



Advancing Authentication through Credible Standards and Robust Next-generation Sequencing Workflows

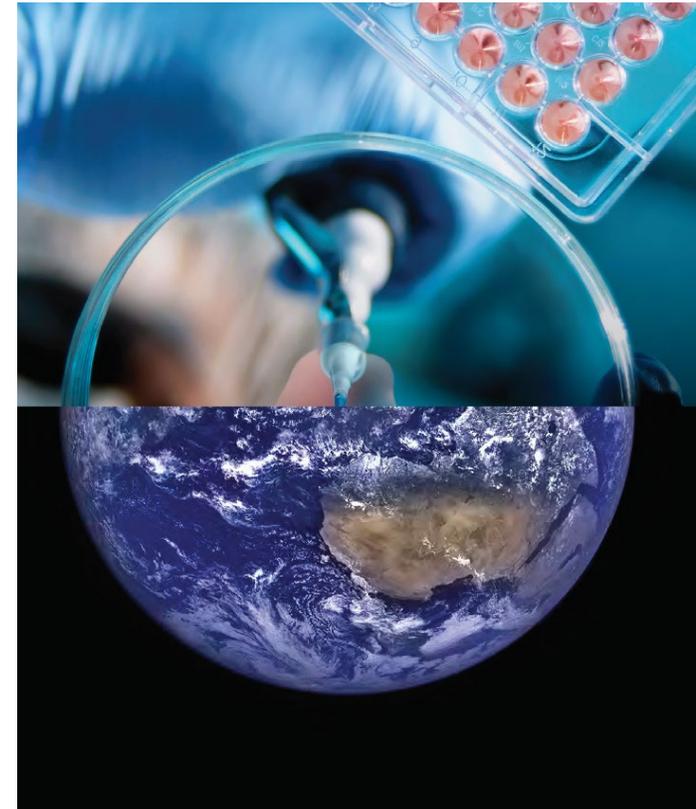
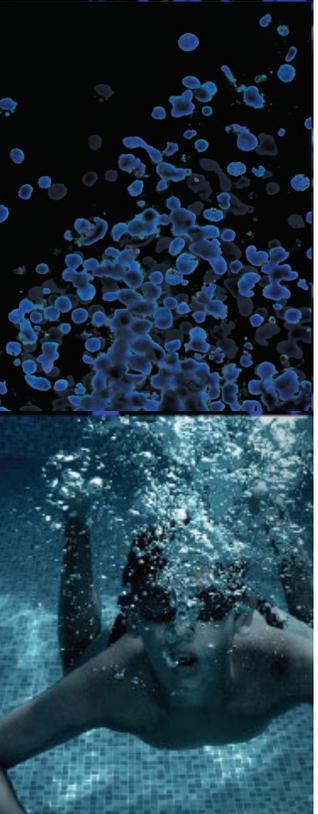
Briana Benton, BS
Technical Manager, ATCC

Sam Minot, PhD
Scientist, Fred Hutchinson Cancer Research Center

Nick Greenfield, MA
Founder and CEO, One Codex

Marco Riojas, PhD
Scientist, ATCC

Credible Leads to InCredible™

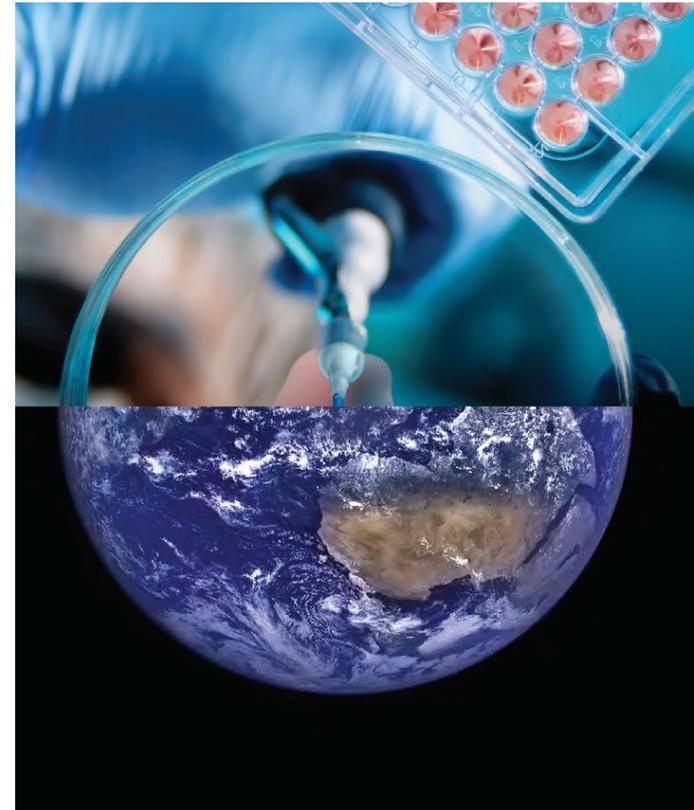
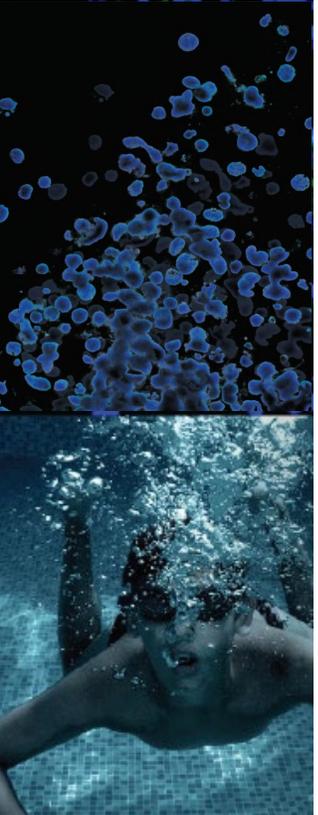




Making Sense Out of Microbiome Data – The Importance of Standards

Briana Benton, BS
Technical Manager, ATCC

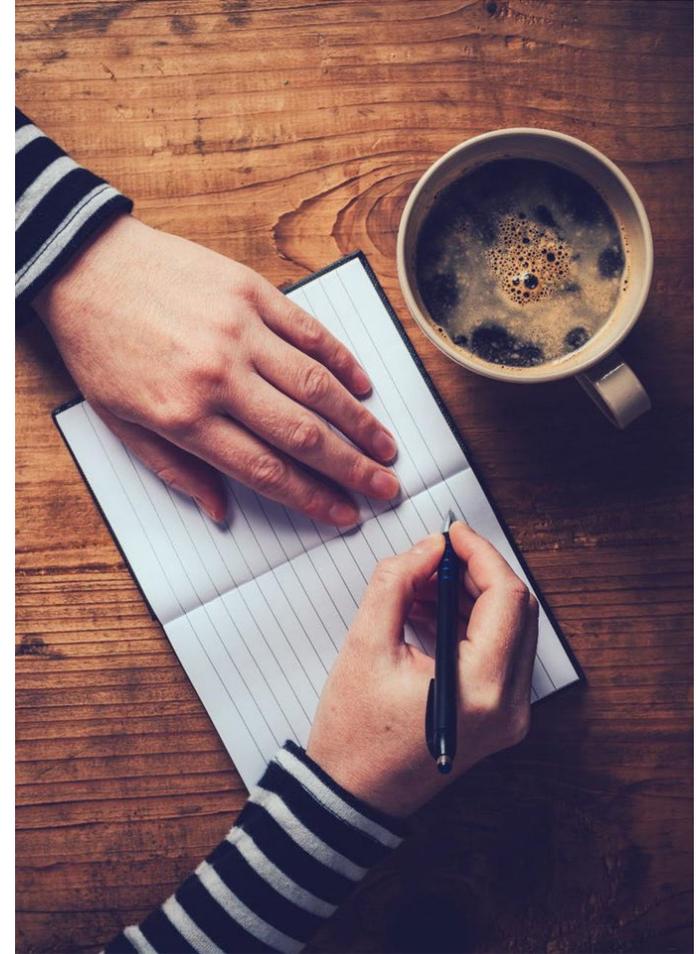
Credible Leads to Incredible™



Agenda

- ✓ Challenges in microbiome analysis and development of standards
- ✓ The ATCC® Microbiome Standards portfolio and upcoming new products
- ✓ Applications of standards in microbiome research

- ✗ Microbiome assay development
- ✗ Show the best data
- ✗ Recommend any specific assay, kit, protocol, or instrument





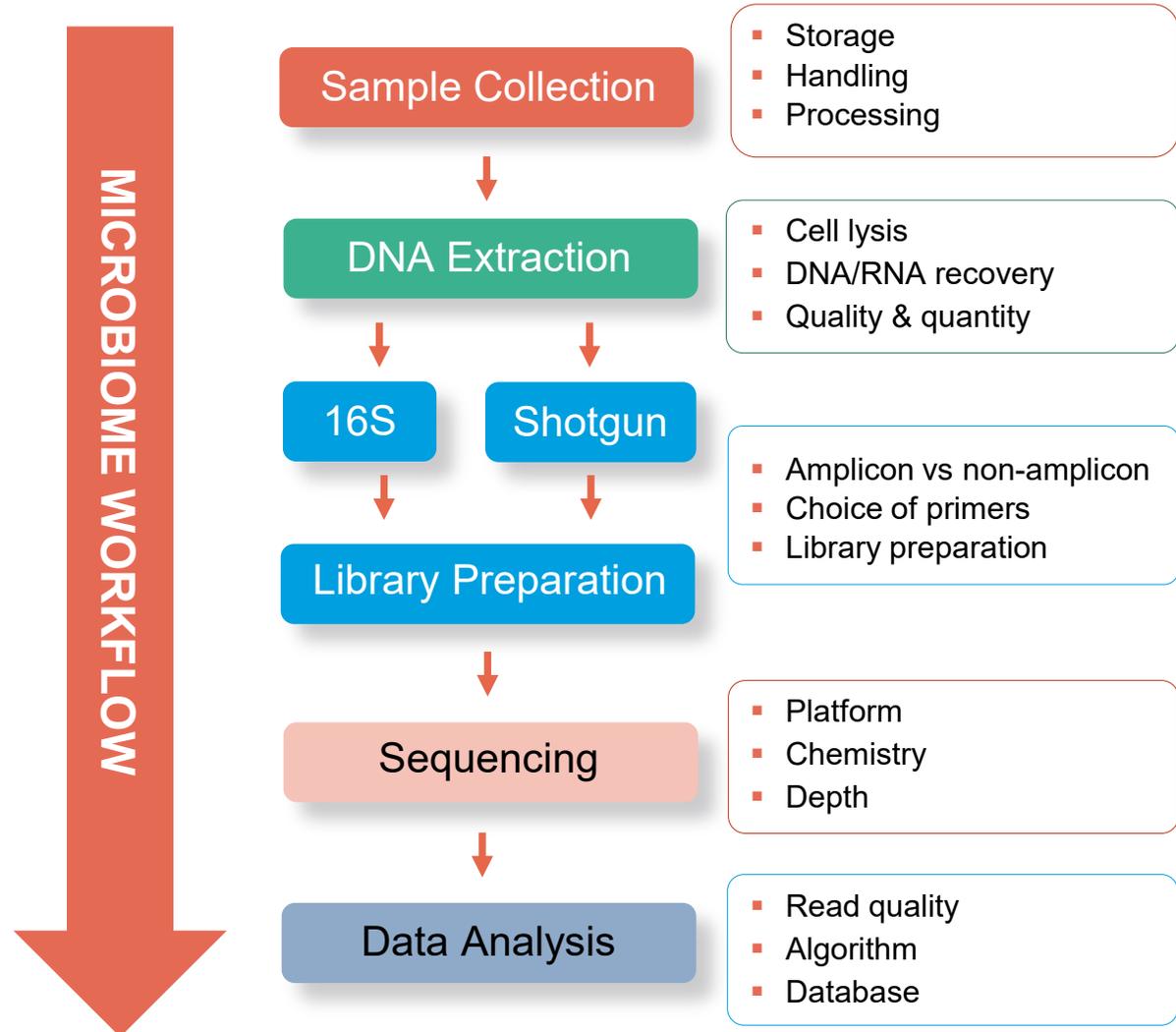
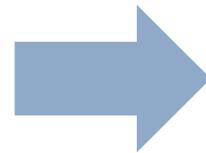
Microbiome Research

The microbiome field is rapidly moving toward translational research pertinent to human health and disease, therapeutics, and personalized medicine

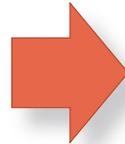
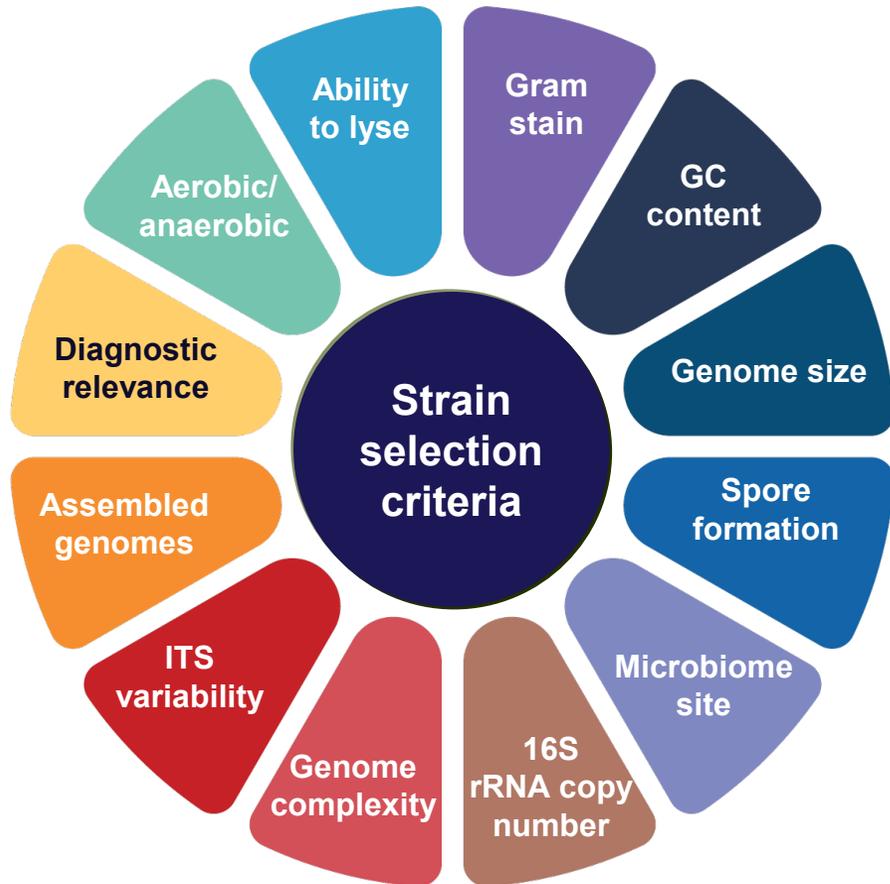
Challenges in Microbiome Research

sequencing viability coverage
extraction technology
amplification
depth bioinformatics

bias



Development of Mock Microbial Communities



Whole Cell Standards

- Authenticated ATCC cultures
- Growth and image cytometry cell counting
- Mixed in even proportion based cell numbers cells
- Storage at 4°C

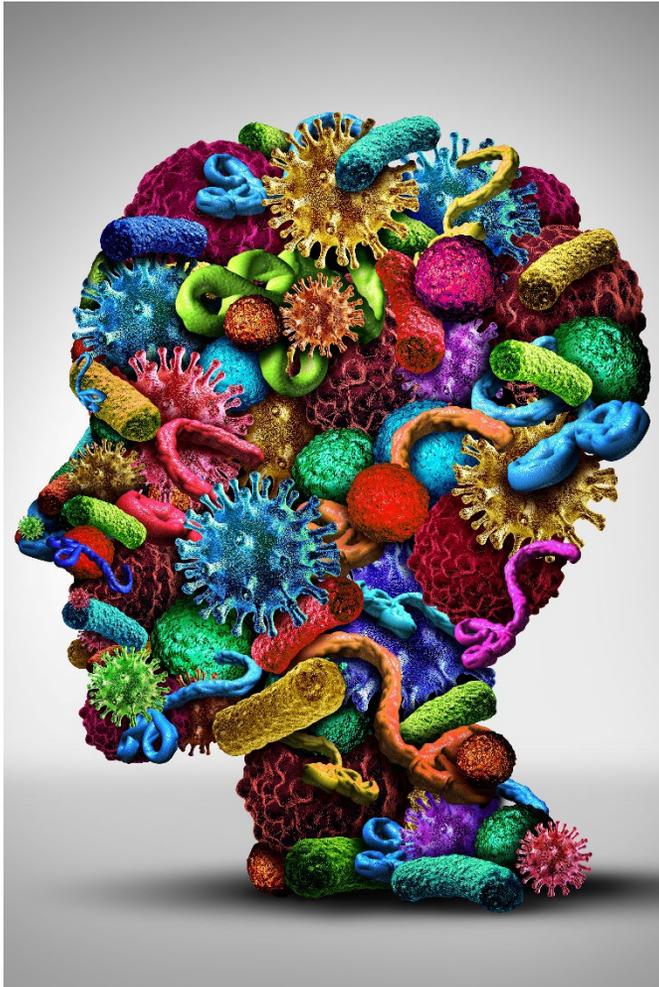


Genomic DNA Standards

- Authenticated ATCC nucleic acids
- Fluorescent dye-based quantification
- Mixed in even proportions based genome copy number
- Storage at -20°C

Assay development, optimization, verification, and quality control

ATCC® Microbiome Standards Portfolio



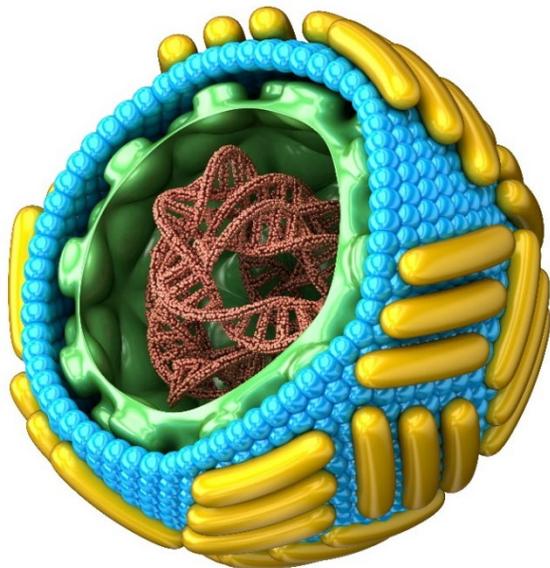
Preparation	ATCC® Catalog No.	Number of Organisms	Composition	Complexity	Importance
Genomic DNA	MSA-1000™	10	Even	Medium	Standards for assay development and optimization
	MSA-1001™	10	Staggered	Medium	
	MSA-1002™	20	Even	High	
	MSA-1003™	20	Staggered	High	
Whole cell	MSA-2003™	10	Even	Medium	
	MSA-2002™	20	Even	High	
Genomic DNA	MSA-4000™	11	Staggered	Medium	NGS-based pathogen detection
Genomic DNA	MSA-3000™	6	Even	Low	Environmental studies
	MSA-3001™	10	Even	Medium	
	MSA-3002™	10	Staggered	Medium	

Site-specific Microbiome Standards



Standard	Preparation	ATCC® Catalog No.	Number of Organisms	Importance
Oral	Whole cell	MSA-2004™	6	<ul style="list-style-type: none"> • Mock microbial communities representing the oral, skin, gut, and vaginal microbiomes • Comprises normal and atypical flora • Anaerobic and aerobic microbial strains • A combination of Gram-positive and -negative bacterial cultures • Even composition
	Genomic DNA	MSA-1004™		
Skin	Whole cell	MSA-2005™	6	
	Genomic DNA	MSA-1005™		
Gut	Whole cell	MSA-2006™	12	
	Genomic DNA	MSA-1006™		
Vaginal	Whole cell	MSA-2007™	6	
	Genomic DNA	MSA-1007™		

ATCC Virome Standards



Composition of Virome Standards

Human herpesvirus 5 strain AD169 (ATCC® VR-538™)

Human mastadenovirus strain F (ATCC® VR-931™)

Influenza B virus strain B/Florida/4/2006 (ATCC® VR-1804™)

Zika virus strain MR 766 (ATCC® VR-1838™)

Reovirus 3 strain Dearing (ATCC® VR-824™)

Human respiratory syncytial virus strain A2 (ATCC® VR-1540™)

Standard	Preparation	ATCC® Catalog No.	Number of Organisms	Specification (ddPCR™)	Applications
Virome	Virus Mix	MSA-2008™	6	2×10^3 genome copies/ μ L per virus	Standards for virome assay development, optimization, verification, and validation; evaluating reproducibility; and use as a daily run quality control
	Nucleic Acid Mix	MSA-1008™	6	2×10^4 genome copies/ μ L per virus	

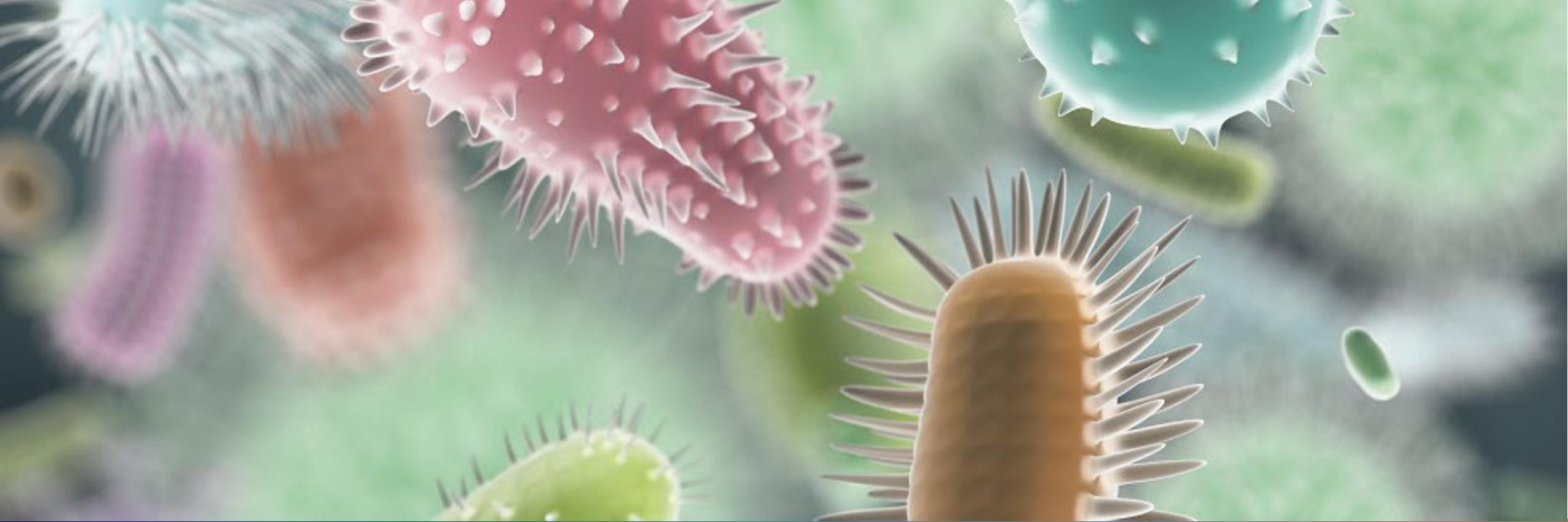
Spike-in and Mycobiome Standards



Standard	Preparation	ATCC® Catalog No.	Number of Organisms	Importance
Spike-in	Whole cell	MSA-2014™	3	<ul style="list-style-type: none"> • Microbiome measurements and data normalization • 16S rRNA and shotgun assay verification, validation, and quality control
	Genomic	MSA-1014™		



Standard	Preparation	ATCC® Catalog No.	Number of Organisms	Importance
Mycobiome	Whole cell	MSA-2010™	10	<ul style="list-style-type: none"> • Fungal mock community standards for assay development, optimization, verification, and validation; evaluating reproducibility; and use as a daily run quality control
	Genomic	MSA-1010™		



Utility and Application of Microbiome Standards

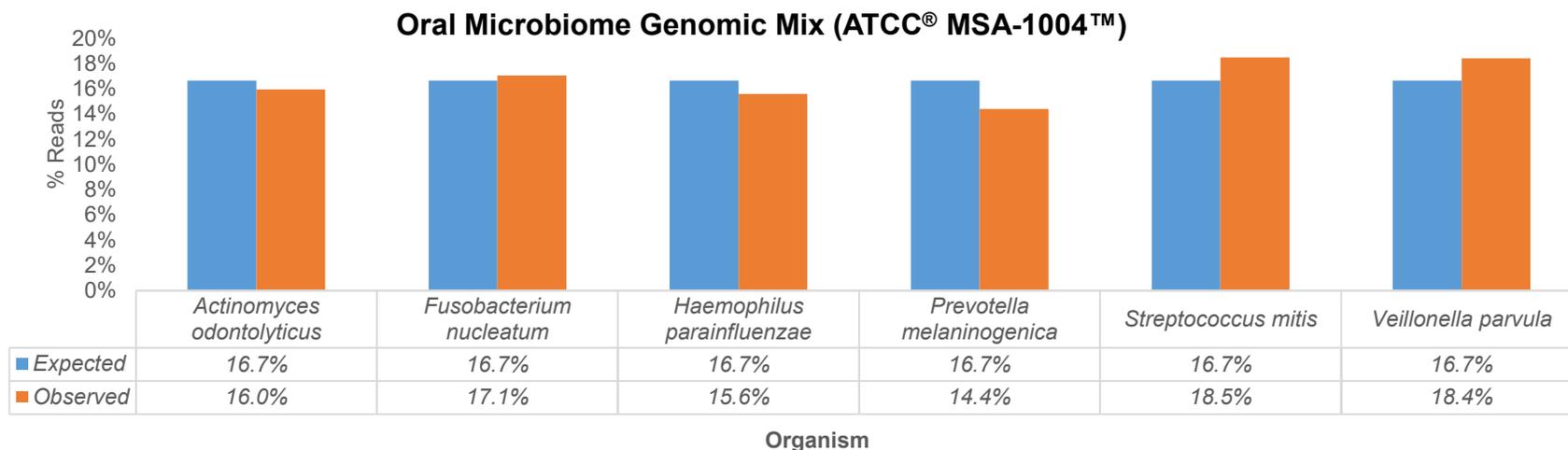


Evaluating DNA Extraction Methods and Kits

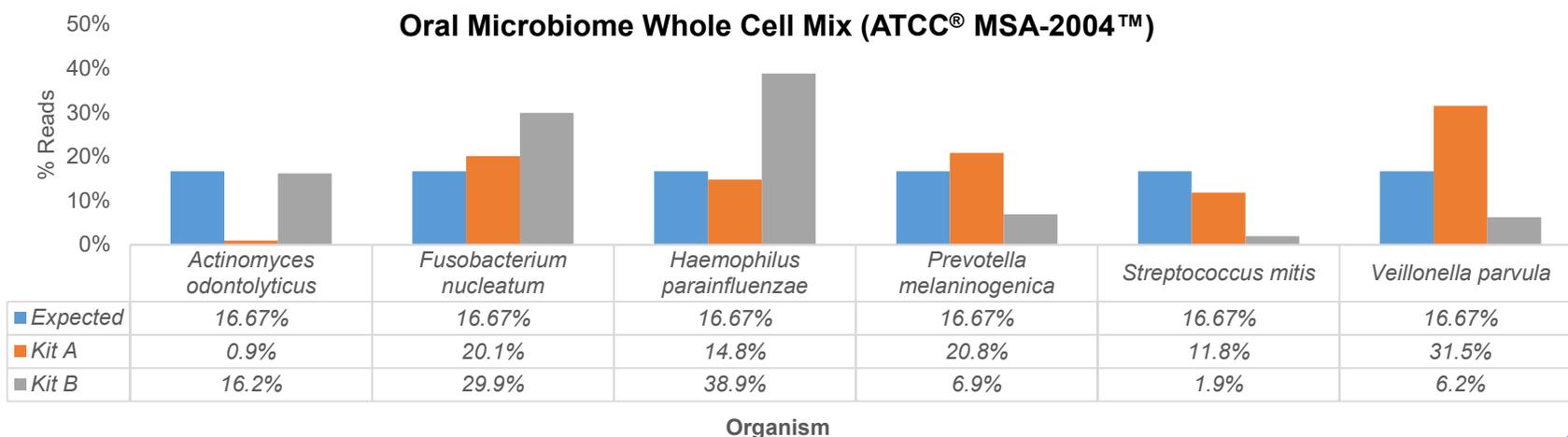
Genomic Versus Whole Cell Standards

DNA extraction methods are not perfect

Shotgun metagenomic analysis of the Oral Microbiome Genomic Mix



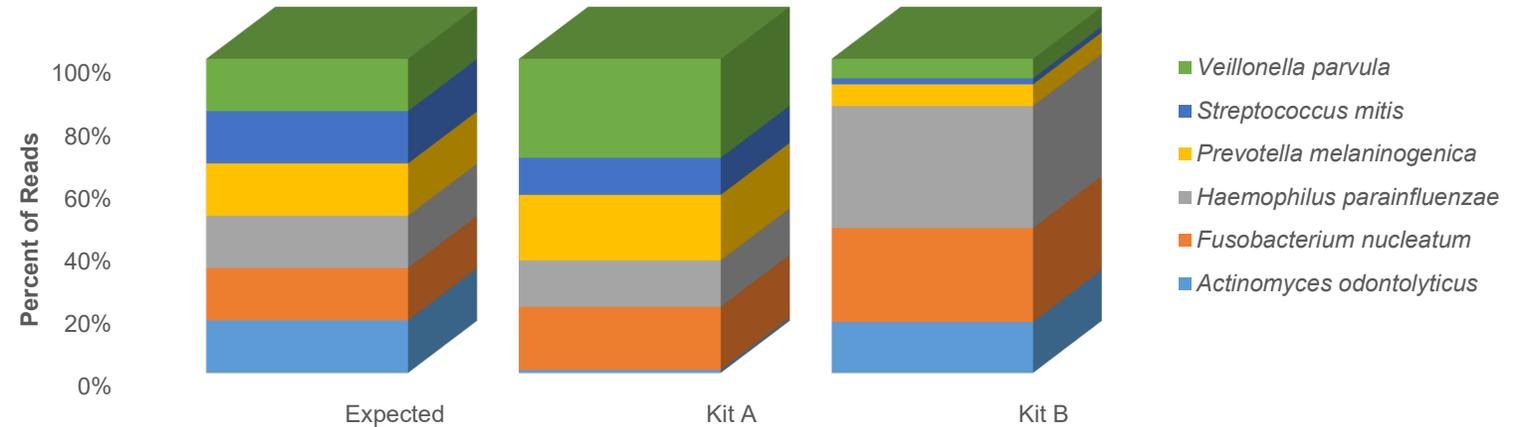
DNA extraction from the Oral Microbiome Whole Cell Mix with two different kits followed by shotgun metagenomic analysis



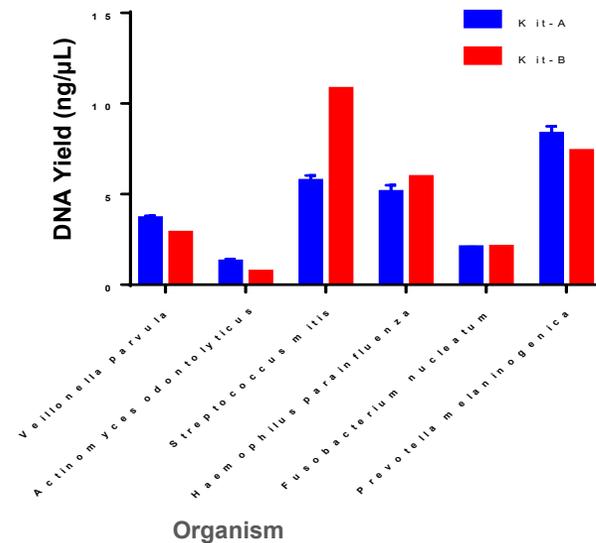
Assess Biases in DNA Extraction

Compare different pre-treatments and extraction methods, optimize protocols, and validate different kits

DNA extraction from the Oral Whole Cell Mix with two different kits followed by shotgun metagenomic analysis



DNA extraction from individual strains that are components of the Oral Whole Cell Mix



Organism	Number of Cells per Component	Gram Stain	Genome size	%GC
<i>Actinomyces odontolyticus</i>	~2x10 ⁷	+	2.39	65.5
<i>Fusobacterium nucleatum</i>		-	2.17	27.2
<i>Haemophilus parainfluenzae</i>		-	2.12	39.3
<i>Prevotella melaninogenica</i>		-	3.17	35.1
<i>Streptococcus mitis</i>		+	1.83	40.5
<i>Veillonella parvula</i>		-	2.16	38.6

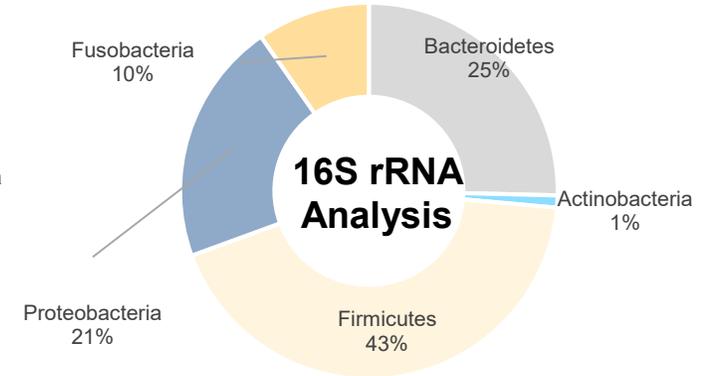
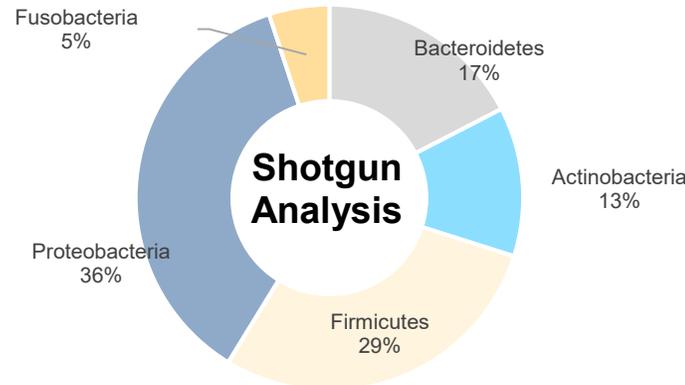
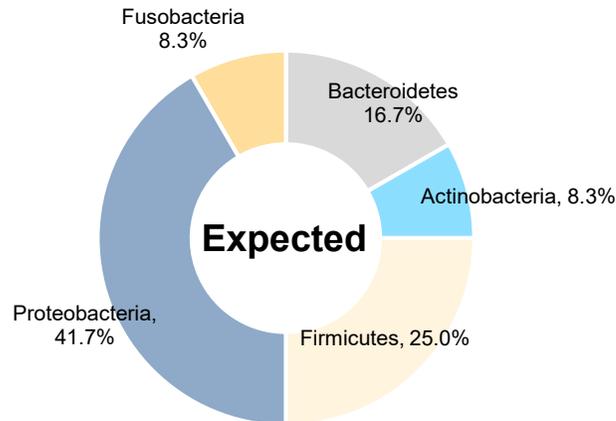
Extraction Kit

Gut Whole Cell Standard

Profiling of gut microbiome standard at the phylum, genus, and species level

The Gut Whole Cell Microbiome Standard (ATCC® MSA-2006™) can be used as a full process control for shotgun and 16S rRNA assays

Phylum	Expected	Strain	Expected	Observed-Shotgun	Observed-16S (V1V2)
Bacteroidetes	16.6%	<i>Bacteroides fragilis</i>	8.3%	12.3%	17.0%
		<i>Bacteroides vulgatus</i>	8.3%	8.6%	8.4%
Actinobacteria	8.33%	<i>Bifidobacterium adolescentis</i>	8.3%	12.0%	1.0%
Firmicutes	25.0%	<i>Clostridium difficile</i>	8.3%	16.5%	29.0%
		<i>Enterococcus faecalis</i>	8.3%	3.1%	1.6%
		<i>Lactobacillus plantarum</i>	8.3%	8.1%	12.3%
Proteobacteria	41.7%	<i>Enterobacter cloacae</i>	8.3%	10.6%	4.0%
		<i>Escherichia coli</i>	8.3%	6.6%	3.3%
		<i>Helicobacter pylori</i>	8.3%	3.8%	7.7%
		<i>Salmonella enterica</i>	8.3%	4.9%	2.2%
		<i>Yersinia enterocolitica</i>	8.3%	8.8%	3.6%
Fusobacteria	8.3%	<i>Fusobacterium nucleatum</i>	8.3%	4.8%	9.7%



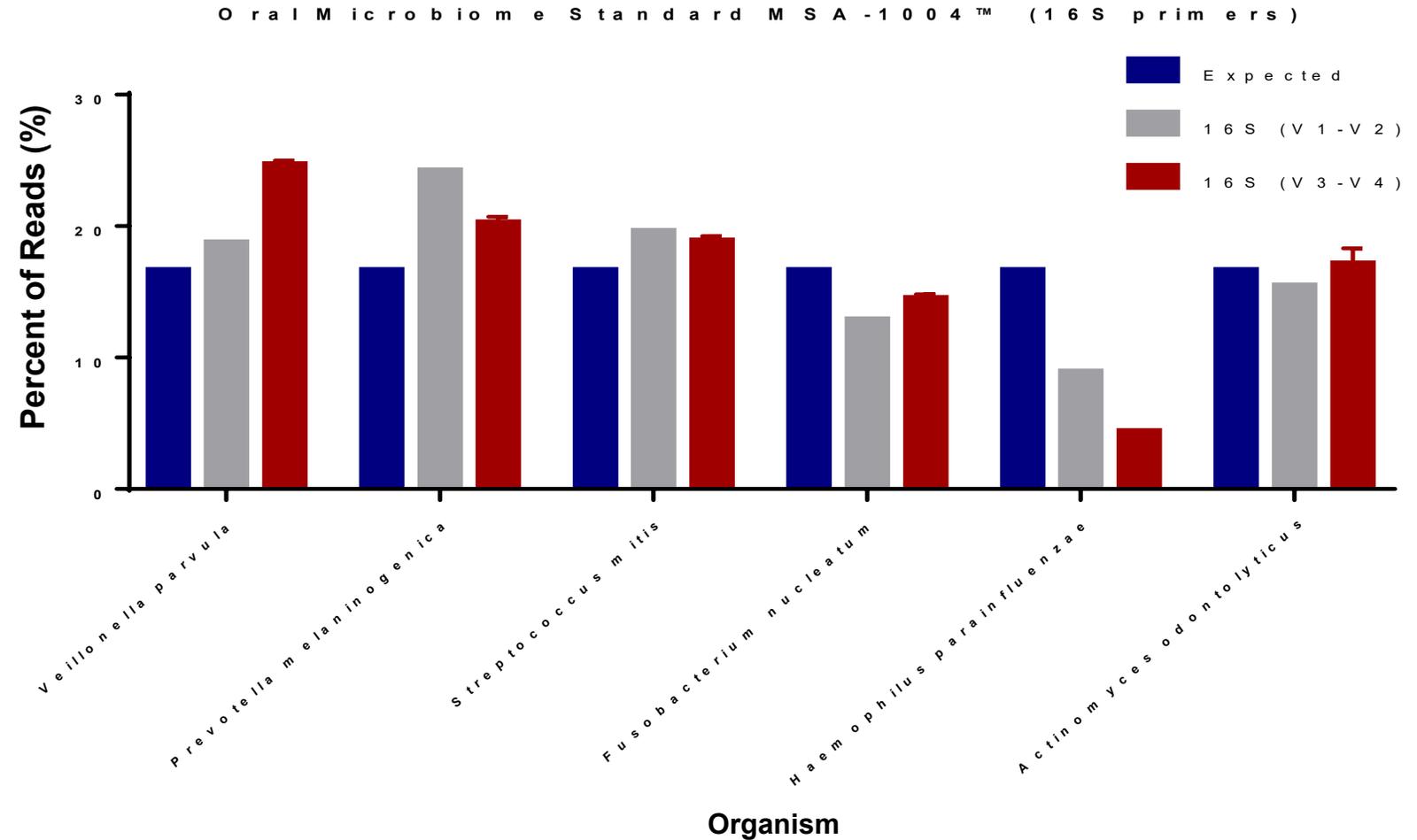


Evaluating 16S rRNA and WGS Library Kits

16S Amplicon-based Analysis: Primer Selection

Compare different primer sets, optimize amplification steps, and validate 16S analysis methods

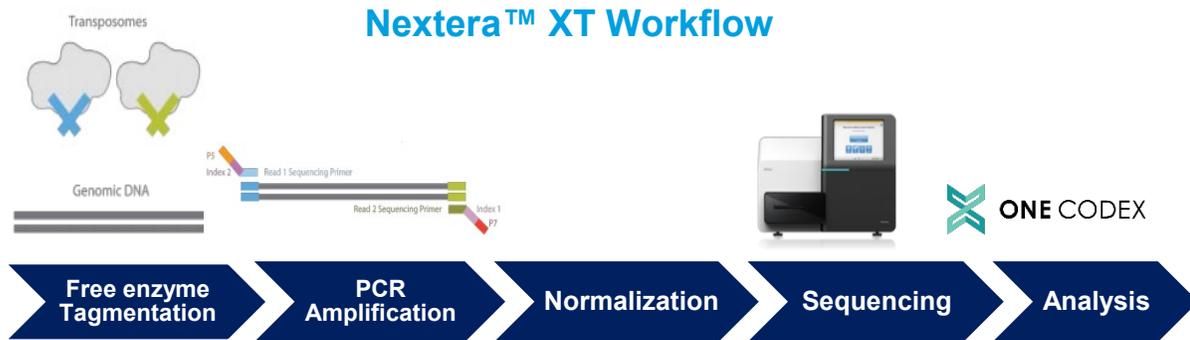
16S rRNA analysis of the Oral Genomic DNA Standard via two primer sets



Comparing Library Preparation Kits

Nextera Flex enables uniform coverage of genomes of low GC content

Oral Microbiome Genomic DNA (ATCC® MSA-1004™)



Sample Composition

Name	Estimated Abundance
Actinomyces odontolyticus	46.07%
Prevotella melaninogenica	16.09%
Streptococcus mitis	16.04%
Veillonella parvula	12.65%
Haemophilus parainfluenzae	7.80%
Fusobacterium nucleatum	1.34%



Sample Composition

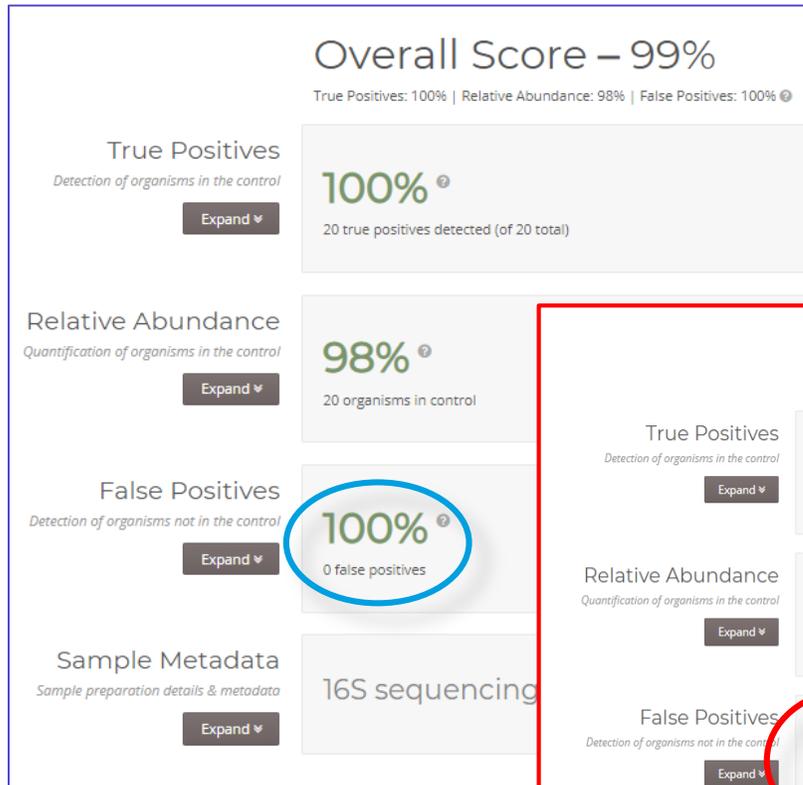
Name	Estimated Abundance
Streptococcus mitis	18.83%
Veillonella parvula	18.42%
Fusobacterium nucleatum	17.13%
Haemophilus parainfluenzae	15.68%
Actinomyces odontolyticus	15.54%
Prevotella melaninogenica	14.40%



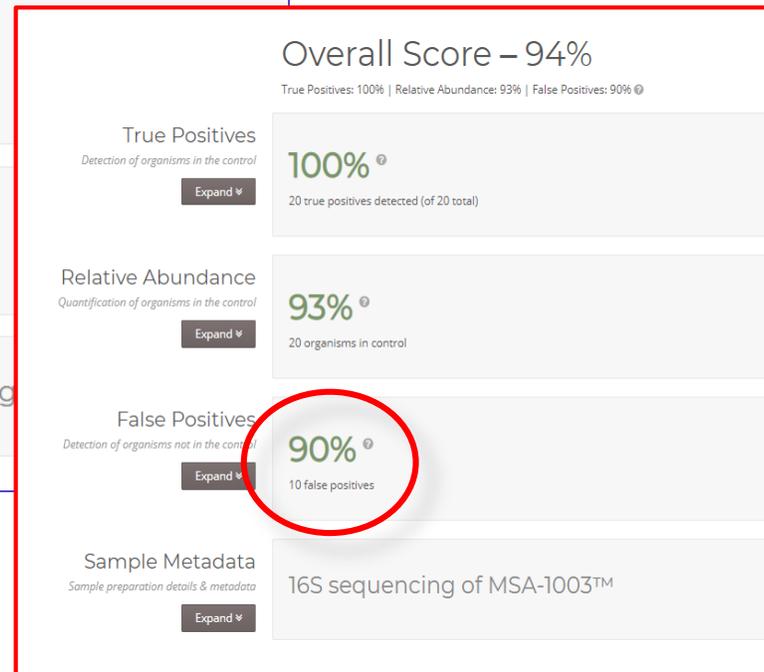
Comparing Library Preparation Kits

The LoopSeq™ 16S rRNA long-read method allows highest sequence accuracy and species-level taxonomy

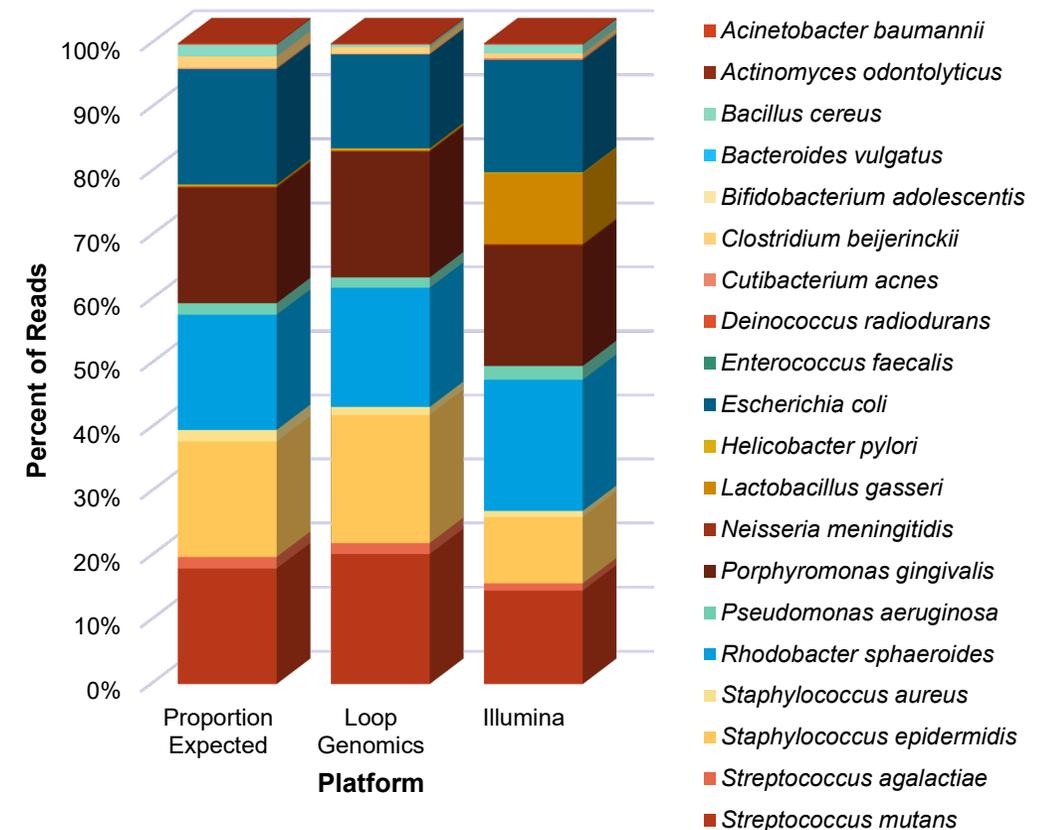
Loop Genomics



Short Reads



Genomic DNA (ATCC® MSA-1003™)



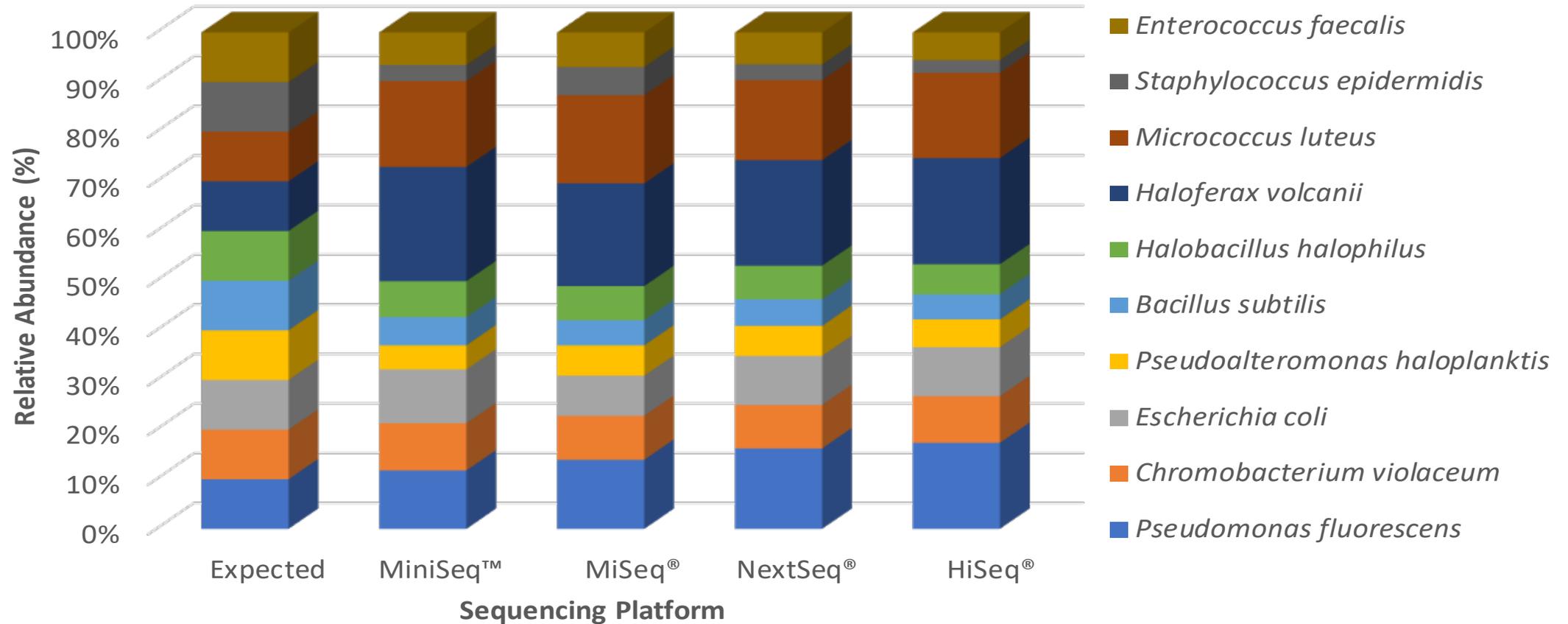


Evaluating NGS Platforms

Short-read Sequencing Platform: Illumina®

Assay reproducibility through different Illumina sequencing platforms

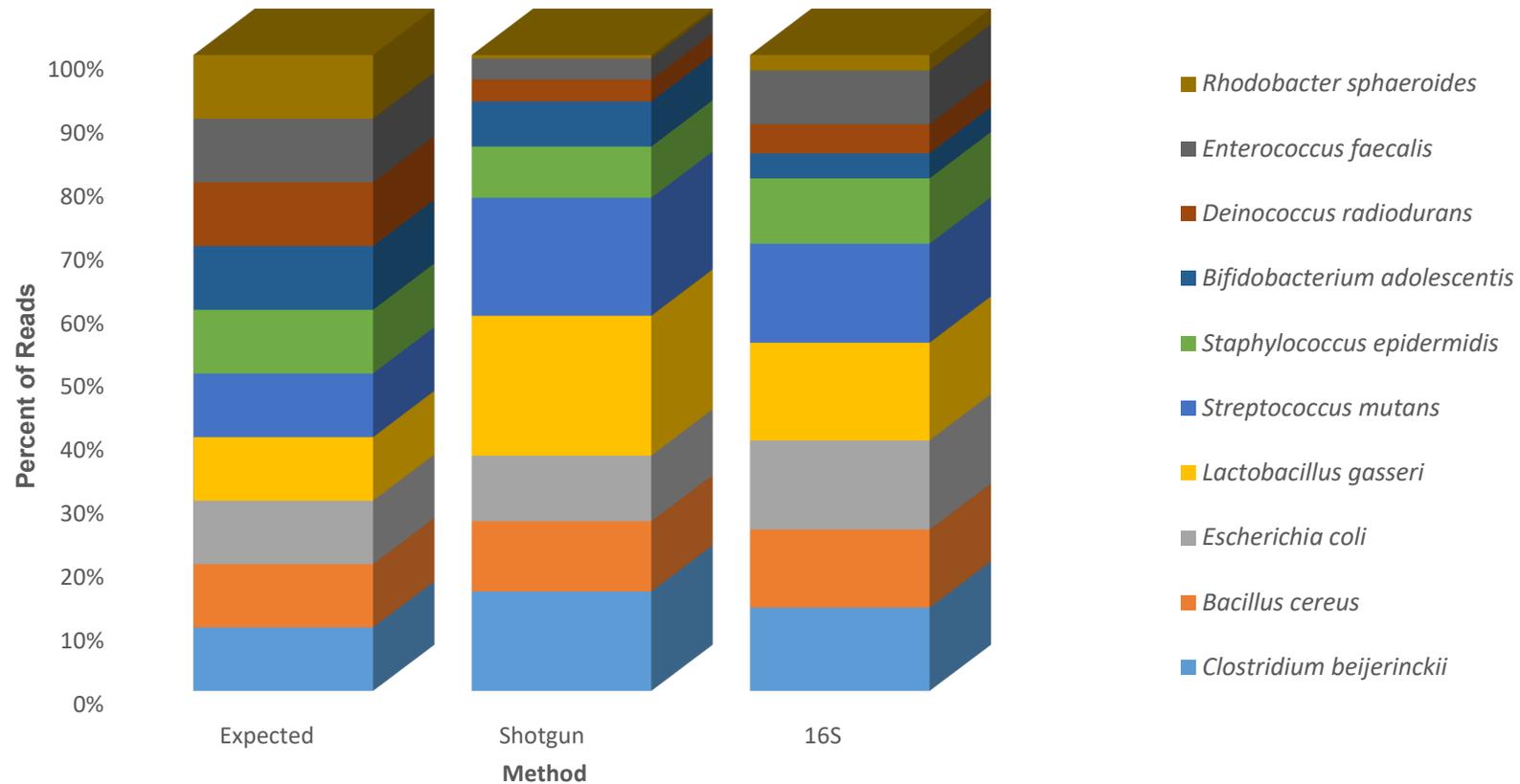
Shotgun Metagenomic Data (ATCC® MSA-3001™)



Short-read Sequencing Platform: Ion Torrent™

16S rRNA and shotgun data on the Ion GPM Platform (ATCC® MSA-1000™)

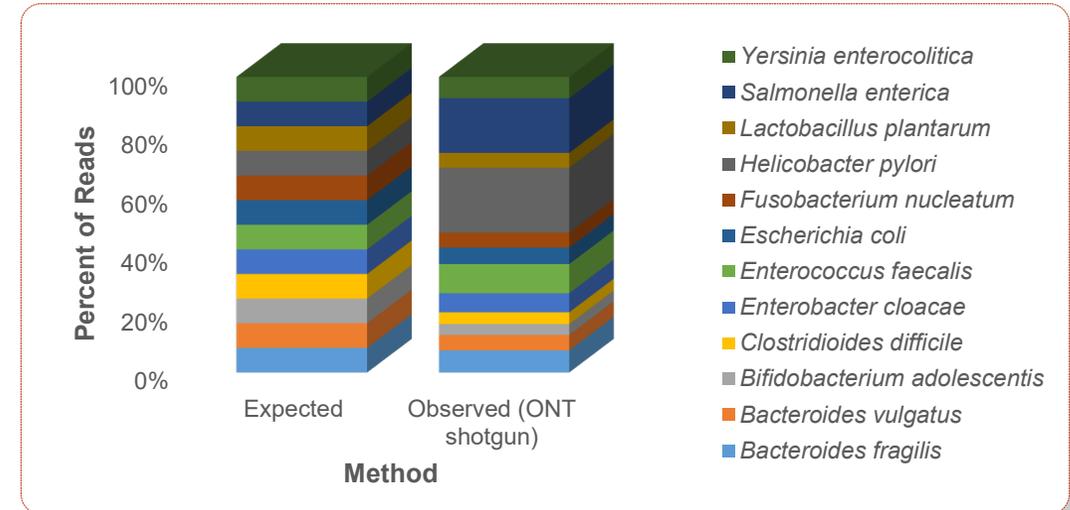
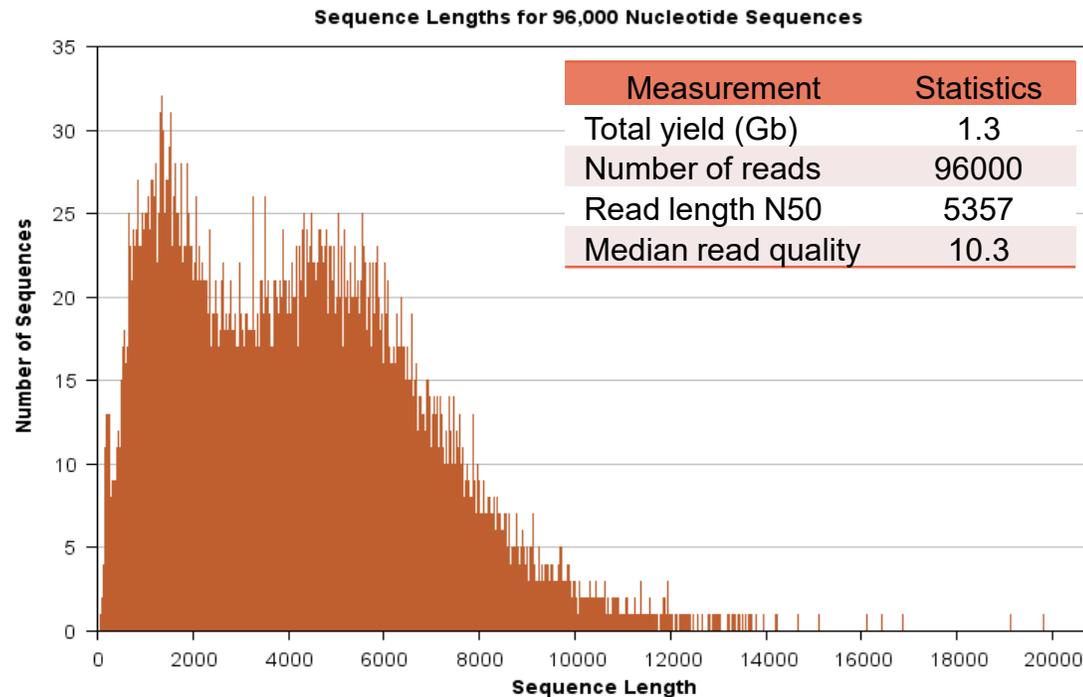
Shotgun vs 16S rRNA assay (V1/V2) (ATCC® MSA-1000™)



Long-read Sequencing Platform: Nanopore®

One hour sequencing coverage was enough to identify all organisms in the mix with sufficient genome coverage

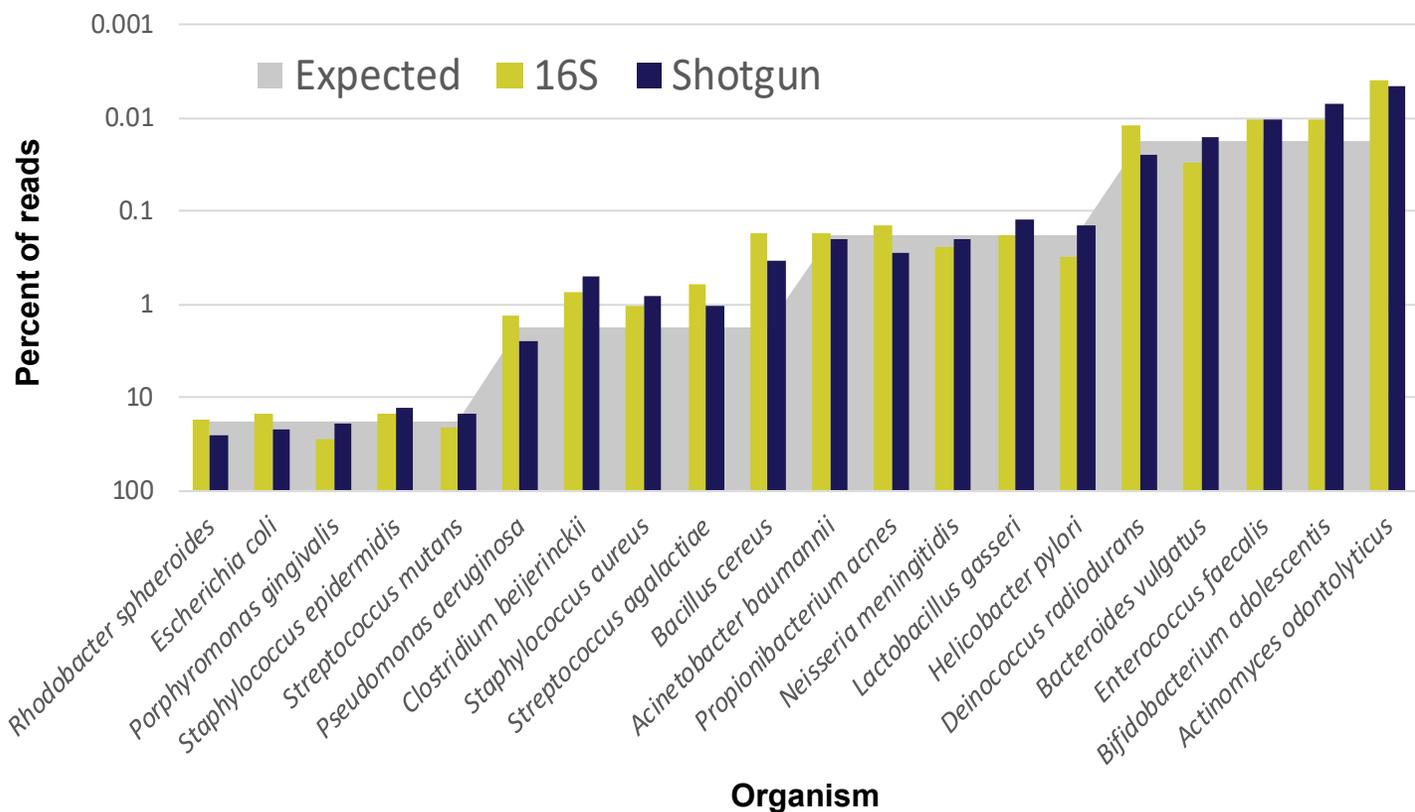
The Gut Microbiome Whole Cell Standard (ATCC® MSA-1006™) was analyzed via shotgun sequencing on the MinION platform



Organism	Genome Coverage (x)
<i>Enterobacter cloacae</i>	9.1
<i>Enterococcus faecalis</i>	14.1
<i>Bacteroides fragilis</i>	8.1
<i>Bacteroides vulgatus</i>	6.6
<i>Bifidobacterium adolescentis</i>	1.8
<i>Clostridioides difficile</i>	7.4
<i>Escherichia coli</i>	6.5
<i>Fusobacterium-nucleatum</i>	4.6
<i>Helicobacter pylori</i>	16.6
<i>Lactobacillus plantarum</i>	6.0
<i>Salmonella enterica</i>	11.1
<i>Yersinia enterocolitica</i>	11.3

Long-read Sequencing Platform: PACBIO®

16S rRNA (full-length) and shotgun data on the PacBio Sequel Platform (ATCC® MSA-1003™)



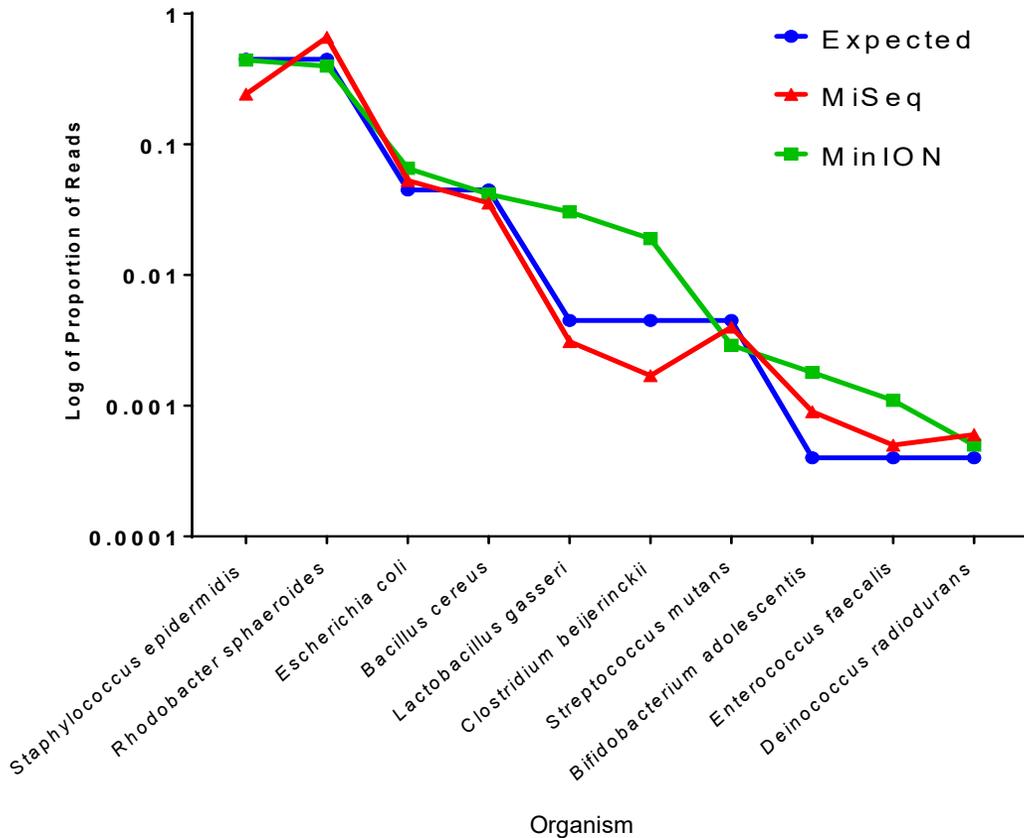
ATCC quality control score (One Codex)

One Codex Analysis	16S rRNA run 1	16S rRNA run 2	Shotgun run 1	Shotgun run 2
True positives	100%	100%	100%	100%
Relative abundance	95%	95%	97%	97%
False positives	100%	100%	88%	84%
Overall score	98%	98%	95%	95%

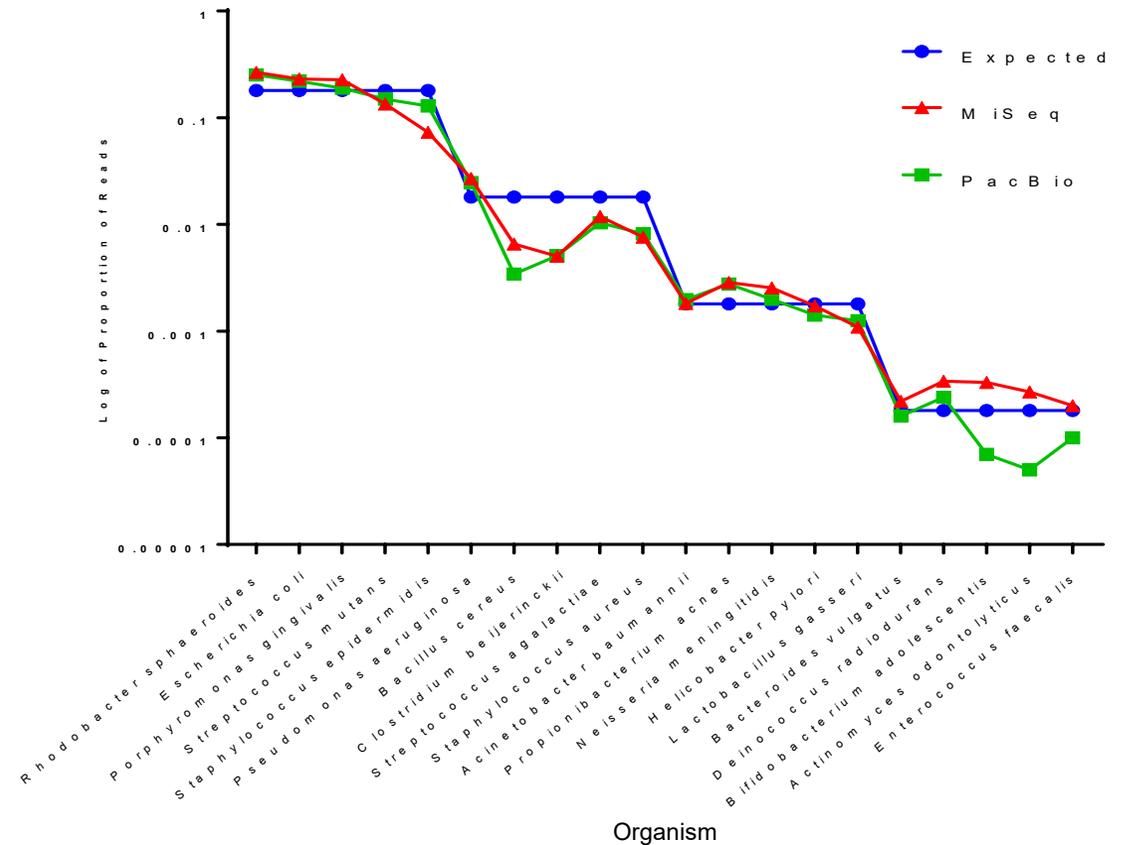
Shotgun Metagenomic Analysis: Short vs Long Reads

ATCC Microbiome Standards are technology agnostic

MSA-1001™ (Staggered 10 Strains)



MSA-1003™ (Staggered 20 Strains)





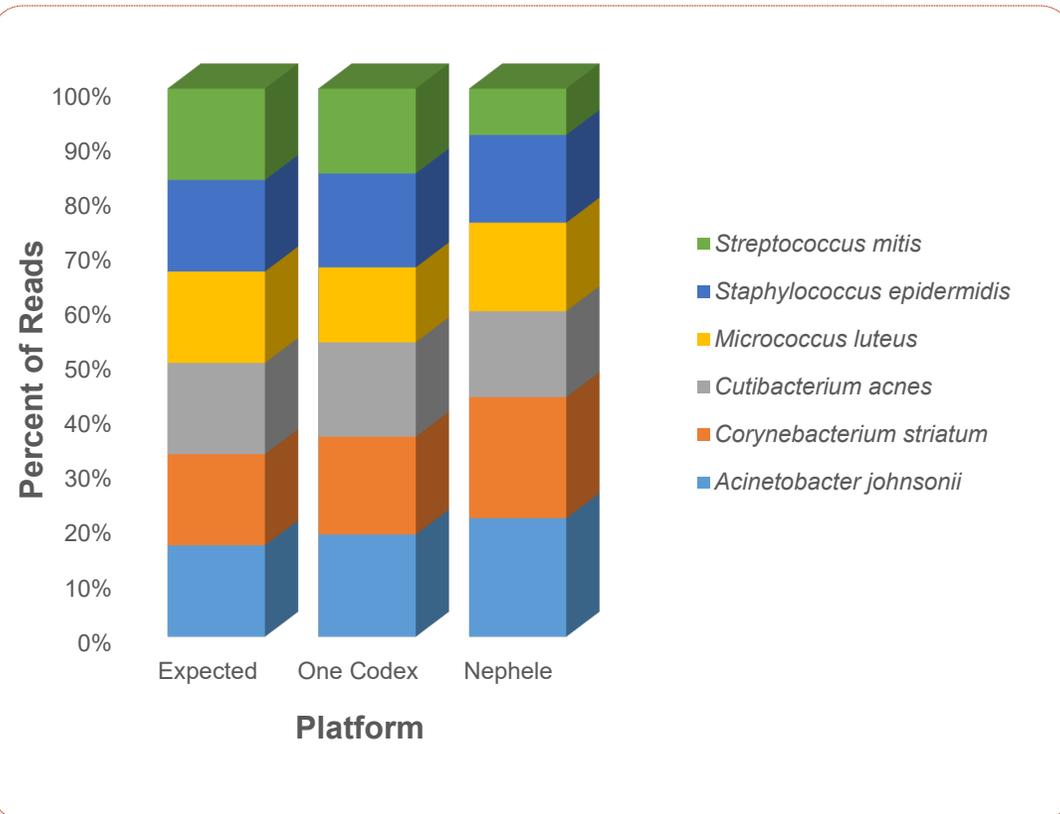
Comparing Bioinformatics and Databases

Data Analysis Using Different Databases

Evaluation of NGS data from microbiome standards in multiple analysis platforms and databases

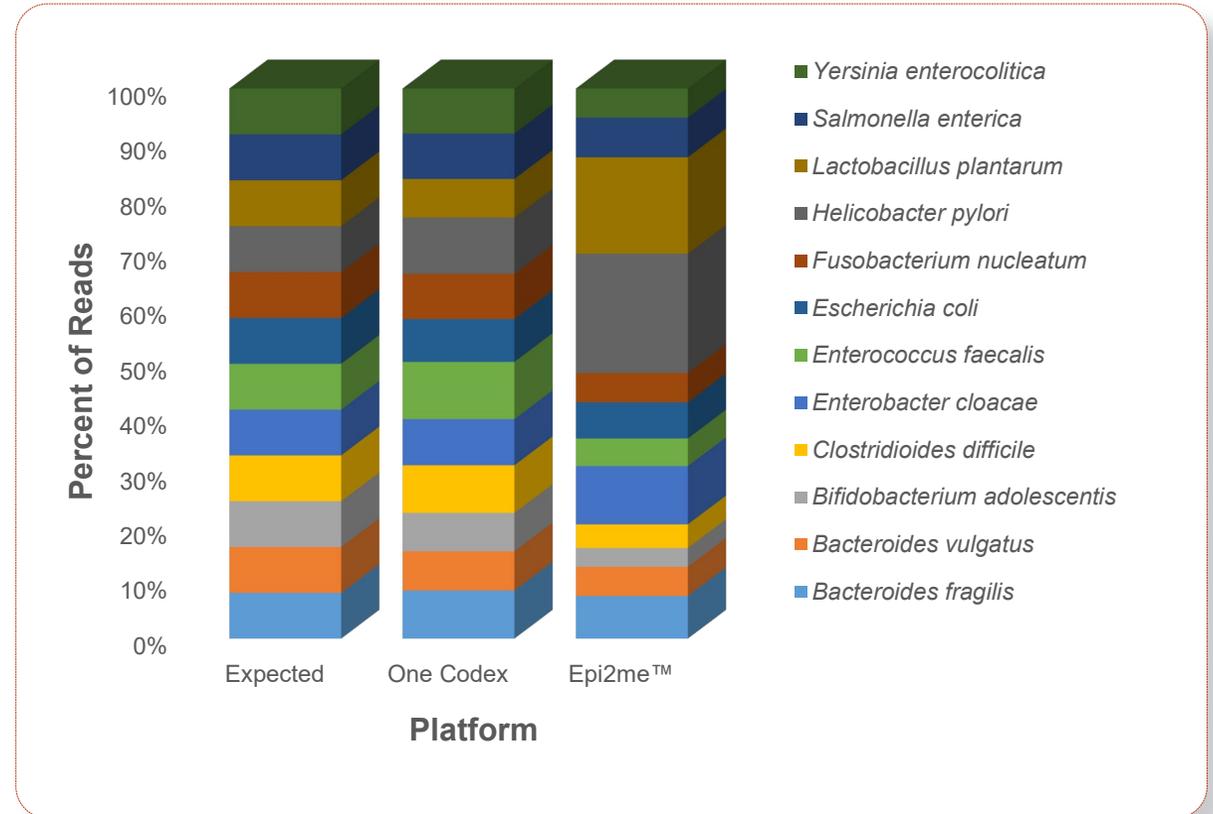
Nephele vs One Codex

Short-read sequencing data from the Skin Genomic DNA Mix (ATCC® MSA-1005™)



Epi2Me vs One Codex

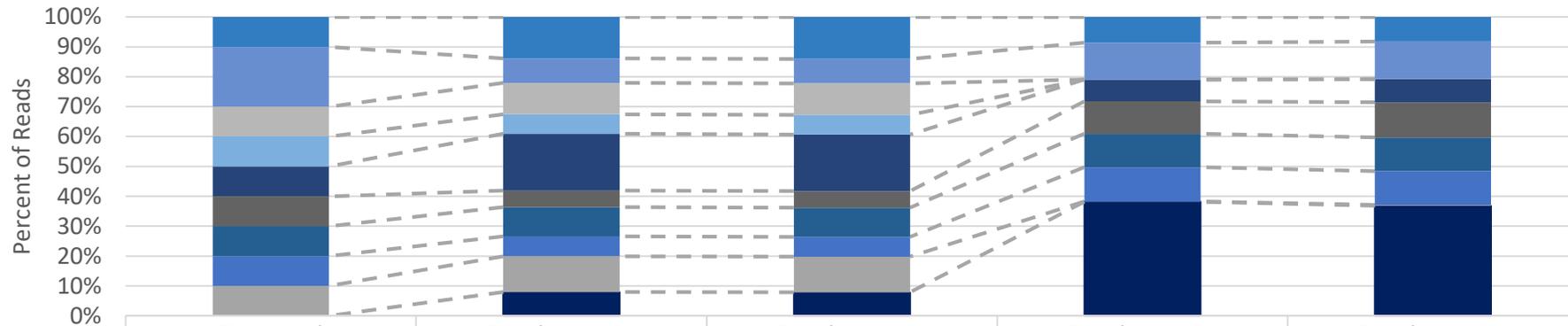
Long-read sequencing data from the Gut Genomic DNA Mix (ATCC® MSA-1006™)



Mycobiome Standards

Data analysis platform impacts strain identification and taxonomic resolution

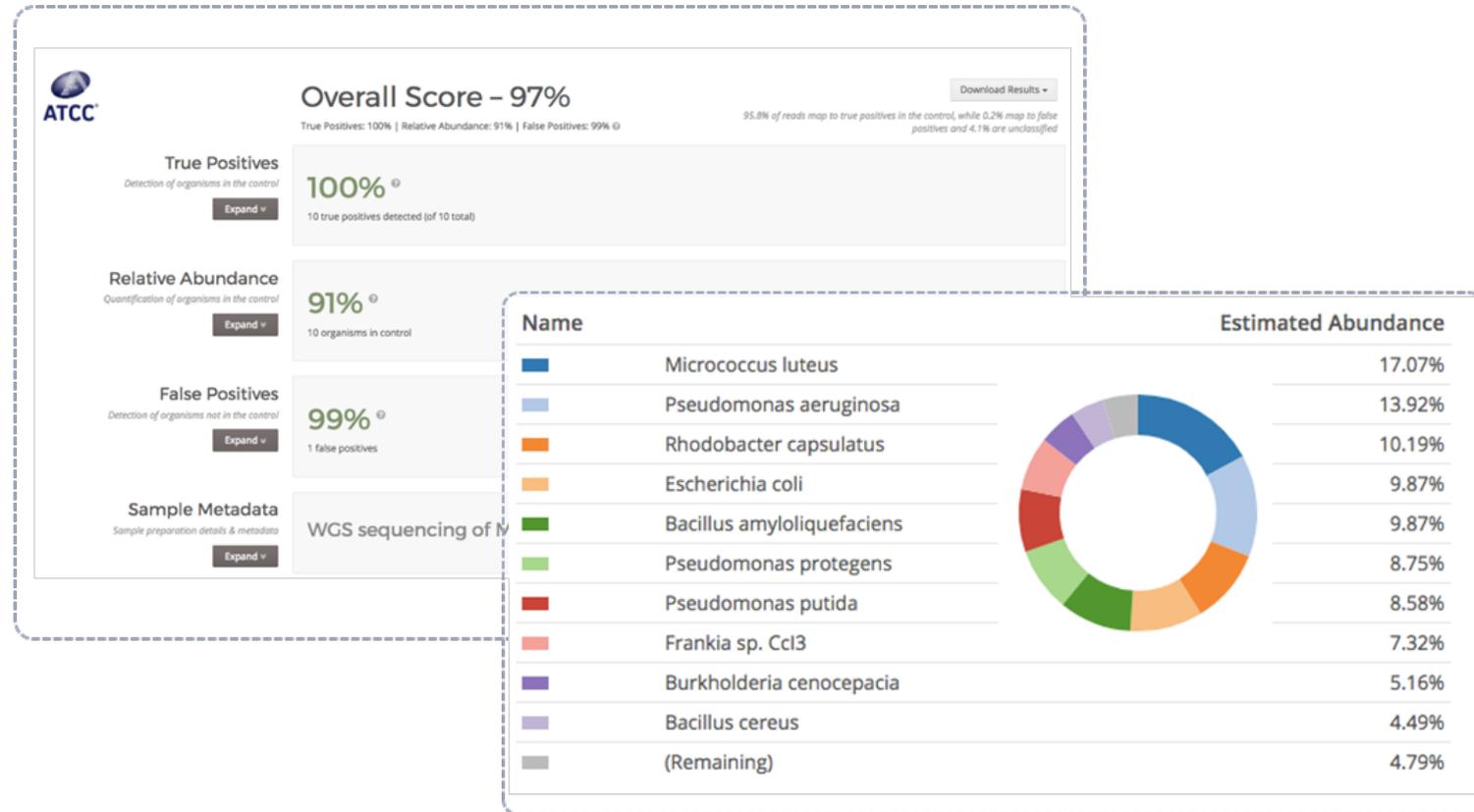
Shotgun Analysis of Genomic DNA standards (ATCC® MSA-1010™)



- Aspergillus
- Candida
- Cryptococcus
- Cutaneotrichosporon
- Fusarium
- Malassezia
- Penicillium
- Saccharomyces
- Trichophyton
- Unclassified

	Expected	Databases A1	Databases A2	Databases B1	Databases B2
Aspergillus	10%	13.88%	13.99%	8.60	8.16
Candida	20%	8.17%	8.22%	12.31	12.61
Cryptococcus	10%	10.52%	10.51%	0.00	0.00
Cutaneotrichosporon	10%	6.53%	6.61%	0.00	0.00
Fusarium	10%	19.00%	18.94%	7.32	7.78
Malassezia	10%	5.53%	5.55%	10.86	11.74
Penicillium	10%	9.81%	9.76%	11.17	11.25
Saccharomyces	10%	6.63%	6.62%	11.48	11.35
Trichophyton	10%	12.01%	11.98%	0.05	0.29
Unclassified	0%	7.93%	7.82%	38.21	36.82

ATCC Data Analysis Solution



WORKFLOW:

1. Drag and drop Fastq files or export via cloud
2. Choose your ATCC product and analysis (16s and shotgun)
3. Download your reports

RESULTS ARE PROVIDED ON A SCORECARD REPORTS:

1. **True positives:** Percentage of organisms detected from the control
2. **False positives:** Detection of organisms not in the control
3. **Relative abundance:** Quantification of organisms in the control

ATCC Data Analysis Solution



Mock Microbial Communities

- Genomic DNA and whole cell standards
- Even and staggered mixtures comprising 10 or 20 strains
- Environmental and pathogen mixtures



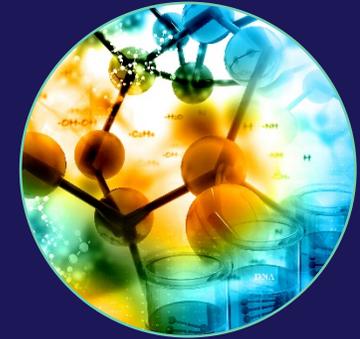
Site-specific Standards

- Genomic DNA and whole cell standards
- Even mixtures of 6-12 strains
- Bacterial strains prevalent in the oral, skin, gut, and vaginal microbiome



Spike-In Standards

- Recombinant strains with a unique DNA tag stably integrated into the chromosome
- Recombinant standards include the Gram-negative and Gram-positive bacteria



New Products

- Genomic DNA and whole cell mock communities representing:
 - Virome
 - Mycobiome

Bundled with data analysis on the One Codex platform

Posters at ASM Microbe 2019

Development of Fungal Mock Community Standards for Mycobiome Studies

Poster Board Number: FRIDAY – MBP-74

Date: Friday, June 21, 2019

Time: 11:00 AM-12:00PM, 4:00 PM-5:00 PM

Utility of Recombinant Bacteria with Unique Tags as Spike-In Controls for Microbiome Studies

Poster Board Number: FRIDAY – CPHM-940

Date: Friday, June 21, 2019

Time: 11:00 AM-12:00PM, 4:00 PM-5:00 PM

Evaluation of ATCC[®] Site-Specific Microbiome Standards on Long-Read Sequencing Platforms

Poster Board Number: SATURDAY – MBP-7

Date: Saturday, June 22, 2019

Time: 11:00 AM-12:00PM, 4:00 PM-5:00 PM

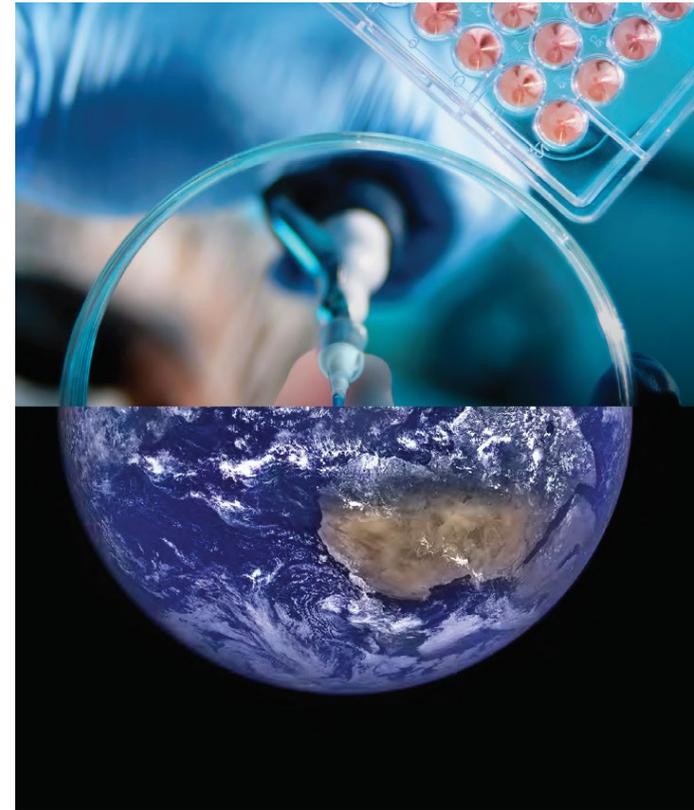
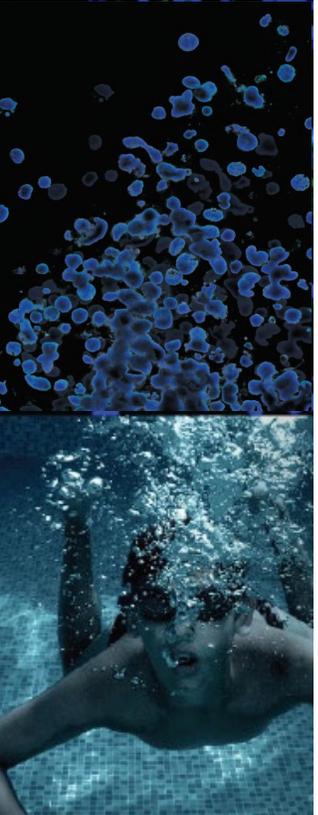
Acknowledgements

- Monique Hunter, MS
- Anna McCluskey, BS
- Stephen King, MS
- Juan Lopera, PhD
- Cara Wilder, PhD
- Dev Mittar, PhD
- Scott Tighe, PhD, UVM
- Denise O'Sullivan, PhD, LGC
- Nick Greenfield, MA, One Codex
- Pat Gillevet, PhD, Microbiome Analysis Center, GMU
- Rohan Patil, Microbiome Analysis Center, GMU
- Stefan Green, PhD, UIC (ABRF-MGRG)
- Joan Wong, PhD, PACBIO®
- Tony Lialin, Loop Genomics



Questions?

Credible Leads to Incredible™



Gene-Level Microbiome Analysis Identifies Culturable Strains Consistently Associated with Human Cancer across Independent Cohorts

Samuel Minot, PhD
Microbiome Research Initiative
Fred Hutch Cancer Research Center
Seattle, Washington, USA

Disclosures

ATCC – Consulting

One Codex – Financial Interest

Collaborators



Amy Willis, PhD
Professor of Biostatistics
University of Washington



Jonathan Golob, MD, PhD
Infectious Disease & Internal Medicine
University of Michigan

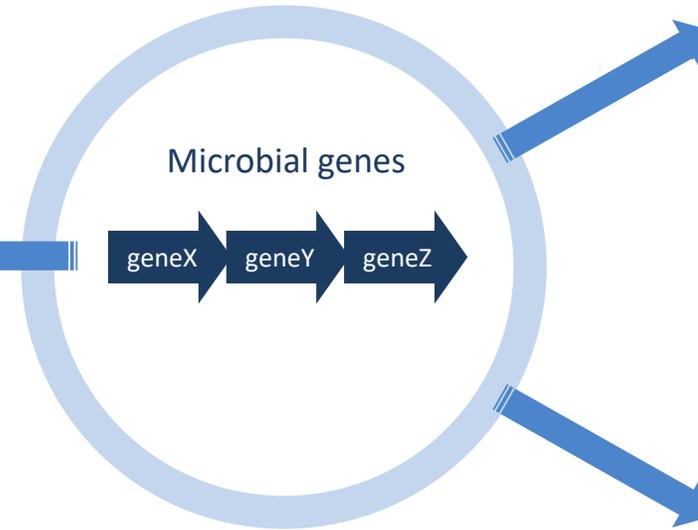
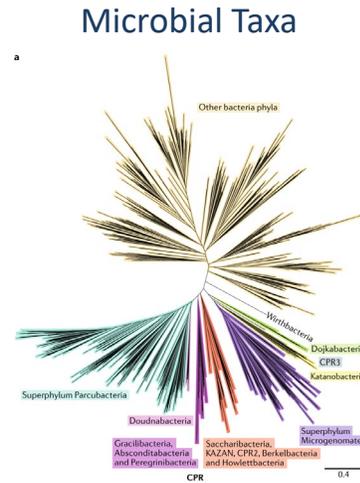
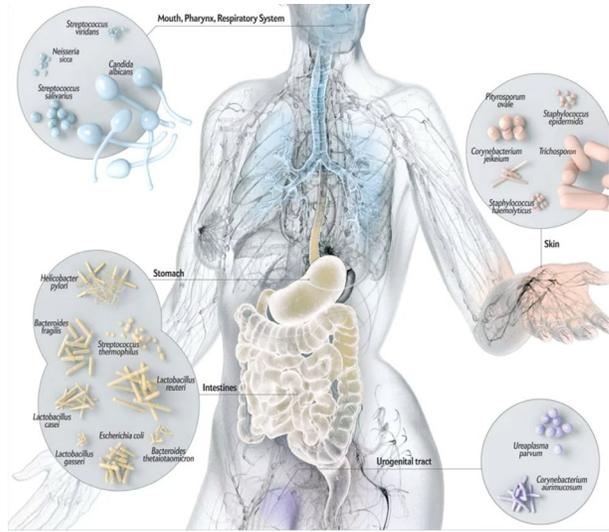
gene-level metagenomics

Applied Microbiology

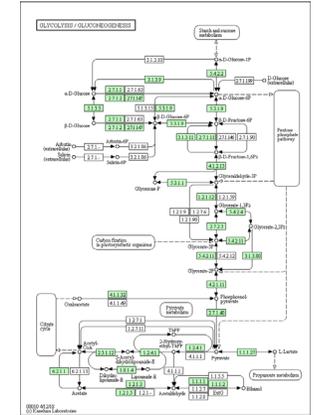


discovery – therapeutics – diagnostics

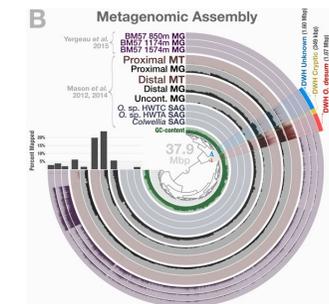
Ontologies of Microbiome Analysis



Metabolic Pathways



Assembled Genomes



Scientific American (2012-05-15)
 Castelle, *et al.* 2018
 Eren, *et al.* 2015
genome.jp/kegg

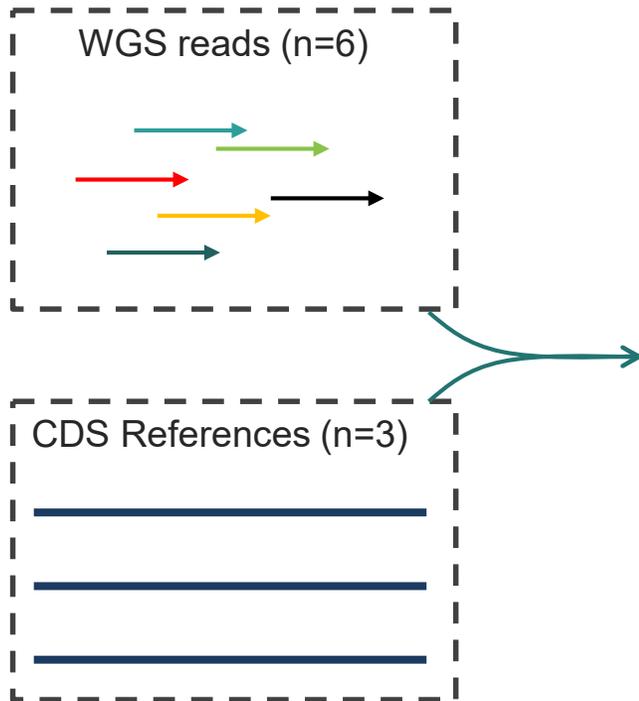
Gene-Level Analysis of the Microbiome

1. Computational methods development
 - Single-sample analysis
 - Cross-sample analysis
2. Application to CRC datasets
3. Validation in mouse model

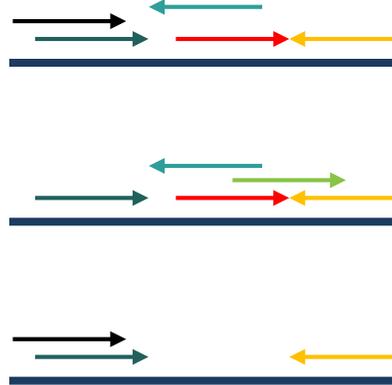
Gene-Level Analysis of the Microbiome

1. Computational methods development
 - **Single-sample analysis**
 - Cross-sample analysis
2. Application to CRC datasets
3. Validation in mouse model

Detecting Genes from WGS Data

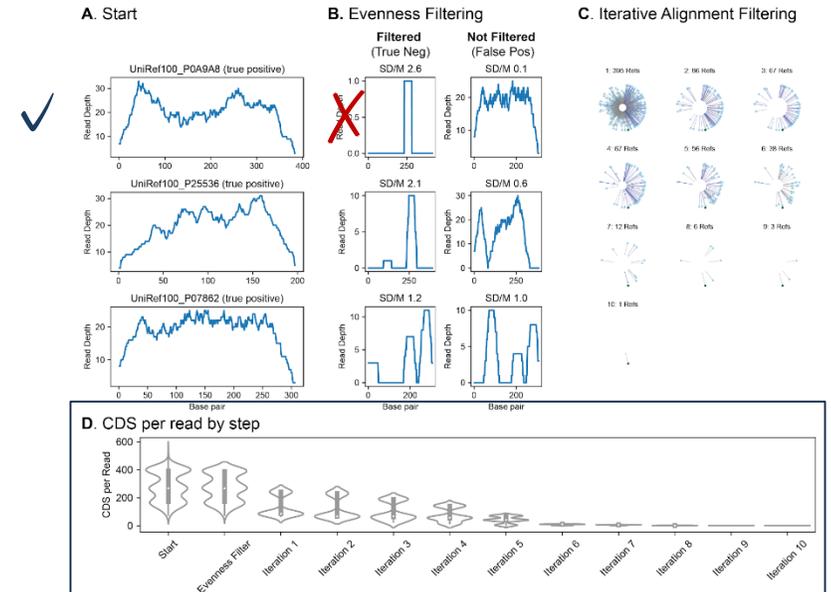


Alignments

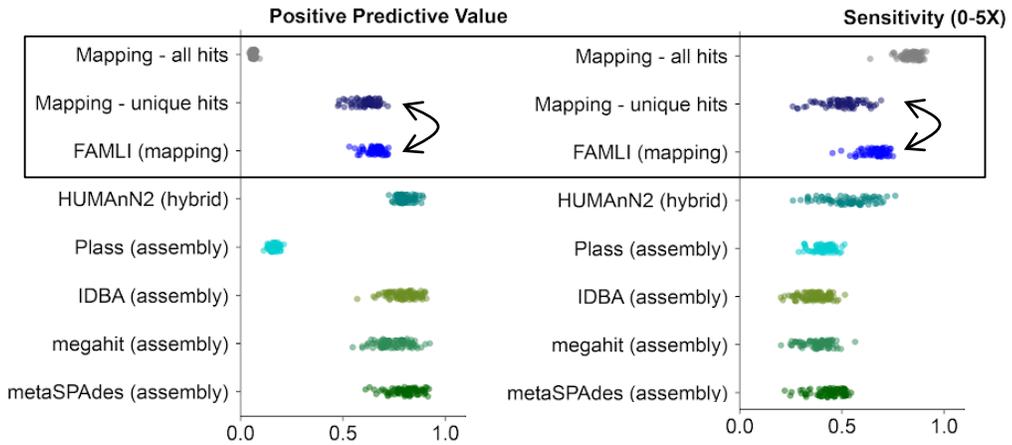
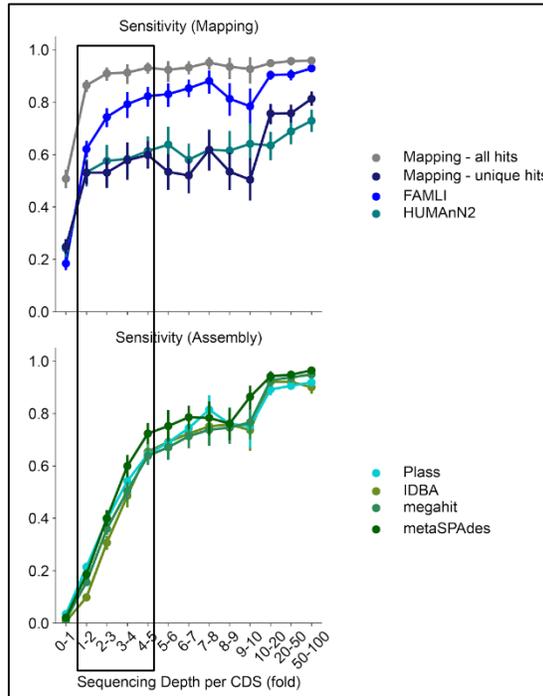


Take all hits?
Take unique hits?

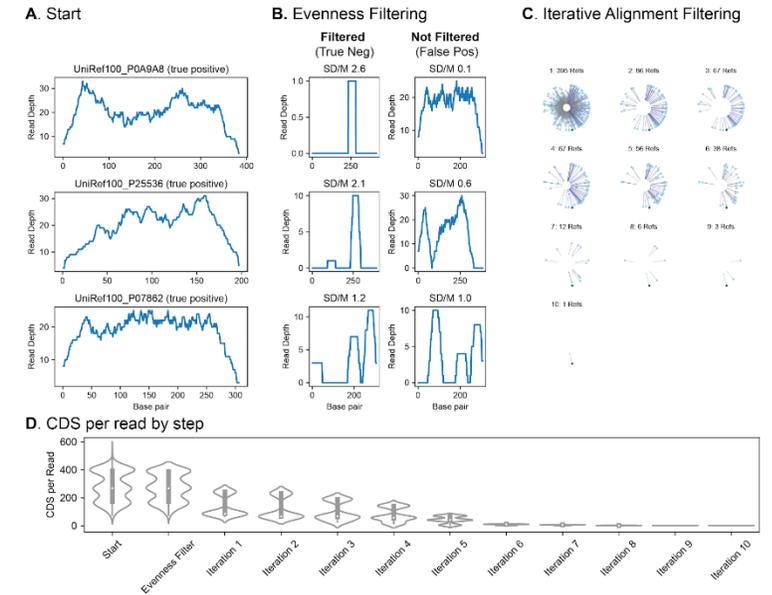
FAMLI Filtering Algorithm



Detecting Genes from WGS Data



FAMLI Filtering Algorithm



Performance Evaluation



Gene-Level Analysis of the Microbiome

1. Computational methods development
 - Single-sample analysis
 - **Cross-sample analysis**
2. Application to CRC datasets
3. Validation in mouse model

Efficiently Comparing Gene-Level Communities

Single stool metagenome:

10-50M reads \rightarrow 100k – 1M genes

Aggregate metagenome:

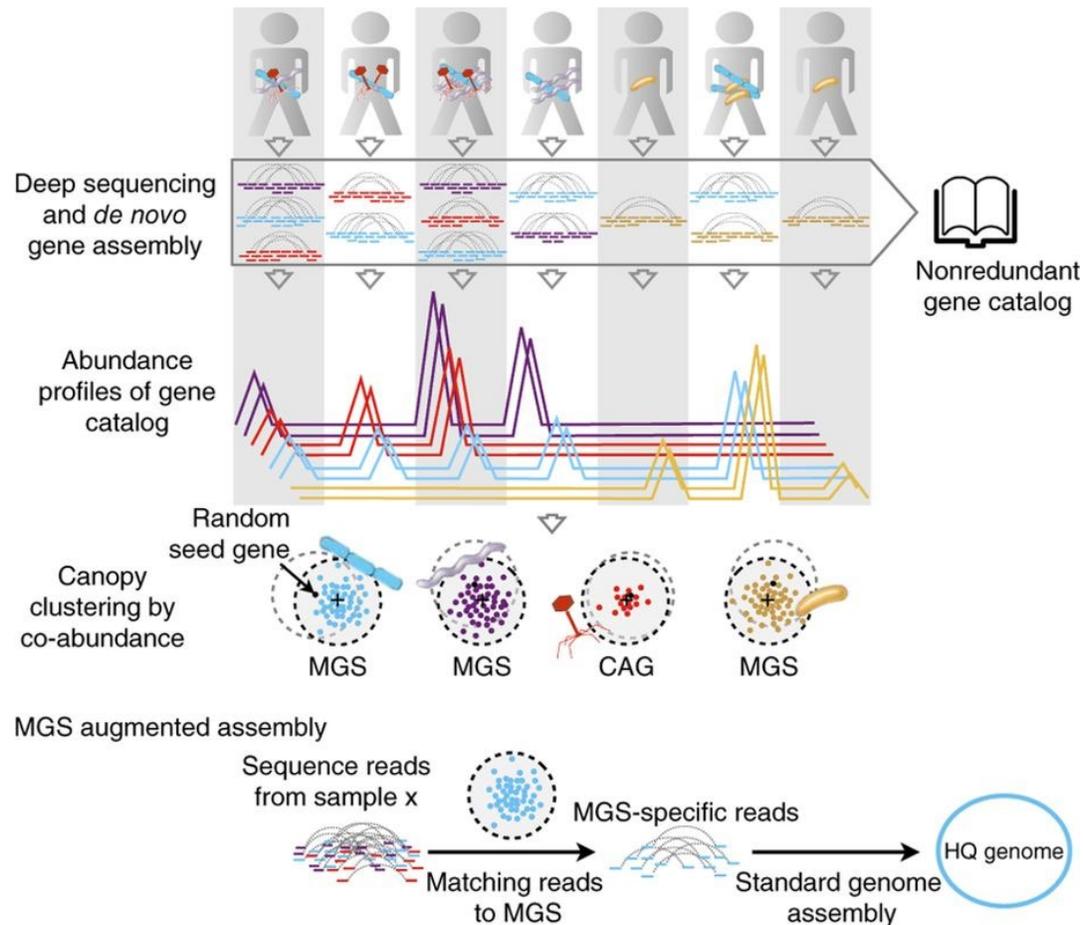
\sim 100 people \rightarrow 1M – 10M genes



Efficiently Comparing Gene-Level Communities

Approach: Co-Abundance Clustering

- Genes → Co-Abundant Gene Groups (CAGs)



Previous state-of-the-art:

Clustering heuristic to identify **species**

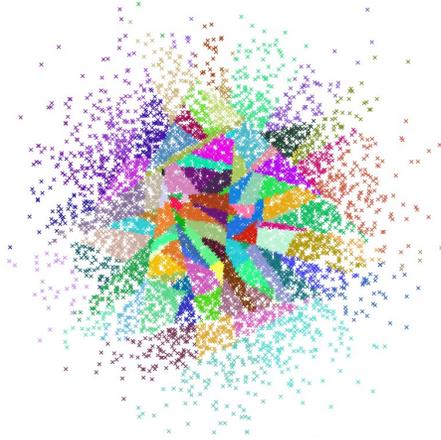
Goal:

Clustering of co-abundant **genes**
(operons, HGT, phages, auxiliary genome)

Enabling Technology:

Approximate Nearest Neighbor Algorithm

Efficiently Comparing Gene-Level Communities



Approximate Nearest Neighbor Algorithm efficiently partitions densely populated high-dimensional space

Previous state-of-the-art:

Clustering heuristic to identify **species**

Goal:

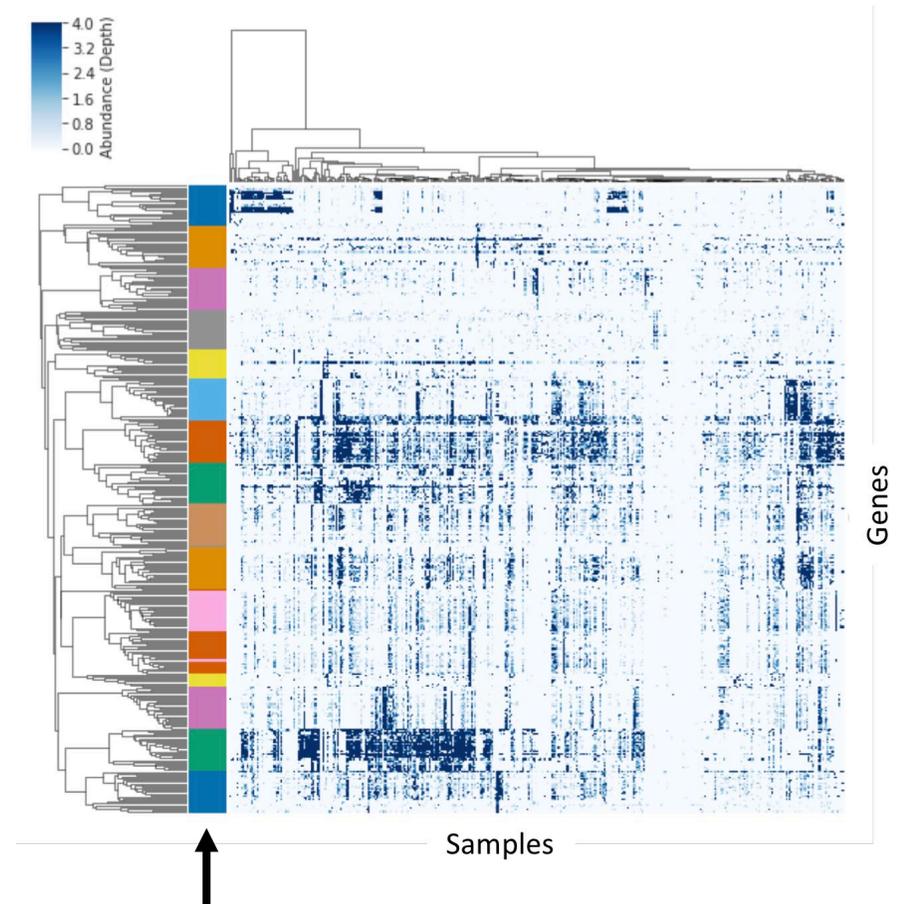
Clustering of co-abundant **genes**
(operons, HGT, phages, auxiliary genome)

Enabling Technology:

Approximate Nearest Neighbor Algorithm

Efficiently Comparing Gene-Level Communities

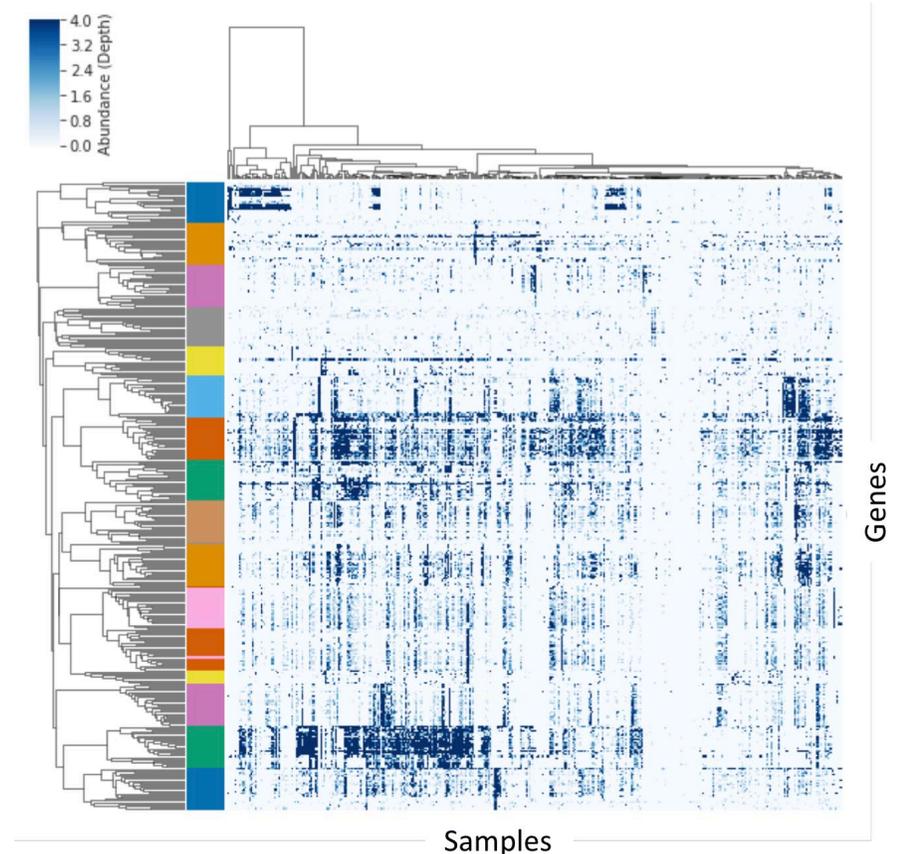
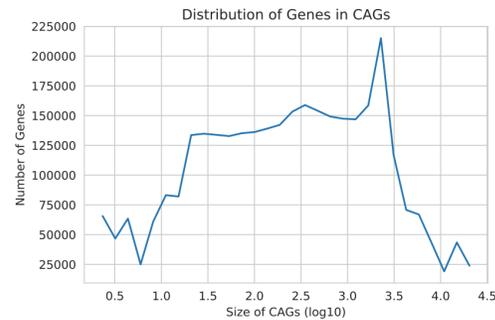
- Source code: <https://github.com/FredHutch/find-cags/>
- Python package: `pip install ann_linkage_clustering`
- Docker container: quay.io/fhcrc-microbiome/find-cags
- Manuscript: [biorxiv.org/content/10.1101/567818v1](https://doi.org/10.1101/567818v1)



Gene Grouping Captures Co-Abundance

Efficiently Comparing Gene-Level Communities

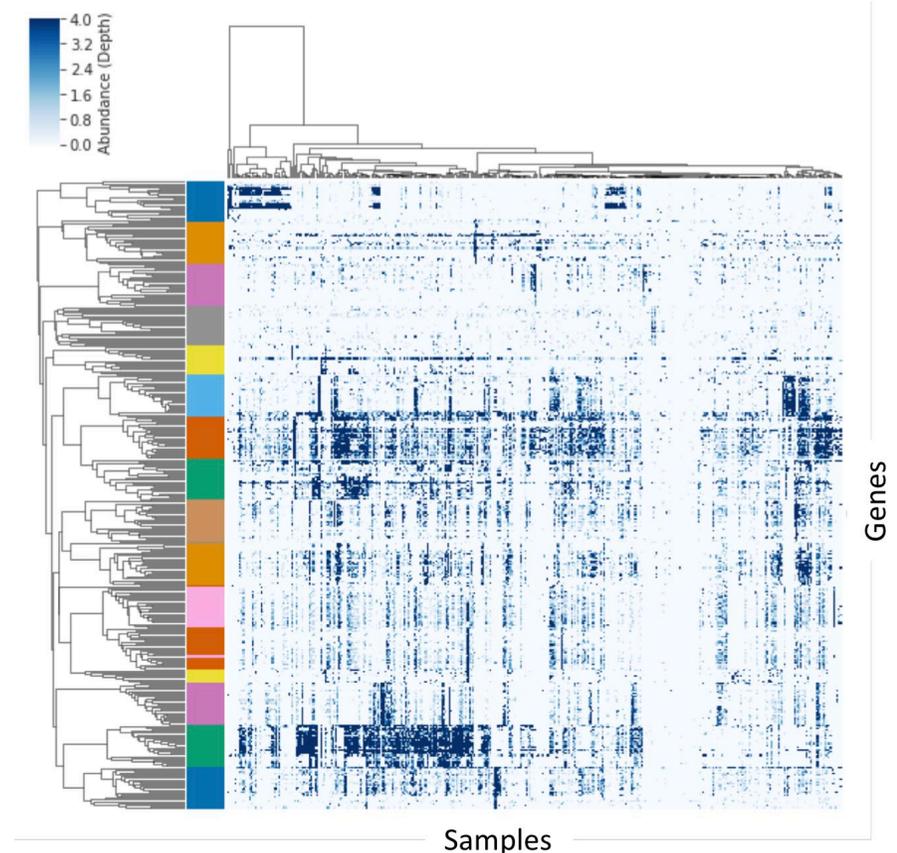
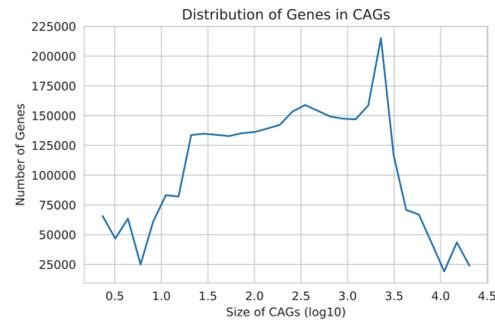
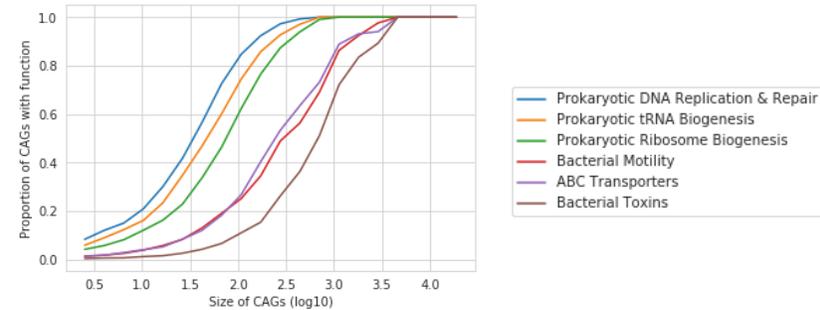
- Source code: <https://github.com/FredHutch/find-cags/>
- Python package: `pip install ann_linkage_clustering`
- Docker container: quay.io/fhcr-microbiome/find-cags
- Manuscript: [biorxiv.org/content/10.1101/567818v1](https://doi.org/10.1101/567818v1)



Co-Abundant Gene Groups (CAGs) range from 20-3,000 genes

Efficiently Comparing Gene-Level Communities

- Source code: <https://github.com/FredHutch/find-cags/>
- Python package: `pip install ann_linkage_clustering`
- Docker container: quay.io/fhcr-microbiome/find-cags
- Manuscript: [biorxiv.org/content/10.1101/567818v1](https://doi.org/10.1101/567818v1)

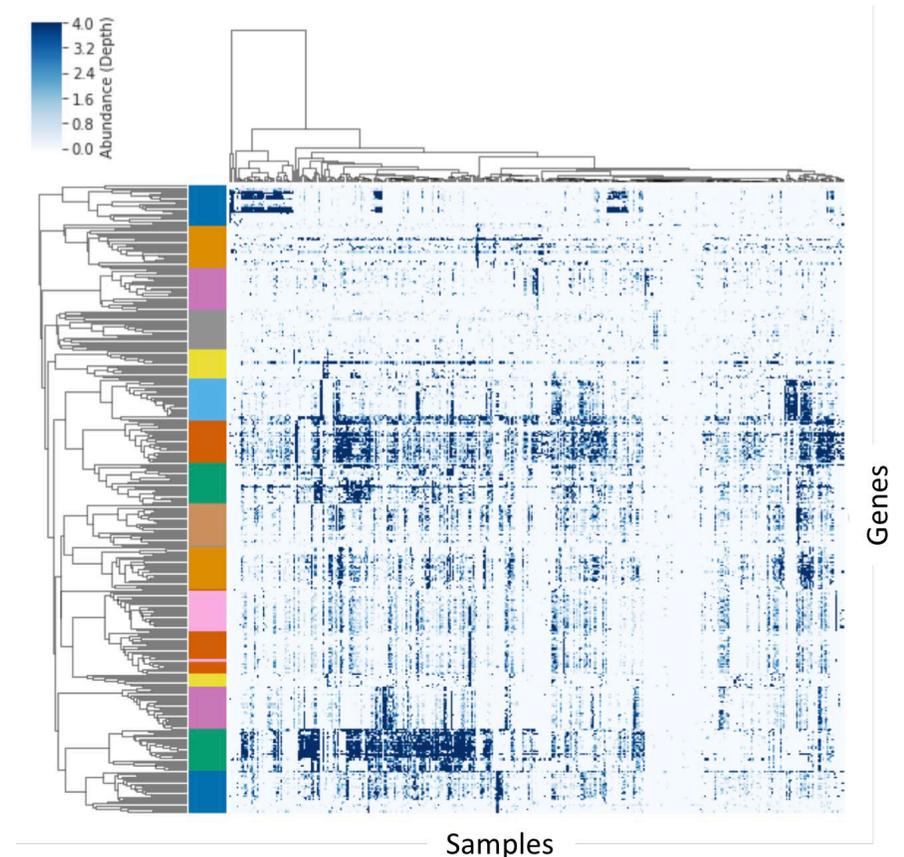
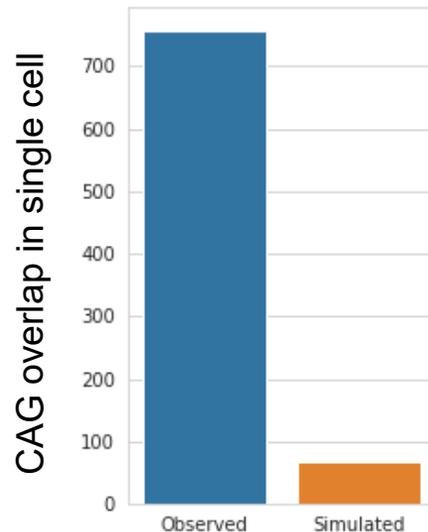


Co-Abundant Gene Groups (CAGs) contain functions required for bacterial life

Efficiently Comparing Gene-Level Communities

- Source code: <https://github.com/FredHutch/find-cags/>
- Python package: `pip install ann_linkage_clustering`
- Docker container: quay.io/fhcr-microbiome/find-cags
- Manuscript: [biorxiv.org/content/10.1101/567818v1](https://doi.org/10.1101/567818v1)

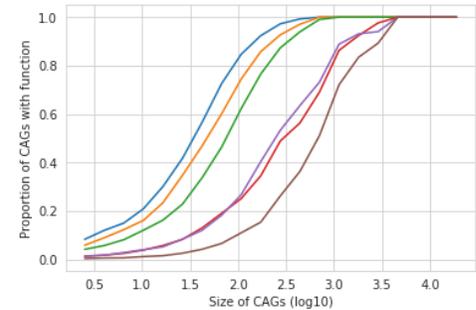
10X Genomics – Microbiome Data



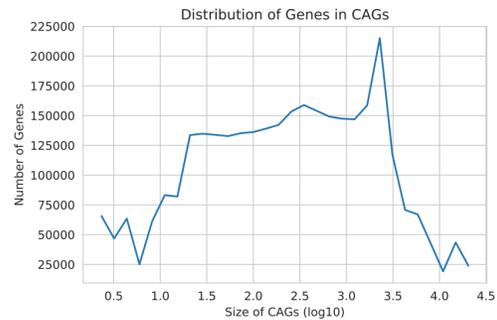
Co-Abundant Gene Groups (CAGs) co-occur in single-cell sequencing datasets

Efficiently Comparing Gene-Level Communities

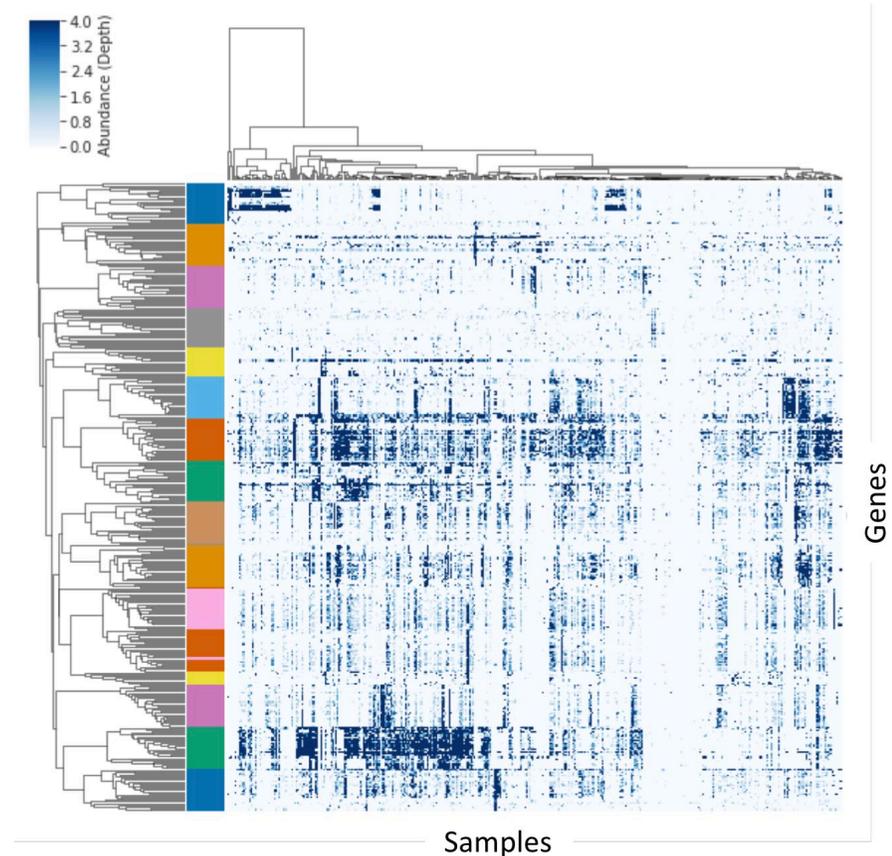
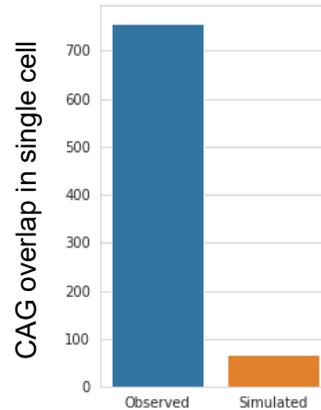
- Source code: <https://github.com/FredHutch/find-cags/>
- Python package: `pip install ann_linkage_clustering`
- Docker container: quay.io/fhcr-microbiome/find-cags
- Manuscript: [biorxiv.org/content/10.1101/567818v1](https://doi.org/10.1101/567818v1)



— Prokaryotic DNA Replication & Repair
— Prokaryotic tRNA Biogenesis
— Prokaryotic Ribosome Biogenesis
— Bacterial Motility
— ABC Transporters
— Bacterial Toxins



10X Genomics – Microbiome Data



Performance: 5M genes ~ 5,000 samples (240GB RAM) → ~10 hours

Efficiently Comparing Gene-Level Communities

- Source code: <https://github.com/FredHutch/find-cags/>
- Python package: `pip install ann_linkage_clustering`
- Docker container: quay.io/fhcr-microbiome/find-cags
- Manuscript: [biorxiv.org/content/10.1101/567818v1](https://doi.org/10.1101/567818v1)

Grouping *millions* of genes → *10,000's* of CAGs

Enables efficient cross-sample comparison

Preserves biological complexity

Gene-Level Analysis of the Microbiome

1. Computational methods development
 - Single-sample analysis
 - Cross-sample analysis
2. **Application to CRC datasets**
3. Validation in mouse model

Identifying Strains Associated with CRC

nature medicine ARTICLES
<https://doi.org/10.1038/nm4191-019-0405-7>

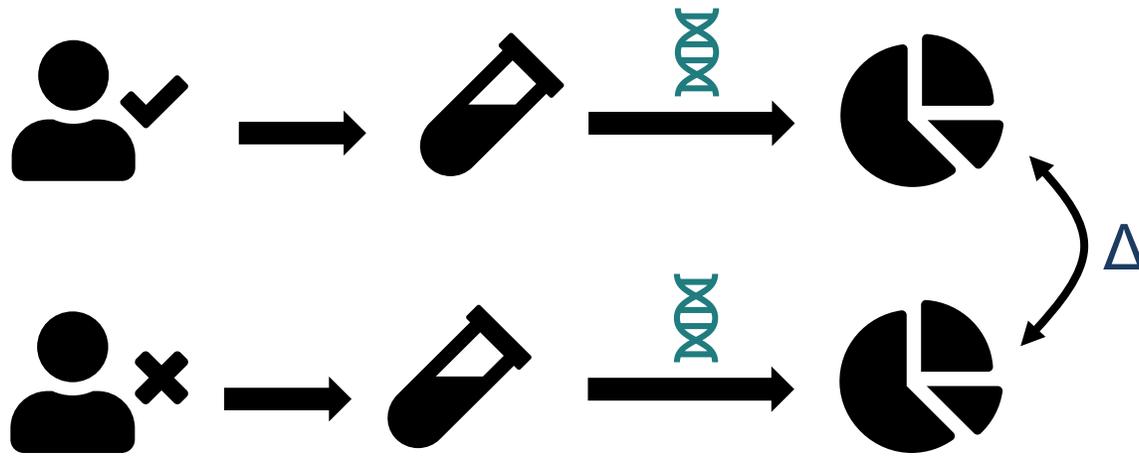
Metagenomic analysis of colorectal cancer datasets identifies cross-cohort microbial diagnostic signatures and a link with choline degradation

Andrew Maltz Thomas^{1,2,3,32}, Paolo Manghi^{1,32}, Francesco Asnicar¹, Edoardo Pasoli¹, Federica Armanini¹, Moreno Zolfo¹, Francesco Beghini¹, Serena Manara¹, Nicolai Karcher¹, Chiara Pozzi⁴, Sara Gandini⁴, Davide Serrano⁴, Sonia Tarallo^{4,5}, Antonio Francavilla^{4,5}, Gaetano Gallo^{4,7}, Mario Trompetto⁴, Giulio Ferrero^{4,8}, Sayaka Mizutani^{9,10}, Hirotugu Shiroma⁷, Satoshi Shiba¹, Tatsuhiro Shibata^{11,12}, Shinichi Yachida^{13,14}, Takuji Yamada^{9,14}, Jakob Wirbel¹⁵, Petra Schrotz-King¹⁶, Cornelia M. Ulrich¹⁷, Hermann Brenner^{18,19,20}, Manimozhiyan Arumugam^{20,21}, Peer Bork^{15,22,23,24}, Georg Zeller¹⁵, Francesca Cordero¹, Emmanuel Dias-Neto^{1,25}, João Carlos Setubal^{1,26}, Adrian Tett¹, Barbara Pardini^{1,27}, Maria Rescigno²⁸, Levi Waldron^{29,30,33}, Alessio Naccarati^{1,31,33} and Nicola Segata^{1,33*}

nature medicine ARTICLES
<https://doi.org/10.1038/nm4191-019-0406-6>

Meta-analysis of fecal metagenomes reveals global microbial signatures that are specific for colorectal cancer

Jakob Wirbel^{1,31}, Paul Theodor Pyl^{2,3,7}, Ece Kartal¹⁴, Konrad Zych¹⁷, Alireza Kashani², Alessio Milanese¹, Jonas S. Fleck¹, Anita Y. Voigt¹⁵, Albert Pallega¹⁶, Ruby Ponnudurai¹, Shinichi Sunagawa¹⁴, Luis Pedro Coelho^{10,19}, Petra Schrotz-King¹⁷, Emily Vogtmann¹, Nina Habermann¹, Emma Nimeus¹⁰, Andrew M. Thomas^{1,12}, Paolo Manghi¹, Sara Gandini¹³, Davide Serrano¹, Sayaka Mizutani^{14,15}, Hirotugu Shiroma¹⁴, Satoshi Shiba¹, Tatsuhiro Shibata^{16,17}, Shinichi Yachida^{18,19}, Takuji Yamada^{14,19}, Levi Waldron^{20,21}, Alessio Naccarati^{22,23}, Nicola Segata¹, Rashmi Sinha¹, Cornelia M. Ulrich¹⁸, Hermann Brenner^{2,25,26}, Manimozhiyan Arumugam^{2,27,32}, Peer Bork^{1,4,28,29,32*} and Georg Zeller^{1,32*}



Discovery cohort:

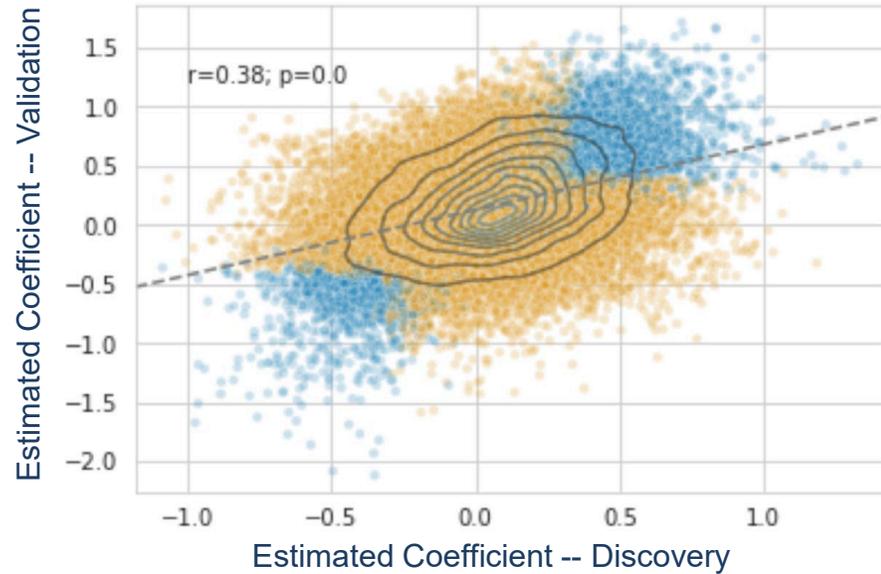
Zeller (2014) Molecular Systems Biology
 156 participants – France

Validation cohort:

Yu (2015) Gut
 128 participants – China



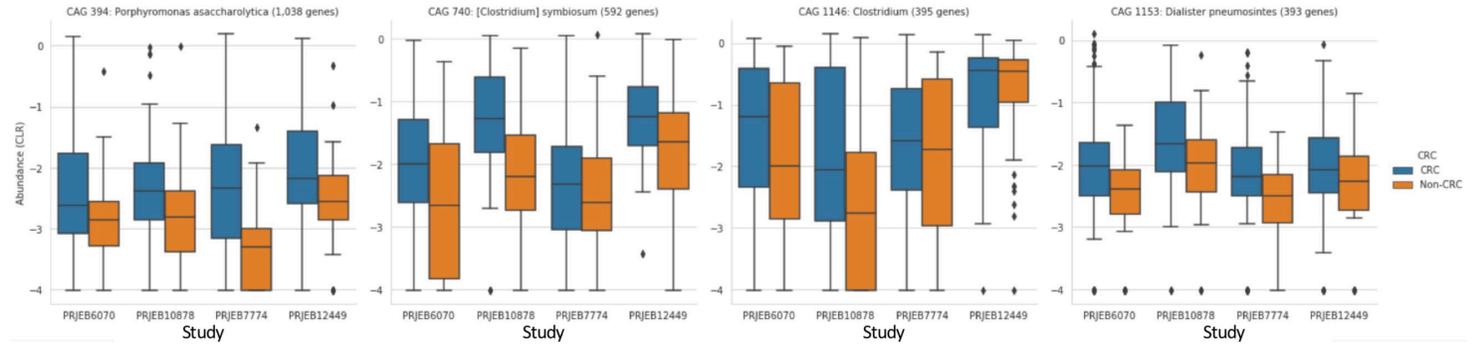
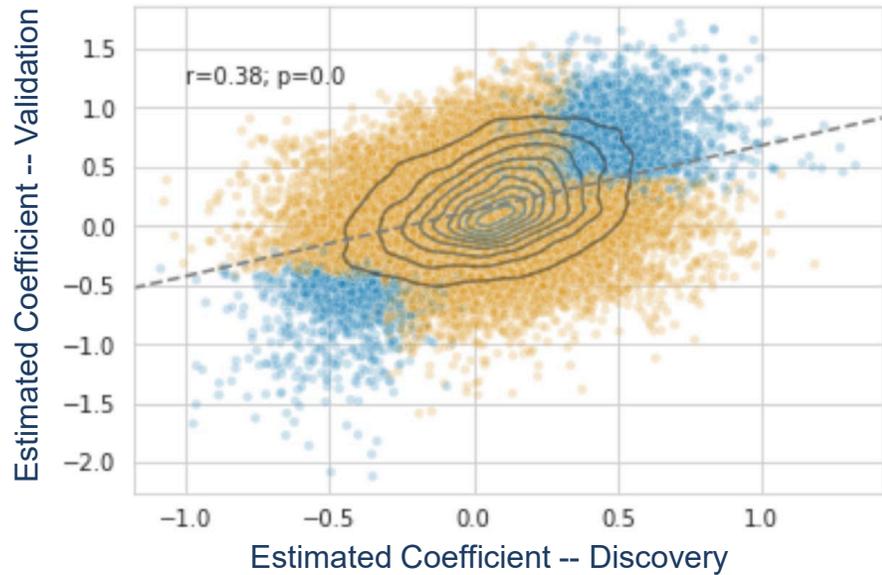
Identifying Strains Associated with CRC



Association of CAG abundance with CRC is **highly reproducible**



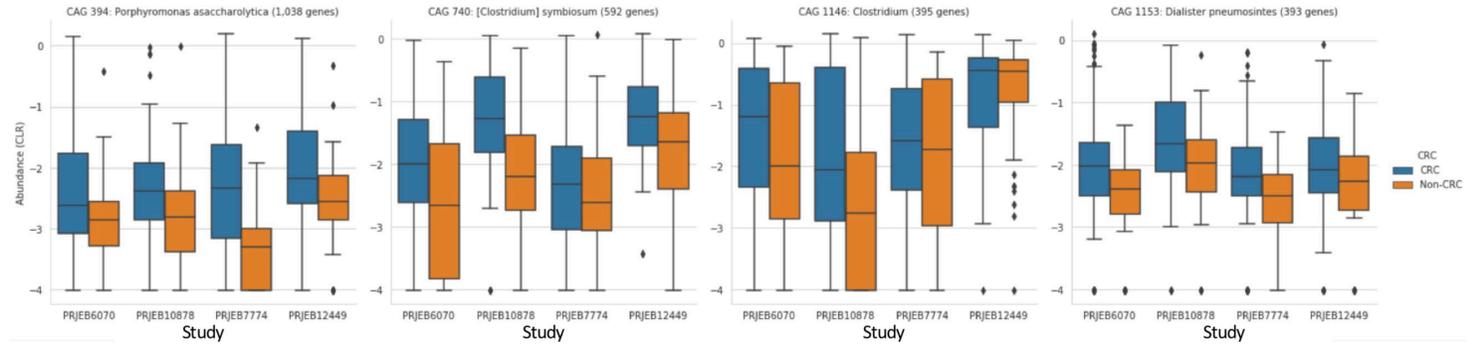
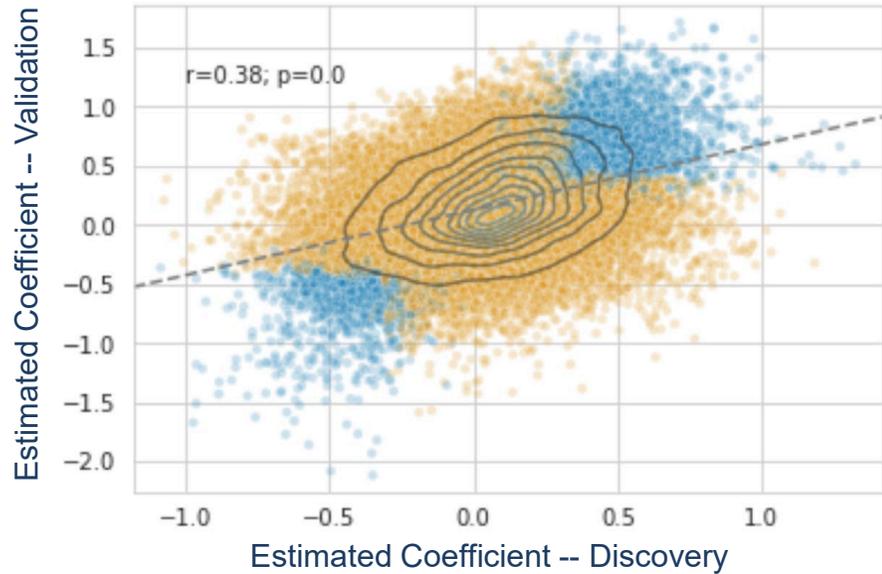
Identifying Strains Associated with CRC



Consistent association with CRC across four independent cohorts



Identifying Strains Associated with CRC



Disease → CAGs

CAGs → Genes

Genes → Function (KEGG)

Genes → Taxonomy (NCBI)

Genes → Genomes (ATCC)

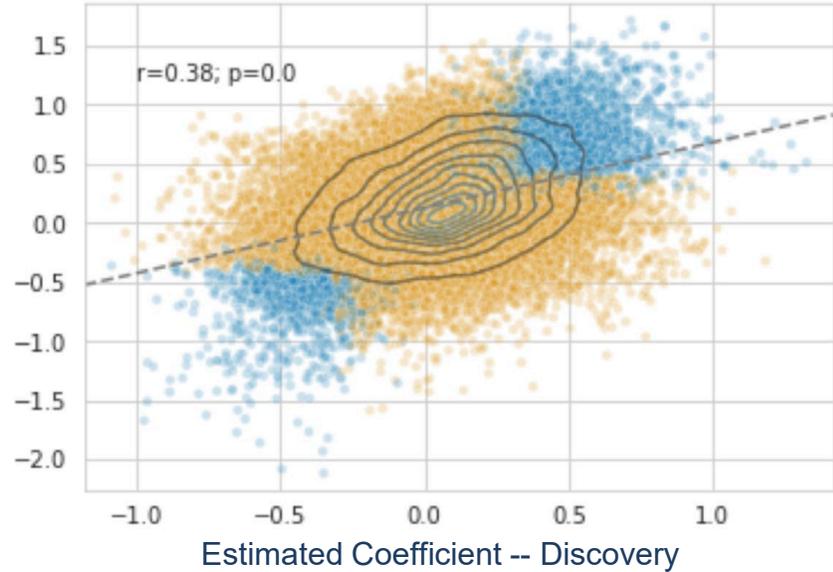


Identifying Strains Associated with CRC

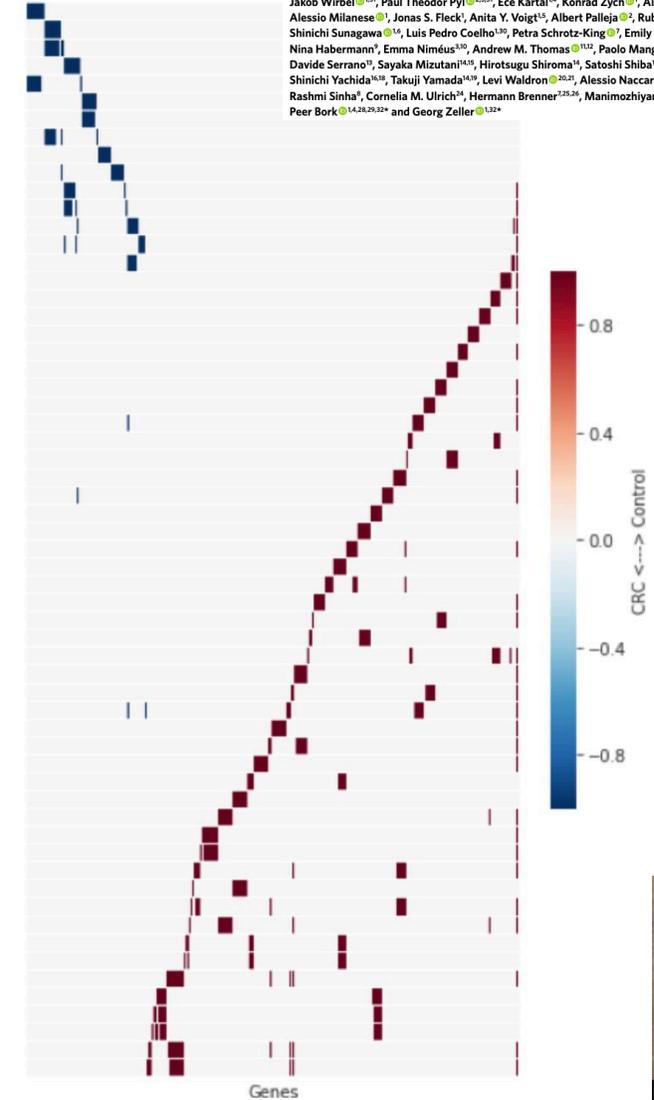
Meta-analysis of fecal metagenomes reveals global microbial signatures that are specific for colorectal cancer

Jakob Wirbel^{1,2,3}, Paul Theodor Pyl^{1,2,3,7}, Ece Kartal^{1,4}, Konrad Zych¹, Alireza Kashani¹, Alessio Milanese¹, Jonas S. Fleck¹, Anita Y. Voigt^{1,5}, Albert Palleja^{1,2}, Ruby Ponnudurai¹, Shinichi Sunagawa^{1,4}, Luis Pedro Coelho^{1,5}, Petra Schrotz-King⁷, Emily Vogtmann¹, Nina Habermann¹, Emma Niménez^{1,2}, Andrew M. Thomas^{1,12}, Paolo Manghi¹¹, Sara Gandini^{1,13}, Davide Serrano¹, Sayaka Mizutani^{1,15}, Hirotosugu Shirota¹, Satoshi Shiba¹, Tatsuhiro Shibata^{1,16,17}, Shinichi Yachida^{1,18}, Takuji Yamada^{1,19}, Levi Waldron^{1,20,21}, Alessio Naccarati^{2,22}, Nicola Segata^{1,10,17}, Rashmi Sinha¹, Cornelia M. Ulrich¹⁴, Hermann Brenner^{1,23,24}, Manimozhyan Arumugam^{2,22,25}, Peer Bork^{14,26,29,32*} and Georg Zeller^{1,32*}

Estimated Coefficient -- Validation



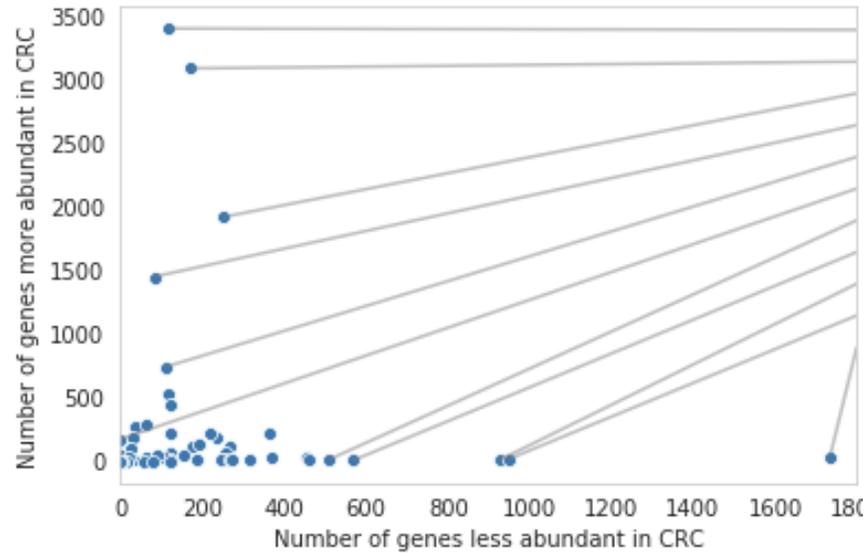
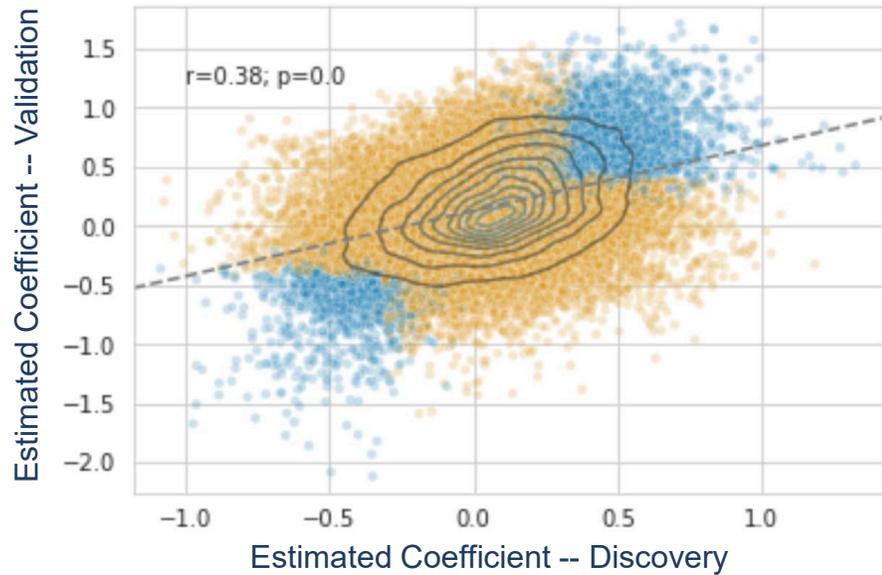
	<i>Bacteroides fragilis</i>	✓ 145	3,568
	<i>Flavonifractor plautii</i>	✓ 37	3,117
→	<i>Lachnospiraceae bacterium 7_1_58FAA</i>	64	3,135
	[<i>Clostridium</i>] <i>clostridioforme</i>	✓ 179	3,184
→	<i>Bacteroides fragilis</i> str. 3397 T10	238	3,082
→	<i>Subdoligranulum</i> sp. OF01-18	21	2,774
→	<i>Ruthenibacterium lactatiformans</i>	33	2,764
→	<i>Clostridiales bacterium VE202-03</i>	17	2,619
→	<i>Phascolarctobacterium faecium</i> DSM 14760	19	2,583
→	<i>Oscillibacter</i> sp. PEA192	23	2,548
	[<i>Clostridium</i>] <i>bolteae</i>	✓ 166	2,598
→	[<i>Clostridium</i>] <i>clostridioforme</i> 90A6	✓ 193	2,349
→	[<i>Ruminococcus</i>] <i>gnavus</i>	✓ 255	1,924
→	[<i>Clostridium</i>] <i>aldenense</i>	✓ 184	1,792
→	[<i>Ruminococcus</i>] <i>gnavus</i>	✓ 416	1,891
→	<i>Firmicutes bacterium</i> AM10-47	2,154	29
→	<i>Ruminococcus</i> sp. AM26-12LB	2,161	32
→	<i>Clostridium</i> sp. AF15-31	2,167	3
→	<i>Eubacterium</i> sp. CAG:251	2,177	3
→	<i>Blautia</i> sp. AF19-10LB	2,211	24
→	<i>Ruminococcus</i> sp. A254.MGS-254	2,204	6
→	<i>Blautia</i> sp. CAG:52	2,238	9
→	<i>Blautia</i> sp. CAG:37	2,242	6
→	<i>Ruminococcus</i> sp. AF32-2AC	2,303	55
→	<i>Ruminococcus</i> sp. CAG:17	2,291	26
→	<i>Ruminococcus</i> sp. CAG:254	2,325	2
→	<i>Eubacterium</i> sp. AM49-13BH	2,336	7
→	<i>Eubacterium ventriosum</i> ATCC 27560	✓ 2,366	4
→	<i>Lachnospiraceae bacterium</i> SNUG30370	2,376	1
→	<i>Firmicutes bacterium</i> AM55-24TS	2,409	2
→	<i>Eubacterium</i> sp. CAG:252	2,410	0
→	<i>Anaerostipes hadrus</i>	✓ 2,423	2
→	<i>Eubacterium</i> sp. CAG:248	2,444	2
→	<i>Clostridium</i> sp. AM32-2	2,570	115
→	<i>Blautia</i> sp. AF19-34	2,554	30
→	<i>Firmicutes bacterium</i> CAG:41	2,541	2
→	<i>Ruminococcus</i> sp. OM08-9BH	2,625	49
→	<i>Clostridium</i> sp. CAG:75	2,719	3
→	<i>Blautia</i> sp. AF14-40	2,816	18
→	[<i>Ruminococcus</i>] <i>torques</i>	✓ 2,891	49
→	<i>Faecalibacterium</i> sp. CAG:74	2,873	2
→	<i>Clostridium</i> sp. OM05-9BH	2,905	9
→	<i>Eubacterium</i> sp. CAG:38	✓ 2,908	4
→	<i>Anaerostipes hadrus</i>	✓ 2,927	7
→	[<i>Eubacterium</i>] <i>hallii</i>	✓ 2,997	11
→	<i>Coprococcus</i> sp. CAG:131	3,014	6
→	<i>Ruminococcus bicirculans</i>	✓ 3,031	1
→	<i>Ruminococcus</i> sp. AM28-13	3,068	1
→	[<i>Eubacterium</i>] <i>eligans</i>	✓ 3,105	8
→	[<i>Eubacterium</i>] <i>hallii</i>	✓ 3,105	7
→	[<i>Eubacterium</i>] <i>eligans</i>	✓ 3,148	17
→	<i>Coprococcus eutactus</i>	✓ 3,191	7
→	<i>Anaerostipes hadrus</i>	✓ 3,243	8
→	<i>Anaerostipes hadrus</i>	✓ 3,441	4
→	<i>Roseburia intestinalis</i> XB6B4	3,650	10
→	<i>Coprobacillus</i> sp. AF31-1BH	3,660	2
→	<i>bacterium</i> LF-3	✓ 3,738	3
→	<i>Coprobacillus</i> sp. AM28-15LB	3,748	10
→	<i>Roseburia intestinalis</i>	✓ 4,017	13
→	<i>Roseburia intestinalis</i>	✓ 4,087	14



Genes identify microbial strains associated with CRC

Identifying Strains Associated with CRC

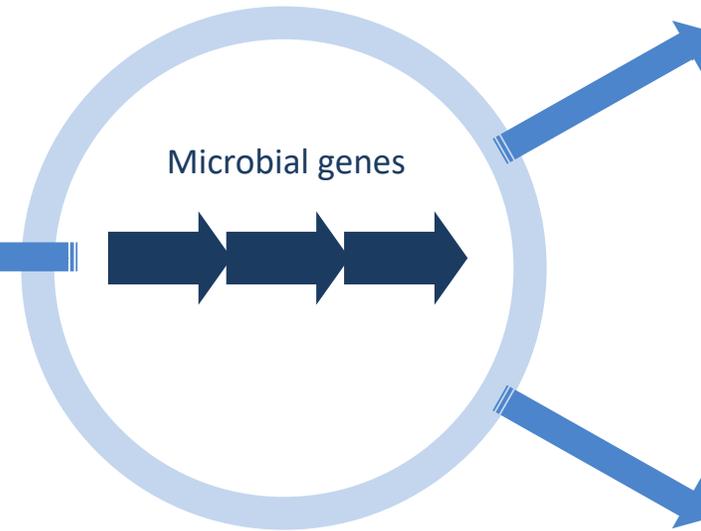
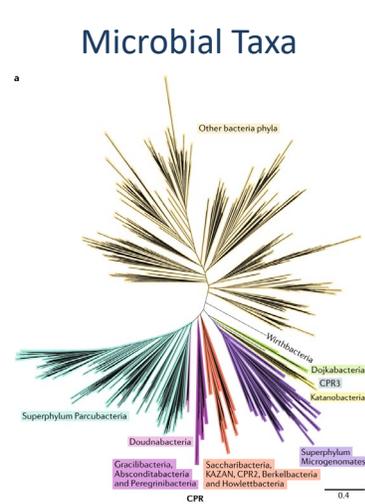
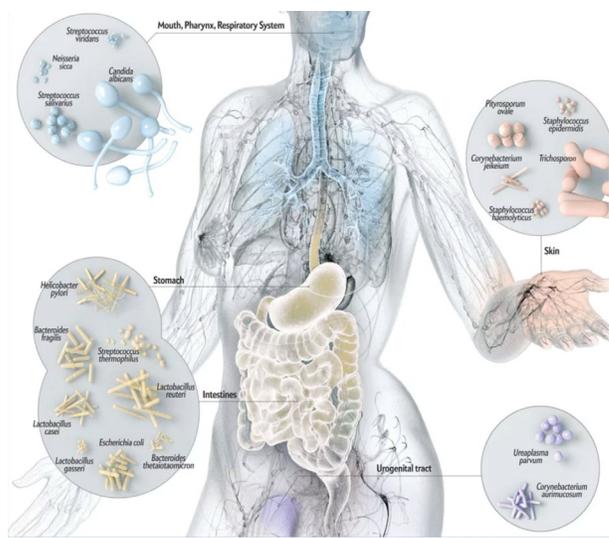
Experimental Validation – Neel Dey MD (Fred Hutch)



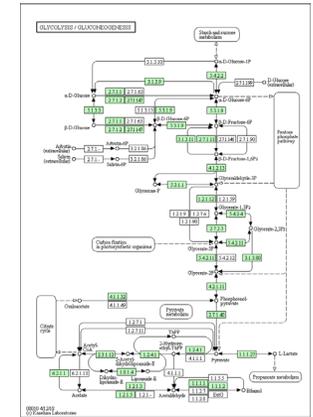
- Bacteroides fragilis* strain A
- Clostridium bolteae* strain B
- Ruminococcus gnavus* strain C
- Clostridium symbiosum* strain D
- Clostridium asparagiforme* strain E
- Fusobacterium nucleatum* strain F
- Dorea longicatena* strain G
- Bifidobacterium pseudocatenulatum* strain H
- Eubacterium rectale* strain I
- Ruminococcus obeum* strain J
- Coprococcus comes* strain K



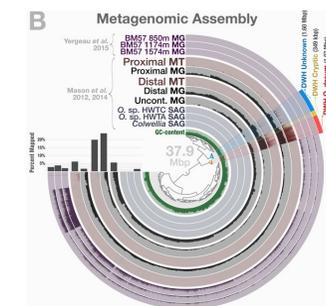
Gene-Level Metagenomics for Microbiome Research



Metabolic Pathways



Assembled Genomes



Scientific American (2012-05-15)
 Castelle, et al. 2018
 Eren, et al. 2015
genome.jp/kegg

THANK YOU



FRED HUTCH
CURES START HERE®

fredhutch.org



CREDIBLE
STANDARDS



INCREDIBLE
RESULTS

Advancing Authentication



Whole genome sequencing has generated data at an unprecedented scale in the biological sciences. However, existing public genomic databases can lack **quality**, **completeness**, **authenticity**, and **traceability**.

Advancing Authentication

Our Enhanced Authentication Initiative aims to enrich the characterization of our biological collections and provide you with the whole genome sequences of the specific, authenticated materials you need to generate credible data

We are giving you the first look at our
ATCC Genome Portal

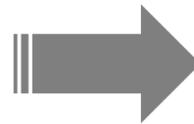
Your resource for reference-quality genomes from
authenticated ATCC products



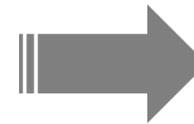
1. Extract DNA from
authenticated
materials



2. Sequence the high-
quality DNA



3. Assemble with short
and long reads



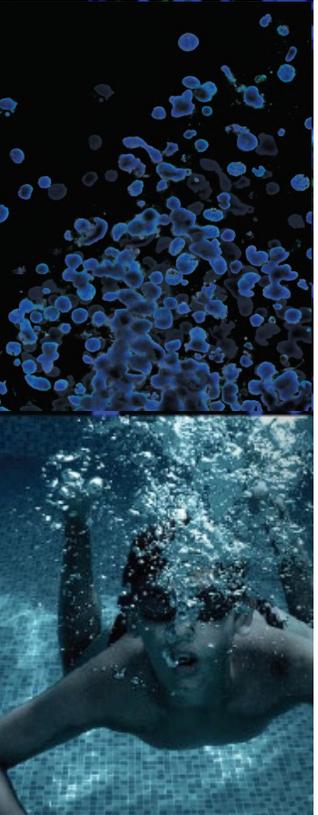
4. Validate strain
designation and
publication



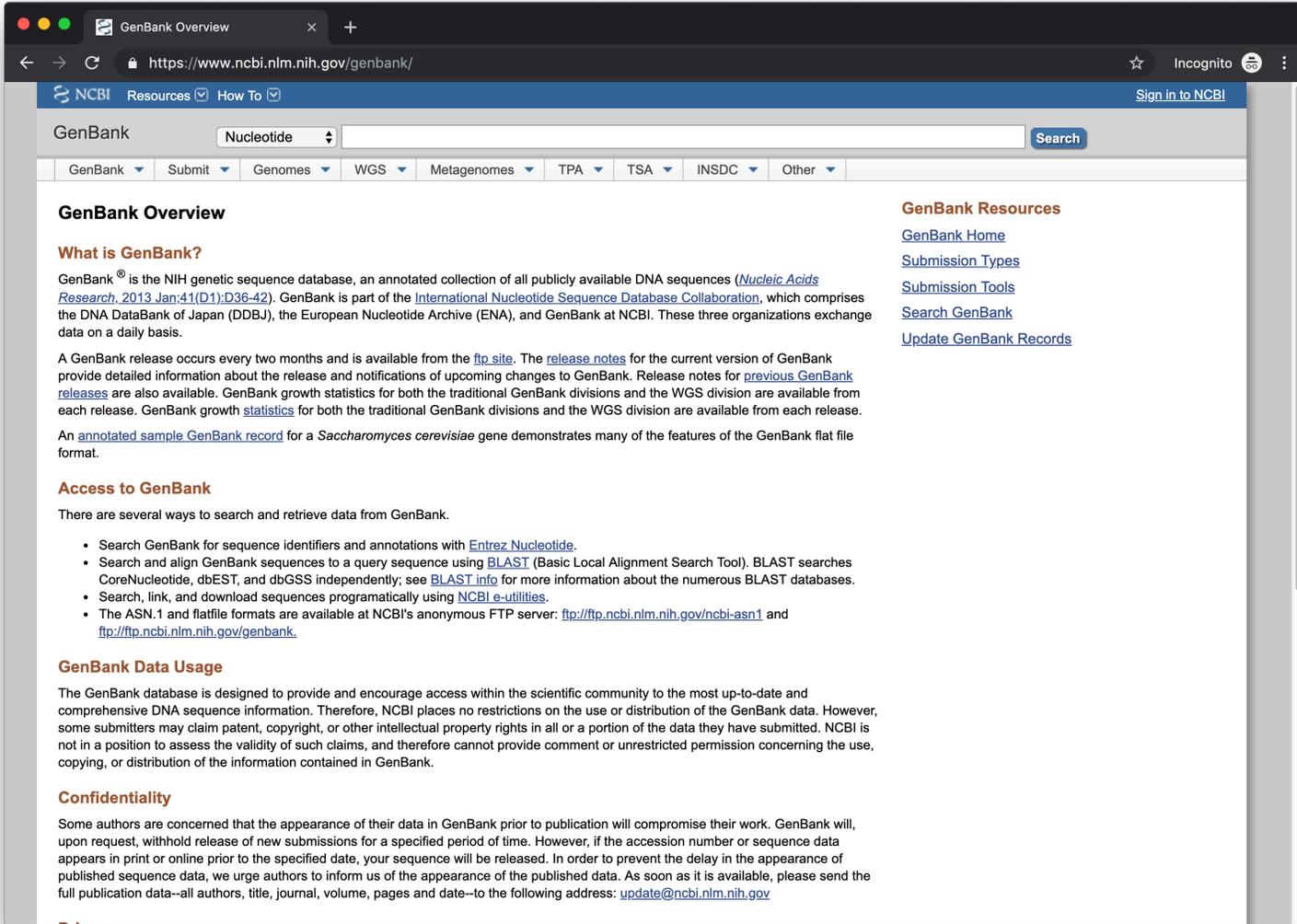
Enhanced Authentication Initiative

Nick Greenfield, MA
Founder and CEO, One Codex

Credible Leads to Incredible™



Current microbial genomics references



The screenshot shows the GenBank Overview page on the NCBI website. The browser address bar displays <https://www.ncbi.nlm.nih.gov/genbank/>. The page features a search bar with a dropdown menu set to "Nucleotide" and a "Search" button. Below the search bar, there are navigation tabs for "GenBank", "Submit", "Genomes", "WGS", "Metagenomes", "TPA", "TSA", "INSDC", and "Other". The main content area is titled "GenBank Overview" and includes sections for "What is GenBank?", "Access to GenBank", "GenBank Data Usage", and "Confidentiality". A sidebar on the right lists "GenBank Resources" with links to "GenBank Home", "Submission Types", "Submission Tools", "Search GenBank", and "Update GenBank Records".

GenBank[®]

- De facto standard
- Lots of genomes
- But relatively little curation
- And highly variable quality

Current microbial genomics references

FDA ARGOS

NCTC

Database for Reference Grade Microbial Sequences (FDA-ARGOS)
Accession: PRJNA231221 ID: 231221

In May 2014, the FDA and collaborators established a publicly available dAtabase for Reference Grade microBial Sequences called FDA-ARGOS. More...

Resource Name	Number of Links
SEQUENCE DATA	
Nucleotide (total)	1997
WGS master	227
Genomic DNA	735
Genomic RNA	22
SRA Experiments	1794
Protein Sequences	2119190

Public Health England reference collections

This project aims to provide annotated and assembled genomes for 3,000 bacteria and 500 viruses as part of a new eResource

The project is split into two parts:

- NCTC 3000:** A joint collaboration between Public Health England, Pacific Biosciences and the Wellcome Sanger Institute to complete the sequencing of 3,000 bacterial strains from PHE's National Collection of Type Cultures (NCTC) using Pacific Biosciences' Single Molecule, Real-Time (SMRT) sequencing technology.
- NCPV 500:** A collaboration between PHE and Sanger to produce 500 viral genomes from PHE's National Collection of Pathogenic Viruses (NCPV) using the Illumina sequencing platform.

Collectively, the data generated will be housed in a publically accessible web-based eResource that integrates metadata and genome sequences for type and reference strains of biomedically important bacterial and viral pathogens. This resource will integrate accession, taxonomy and authentication information with publications, genome sequences, comparative analysis databases and other resources at EMBL and NCBI.

This is a community resource project. Data will be available from here, and from the NCTC. We will submit assembled, annotated sequences to the International Sequence Databases as they become available. We request that you cite this webpage in any publication using the data, and would appreciate it if you contact us to discuss the use of this data.

Data Downloads

- Download annotated assemblies
- BLAST server

Please note: these are pre-submission assemblies that should not be treated as final versions. Assemblies contain both chromosomal and plasmid contigs.

Background

... and numerous other specialty collections

Current database challenges

#1 Quality

Mukherjee et al. *Standards in Genomic Sciences* 2015, **10**:18
<http://www.standardsingenomics.com/content/10/1/18>



COMMENTARY **Open Access**

Large-scale contamination of microbial isolate genomes by Illumina PhiX control

Supratim Mukherjee^{1*}, Marcel Huntemann¹, Natalia Ivanova¹, Nikos C Kyrpides^{1,2} and Amrita Pati¹

Abstract

With the rapid growth and development of sequencing technologies, genomes have become the new go-to for exploring solutions to some of the world's biggest challenges such as searching for alternative energy sources and exploration of genomic dark matter. However, progress in sequencing has been accompanied by its share of errors that can occur during template or library preparation, sequencing, imaging or data analysis. In this study we screened over 18,000 publicly available microbial isolate genome sequences in the Integrated Microbial Genomes database and identified more than 1000 genomes that are contaminated with PhiX, a control frequently used during Illumina sequencing runs. Approximately 10% of these genomes have been published in literature and 129 contaminated genomes were sequenced under the Human Microbiome Project. Raw sequence reads are prone to contamination from various sources and are usually eliminated during downstream quality control steps. Detection of PhiX contaminated genomes indicates a lapse in either the application or effectiveness of proper quality control measures. The presence of PhiX contamination in several publicly available isolate genomes can result in additional errors when such data are used in comparative genomics analyses. Such contamination of public databases have far-reaching consequences in the form of erroneous data interpretation and analyses, and necessitates better measures to proofread raw sequences before releasing them to the broader scientific community.

Keywords: Next-generation sequencing, PhiX, Contamination, Comparative genomics

Background

The ability to produce large numbers of high-quality, low-cost reads has revolutionized the field of microbiology [1-3]. Starting from a meager 1575 registered projects in September 2005, there has been a steady increase in the number of sequencing projects according to the Genomes OnLine Database [4]. As of November 17th 2014, there were 41,553 bacterial and archaeal isolate genome sequencing projects reported in GOLD [4,5]. This explosion of genome sequencing projects especially during the last 5 years has been largely catalyzed by the development of several next-generation sequencing platforms offering rapid and accurate genome information at a low cost. Among the different NGS technologies available commercially, the sequencing by synthesis technology [6] championed by Illumina [7] is the most widely used.

Despite its high accuracy, the Illumina sequencing platform does come with its share of challenges [8] that need to be addressed by the users of this technology. One such challenge is the protocol in which PhiX is used as a quality and calibration control for sequencing runs. PhiX is an icosahedral, nontailed bacteriophage with a single-stranded DNA. It has a tiny genome with 5386 nucleotides and was the first DNA genome to be sequenced by Fred Sanger [9]. Due to its small, well-defined genome sequence, PhiX has been commonly used as a control for Illumina sequencing runs. For the majority of its library preparations Illumina recommends using PhiX at a low concentration of 1%, which can be raised up to 40% for low diversity samples. Depending on the concentration of PhiX used, it can be spiked in the same lane along with the sample or used as a separate lane. Addition of PhiX as a sequencing control necessitates subsequent quality control steps to remove the sequences such that they do not get integrated as part of the target genome.

* Correspondence: supratim@mukherjee@lbl.gov
¹DOE Joint Genome Institute, Walnut Creek, CA, USA
Full list of author information is available at the end of the article

 © 2015 Mukherjee et al.; licensee BioMed Central. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly credited. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated.

“...[we] identified more than 1000 genomes that are contaminated with PhiX...”

Current database challenges

#1 Quality

Downloaded from genome.cshlp.org on June 17, 2019 - Published by Cold Spring Harbor Laboratory Press

Research

Human contamination in bacterial genomes has created thousands of spurious proteins

Florian P. Breitwieser,¹ Mihaela Perlea,^{1,2} Aleksey V. Zimin,^{1,3} and Steven L. Salzberg^{1,2,3,4}

¹Center for Computational Biology, McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins School of Medicine, Baltimore, Maryland 21205, USA; ²Department of Computer Science, Whiting School of Engineering, Johns Hopkins University, Baltimore, Maryland 21218, USA; ³Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21218, USA; ⁴Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, Maryland 21205, USA

Contaminant sequences that appear in published genomes can cause numerous problems for downstream analyses, particularly for evolutionary studies and metagenomics projects. Our large-scale scan of complete and draft bacterial and archaeal genomes in the NCBI RefSeq database reveals that 2250 genomes are contaminated by human sequence. The contaminant sequences derive primarily from high-copy human repeat regions, which themselves are not adequately represented in the current human reference genome, GRCh38. The absence of the sequences from the human assembly offers a likely explanation for their presence in bacterial assemblies. In some cases, the contaminating contigs have been erroneously annotated as containing protein-coding sequences, which over time have propagated to create spurious protein “families” across multiple prokaryotic and eukaryotic genomes. As a result, 2437 spurious protein entries are currently present in the widely used nr and TrEMBL protein databases. We report here an extensive list of contaminant sequences in bacterial genome assemblies and the proteins associated with them. We found that nearly all contaminants occurred in small contigs in draft genomes, which suggests that filtering out small contigs from draft genome assemblies may mitigate the issue of contamination while still keeping nearly all of the genuine genomic sequences.

[Supplemental material is available for this article.]

Over the past two decades, the number of publicly available genomes has grown from just a handful of species to well over 100,000 genomes today. These genomes are pivotal resources for countless biomedical research questions, including microbiome studies that use them to identify species in complex samples (Breitwieser et al. 2017). Ideally, all genomes in reference databases would be complete and accurate (Fraser et al. 2002), but for practical reasons, the vast majority of genomes available today are still “drafts.” A draft genome consists of multiple contigs or scaffolds that are typically unordered and not assigned into chromosomes (Ghurye et al. 2016). A genome is not truly complete or “finished” until every base pair has been determined for every chromosome and organelle, end-to-end, with no gaps. Even the human genome, although far more complete than most other animal genomes, is still unfinished: The current human assembly, GRCh38.p13 (released Feb. 28, 2019), has 473 scaffolds that contain 875 internal gaps. While most of the human sequence has been placed on chromosomes, some highly repetitive regions are underrepresented (Altmeppen et al. 2014), leading to problems that we discuss below. Draft genomes of other species vary widely in quality as well as contiguity, with some having thousands of contigs and others having a much smaller number.

Contamination of genome assemblies with sequences from other species is not uncommon, especially in draft genomes (Longo et al. 2011; Merchant et al. 2014; Delmont and Eren 2016; Kryukov and Imanishi 2016; Lu and Salzberg 2018). In 2011, researchers reported that over 10% of selected nonprimate assemblies in the NCBI and UCSC Genome Browser databases were contaminated with the primate-specific *AluY* repeats (Longo et al. 2011). Although validation pipelines have improved substantially since then (Tatusova et al. 2016; Haft et al. 2018), some contaminants still remain, as we describe below. Furthermore, when open reading frames (ORFs) in the contaminated contigs get annotated as protein-coding genes, their protein sequence may be added to other databases. Once in those databases, these spurious proteins may in turn be used in future annotation, leading to the so-called “transitive catastrophe” problem where errors are propagated widely (Karp 1998; Salzberg 2007; Danchin et al. 2018). Indeed, one study found that the percentage of misannotated entries in the NCBI nonredundant (nr) protein collection, which is used for thousands of BLAST searches every day, has been increasing over time (Schnoes et al. 2009).

Contamination of genomic sequences can be particularly problematic for metagenomic studies. For example, if a genome labeled as species X contains fragments of the human genome, then any sample containing human DNA might erroneously be identified as also containing species X. Since human DNA is virtually always present in the environment of sequencing laboratories, human contamination is very common in sequencing experiments of all types. Contamination of laboratory reagents with DNA from other organisms can also lead to serious misinterpretations, such as the supposed detection of the novel virus NIH-CQV

Corresponding authors: florian.bw@gmail.com, salzberg@jhu.edu
Articles published online before print. Article, supplemental material, and publication date are at <http://www.genome.org/cgi/doi/10.1101/gr.243373.118>.

© 2019 Breitwieser et al. This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <http://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

954 Genome Research 29:954–960 Published by Cold Spring Harbor Laboratory Press; ISSN 1088-9051/19; www.genome.org
www.genome.org

“...2250 [microbial] genomes are contaminated by human sequence...”

Current database challenges

#1 Quality

Mukherjee et al. Standards in Genomic Sciences 2015, 10:18
<http://www.standardsin-genomics.com/content/10/1/18>



**Zamin Iqbal**
@ZaminIqbal

Following

Urgh. "2250 genomes are contaminated by human sequence. The contaminant sequences derive primarily from high-copy human repeat regions, which themselves are not adequately represented in the current human reference genome, GRCh38." which have made their way into databases

Correspondence: zamin@stanford.edu
© 2015 Mukherjee et al.; licensee BioMed Central. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly credited. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated.



" 2250 [microbial] are contaminated by human sequence..."

Current database challenges

#2 Completeness

Land et al. *Standards in Genomic Sciences* 2014, **9**:20
<http://www.standardsingenomics.com/content/9/1/20>

 Standards in Genomic Sciences

RESEARCH **Open Access**

Quality scores for 32,000 genomes

Miriam L Land^{1*}, Doug Hyatt^{1,2}, Se-Ran Jun¹, Guruprasad H Kora³, Loren J Hauser^{1,2,4}, Oksana Lukjancenko⁶ and David W Ussey^{1,2,5}

Abstract

Background: More than 80% of the microbial genomes in GenBank are of 'draft' quality (12,553 draft vs. 2,679 finished, as of October, 2013). We have examined all the microbial DNA sequences available for complete, draft, and Sequence Read Archive genomes in GenBank as well as three other major public databases, and assigned quality scores for more than 30,000 prokaryotic genome sequences.

Results: Scores were assigned using four categories: the completeness of the assembly, the presence of full-length rRNA genes, tRNA composition and the presence of a set of 102 conserved genes in prokaryotes. Most (~88%) of the genomes had quality scores of 0.8 or better and can be safely used for standard comparative genomics analysis. We compared genomes across factors that may influence the score. We found that although sequencing depth coverage of over 100x did not ensure a better score, sequencing read length was a better indicator of sequencing quality. With few exceptions, most of the 30,000 genomes have nearly all the 102 essential genes.

Conclusions: The score can be used to set thresholds for screening data when analyzing "all published genomes" and reference data is either not available or not applicable. The scores highlighted organisms for which commonly used tools do not perform well. This information can be used to improve tools and to serve a broad group of users as more diverse organisms are sequenced. Unexpectedly, the comparison of predicted tRNAs across 15,000 high quality genomes showed that anticodons beginning with an 'A' (codons ending with a 'U') are almost non-existent, with the exception of one arginine codon (CGU); this has been noted previously in the literature for a few genomes, but not with the depth found here.

Keywords: DNA, Sequencing, Database, Quality, Evaluation, Status

Background

The introduction of second-generation sequencing began an exponential growth in sequencing data [1-4] and in the number of genomes submitted to public repositories. The drop in sequencing cost that came with this technology, however, had little effect the mostly manual cost of finishing genomes. Finishing second-generation sequenced genomes continues to be expensive and many researchers have no plans to finish their draft genomes [5]. There is still an open question of whether whole genome sequencing projects with less than 5% of the genes missing is adequate quality for most purposes [6] or if there continues to be value in finishing most microbial genomes [7]. Even though single molecule, or 'third-generation' sequencing will facilitate the generation of closed genomes, currently most of the genomes in the database are of varying levels of draft quality.

The establishment of a quality nomenclature by Chain *et al.* in 2009 [8] provides a mechanism for comparing draft sequences and understanding the qualifiers associated with a single genome sequence. It does not, however, shine any light on the impact that predominately draft genomes have on the quality of the repository databases. With more than 30,000 unique publicly available genome sequences of varying qualities, there is enough data to score genomes on the basis of completeness and compare quality among data sources.

DNA sequences were obtained from two sources at GenBank and the National Center for Biotechnology Information [9]: draft genomes ('WGS' or 'draft') and complete finished genomes ('complete'). An assembled version of the GenBank Sequence Read Archive was obtained for analysis [10]. Despite major overlaps, three additional

* Correspondence: landml@ornl.gov
¹Comparative Genomics Group, Biosciences Division, Oak Ridge National Laboratory, P.O. Box 2008, MS 6420, Oak Ridge, TN 37831-6420, USA
Full list of author information is available at the end of the article

 BioMed Central

© 2014 Land et al.; licensee BioMed Central Ltd. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly credited. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated.

“More than 80% of the microbial genomes in GenBank are of ‘draft’ quality...”

Current database challenges

#3 Authenticity

 PLOS ONE

RESEARCH ARTICLE

Strategies to Avoid Wrongly Labelled Genomes Using as Example the Detected Wrong Taxonomic Affiliation for *Aeromonas* Genomes in the GenBank Database

Roxana Beaz-Hidalgo¹, Mohammad J. Hossain², Mark R. Liles², Maria-Jose Figueras^{1*}

¹ Unitat de Microbiologia, Departament de Ciències Mèdiques Bàsiques, Facultat de Medicina i Ciències de la Salut, IISPV, Universitat Rovira i Virgili, Reus, Spain, ² Department of Biological Sciences, Auburn University, Auburn, Alabama, United States of America

* mariajose.figueras@urv.cat

 CrossMark

OPEN ACCESS

Citation: Beaz-Hidalgo R, Hossain MJ, Liles MR, Figueras M-J (2015) Strategies to Avoid Wrongly Labelled Genomes Using as Example the Detected Wrong Taxonomic Affiliation for *Aeromonas* Genomes in the GenBank Database. PLOS ONE 10(1): e0115813. doi:10.1371/journal.pone.0115813

Academic Editor: Turgay Unver, Cankiri Karatekin University, TURKEY

Received: August 14, 2014

Accepted: December 1, 2014

Published: January 21, 2015

Copyright: © 2015 Beaz-Hidalgo et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All genome files are available from the NCBI database, at the web site <http://www.ncbi.nlm.nih.gov/genomes/Bacteria/Aeromonas>

Funding: This work was supported in part by the project with reference AGL2011-30461-C02-02 by the "Ministerio de Ciencia e Innovación" (Spain) and by funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement no. 311346. The authors are solely responsible for the content of this publication. It does not represent the opinion of the European

Abstract

Around 27,000 prokaryote genomes are presently deposited in the Genome database of GenBank at the National Center for Biotechnology Information (NCBI) and this number is exponentially growing. However, it is not known how many of these genomes correspond correctly to their designated taxon. The taxonomic affiliation of 44 *Aeromonas* genomes (only five of these are type strains) deposited at the NCBI was determined by a multilocus phylogenetic analysis (MLPA) and by pairwise average nucleotide identity (ANI). Discordant results in relation to taxa assignment were found for 14 (35.9%) of the 39 non-type strain genomes on the basis of both the MLPA and ANI results. Data presented in this study also demonstrated that if the genome of the type strain is not available, a genome of the same species correctly identified can be used as a reference for ANI calculations. Of the three ANI calculating tools compared (ANI calculator, EzGenome and JSpecies), EzGenome and JSpecies provided very similar results. However, the ANI calculator provided higher intra- and inter-species values than the other two tools (differences within the ranges 0.06–0.82% and 0.92–3.38%, respectively). Nevertheless each of these tools produced the same species classification for the studied *Aeromonas* genomes. To avoid possible misinterpretations with the ANI calculator, particularly when values are at the borderline of the 95% cutoff, one of the other calculation tools (EzGenome or JSpecies) should be used in combination. It is recommended that once a genome sequence is obtained the correct taxonomic affiliation is verified using ANI or a MLPA before it is submitted to the NCBI and that researchers should amend the existing taxonomic errors present in databases.

Introduction

Members of the genus *Aeromonas* are found in aquatic environments worldwide and have been implicated in human and fish diseases [3, 12]. The genus now accounts for 27 species and

PLOS ONE | DOI:10.1371/journal.pone.0115813 January 21, 2015 1/13

“Discordant results in relation to taxa assignment were found for 14 (35.9%) of the 39 non-type strain genomes...”

Current database challenges

#4 Traceability

 **Microbiology**
Resource Announcements

LETTER TO THE EDITOR



Beware of False "Type Strain" Genome Sequences

Francisco Salvà-Serra,^{a,b,c,d,e} Daniel Jaén-Luchoro,^{a,b,c,d} Roger Karlsson,^{a,b,d,f} Antoni Bennasar-Figueras,^{a,g} Hedvig E. Jakobsson,^{h,i} Edward R. B. Moore^{a,h,c,d}

^aDepartment of Infectious Diseases, Institute of Biomedicine, Sahlgrenska Academy, University of Gothenburg, Gothenburg, Sweden
^bClinical Microbiology, Sahlgrenska University Hospital, Gothenburg, Sweden
^cCulture Collection University of Gothenburg, Sahlgrenska Academy, University of Gothenburg, Gothenburg, Sweden
^dCentre for Antibiotic Resistance Research (CARe), University of Gothenburg, Gothenburg, Sweden
^eMicrobiology, Department of Biology, University of the Balearic Islands, Palma de Mallorca, Spain
^fNanox Consulting AB, Gothenburg, Sweden
^gArea of Infectious Diseases, Research Institute of Health Sciences (IUNICS-UIB), University of the Balearic Islands, Palma de Mallorca, Spain

With this letter, we warn users of bacterial DNA sequence data about recent cases misusing the term "type strain" in bacterial genome sequence reports and highlight the importance that the term is used in the correct context.

In recent articles published in the journal *Genome Announcements* (GA), (now *Microbiology Resource Announcements* [MRA]) (1–5), the complete genome sequences of five strains of five species of bacteria were reported. The titles of the articles stated that the strains represent the respective "type strains" of the five species, although the articles do not present details about the processes of defining the type strains. At this point, it is important to point out that the concept and description of "type strain" are not arbitrary; the type strains of bacterial species are defined by rule 18a of the International Code of Nomenclature of Prokaryotes as follows: "The type strain is made up of living cultures of an organism, which are descended from a strain designated as the nomenclatural type" (6). Strains serving as nomenclatural type material are designated precisely, and this may be done in various ways (7); typically, a "holotype strain" is designated at the time of valid publication of a new species name (rule 18b), in which the type strain must conform to defined conditions (rule 30). The five species described in the GA publications had been validly published previously, and the type strains for those species were already defined, preserved, and publicly available in numerous culture collections. Meanwhile, the strains reported in the five GA articles were recently isolated independently. Therefore, since they are not descended from the already defined type strains, they cannot represent the authentic type strains of the species.

To confirm these observations, we performed average nucleotide identity (ANI) analyses using *iSpeciesWS* (8, 9) between the five genome sequences reported in the GA articles and the genome sequences of the documented type strains of the five species (Table 1). The ANI values range from 83.10 to 98.96% (percentages aligned, 75.56 to 91.18%), confirming that the described strains are different and, furthermore, confirming that one of them, *Lelliottia nimpresuralis* SGAir0187, is misclassified and is not a strain of this species.

Representing genome sequences from false type strains of given species has the potential to lead to erroneous conclusions in future studies that may rely on the published misinformation, as they may be used as the wrong reference points. We encourage the authors of the five GA publications and the journal to publish corrigenda and remove the words "type strain" from the titles of the publications. We also warn users of publicly available genome sequence data to be cautious in accepting the metadata associated with genome sequences; recent studies have clearly demonstrated the presence of high numbers of misclassified genome se-

Citation Salvà-Serra F, Jaén-Luchoro D, Karlsson R, Bennasar-Figueras A, Jakobsson HE, Moore ERB. 2019. Beware of false "type strain" genome sequences. *Microbiol Resour Announc* 8:e03069-19. <https://doi.org/10.1128/MRA.02919-19>

Editor Irene L. G. Newton, Indiana University, Bloomington

Copyright © 2019 Salvà-Serra et al. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license.

Address correspondence to Francisco Salvà-Serra, francisco.salva@gu.se, or Edward R. B. Moore, erbmoores@cug.ac.uk.

Ed Note: The authors of the published articles did not respond.

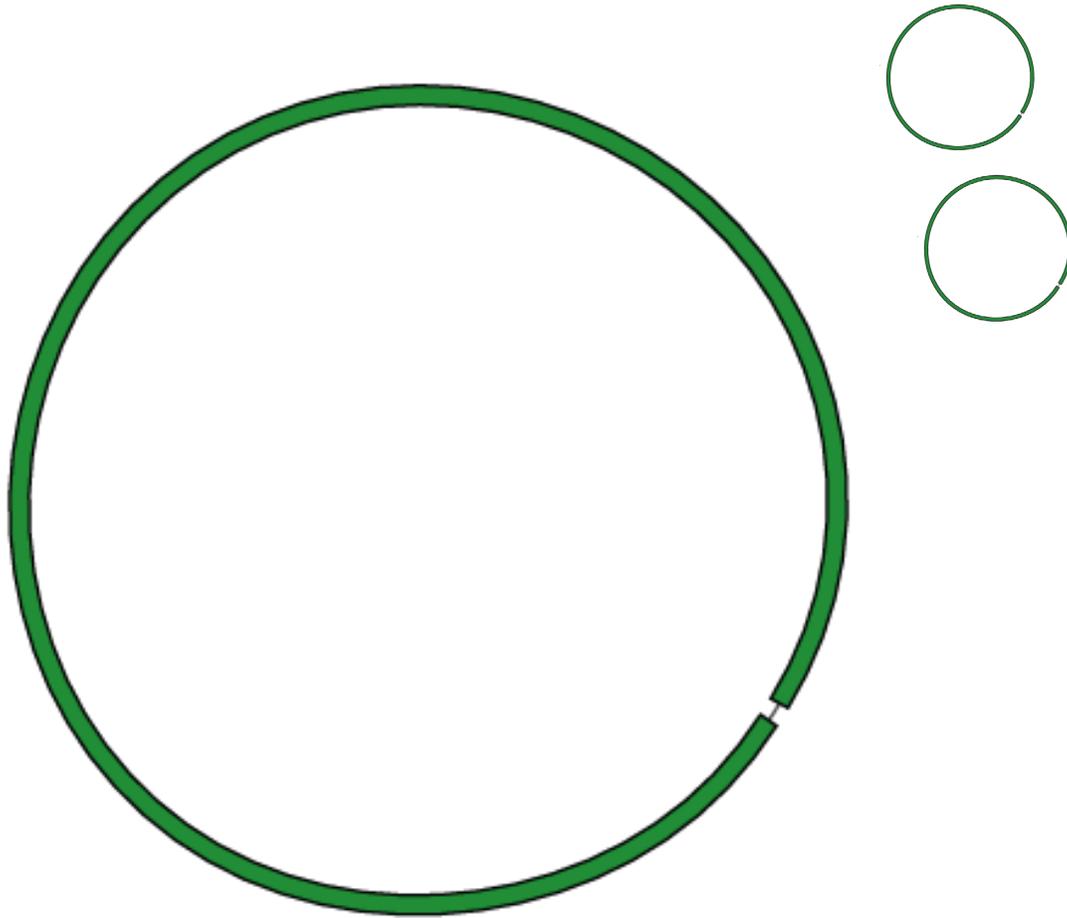
Published 10 May 2019

Downloaded from <http://mra.asm.org/> on June 13, 2019 by guest

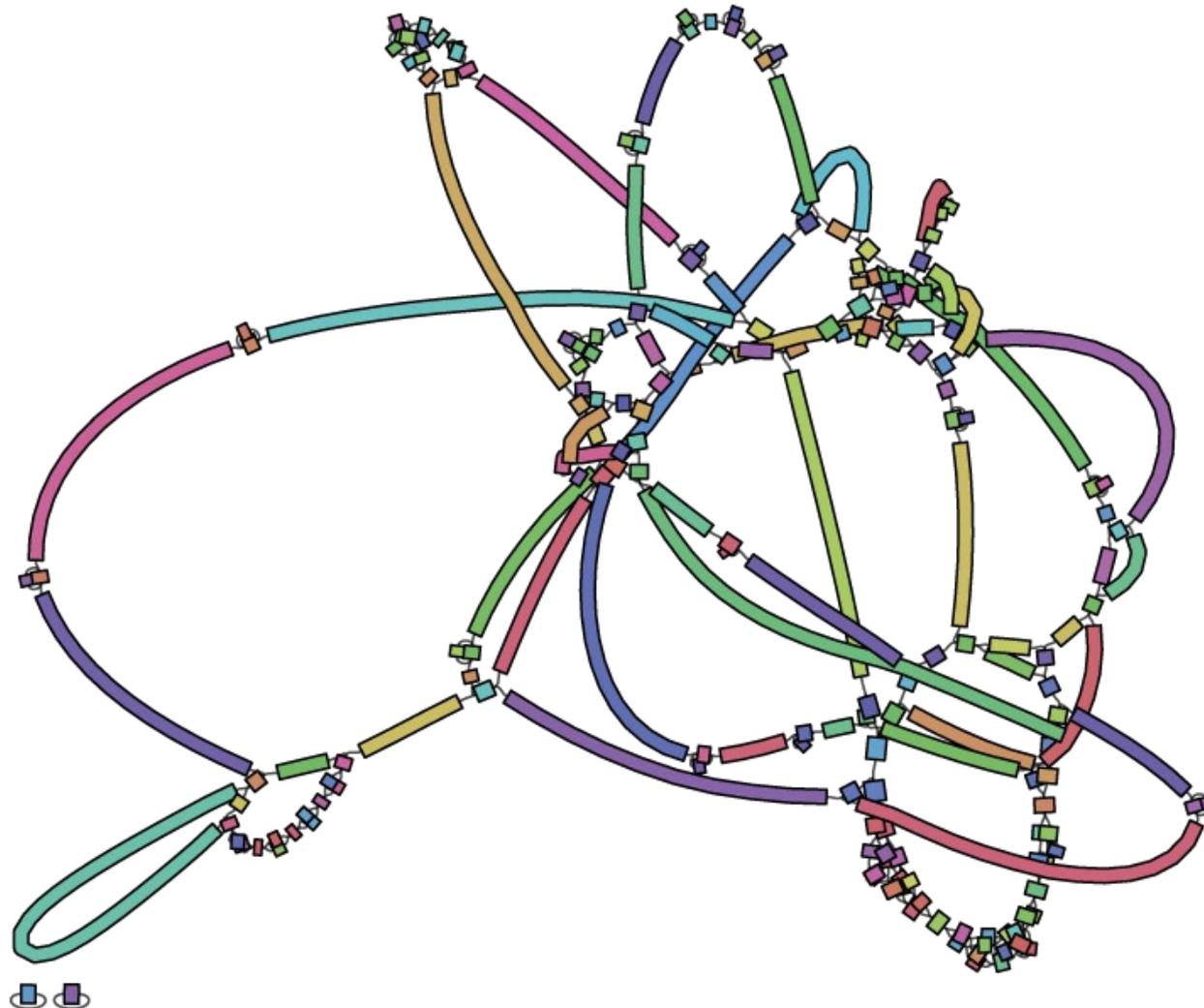
Volume 8 | Issue 22 | e03069-19  mra.asm.org 1

“We also warn users of publicly available genome sequence data to be cautious in accepting the metadata associated with genome sequences...”

What does a good genome assembly look like?



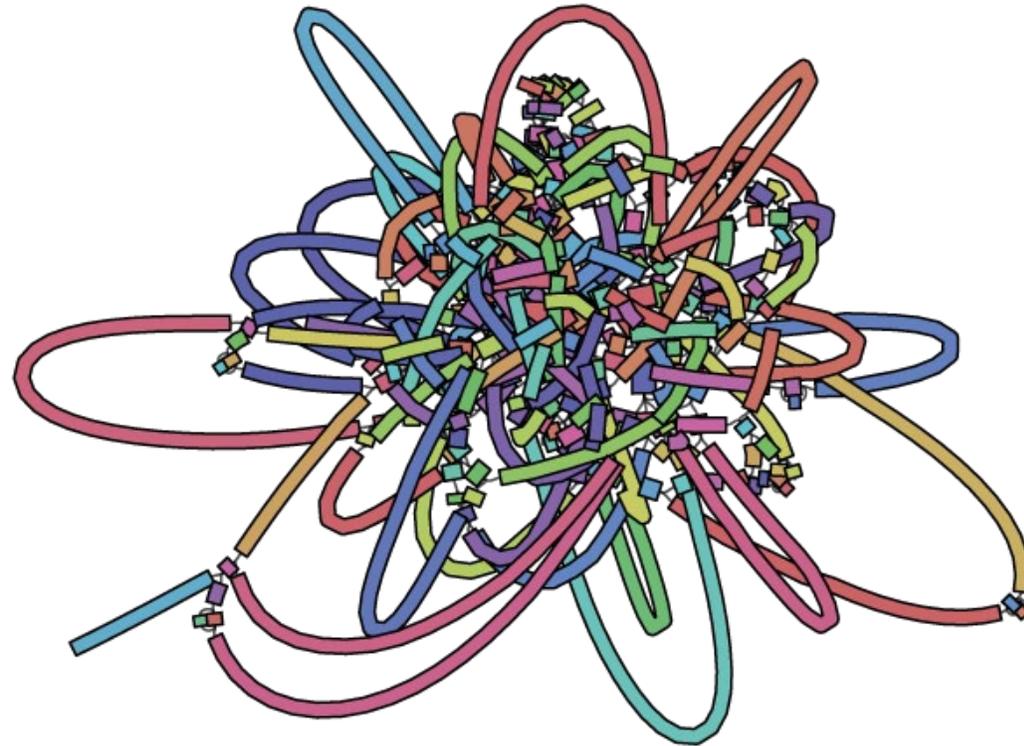
What does a bad (less good) assembly look like?



*Good Illumina-
only assembly*

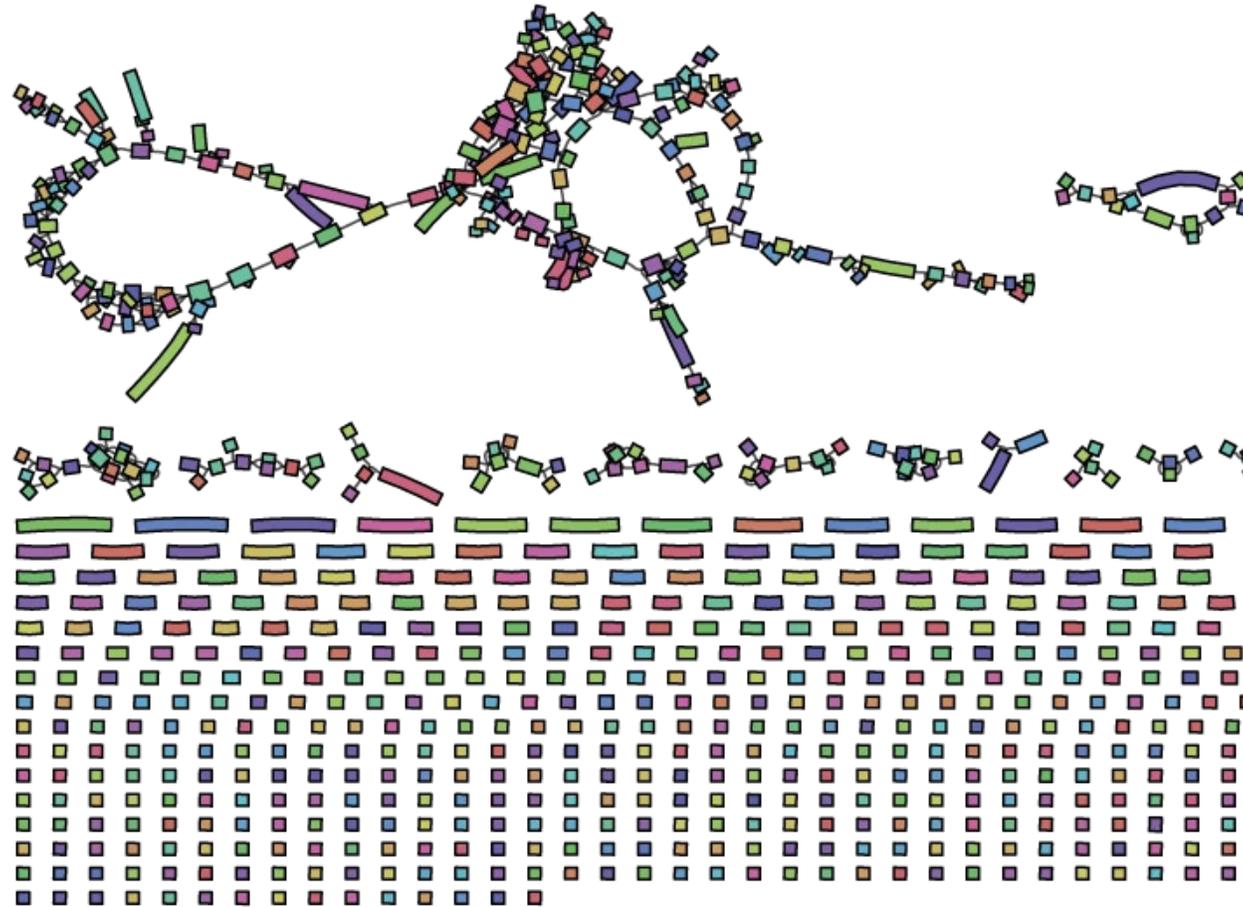


What does a bad (less good) assembly look like?



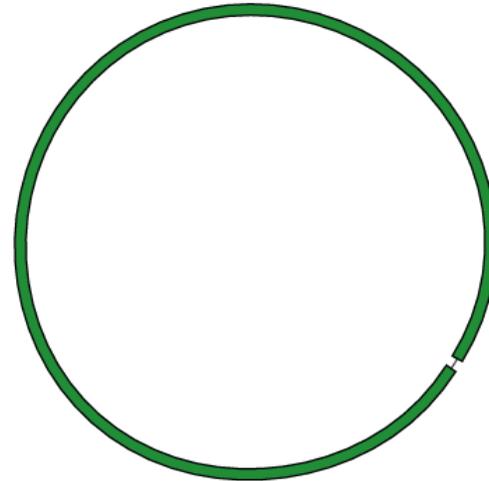
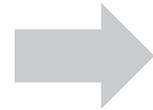
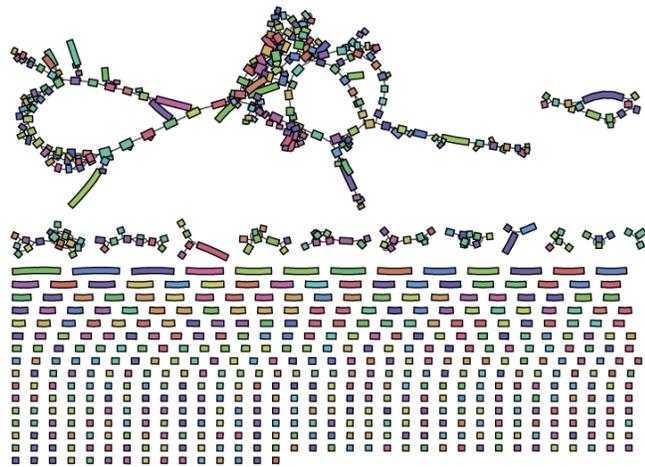
*OK Illumina-
only assembly*

What does a bad (less good) assembly look like?



Bad Illumina-only assembly

Enhanced Authentication Initiative – Goals



Improve

- **Quality**
- **Completeness**
- **Authenticity**
- **Traceability**

Enhanced Authentication Initiative – Overview



1. Extract DNA from authenticated materials



2. Sequence the high-quality DNA



3. Assemble with short and long reads



4. Validate strain designation and publication

Enhanced Authentication Initiative – Overview

The screenshot shows a web browser window with the following elements:

- Browser Tab:** ATCC - Genomes
- Address Bar:** <https://genomes.atcc.org>
- Page Header:** ATCC logo, GENOMES, SEQUENCE SEARCH, and a user profile icon labeled DEMO@ONECODEX.COM.
- Main Content:**
 - ## Welcome to the ATCC Genome Portal
 - A comprehensive collection of high-quality microbial genomics reference data
 - [VIEW ALL GENOMES >](#)
 - ### Search for a genome
 -
 -
 - ### Recently updated
 -  **Acinetobacter johnsonii (ATCC® 17909™)**
Added 05/13/2019
- Footer:** Powered by ONE CODEX logo.

The background of the page features a stylized, glowing DNA double helix structure in shades of green and yellow.

Enhanced Authentication Initiative – Quick Tour

The screenshot shows the ATCC Genomes website interface. The browser address bar displays 'https://genomes.atcc.org/genomes'. The page header includes the ATCC logo, navigation links for 'GENOMES' and 'SEQUENCE SEARCH', and a user profile 'DEMO@ONECODEX.COM'. Below the header, the 'Genomes' section is active, with tabs for 'All Genomes' and 'My Genomes'. A sorting dropdown is set to 'Taxonomic Name' with up and down arrows, and a search bar is present. The main content is a table listing various bacterial genomes.

Taxonomic name	ATCC Product Name	Date Published	Length	Genomic Data	Download
<i>Acinetobacter baumannii</i>	ATCC® BAA-1710™	May 14, 2019	4.0 Mb	View	Download
<i>Acinetobacter baumannii</i>	ATCC® 19606™	May 14, 2019	4.0 Mb	View	Download
<i>Acinetobacter baumannii</i>	ATCC® BAA-1605™	May 14, 2019	4.1 Mb	View	Download
<i>Acinetobacter johnsonii</i>	ATCC® 17909™	May 14, 2019	3.6 Mb	View	Download
<i>Actinobacillus pleuropneumoniae</i>	ATCC® 27088™	May 14, 2019	2.3 Mb	View	Download
<i>Alcanivorax borkumensis</i>	ATCC® 700651™	May 14, 2019	3.1 Mb	View	Download
<i>Bacillus cereus</i>	ATCC® 10702™	May 14, 2019	5.6 Mb	View	Download
<i>Bacillus subtilis</i>	ATCC® 6633™	May 14, 2019	4.0 Mb	View	Download
<i>Bartonella henselae</i>	ATCC® 49882™	May 14, 2019	2.0 Mb	View	Download

Enhanced Authentication Initiative – Quick Tour

ATCC GENOMES SEQUENCE SEARCH DEMO@ONECODEX.COM

Acinetobacter baumannii (ATCC® BAA-1710™) *Acinetobacter calcoaceticus/baumannii* complex *Acinetobacter baumannii*

[Overview](#) [Genome Browser](#) [Related Genomes](#) [Quality Control](#)

[DOWNLOAD ASSEMBLY](#) [DOWNLOAD ANNOTATIONS](#)

Assembly Summary	
Date Published	May 13, 2019
Length	3,960,239 bp
Sequencing Technology	Illumina + Oxford Nanopore Hybrid Assembly
Number of Contigs	4 (All Circularized)
N50	3,942,209 bp
%GC	39.35%

Genome Summary	
Name	ATCC® BAA-1710™
Isolation	-
Applications	-
Biosafety Level	2
Annotated Genes	-
MIC Range	-
Type Strain	-
Preceptrol	-
Genotype	-
Antigenic Properties	-

Enhanced Authentication Initiative – Quick Tour

ATCC GENOMES SEQUENCE SEARCH DEMO@ONECODEX.COM

Acinetobacter baumannii (ATCC® BAA-1710™) *Acinetobacter calcoaceticus/baumannii* complex *Acinetobacter baumannii*

Overview Genome Browser Related Genomes Quality Control Annotation Legend

Display Hypothetical Proteins Search Download Table CSV

Contig	Start	End	Name	Product	EC Number	Type	Uniprot ID	Jump
1	1	1399	<i>dnaA</i>	Chromosomal replication initiator protein DnaA		CDS	P03004	↗
1	1496	2645	<i>dnaN</i>	Beta sliding clamp		CDS	Q9I7C4	↗
1	2659	3742	<i>recF_1</i>	DNA replication and repair protein RecF		CDS	POA7H0	↗
1	3794	6263	<i>gyrB</i>	DNA gyrase subunit B	5.99.1.3	CDS	POAES6	↗
1	6300	6693	<i>cybC</i>	Soluble cytochrome b562		CDS	POABE7	↗
1	6776	7334	<i>dedA</i>	Protein DedA		CDS	POABP6	↗
1	7584	9516	<i>yheS</i>	putative ABC transporter ATP-binding protein YheS		CDS	P63389	↗

Enhanced Authentication Initiative – Quick Tour

ATCC GENOMES SEQUENCE SEARCH DEMO@ONECODEX.COM

Acinetobacter baumannii (ATCC® BAA-1710™) > *Acinetobacter calcoaceticus/baumannii* complex > *Acinetobacter baumannii*

Overview Genome Browser Related Genomes Quality Control

Most similar Genomes

The following genomes have the greatest genomic similarity to this one (>95% average nucleotide identity).

Genome Name	Similarity	Contigs	Size	Action
<i>Acinetobacter baumannii</i> (ATCC® BAA-1605™)	99.6% similar	3 contigs	4.1 Mb	View Genome
<i>Acinetobacter baumannii</i> (ATCC® 19606™)	97.9% similar	3 contigs	4.0 Mb	View Genome

Other members of this genus

The following genomes share the same genus according to the NCBI taxonomy.

Genome Name	Contigs	Size
<i>Acinetobacter baumannii</i> (ATCC® BAA-1605™)	3 contigs	4.1 Mb
<i>Acinetobacter baumannii</i> (ATCC® 19606™)	3 contigs	4.0 Mb
<i>Acinetobacter johnsonii</i> (ATCC® 17909™)	3 contigs	4.0 Mb

Enhanced Authentication Initiative – Quick Tour

The screenshot displays the ATCC Genomes website interface. The browser address bar shows the URL <https://genomes.atcc.org/genomes/65192d5a944b42d5>. The user is logged in as DEMO@ONECODEX.COM. The main heading is *Acinetobacter baumannii* (ATCC® BAA-1710™). Below the heading are navigation tabs: Overview, Genome Browser, Related Genomes, and Quality Control (which is selected). The page is divided into two columns: Sequencing Quality Control and Assembly Quality Control. Both columns show a green circular progress indicator with '3/3' and the text '3 out of 3 passed'. The Sequencing Quality Control section lists three metrics: Number of trimmed reads (5,442,598), Median Q score, all bases (38), and Ambiguous content (% N bases) (0). The Assembly Quality Control section lists three metrics: Estimated genome completeness (100%), Estimated genome contamination (0.27%), and Average depth of coverage (368.772x). Each metric is accompanied by a green 'Passed' badge.

ATCC GENOMES SEQUENCE SEARCH DEMO@ONECODEX.COM

Acinetobacter baumannii (ATCC® BAA-1710™) *Acinetobacter calcoaceticus/baumannii* complex *Acinetobacter baumannii*

Overview Genome Browser Related Genomes Quality Control

Sequencing Quality Control

Quality control statistics on Illumina sequencing data.

3/3 3 out of 3 passed

Passed	Number of trimmed reads	5,442,598
Passed	Median Q score, all bases	38
Passed	Ambiguous content (% N bases)	0

Assembly Quality Control

Metrics assessing the assembly quality (from CheckM).

3/3 3 out of 3 passed

Passed	Estimated genome completeness	100%
Passed	Estimated genome contamination	0.27%
Passed	Average depth of coverage	368.772x

Enhanced Authentication Initiative – Quick Tour

ATCC GENOMES SEQUENCE SEARCH DEMO@ONECODEX.COM

Acinetobacter baumannii (ATCC® BAA-1710™) *Acinetobacter calcoaceticus/baumannii* complex *Acinetobacter baumannii*

Overview Genome Browser Related Genomes Quality Control Annotation Legend

Display Hypothetical Proteins Download Table CSV

Contig	Start	End	Name	Product	EC Number	Type	Uniprot ID	Jump
1	391861	393046	<i>tetA</i>	Tetracycline resistance protein, class C		CDS	P02981	↗
1	3647608	3648883	<i>tetA</i>	Tetracycline resistance protein, class C		CDS	P02981	↗
1	3677181	3678357	<i>tetA</i>	Tetracycline resistance protein, class C		CDS	P02981	↗

Powered by ONE CODEX

Enhanced Authentication Initiative – Quick Tour

The screenshot shows a web browser window with the URL `https://genomes.atcc.org/sequence-search`. The page features the ATCC logo and navigation links for "GENOMES" and "SEQUENCE SEARCH". A user profile is logged in as "DEMO@ONECODEX.COM". The main heading is "Search for a genome". Below this, instructions state: "Enter a nucleotide sequence (at least 40 bases) to find genomes that match >80% of the sequence". A large text input box contains the placeholder text: "Enter or paste a nucleotide sequence here to find all of the genomes that match your query". To the right of the input box is a "Results" section with the instruction: "Use the search box on the left to enter your nucleotide sequence of interest." At the bottom right of the input area is a "Search" button with a magnifying glass icon.

Enhanced Authentication Initiative – Quick Tour

The screenshot shows a web browser window with the URL `https://genomes.atcc.org/sequence-search`. The page header includes the ATCC logo, navigation links for "GENOMES" and "SEQUENCE SEARCH", and a user profile icon labeled "DEMO@ONECODEX.COM".

Search for a genome

Enter a nucleotide sequence (at least 40 bases) to find genomes that match >80% of the sequence

Results

Use the search box on the left to enter your nucleotide sequence of interest.

```
CTTTTTCACATAGCCTAATTTTATTGCTAGTAGGTCGCATTATGCTGGAATC
ACCAAGTCCCAACATGGCTGTTGCAAGTGCCTTATATTGTCGATGTTTGCACGA
AAATAACCGCGCAAAATATTTTGGTTAATCAATGCTATGTTTGGTGCAGGCT
TCATTATTGGCCCTGTATTGGCGGATTCCTGAGTGAATATGGATTAAGACTT
CCTTCTTTGCGGCAGCTATATTGACAGGGCTTAATCTTTTATCTGCCTATTT
TGTTTTCCTGAATCCGAAAAGTGAATTTGGAGAACAACAATTATCTACAT
TAAACCCCTTTAAAAATATTGCTGGTATTAGCTCTATTTCGTGGTACTTCCA
CTTATTAGACCTTTTTTATCTTTAGTGCATAGGGGAGGTATATGGAGTCTG
CTGGGCATTATGGGGACATGACACATTCAGTGGAGCGTTTCTGGGTAGGTC
TTTCTCTAGGCGCATTGGTCTATGTCAAATGCTGGTACAGGCTCTTATTCCG
AGTCATGCTTCAAGATTGCTGGGTAATCGTAATGCTGTGCTGGCTGGTATTGC
TTGTTCTTGTTTTGGCTTTAGCAGTAATGGCTTTCGCCCAAAGTGGTTGGATGA
TTTTTGCTATTATGCCTATTTTGGCTAGGGAGTATGGGGACACCTTCATTA
CAAGCCTTAGCTTCTCAAAGGTTTCTGCTGACCAGCAAGGACAGTTTCAGGG
AGTGATAGCATCTACGGTAAGTATGGCCCTATGATTGCCCTATGTTTTTCT
CCACTCTTATTTTTCAGTTTCAGGAAAAATGGCCGGGAGCAATTGGTTGAGC
GTGATTTGATTTACCTGCTCACCTTACCGATCATTFTGTATAGTACCCGACC
AGTCGTACAACAAAGATAG
```

Search

Enhanced Authentication Initiative – Quick Tour

The screenshot shows the ATCC Genomes Sequence Search interface. The browser address bar displays `https://genomes.atcc.org/sequence-search`. The page title is "Search for a genome". Below the title, there is a search instruction: "Enter a nucleotide sequence (at least 40 bases) to find genomes that match >80% of the sequence". A search input field contains a long nucleotide sequence, and a "Search" button is located below it. The results section, titled "Results on 1185 bases", lists three matches for *Acinetobacter baumannii* strains:

- Acinetobacter baumannii (ATCC® BAA-1710™)**: 1185 bases matched (100%). 4 contigs, 4.0 Mb. [View Genome](#)
- Acinetobacter baumannii (ATCC® BAA-1605™)**: 1185 bases matched (100%). 3 contigs, 4.1 Mb. [View Genome](#)
- Acinetobacter baumannii (ATCC® 19606™)**: 990 bases matched (83%). 3 contigs, 4.0 Mb. [View Genome](#)

Enhanced Authentication Initiative – Quick Tour

The screenshot displays the ATCC Genomes Sequence Search web application. The browser address bar shows the URL `https://genomes.atcc.org/sequence-search`. The page header includes the ATCC logo, navigation links for "GENOMES" and "SEQUENCE SEARCH", and a user profile icon labeled "DEMO@ONECODEX.COM".

Search for a genome

Enter a nucleotide sequence (at least 40 bases) to find genomes that match >80% of the sequence

```
ATGAATCGATCTCTATTTATTATCTTTGCAACTATTGCTTTAGATGCTATTGG
TATCGGCTTATTTTTCCGATTCTTCCTTTATTATTACAAGATATGACGCATA
GCACTCATATTTCTATATATATGGGTATATGGCCAGTCTCTATGCGGCCATG
CAATTTATCTTCTCCTTTATTAGGTGCGTTAAGTGACAGATGGGGCGTAG
ACCGGTCTTGCTTATTTCACTGGCTGGGTGAGCAGTTAATTATCTCTTTCTAA
CTTTTTCACATAGCCTAATTTTATTGCTAGTAGGTCGCATTATGCTGGAATC
ACCAGTGCCAACATGGCTGTGCAAGTGCCTATATTGTCGATGTTTGCACGA
AAATAACCGCGCAAAATATTTGGTTTAAATCAATGCTATGTTTGGTGCAGGCT
TCATTATTGGCCCTGTATTGGGCGGATTCTTGAGTGAATATGGATTAAGACTT
CCTTTCTTTGCGGCAGCTATATTGACAGGCTTAATCTTTTATCTGCCATTTT
TGTTTTGCTTGAATCCCGAAAAGTGACTTTGGAGAACAAACAATTATCTACAT
TAAACCCCTTTAAAATATTGCTGGTATTAGCTCTATTCGTGGTGTACTTCCA
CTTATTACGACCTTTTTTATCTTTAGTGCCATAGGGGAGGTATATGGAGTCTG
CTGGGCATTATGGGACATGACACATTTCACTGGAGCGGTTTCTGGGTAGTCT
TTTCTTAGGCGCATTTGGTCTATGTCAAATGCTGGTACAGGCTCTTATTCCG
AGTCATGCTTCAAGATTGCTGGGTAATCGTAATGCTGTGCTGGCTGGTATTGC
TTGTTCTTGTTTGCTTTAGCAGTAATGGCTTTCGCCCAAAGTGGTTGGATGA
```

Results on 1185 bases

- Acinetobacter baumannii (ATCC® BAA-1710™)** 1185 bases matched (100%)
4 contigs, 4.0 Mb [View Genome](#)
- Acinetobacter baumannii (ATCC® BAA-1605™)** 1185 bases matched (100%)
3 contigs, 4.1 Mb [View Genome](#)
- Acinetobacter baumannii (ATCC® 19606™)** 990 bases matched (83%)
3 contigs, 4.0 Mb [View Genome](#)

[Search](#)

Enhanced Authentication Initiative – Quick Tour

The screenshot shows the ATCC Genomes Sequence Search interface. The browser address bar displays `https://genomes.atcc.org/sequence-search`. The page title is "ATCC GENOMES SEQUENCE SEARCH" and the user is logged in as "DEMO@ONECODEX.COM".

Search for a genome

Enter a nucleotide sequence (at least 40 bases) to find genomes that match >80% of the sequence

```
ATGAATCGATCTCTATTTATTATCTTTGCAACTATTGCTTTAGATGCTATTGG
TATCGGCTCTATTTTCCGATTCTTCCTTTATTATTACAAGATATGACGCATA
GCACTCATATTTCTATATATATGGGTATATGGCCAGTCTCTATGCGGCCATG
CAATTTATCTTCTCCTTTATTAGGTGCGTTAAGTGACAGATGGGGCGTAG
ACCGGTCTTGCTTATTTCACTGGCTGGGTGAGCAGTTAATTATCTCTTTCTAA
CTTTTTCACATAGCCTAATTTTATTGCTAGTAGGTCGCATTATTGCTGGAATC
ACCAGTGCCAACATGGCTGTGCAAGTGCTTATATTGTCGATGTTTGCACGA
AAATAACCGCGCAAAATATTTGGTTAAATCAATGCTATGTTTGGTGCAGGCT
TCATTTTGGCCCTGTATTGGGCGGATTCTTGAGTGAATATGGATTAAGACTT
CCTTTCTTTGCGGCAGCTATATTGACAGGCTTAATCTTTTATCTGCCATTTT
TGTTTTGCTGAAATCCGAAAAGTGACTTTGGAGAACAAACAATTATCTACAT
TAAACCCCTTTAAAATATTGCTGGTATTAGCTCTATTCGTGGTGTACTTCCA
CTTATTACGACCTTTTTTATCTTTAGTGCCATAGGGGAGGTATATGGAGTCTG
CTGGGCATTATGGGACATGACACATTTCACTGGAGCGGTTTCTGGGTAGTCT
TTTCTTAGGCGCATTTGGTCTATGTCAAATGCTGGTACAGGCTCTTATTCCG
AGTCATGCTTCAAGATTGCTGGGTAATCGTAATGCTGTGCTGGCTGGTATTGC
TTGTTCTTGTTTTGCCTTAGCAGTAATGGCTTTCGCCCAAAGTGGTTGGATGA
```

Results on 1185 bases

- Acinetobacter baumannii (ATCC® BAA-1710™)** 1185 bases matched (100%)
4 contigs, 4.0 Mb [View Genome](#)
- Acinetobacter baumannii (ATCC® BAA-1605™)** 1185 bases matched (100%)
3 contigs, 4.1 Mb [View Genome](#)
- Acinetobacter baumannii (ATCC® 19606™)** 990 bases matched (83%)
3 contigs, 4.0 Mb [View Genome](#)

[Search](#)

Enhanced Authentication Initiative – Quick Tour

The screenshot displays the ATCC Genomes Sequence Search web application. The browser address bar shows the URL `https://genomes.atcc.org/sequence-search`. The page header includes the ATCC logo, navigation links for "GENOMES" and "SEQUENCE SEARCH", and a user profile for "DEMO@ONECODEX.COM".

The main heading is "Search for a genome". Below it, instructions state: "Enter a nucleotide sequence (at least 40 bases) to find genomes that match >80% of the sequence".

The search results are displayed on the right side, titled "Results on 1185 bases". Three results are shown:

- Acinetobacter baumannii (ATCC® BAA-1710™)**: 1185 bases matched (100%). 4 contigs, 4.0 Mb. [View Genome](#)
- Acinetobacter baumannii (ATCC® BAA-1605™)**: 1185 bases matched (100%). 3 contigs, 4.1 Mb. [View Genome](#)
- Acinetobacter baumannii (ATCC® 19606™)**: 990 bases matched (83%). 3 contigs, 4.0 Mb. [View Genome](#)

The search input field on the left contains a long nucleotide sequence with a "Search" button below it.

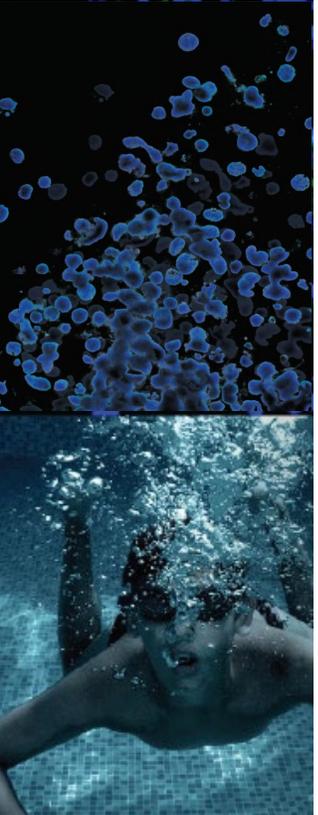
Enhanced Authentication Initiative – Outcomes

- Modern, easy-to-use genome browser and data platform
- Improved **quality, completeness, authenticity, traceability**
- Publicly available September 2019



Questions?

Credible Leads to Incredible™

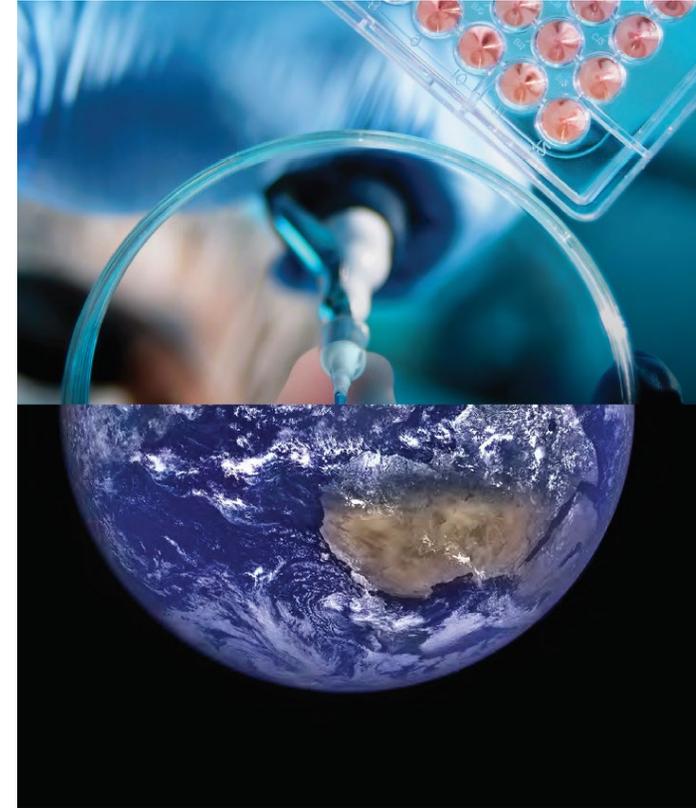
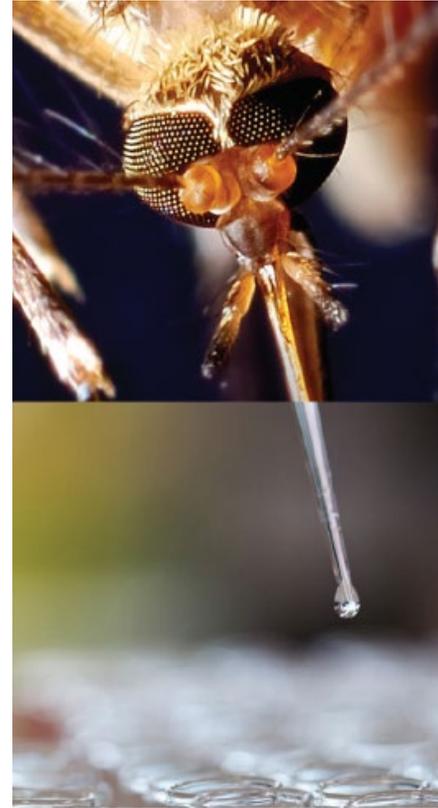
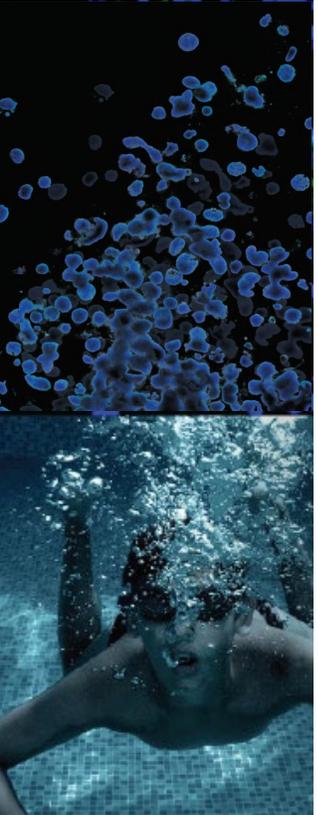




Using Whole-Genome Sequencing for the Enhanced Authentication of ATCC's *Bacillus cereus* Group Strains

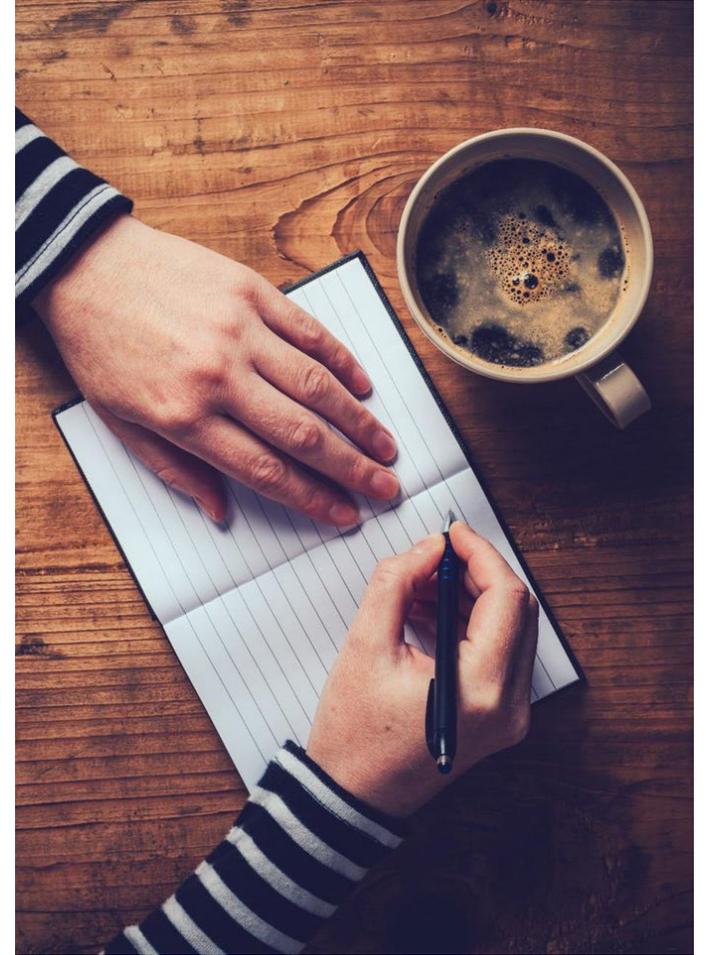
Marco A. Riojas, PhD
Scientist, ATCC

Credible Leads to Incredible™



Agenda

- *Bacillus* and *Bacillus cereus* Group (BcG)
- Type strains and species definitions
- BcG strain analysis





Bacillus and *Bacillus cereus* Group (BcG)

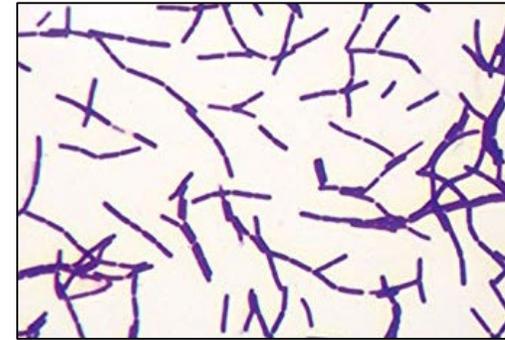
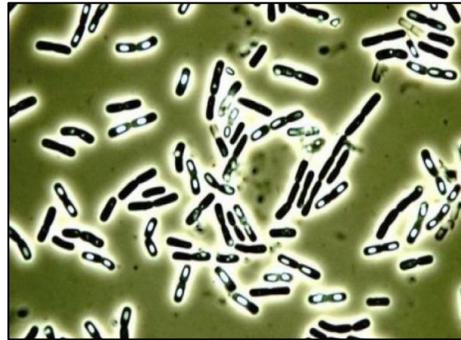
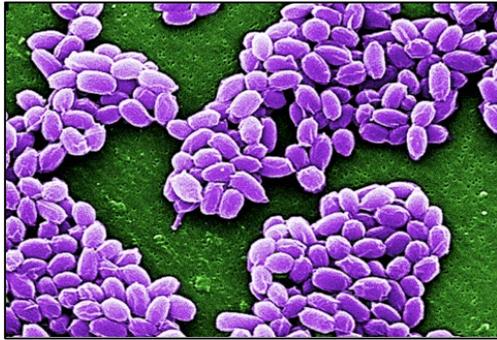
The Genus *Bacillus*

- Composed of Gram-positive, aerobic, endospore-forming bacteria
- Formed in 1872 with the description of *Bacillus subtilis* (the type species) and *B. anthracis*
- Currently a total of 281 species and subspecies



The *Bacillus cereus* Group (BcG)

Currently 18 species, some of which are pathogenic to various plant and animal species



B. anthracis



B. cereus



B. thuringiensis

The *Bacillus cereus* Group (BcG)

A group of closely related species, including some important to health and biotechnology

Species	Year Identified
<i>B. anthracis</i>	1872
<i>B. mycoides</i>	1886
<i>B. cereus</i>	1887
<i>B. thuringiensis</i>	1915
<i>B. pseudomycoides</i>	1998
<i>B. weihenstephanensis</i> *	1998
<i>B. cytotoxicus</i>	2013
<i>B. toyonensis</i>	2014
<i>B. weidmannii</i>	2016

Species	Year Identified
<i>B. paranthracis</i>	2017
<i>B. pacificus</i>	2017
<i>B. tropicus</i>	2017
<i>B. albus</i>	2017
<i>B. mobilis</i>	2017
<i>B. luti</i>	2017
<i>B. proteolyticus</i>	2017
<i>B. nitratireductans</i>	2017
<i>B. paramycoides</i>	2017



Type Strains and Species Definitions

How Are Species Defined?

Background

- Each species is represented by a type strain and a description of that strain
 - Usually the first strain identified
 - Not necessarily the most typical or representative of the species
 - **The type strain is essentially the “definition” of a species**
- If the strain upon which the description was based upon cannot be found or is no longer appropriate, a neotype strain may be proposed
- *Mycobacterium tuberculosis*
 - Identified in 1882 by Robert Koch
 - Strain H37Rv was proposed as the neotype strain in 1972 (Kubica, *et al.*)
 - Now *Mycobacterium tuberculosis* H37Rv^T is the type strain

How is a New Strain Assigned to a Species?

Background

- The characteristics of the new strain are compared to the characteristics of species type strains
 - Historically: phenotypic characteristics
- A strain that shares enough of the essential characteristics of a type strain is said to be within the circumscription of that species/type strain
 - Therefore, it belongs to that species
- Recognizing that phenotypes can be quite unreliable, today we rely more heavily on genotypic comparisons
 - 16S rRNA genes, *hsp65*, *rpoB*
- Small numbers of genes can still provide misleading results
 - Most accurate comparison would be between whole genomes

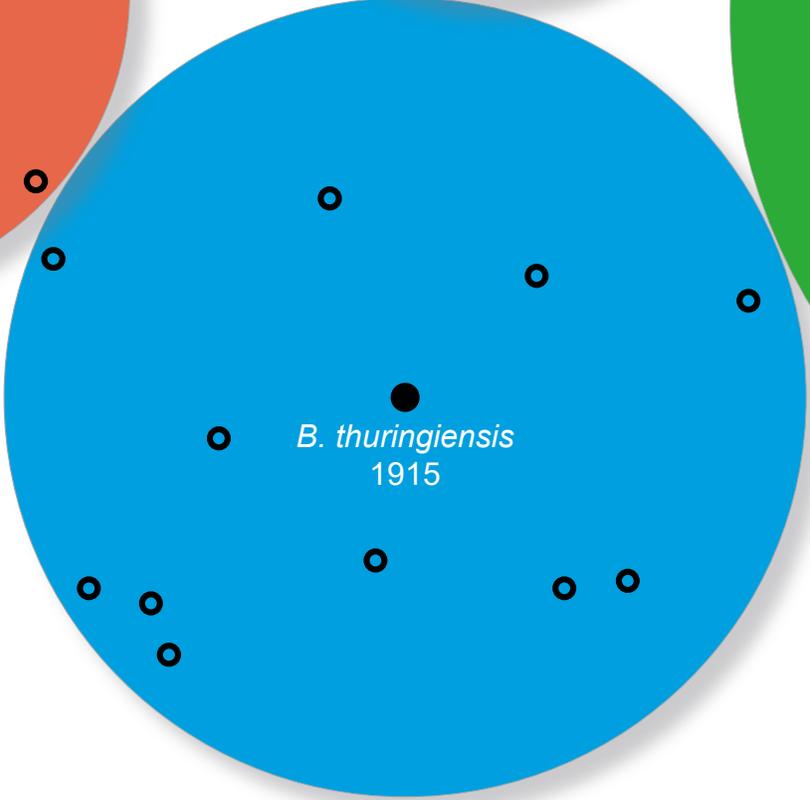
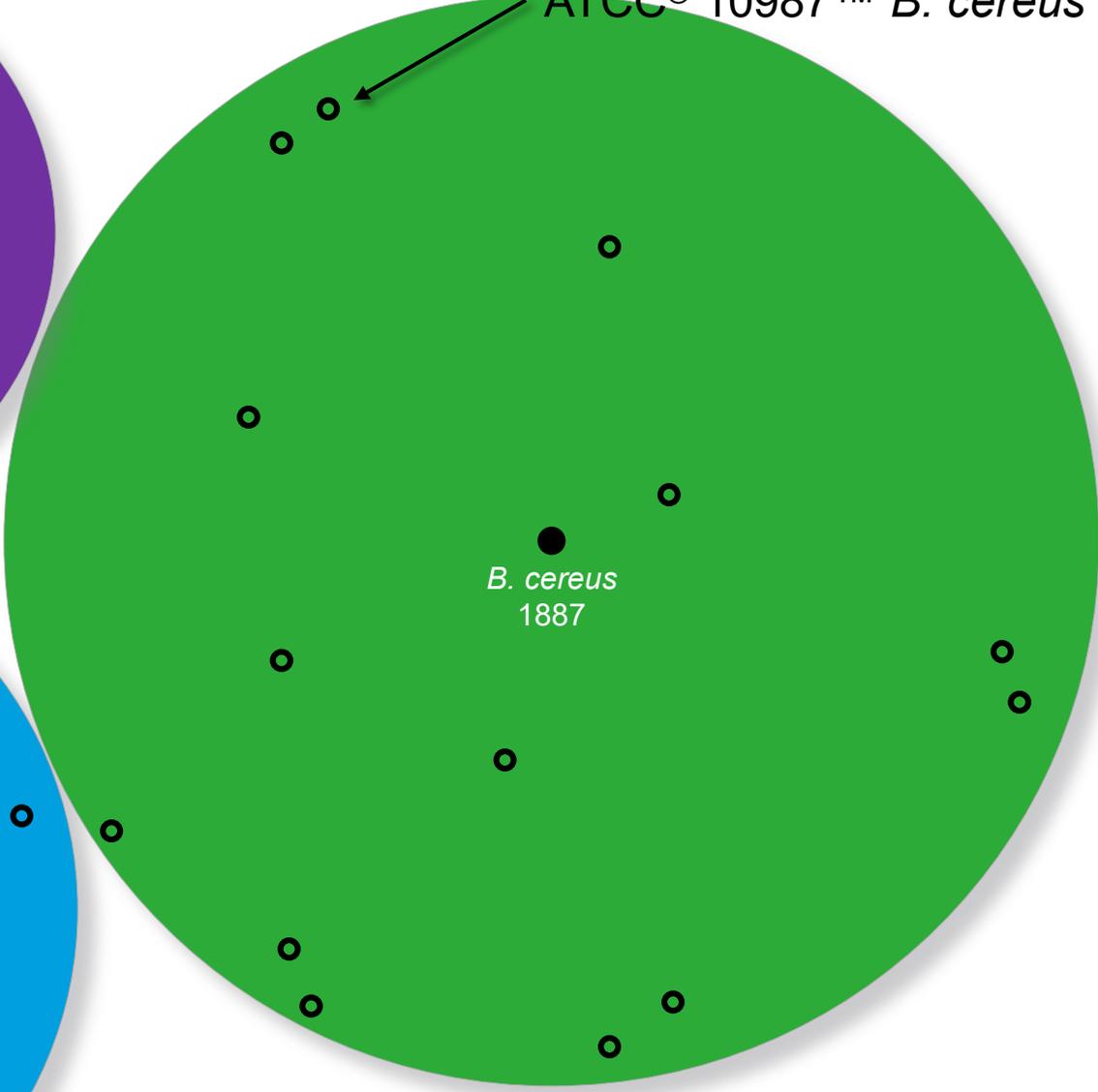
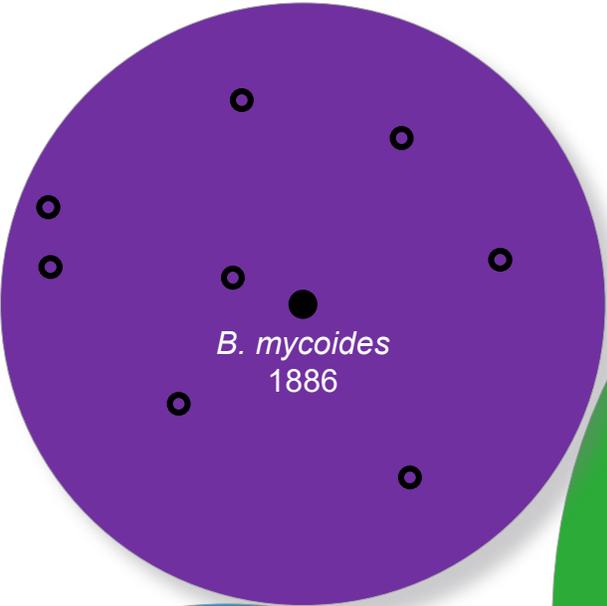
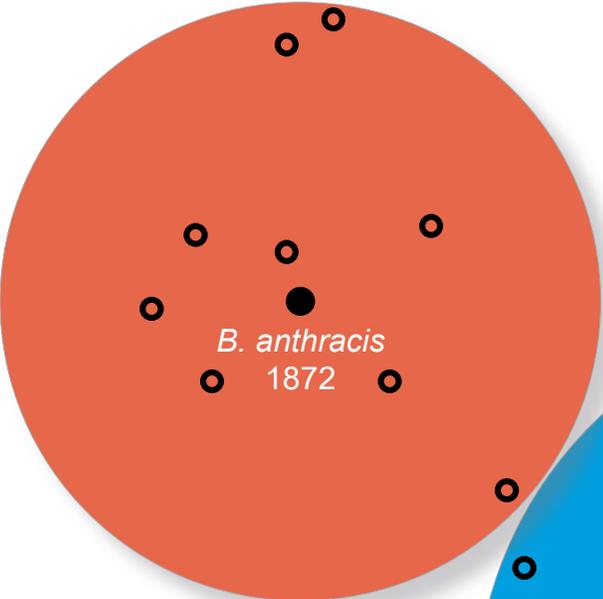
Bacillus cereus Group (BcG) Strains at ATCC

Background

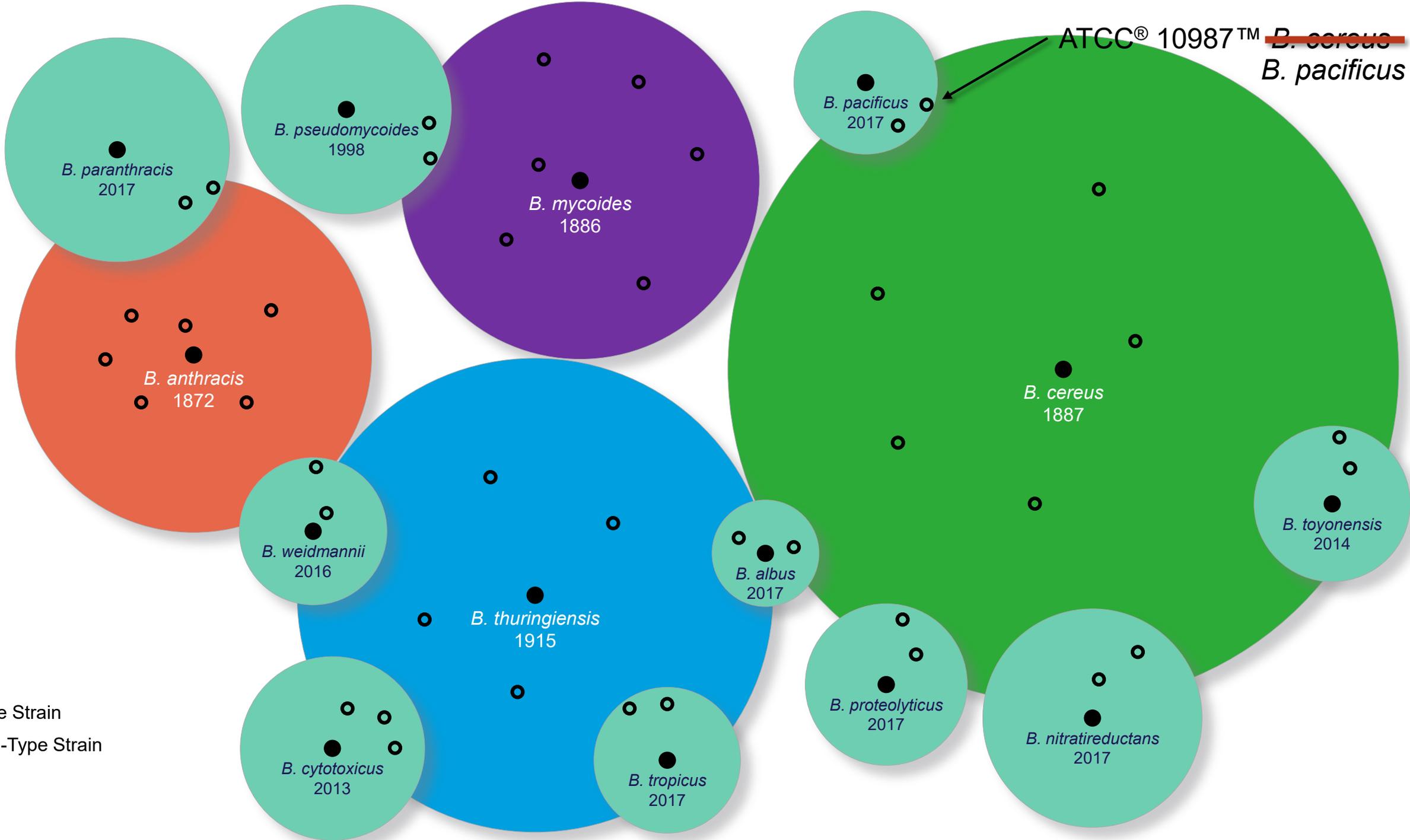
- Many of our strains were deposited many decades ago
 - E.g., ATCC® 246™ *Bacillus cereus* was deposited in 1925
 - New species have been discovered
 - Previous classifications of our strains were based on less accurate/comprehensive methods, such as phenotypic observations or biochemical testing
- Do our *Bacillus* strains match up with current taxonomy?
 - GOAL: Identify our strains using the most current techniques and definitions
- Whole-genome sequencing (WGS) of expertly authenticated material
 - Primary focus on BcG strains from ATCC and BEI Resources*
 - Secondary focus on the BcG species at large

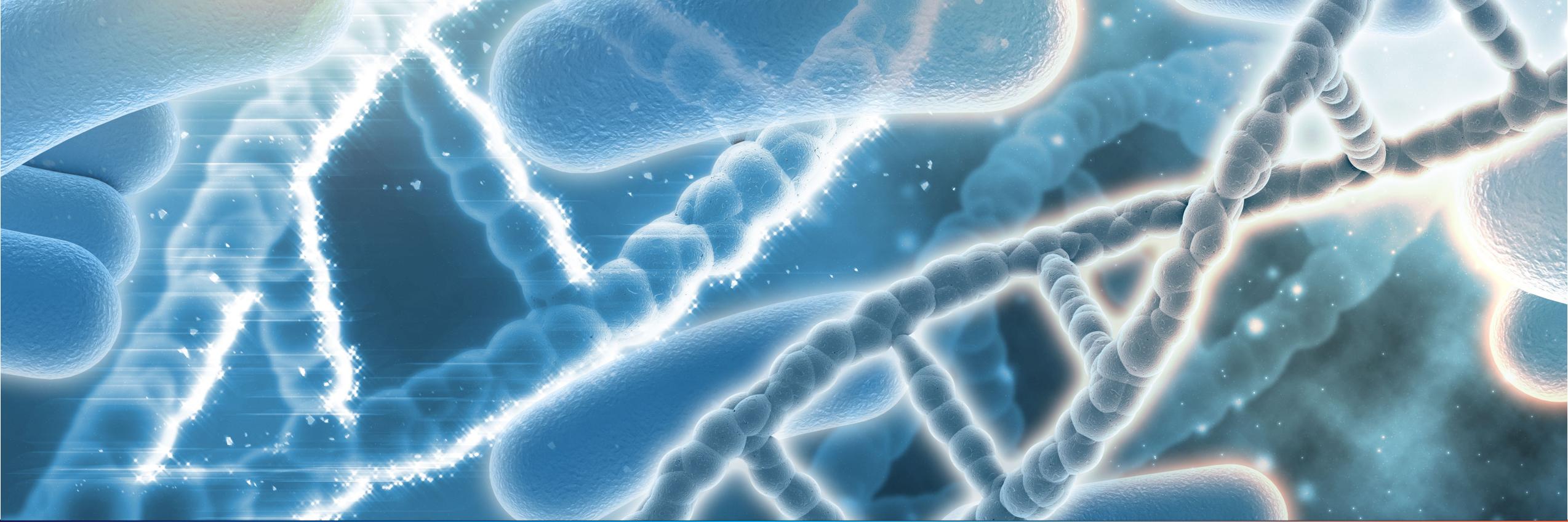
* BEI Resources material is the property of the National Institute of Allergy and Infectious Diseases (NIAID), NIH

ATCC® 10987™ *B. cereus*



- Type Strain
- Non-Type Strain





BcG Strain Analysis

Bacillus cereus Group (BcG) Strains

Methodology

- Selected a subset of BcG and *Bacillus* strains
 - Priority given to lower accession numbers (e.g., ATCC[®] 246[™], 4342[™], 6463[™], etc.)
 - Older deposits more likely to be misclassified
 - Also obtained strains from BEI Resources, a NIAID funded repository
- Whole-Genome Sequencing (WGS)
 - Illumina[®] MiSeq[®] v3 (2×300)
- Genome Assembly
 - SPAdes
- Genomic Comparison of Taxonomy
 - Digital DNA-DNA Hybridization (dDDH)
 - Genome-to-Genome Distance Calculator (GGDC)
 - Average Nucleotide Identity (ANI)
 - OrthoANIb

dDDH Range	Interpretation	
≥ 80%	Same species	Same subspecies
70 – 80%	Same species	Different subspecies
< 70%	Different species	

Meier-Kolthoff JP, et al. (2014); Meier-Kolthoff JP, et al. (2013); Auch, et al. (2013a); Auch, et al. (2013b).

B. anthracis and B. mycooides

Species/Subspecies	Strain	Source	A0488 ^T	NR-1041	NR-1202	NR-1389	NR-1408	NR-36073	NR-36091	NR-3838	NR-411	NR-41	NR-46	NR-51483	NR-51484	NR-51485	NR-51486	ATCC 6462 ^T	ATCC 6462 ^T	NR-613	NR-22171	NBRC 101238 ^T	KBAB4	DSM 20231 ^T
<i>B. anthracis</i>	A0488 ^T	ABJC01	100	99.8	99.7	99.7	99.6	99.7	99.6	99.7	99.8	99.5	99.6	99.6	99.8	99.8	99.8	98.2	98	99.7	97.8	97.9	98.2	28
<i>B. anthracis</i>	NR-1041	This Work	99.92	100	99.8	99.8	99.7	99.7	99.7	99.8	99.8	99.5	99.7	99.7	99.8	99.9	99.8	98.3	98.1	99.8	98	98	98.2	23
<i>B. anthracis</i>	NR-1202	This Work	99.93	99.92	100	99.8	99.8	99.8	99.7	99.9	99.9	99.6	99.8	99.9	99.8	99.8	99.8	98.2	98.1	99.8	97.9	97.9	98.2	23
<i>B. anthracis</i>	NR-1389	This Work	99.89	99.91	99.92	100	99.6	99.9	99.8	99.8	99.8	99.7	99.7	99.7	99.8	99.8	99.8	98.3	98.1	99.8	97.9	98	98.2	23
<i>B. anthracis</i>	NR-1408	This Work	99.87	99.89	99.92	99.92	100	99.7	99.6	99.8	99.8	99.5	99.7	99.7	99.6	99.7	99.6	98.8	98.6	99.9	98.4	98.5	98.8	23
<i>B. anthracis</i>	NR-36073	This Work	99.90	99.93	99.90	99.93	99.87	100	99.8	99.8	99.7	99.7	99.7	99.7	99.7	99.7	99.7	98.2	98.1	99.8	97.9	98	98.2	22
<i>B. anthracis</i>	NR-36091	This Work	99.89	99.88	99.88	99.90	99.88	99.91	100	99.7	99.7	99.6	99.7	99.6	99.7	99.7	99.6	98.3	98.2	99.8	98	98.1	98.3	23
<i>B. anthracis</i>	NR-3838	This Work	99.93	99.91	99.97	99.91	99.91	99.92	99.89	100	100	99.6	99.8	99.8	99.8	99.8	99.8	98.2	98.1	99.8	97.9	98	98.2	22
<i>B. anthracis</i>	NR-411	This Work	99.91	99.91	99.93	99.93	99.94	99.92	99.90	99.94	100	99.6	99.8	99.8	99.8	99.8	99.8	98.2	98.1	99.8	97.9	98	98.2	23
<i>B. anthracis</i>	NR-41	This Work	99.92	99.91	99.93	99.92	99.91	99.93	99.88	99.93	99.91	100	99.5	99.6	99.5	99.5	99.5	98.3	98.1	99.8	98	98	98.2	23
<i>B. anthracis</i>	NR-46	This Work	99.89	99.85	99.89	99.90	99.90	99.89	99.87	99.91	99.88	99.87	100	99.7	99.7	99.7	99.7	98.6	98.5	99.9	98.2	98.3	98.6	23
<i>B. anthracis</i>	NR-51483	This Work	99.91	99.91	99.95	99.91	99.93	99.92	99.90	99.93	99.95	99.94	99.88	100	99.7	99.7	99.7	98.2	98.1	99.8	97.9	98	98.2	23
<i>B. anthracis</i>	NR-51484	This Work	99.92	99.95	99.94	99.91	99.89	99.91	99.89	99.92	99.93	99.92	99.89	99.90	100	99.9	100	98.2	98.1	99.8	97.9	98	98.2	23
<i>B. anthracis</i>	NR-51485	This Work	99.91	99.95	99.93	99.91	99.87	99.94	99.91	99.92	99.91	99.91	99.86	99.92	99.96	100	99.9	98.2	98.1	99.8	97.9	98	98.2	23
<i>B. anthracis</i>	NR-51486	This Work	99.92	99.94	99.93	99.92	99.86	99.91	99.89	99.94	99.89	99.91	99.85	99.93	99.98	99.94	100	98.2	98.1	99.8	97.9	98	98.2	23
<i>B. mycooides</i>	ATCC 6462 ^T	CP009692.1	89.35	89.47	89.37	89.37	89.49	89.34	89.27	89.39	89.40	89.45	89.45	89.38	89.31	89.34	89.41	100	96.8	93.2	89.9	78.7	78.8	27
<i>B. mycooides</i>	ATCC 6462 ^T	This Work	89.24	89.47	89.38	89.37	89.53	89.40	89.31	89.33	89.35	89.40	89.47	89.38	89.38	89.41	89.36	99.94	100	92.3	88.8	78	78.4	25
<i>B. mycooides</i>	NR-613	This Work	89.73	89.86	89.69	89.88	89.89	89.81	89.88	89.84	89.83	89.88	89.85	89.87	89.86	89.86	89.71	99.08	99.13	100	87.7	79.7	80.7	23
<i>B. cereus</i>	NR-22171	This Work	89.28	89.46	89.45	89.38	89.42	89.33	89.32	89.34	89.38	89.37	89.47	89.29	89.32	89.31	89.32	98.81	98.65	98.33	100	79.8	79	23
<i>B. weihenstephanensis</i>	NBRC 101238 ^T	BAUY01	89.27	89.38	89.31	89.32	89.53	89.30	89.36	89.38	89.31	89.37	89.45	89.28	89.35	89.31	89.29	97.71	97.62	97.54	97.78	100	84.4	22
<i>B. weihenstephanensis</i>	KBAB4	CP000903.1	89.31	89.37	89.40	89.45	89.51	89.48	89.39	89.41	89.37	89.51	89.53	89.36	89.35	89.34	89.34	97.53	97.52	97.58	97.75	98.33	100	27
<i>Staphylococcus aureus</i> subsp. <i>aureus</i>	DSM 20231 ^T	CP011526.1	67.29	67.62	67.65	67.06	67.27	67.22	67.26	67.61	67.47	67.33	67.62	67.38	67.49	67.58	67.50	68.13	67.27	67.74	67.52	67.80	68.00	100

ANI	Interpretation	dDDH
98.0 - 100	Same species and subspecies	80.0 - 100
96.5 - 97.999	Same species, different subspecies	70.0 - 79.9
92.0 - 96.499	Different species	50.0 - 69.9
85.0 - 91.999		30.0 - 49.9
0.0 - 84.999		0.0 - 29.9

Type Strain
Non-Type Strain
Subspecies Circumscription
Species Circumscription

B. anthracis and B. mycooides

Species/Subspecies	Strain	Source	A0488 ^T	NR-1041	NR-1202	NR-1389	NR-1408	NR-36073	NR-36091	NR-3838	NR-411	NR-41	NR-46	NR-51483	NR-51484	NR-51485	NR-51486	ATCC 6462 ^T	ATCC 6462 ^T	NR-613	NR-22171	NBRC 101238 ^T	KBAB4	DSM 20231 ^T
<i>B. anthracis</i>	A0488 ^T	ABJC01	100	99.8	99.7	99.7	99.6	99.7	99.6	99.7	99.8	99.5	99.6	99.6	99.8	99.8	99.8	38.2	38	39.7	37.8	37.9	38.2	28
<i>B. anthracis</i>	NR-1041	This Work	100	100	99.8	99.8	99.7	99.7	99.7	99.8	99.8	99.5	99.7	99.7	99.8	99.9	99.8	38.3	38.1	39.8	38	38	38.2	23
<i>B. anthracis</i>	NR-1202	This Work			100	99.8	99.8	99.8	99.7	99.9	99.9	99.6	99.8	99.9	99.8	99.8	99.8	38.2	38.1	39.8	37.9	37.9	38.2	23
<i>B. anthracis</i>	NR-1389	This Work				100	99.6	99.9	99.8	99.8	99.8	99.7	99.7	99.7	99.8	99.8	99.8	38.3	38.1	39.8	37.9	38	38.2	23
<i>B. anthracis</i>	NR-1408	This Work					100	99.7	99.6	99.8	99.8	99.5	99.7	99.7	99.8	99.8	99.8	38.3	38.1	39.8	37.9	38.5	38.8	23
<i>B. anthracis</i>	NR-36073	This Work						100	99.8	99.8	99.7	99.7	99.7	99.7	99.8	99.8	99.8	38.3	38.1	39.8	37.9	38.1	38.3	22
<i>B. anthracis</i>	NR-36091	This Work							100	99.7	99.7	99.6	99.7	99.6	99.8	99.8	99.8	38.2	38.1	39.8	37.9	38.1	38.3	23
<i>B. anthracis</i>	NR-3838	This Work								100	100	99.6	99.8	99.8	99.8	99.8	99.8	38.2	38.1	39.8	37.9	38	38.2	22
<i>B. anthracis</i>	NR-411	This Work									100	99.6	99.8	99.8	99.8	99.8	99.8	38.2	38.1	39.8	37.9	38	38.2	23
<i>B. anthracis</i>	NR-41	This Work										100	99.5	99.6	99.7	99.7	99.7	38.2	38.1	39.8	37.9	38	38.2	23
<i>B. anthracis</i>	NR-46	This Work											100	99.7	99.7	99.7	99.7	38.6	38.5	39.9	38.2	38.3	38.6	23
<i>B. anthracis</i>	NR-51483	This Work												100	99.7	99.7	99.7	38.2	38.1	39.8	37.9	38	38.2	23
<i>B. anthracis</i>	NR-51484	This Work													100	99.9	100	38.2	38.1	39.8	37.9	38	38.2	23
<i>B. anthracis</i>	NR-51485	This Work														100	99.9	38.2	38.1	39.8	37.9	38	38.2	23
<i>B. anthracis</i>	NR-51486	This Work															100	38.2	38.1	39.8	37.9	38	38.2	23
<i>B. mycooides</i>	ATCC 6462 ^T	CP009692.1																100	96.8	93.2	89.9	78.7	78.8	27
<i>B. mycooides</i>	ATCC 6462 ^T	This Work																	100	92.3	88.8	78	78.4	25
<i>B. mycooides</i>	NR-613	This Work																		100	87.7	79.7	80.7	23
<i>B. cereus</i>	NR-22171	This Work																			100	79.8	79	23
<i>B. weihenstephanensis</i>	NBRC 101238 ^T	BAUY01																				100	84.4	22
<i>B. weihenstephanensis</i>	KBAB4	CP000903.1																					100	27
<i>Staphylococcus aureus</i> subsp. <i>aureus</i>	DSM 20231 ^T	CP011526.1																						100

dDDH

ANI	Interpretation	dDDH
98.0 - 100	Same species and subspecies	80.0 - 100
96.5 - 97.999	Same species, different subspecies	70.0 - 79.9
92.0 - 96.499	Different species	50.0 - 69.9
85.0 - 91.999		30.0 - 49.9
0.0 - 84.999		0.0 - 29.9

- Type Strain
- Non-Type Strain
- Subspecies Circumscription
- Species Circumscription



B. anthracis and B. mycooides

Species/Subspecies	Strain	Source	A0488 ^T	NR-1041	NR-1202	NR-1389	NR-1408	NR-36073	NR-36091	NR-3838	NR-411	NR-41	NR-46	NR-51483	NR-51484	NR-51485	NR-51486	ATCC 6462 ^T	ATCC 6462 ^T	NR-613	NR-22171	NBRC 101238 ^T	KBAB4	DSM 20231 ^T
<i>B. anthracis</i>	A0488 ^T	ABJC01	100																					
<i>B. anthracis</i>	NR-1041	This Work	99.92	100																				
<i>B. anthracis</i>	NR-1202	This Work	99.93	99.92	100																			
<i>B. anthracis</i>	NR-1389	This Work	99.89	99.91	99.92	100																		
<i>B. anthracis</i>	NR-1408	This Work	99.87	99.89	99.92	99.92	100																	
<i>B. anthracis</i>	NR-36073	This Work	99.90	99.93	99.90	99.93	99.87	100																
<i>B. anthracis</i>	NR-36091	This Work	99.89	99.88	99.88	99.90	99.88	99.91	100															
<i>B. anthracis</i>	NR-3838	This Work	99.93	99.91	99.97	99.91	99.91	99.92	99.89	100														
<i>B. anthracis</i>	NR-411	This Work	99.91	99.91	99.93	99.93	99.94	99.92	99.90	99.94	100													
<i>B. anthracis</i>	NR-41	This Work	99.92	99.91	99.93	99.92	99.91	99.93	99.88	99.93	99.91	100												
<i>B. anthracis</i>	NR-46	This Work	99.89	99.85	99.89	99.90	99.90	99.89	99.87	99.91	99.88	99.87	100											
<i>B. anthracis</i>	NR-51483	This Work	99.91	99.91	99.95	99.95	99.95	99.95	99.95	99.95	99.95	99.95	99.95	100										
<i>B. anthracis</i>	NR-51484	This Work	99.92	99.95	99.94	99.94	99.94	99.94	99.94	99.94	99.94	99.94	99.94	99.94	100									
<i>B. anthracis</i>	NR-51485	This Work	99.91	99.95	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	100								
<i>B. anthracis</i>	NR-51486	This Work	99.92	99.94	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	99.93	100							
<i>B. mycooides</i>	ATCC 6462 ^T	CP009692.1	89.35	89.47	89.37	89.37	89.37	89.37	89.37	89.37	89.37	89.37	89.37	89.37	89.37	89.37	89.37	100						
<i>B. mycooides</i>	ATCC 6462 ^T	This Work	89.24	89.47	89.38	89.38	89.38	89.38	89.38	89.38	89.38	89.38	89.38	89.38	89.38	89.38	89.38	99.94	100					
<i>B. mycooides</i>	NR-613	This Work	89.73	89.86	89.69	89.88	89.89	89.81	89.88	89.84	89.83	89.88	89.85	89.87	89.86	89.86	89.71	99.08	99.13	100				
<i>B. cereus</i>	NR-22171	This Work	89.28	89.46	89.45	89.38	89.42	89.33	89.32	89.34	89.38	89.37	89.47	89.29	89.32	89.31	89.32	98.81	98.65	98.33	100			
<i>B. weihenstephanensis</i>	NBRC 101238 ^T	BAUY01	89.27	89.38	89.31	89.32	89.53	89.30	89.36	89.38	89.31	89.37	89.45	89.28	89.35	89.31	89.29	97.71	97.62	97.54	97.78	100		
<i>B. weihenstephanensis</i>	KBAB4	CP000903.1	89.31	89.37	89.40	89.45	89.51	89.48	89.39	89.41	89.37	89.51	89.53	89.36	89.35	89.34	89.34	97.53	97.52	97.58	97.75	98.33	100	
<i>Staphylococcus aureus</i> subsp. <i>aureus</i>	DSM 20231 ^T	CP011526.1	67.29	67.62	67.65	67.06	67.27	67.22	67.26	67.61	67.47	67.33	67.62	67.38	67.49	67.58	67.50	68.13	67.27	67.74	67.52	67.80	68.00	100

ANI

ANI	Interpretation	dDDH
98.0 - 100	Same species and subspecies	80.0 - 100
96.5 - 97.999	Same species, different subspecies	70.0 - 79.9
92.0 - 96.499	Different species	50.0 - 69.9
85.0 - 91.999		30.0 - 49.9
0.0 - 84.999		0.0 - 29.9

Type Strain
Non-Type Strain
Subspecies Circumscription
Species Circumscription



Bacillus cereus and Bacillus thuringiensis

Species/Subspecies	Strain	Source	ATCC 14579 ^T	ATCC 13367	ATCC 11778	ATCC 33019	ATCC 10792 ^T	ATCC 10792 ^T	NR-28583	ATCC 246	NR-610	ATCC 33679	ATCC 10876	ATCC 7039	ATCC 9592	ATCC 700872	ATCC 35646	ATCC 19266
<i>B. cereus</i>	ATCC 14579 ^T	NC_004722.1	100	82.8	82.5	81.8	71.1	71.3	71	74.1	73.2	73.2	72.5	65.2	65.2	65.8	65.7	65.2
<i>B. thuringiensis</i>	ATCC 13367	This Work	98.15	100	89.4	87.3	69.8	69.5	69.6	75.2	73.2	73.3	73.6	64.	Subspecies 1			64.5
<i>B. cereus</i>	ATCC 11778	This Work	98.06	98.84	100	92.1	71.1	71.2	71	77.7	76.8	76.7	75.7	66.	Subspecies 1			66.9
<i>B. cereus</i>	ATCC 33019	This Work	97.96	98.67	99.10	100	70.7	70.8	70.5	77.5	75.7	75.5	74.8	65.9	66	66.3	66	66.1
<i>B. thuringiensis</i>	ATCC 10792 ^T	CP020754.1	96.65	96.46	96.58	96.69	100	99.1	82.6	68.1	67	67.1	67.2	68.9	69	69.3	68.9	68.8
<i>B. thuringiensis</i>	ATCC 10792 ^T	This Work	96.60	96.44	96.73	96.75	99.83	100	82.2	67.9	67.2	67.3	66.8	68.	Subspecies 2			68
<i>B. thuringiensis</i>	NR-28583	This Work	96.62	96.38	96.66	96.65	98.02	98.01	100	68.7	68	68	67.9	70	70.2	70.3	69.8	69.9
<i>B. cereus</i>	ATCC 246	This Work	97.13	97.28	97.51	97.54	96.21	96.14	96.35	100	87.9	87.7	81.6	65.1	65.2	65.7	65.4	64.7
<i>B. thuringiensis</i>	NR-610	This Work	97.04	97.16	97.42	97.32	96.15	96.21	96.20	98.80	100	96.7	79.5	63.	Subspecies 3			63.6
<i>B. thuringiensis</i>	ATCC 33679	This Work	97.04	97.13	97.41	97.34	96.17	96.23	96.33	98.78	99.92	100	79.6	63.	Subspecies 3			63.5
<i>B. cereus</i>	ATCC 10876	This Work	96.91	97.11	97.32	97.26	96.10	96.04	96.31	98.04	97.85	97.86	100	64.2	64.3	65.2	64.8	63.5
<i>B. cereus</i>	ATCC 7039	This Work	95.88	95.78	96.04	96.11	96.42	96.25	96.49	95.82	95.66	95.63	95.69	100	95.6	93.6	92.6	92.5
<i>B. cereus</i>	ATCC 9592	This Work	95.83	95.82	96.06	96.04	96.42	96.27	96.59	95.80	95.65	95.62	95.70	99.92	100	94	93	92.8
<i>B. thuringiensis</i>	ATCC 700872	This Work	95.97	95.91	96.10	96.04	96.42	96.32	96.50	95.93	95.74	95.67	95.89	99.27	99.32	100	99.7	91.8
<i>B. thuringiensis</i>	ATCC 35646	This Work	95.91	95.80	96.11	96.02	96.38	96.21	96.53	95.83	95.65	95.60	95.75	99.20	99.16	99.93	100	91.1
<i>B. thuringiensis</i>	ATCC 19266	This Work	95.83	95.73	96.15	96.00	96.33	96.28	96.53	95.83	95.61	95.61	95.57	99.12	99.11	99.06	99.02	100

ANI	Interpretation	dDDH
98.0 - 100	Same species and subspecies	80.0 - 100
96.5 - 97.999	Same species, different subspecies	70.0 - 79.9
92.0 - 96.499	Different species	50.0 - 69.9
85.0 - 91.999		30.0 - 49.9
0.0 - 84.999		0.0 - 29.9

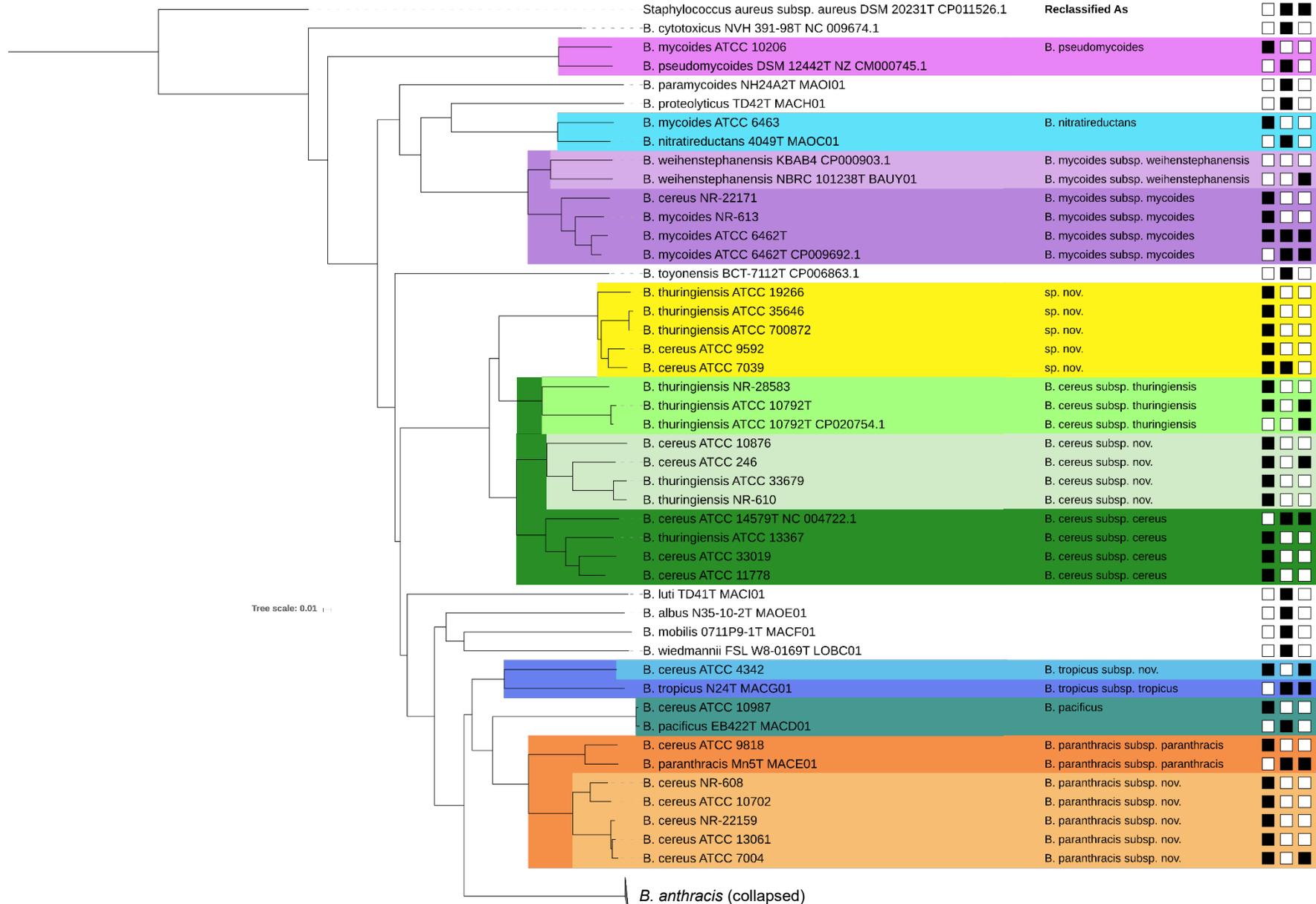
Novel Species

Other BcG Strains

Species/Subspecies	Strain	Source	BCT-7112 ^T	FSL W8-0169 ^T	Mn5 ^T	ATCC 9818	ATCC 10702	NR-608	NR-22159	ATCC 13061	ATCC 7004	EB422 ^T	ATCC 10987	N24 ^T	ATCC 4342	N35-10-2 ^T	0711P9-1 ^T	TD41 ^T	TD42 ^T	4049 ^T	ATCC 6463	NH24A2 ^T	DSM 12442 ^T	ATCC 10206	NVH 391-98 ^T
<i>B. toyonensis</i>	BCT-7112 ^T	CP006863.1	100	44	42.8	42.8	42.7	43.2	42.7	42.8	42.8	42.5	42.5	43.2	43.1	42.6	42.7	43.6	42.4	42	42	39.3	27.3	27.3	26.2
<i>B. wiedmannii</i>	FSL W8-0169 ^T	LOBC01	91.30	100	50.9	50.8	51.3	51.8	51.1	51.1	51.1	51.6	51.7	52	51.6	54.1	58.1	45.2	40.9	42.2	42.2	37.7	27	27	25.5
<i>B. paranthracis</i>	Mn5 ^T	MACE01	90.98	93.18	100	91.7	78.6	79	77.4	77.2	77.4	65.4	65.3	61.9	64.3	52.7	49.9	44.6	39.1	39.7	39.6	36.6	26.8	26.9	25.4
<i>B. cereus</i>	ATCC 9818	This Work	90.92	93.14	99.08	100	78.5	79	77.8	77.5	77.2	65.7	65.5	61.6	64.4	52.8	49.9	44.7	39.2	39.6	39.6	36.7	26.8	26.9	25.6
<i>B. cereus</i>	ATCC 10702	This Work	90.81	93.15	97.58	97.60	100	95.2	90	89.9	89.4	66	67	62.6	65.4	52.1	51	44.3	39.3	39.6	39.6	36.8	27	27.4	25.5
<i>B. cereus</i>	NR-608	This Work	90.95	93.33	97.69	97.67	99.93	100	90.3	90.2	89.8	66.6	67.6	63.1	65.8	52.6	51.6	44.9	39.9	40.1	40.1	37.3	27.2	27.4	25.8
<i>B. cereus</i>	NR-22159	This Work	90.88	93.24	97.47	97.55	98.90	98.86	100	99	98.6	66.7	66.8	62.6	65.3	51.9	50.7	44.3	39.2	39.6	39.5	36.6	26.8	27	25.5
<i>B. cereus</i>	ATCC 13061	This Work	90.84	93.23	97.46	97.54	98.87	98.82	99.78	100	98.9	66.4	66.7	62.6	65.1	52	50.2	44.3	39.2	39.6	39.7	36.7	26.8	27	25.5
<i>B. cereus</i>	ATCC 7004	This Work	90.89	93.16	97.48	97.36	98.79	98.79	99.80	99.82	100	66.4	66.5	62.4	65.1	52	50	44.3	39.2	39.6	39.6	36.7	26.9	27.1	25.5
<i>B. pacificus</i>	EB422 ^T	MACD01	90.77	93.23	95.94	95.82	95.88	96.01	96.01	96.00	95.93	100	99.3	60	59.7	52.4	51.7	44.3	39.2	39.9	39.8	36.7	26.8	26.8	25.6
<i>B. cereus</i>	ATCC 10987	This Work	90.77	93.24	95.82	95.87	96.08	96.15	96.06	96.07	96.00	99.86	100	60.3	59.9	52.6	52	44.4	39.2	39.9	39.9	36.7	26.8	26.9	25.8
<i>B. tropicus</i>	N24 ^T	MACG01	91.09	93.36	95.33	95.27	95.37	95.51	95.48	95.44	95.40	94.97	95.04	100	70.7	52.2	50.7	45.3	39.8	39.9	40	37	26.9	26.9	25.5
<i>B. cereus</i>	ATCC 4342	This Work	90.96	93.30	95.66	95.68	95.76	95.90	95.84	95.85	95.80	94.89	94.94	96.70	100	53.1	50.5	44.8	39.6	40	39.9	36.9	27	27.3	25.6
<i>B. albus</i>	N35-10-2 ^T	MAOE01	90.95	93.77	93.45	93.46	93.32	93.54	93.35	93.42	93.28	93.45	93.44	93.37	93.58	100	54.4	44.1	39.1	40.1	40.1	36.6	26.9	26.8	25.4
<i>B. mobilis</i>	0711P9-1 ^T	MACF01	90.90	94.65	92.96	92.93	93.18	93.37	93.16	92.95	92.97	93.28	93.39	93.10	93.07	93.94	100	44.1	39.9	41.3	41.3	37.3	26.7	26.7	25.5
<i>B. luti</i>	TD41 ^T	MACI01	91.13	91.62	91.46	91.50	91.51	91.66	91.50	91.43	91.45	91.44	91.43	91.84	91.64	91.35	91.41	100	40.7	40	40.1	38	26.9	26.9	25.5
<i>B. proteolyticus</i>	TD42 ^T	MACH01	90.71	90.26	89.72	89.66	89.73	89.89	89.75	89.67	89.67	89.66	89.80	89.81	89.74	89.55	89.86	90.15	100	58.4	57.6	43.8	28.1	28	25.6
<i>B. nitratireductans</i>	4049 ^T	MAOC01	90.66	90.74	90.03	89.95	89.85	90.10	89.99	89.91	89.91	90.03	90.06	90.07	90.06	90.07	90.47	90.10	94.75	100	86.3	43.7	27.7	27.6	25.6
<i>B. mycooides</i>	ATCC 6463	This Work	90.64	90.70	89.80	89.81	89.85	90.04	89.80	89.85	89.93	89.93	89.98	90.05	89.93	90.03	90.42	90.11	94.59	98.55	100	43.6	28	27.7	25.7
<i>B. paramycooides</i>	NH24A2 ^T	MAOI01	89.66	89.17	88.73	88.84	88.79	88.91	88.94	88.89	88.76	88.74	88.86	89.02	88.69	88.73	89.05	89.15	91.20	91.25	91.25	100	28.1	28	25.9
<i>B. pseudomycooides</i>	DSM 12442 ^T	NZ_CM000745.1	82.65	82.50	82.32	82.36	82.23	82.38	82.37	82.39	82.40	82.32	82.30	82.37	82.47	82.33	82.18	82.53	83.17	82.95	83.11	83.41	100	86.5	28
<i>B. mycooides</i>	ATCC 10206	This Work	82.61	82.21	82.23	82.34	82.16	82.49	82.29	82.29	82.22	82.33	82.32	82.24	82.29	82.24	82.22	82.45	82.92	82.80	82.88	83.14	98.59	100	27.9
<i>B. cytotoxicus</i>	NVH 391-98 ^T	NC_009674.1	81.44	81.29	81.15	81.15	81.24	81.30	81.38	81.22	81.23	81.32	81.39	81.27	81.16	81.31	81.27	81.37	81.31	81.58	81.48	81.86	83.20	83.07	100

ANI	Interpretation	dDDH
98.0 - 100	Same species and subspecies	80.0 - 100
96.5 - 97.999	Same species, different subspecies	70.0 - 79.9
92.0 - 96.499		50.0 - 69.9
85.0 - 91.999	Different species	30.0 - 49.9
0.0 - 84.999		0.0 - 29.9

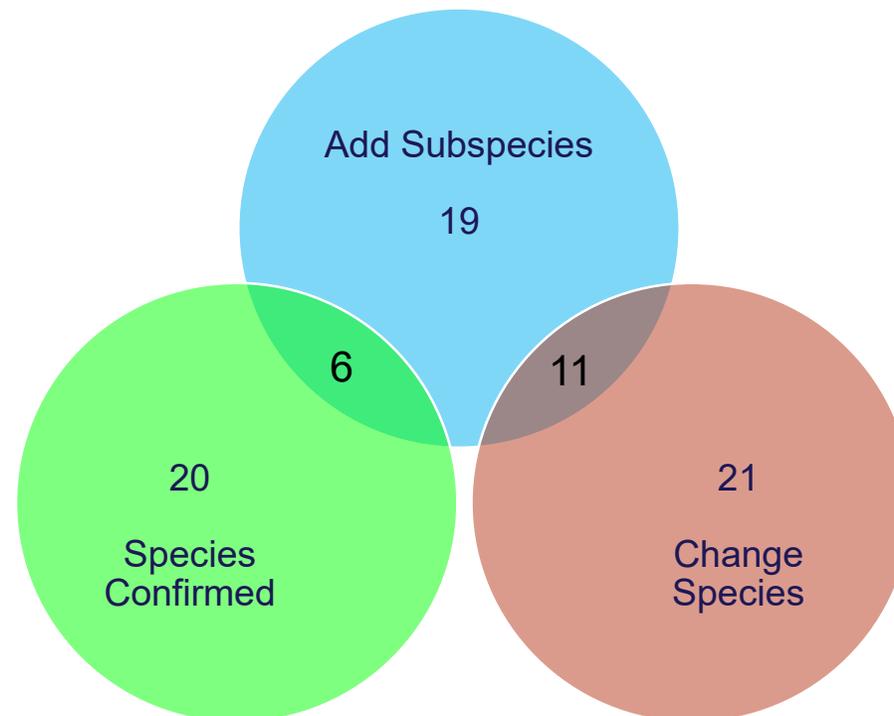
BcG Phylogenomic Tree



Summary of Strains from ATCC and BEI Resources

41 in-house strains sequenced

- 27 items should receive some form of name change
- 23 transferred to a different species (e.g., *B. cereus* to *B. pacificus*)
- 19 added to a subspecies for additional specificity



Publication of Results

- **Poster: ASM Microbe 2019, San Francisco, CA**

- EEB02 - Microbial Evolution and Comparative Genomics 1
- Poster Board Number: **FRIDAY - EEB-488**
- Friday, June 21, 2019
 - 11 AM - 12 PM PT
 - 4 PM - 5 PM PT

- **Scientific Manuscript Submitted: *Int J Syst Evol Microbiol***

- Phylogenomic Reclassification of ATCC *Bacillus* Strains and Various Taxa within the Genus *Bacillus*
 - Marco A. Riojas, Andrew M. Frank, Samantha L. Fenn, Stephen King, Sonia Brower, Manzour Hernando Hazbón

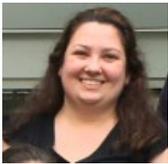
Acknowledgements



■ Joe Leonelli, PhD



■ Briana Benton



■ Samantha Fenn



■ Andrew Frank



■ Manzour Hazbón, PhD



■ Stephen King



■ Juan Lopera, PhD



■ Anna McCluskey



National Institute of Allergy and Infectious Diseases (NIAID)

The following reagents were obtained through BEI Resources, NIAID, NIH: NR-41, NR-46, NR-411, NR-608, NR-610, NR-613, NR-1041, NR-1202, NR-1389, NR-1408, NR-3838, NR-22159, NR-22171, NR-28583, NR-36073, NR-36091, NR-51483, NR-51484, NR-51485, and NR-51486.



Questions?

Credible Leads to Incredible™

