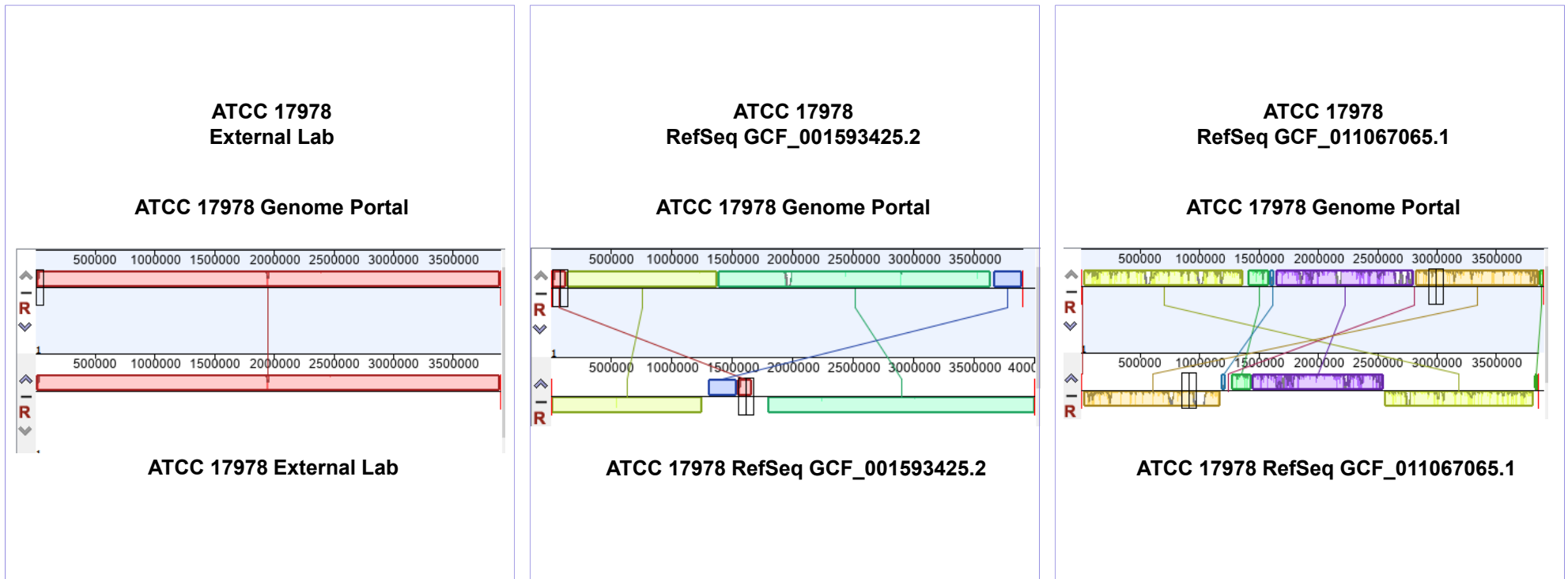
MUMmer alignment with the de novo ATCC 17978 versus GenBank RefSeq genome assemblies GCF_001593425.2 and GCF_011067065.1

Acinetobacter baumannii strain 5377 (ATCC 17978)



Evaluation of public sequences for ATCC 17978

Acinetobacter baumannii strain 5377 (ATCC 17978)

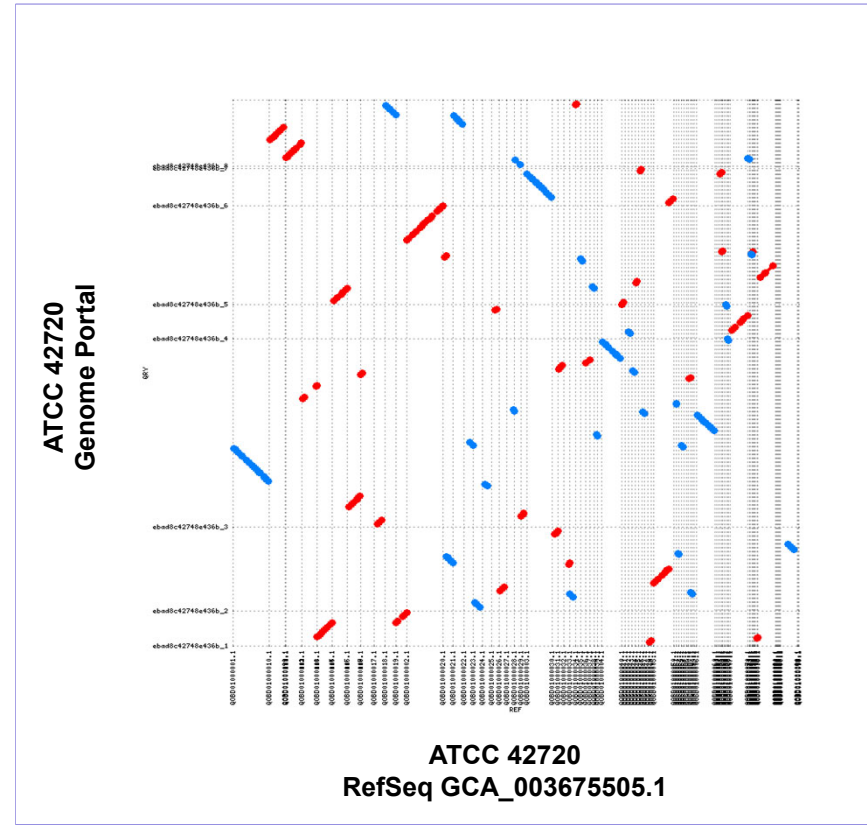
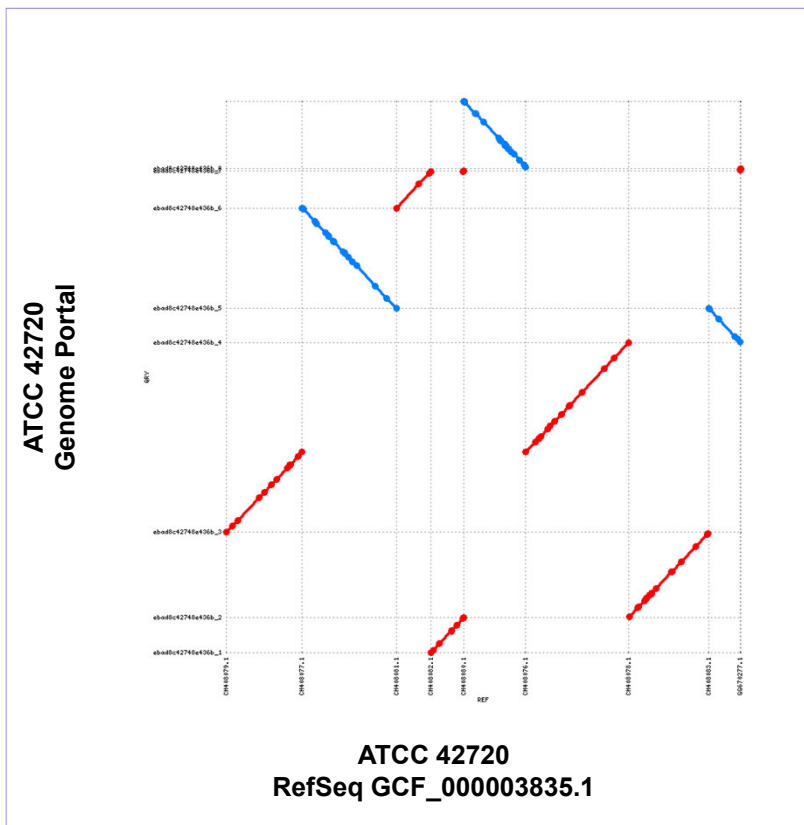


Evaluation of genome sequences from public databases

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 lusitanae</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

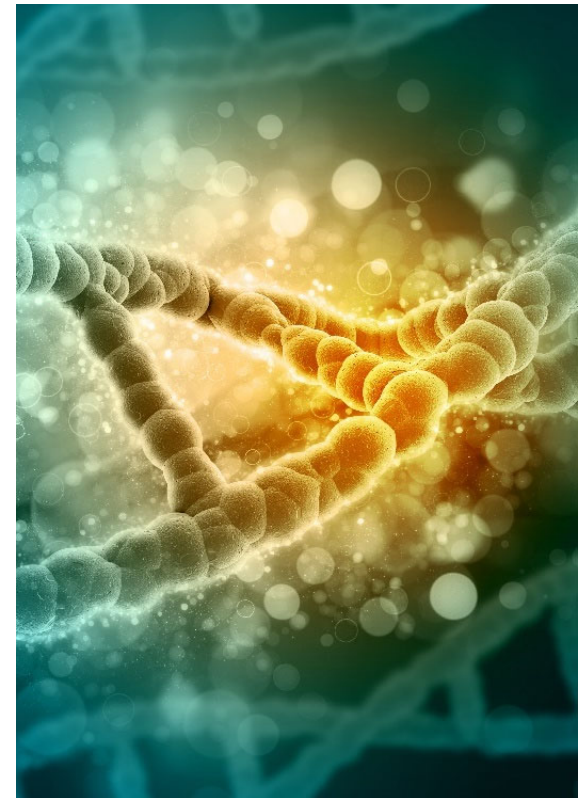
Evaluation of public sequences for ATCC 42720

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
- **Importance of authenticated reference genomes**
- Development roadmap preview



Selected timeline for (microbial) genomics standards



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

Recognition of the importance of traceability to biomaterials

nature
biotechnology

PERSPECTIVE

“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 ...

The minimum information about a genome sequence (MIGS) specification

Dawn Field^{*1}, George Garrity², Tanya Gray¹, Norman Morrison^{3,4}, Jeremy Selengut⁵, Peter Sterk⁶, Tatiana Tatusova⁷, Nicholas Thomson⁸, Michael J Allen⁹, Samuel V Angiuoli^{5,10}, Michael Ashburner¹¹, Nelson Axelrod⁵, Sandra Baldauf¹², Stuart Ballard¹³, Jeffrey Boore¹⁴, Guy Cochrane⁶, James Cole², Peter Dawyndt¹⁵, Paul De Vos^{16,17}, Claude dePamphilis¹⁸, Robert Edwards^{19,20}, Nadeem Faruque⁶, Robert Feldman²¹, Jack Gilbert⁹, Paul Gilna²², Frank Oliver Glöckner²³, Philip Goldstein²⁴, Robert Guralnick²⁴, Dan Haft⁵, David Hancock^{3,4}, Henning Hermjakob⁶, Christiane Hertz-Fowler⁸, Phil Hugenholtz²⁵, Ian Joint⁹, Leonid Kagan⁵, Matthew Kane²⁶, Jessie Kennedy²⁷, George Kowalchuk²⁸, Renzo Kottmann²³, Eugene Kolker^{29–31}, Saul Kravitz⁵, Nikos Kyrpides³², Jim Leebens-Mack³³, Suzanna E Lewis³⁴, Kelvin Li⁵, Allyson L Lister^{35,36}, Phillip Lord³⁵, Natalia Maltsev²⁰, Victor Markowitz³⁷, Jennifer Martiny³⁸, Barbara Methé⁵, Ilene Mizrahi⁷, Richard Moxon³⁹, Karen Nelson^{5,40}, Julian Parkhill⁸, Lita Proctor²⁶, Owen White¹⁰, Susanna-Assunta Sansone⁶, Andrew Spiers⁴², Robert Stevens³, Paul Swift¹, Chris Taylor⁶, Yoshio Tateno⁴³, Adrian Tett¹, Sarah Turner¹, David Ussery⁴⁴, Bob Vaughan⁶, Naomi Ward⁴⁵, Trish Whetzel⁴⁶, Ingio San Gil⁴¹, Gareth Wilson¹ & Anil Wipat^{35,36}

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

This was in 2008.

We agree.

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

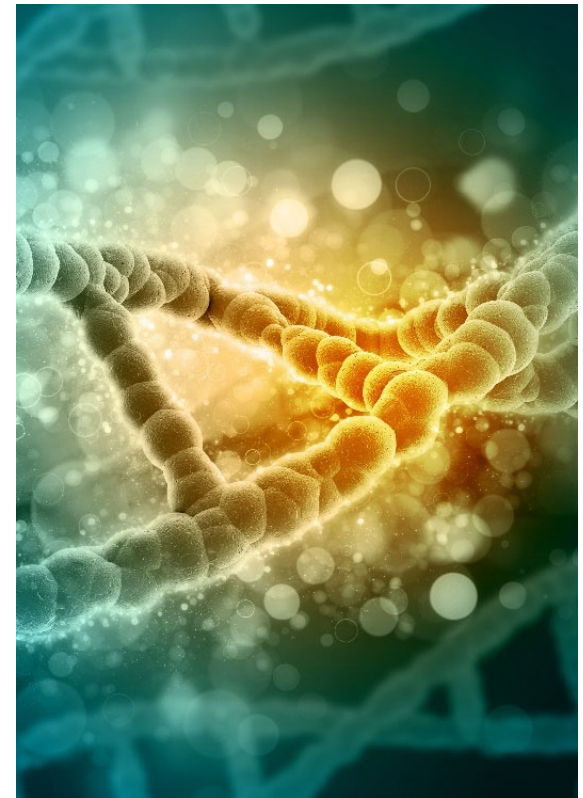
Chain of custody of biomaterials is rarely or poorly documented.

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).

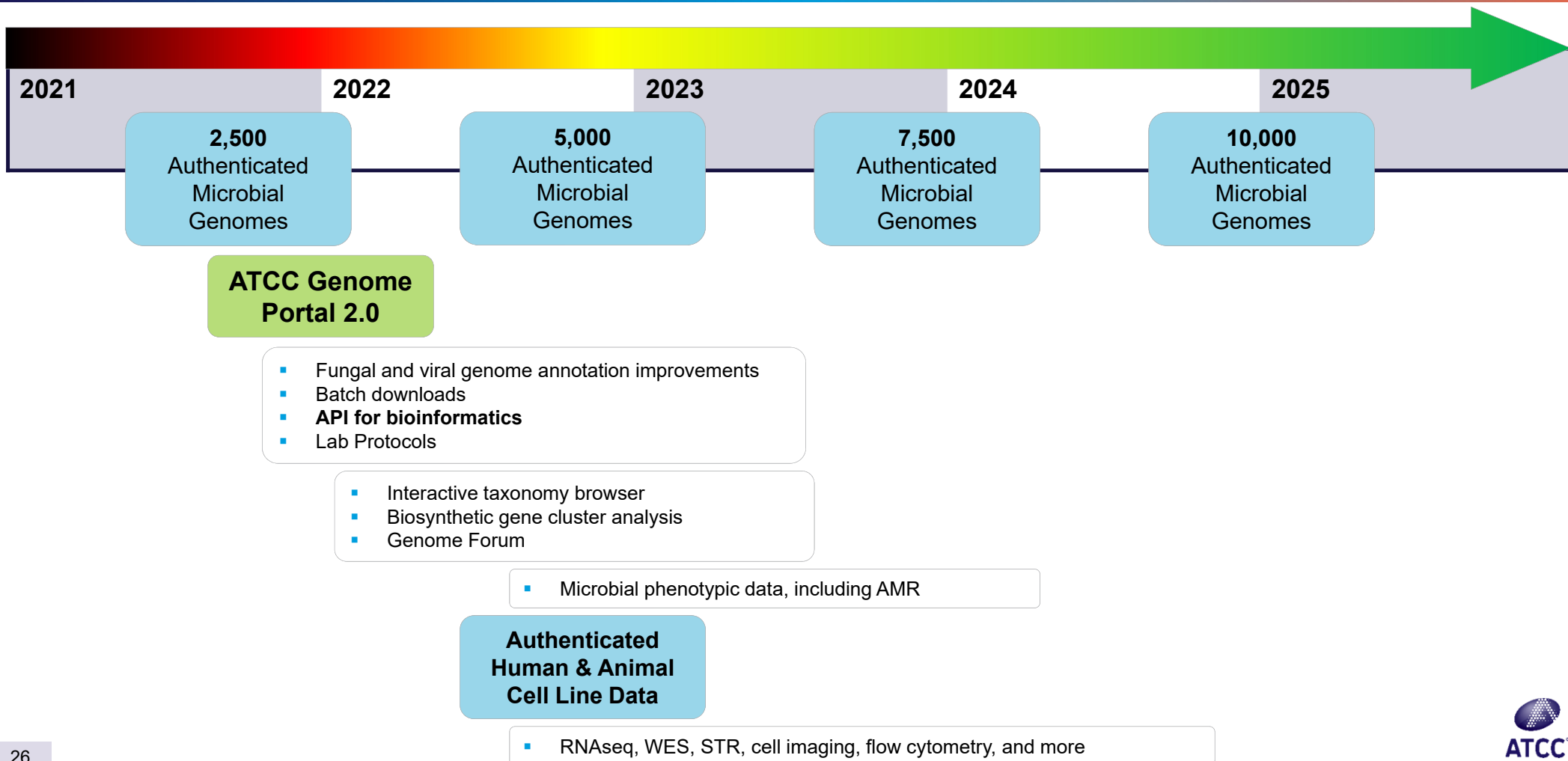


Overview

- 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

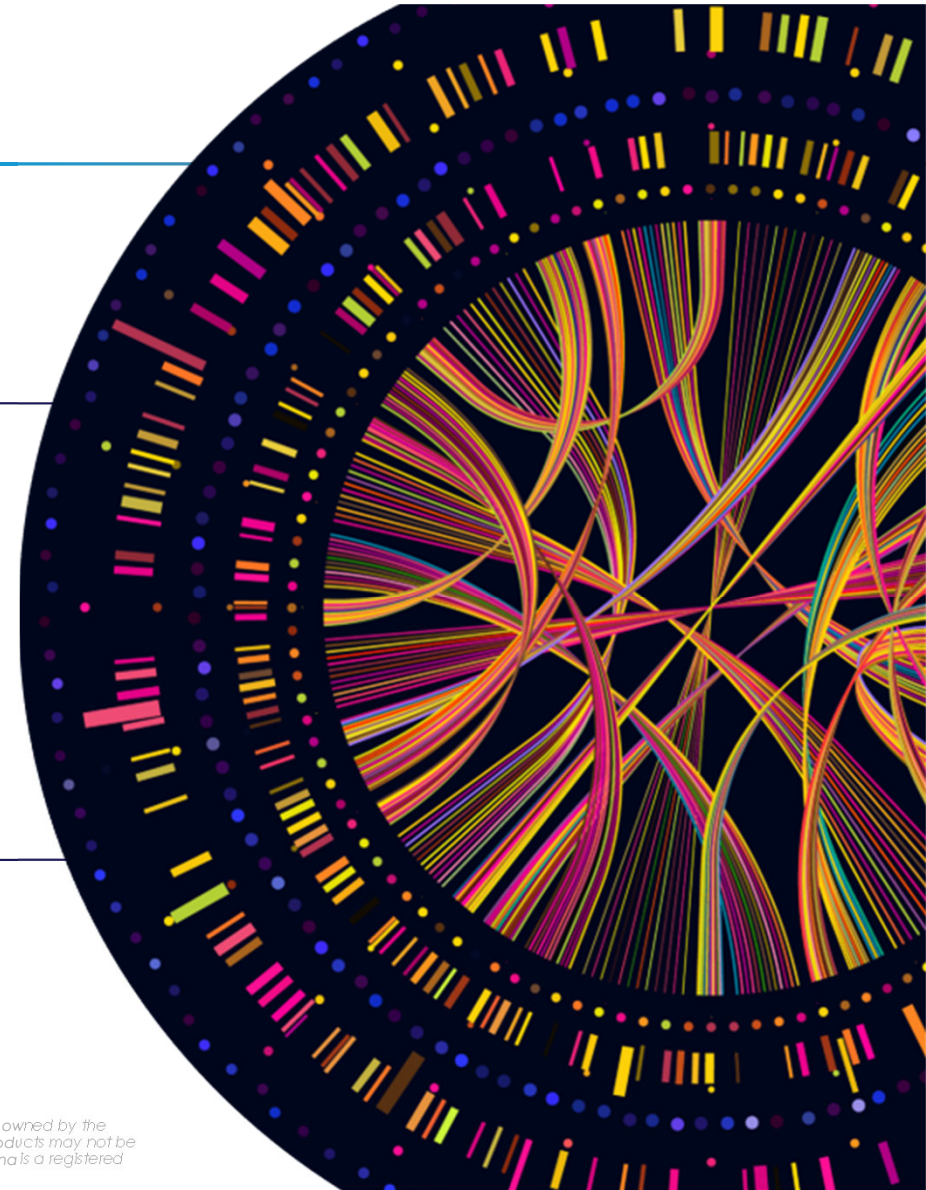
Juan Lopera, PhD

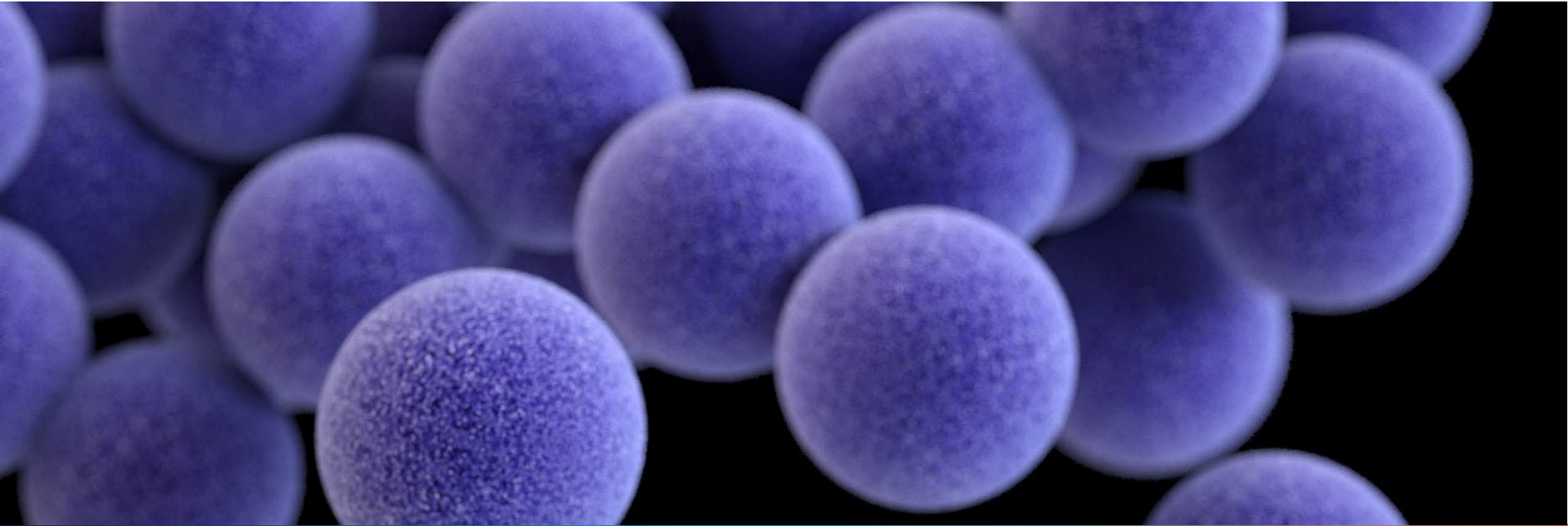
Marco Riojas, PhD

... and One Codex!

JOIN OUR TEAM! We're hiring!

© 2021 American Type Culture Collection. The ATCC trademark and trade name, and any other trademarks listed in this publication are trademarks owned by the American Type Culture Collection unless indicated otherwise. These products are for laboratory use only. Not for human or diagnostic use. ATCC products may not be resold, modified for resale, used to provide commercial services, or to manufacture commercial products without prior ACC written approval. Illumina is a registered trademark of Illumina, Inc. Oxford Nanopore is a registered trademark of Oxford Nanopore Technologies Limited.





Thank you