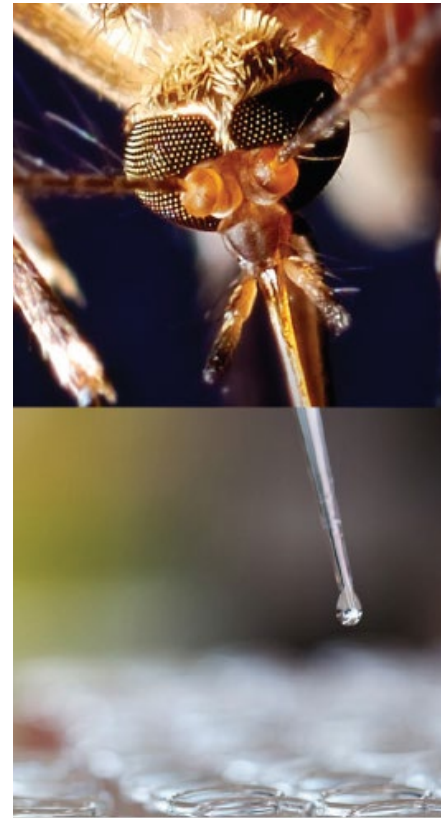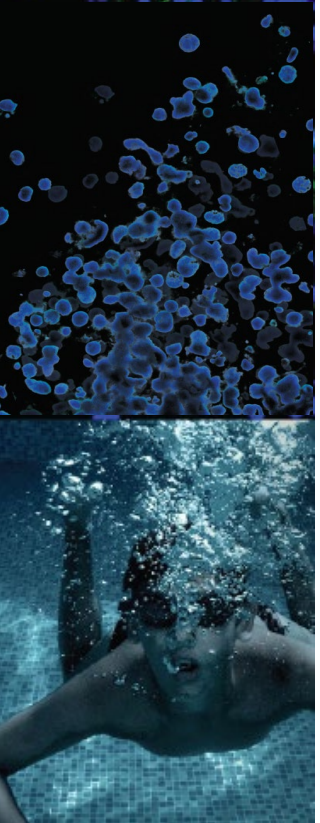# Genomic Data Quality

Connecting the Dots Between Bioinformatics and Physical Materials

**Jonathan Jacobs, PhD**
Senior Director, Bioinformatics
Sequencing & Bioinformatics Center
ATCC

Credible Leads to Incredible™

# About ATCC

- Founded in 1925

- 501(c)(3) not-for-profit organization

- World's largest, most diverse biorepository

- Quality Accreditation by multiple industry standards
  – ISO 9001 Certified
  – ISO 13485 Certified
  – ISO/IEC 17025 Accredited
  – ISO 17034 Accredited

- Standards development partner with multiple industry working groups
  – ANSI Standards Working Groups
  – AOAC International Working Group
  – IMMSA/NIST Microbiome Standards

- Global supplier of authenticated cell lines, microorganisms, and molecular standards

- Sales and Distribution to 150+ countries

- Full Talented team of 500+ employees

Thousands of authenticated biomaterials
- 5,000+ cell lines & primary tissue
- 2,500+ viruses
- 9,500+ bacteria
- 38,000+ fungi and protists

Visit **atcc.org**

# Genomics data quality

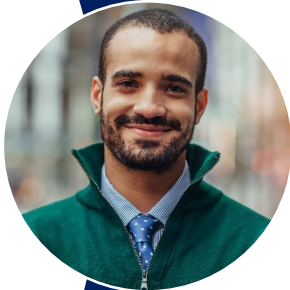*Connecting the dots between bioinformatics and physical materials*

- Review challenges associated with genomics data quality and authenticity

- Discuss *why* ATCC is committed to **providing reference-quality genomes** for our materials

- Discuss our current efforts to produce standardized genomics reference data

- Explore the ATCC Genome Portal

- Explore the ATCC Cell Line Land
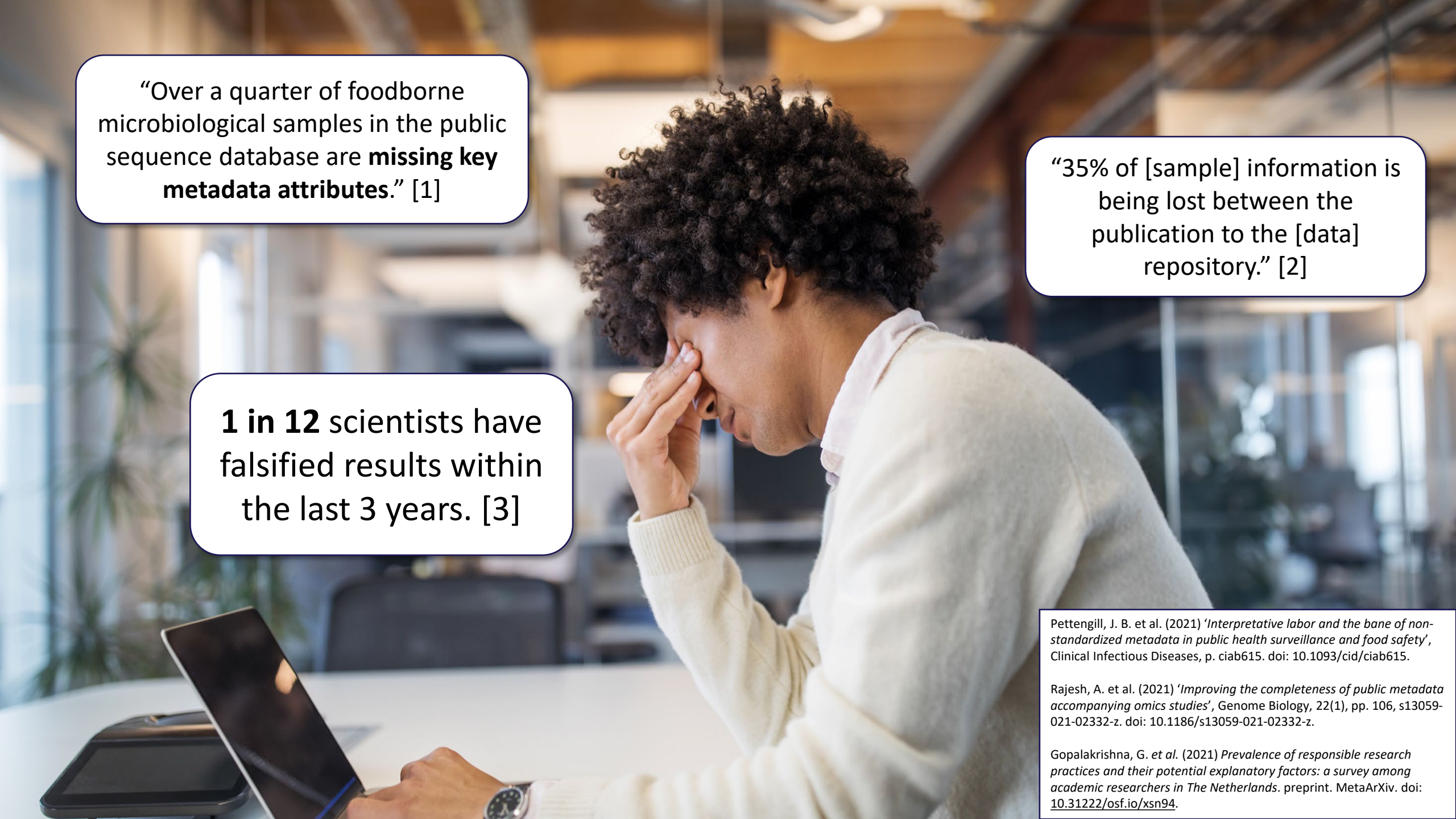
ATCC

# Challenges stemming from poor data quality...

"***Finding the right cell lines*** for my research is a challenge."

"Many cell types are ***not good models*** for the disease I'm studying."

"Pre-existing results are difficult to reproduce and often ***not reproducible***."

"Over a quarter of foodborne microbiological samples in the public sequence database are **missing key metadata attributes**." [1]

"35% of [sample] information is being lost between the publication to the [data] repository." [2]

**1 in 12** scientists have falsified results within the last 3 years. [3]

Pettengill, J. B. et al. (2021) 'Interpretative labor and the bane of non-standardized metadata in public health surveillance and food safety', Clinical Infectious Diseases, p. ciab615. doi: 10.1093/cid/ciab615.

Rajesh, A. et al. (2021) 'Improving the completeness of public metadata accompanying omics studies', Genome Biology, 22(1), pp. 106, s13059-021-02332-z. doi: 10.1186/s13059-021-02332-z.

Gopalakrishna, G. et al. (2021) Prevalence of responsible research practices and their potential explanatory factors: a survey among academic researchers in The Netherlands. preprint. MetaArXiv. doi: 10.31222/osf.io/xsn94.

# Fake data was first discovered in GenBank in 1997

> *"Mr. Hajra, former graduate student, University of Michigan, engaged in scientific misconduct by falsifying and fabricating research data in five published research papers, two published review articles, one submitted but unpublished paper, in his doctoral dissertation, and in a **submission to the GenBank computer data base.**" – The Federal Register, v62, n135 (**1997**)*

Federal Register / Vol. 62, No. 135 / Tuesday, July 15, 1997 / Notices    37921

author of the application is identified and that person's role in the project is ...

years. In the event a consortium of applicants is proposed, the project history of prior joint work should be provided. The previous Federal assistance is identified by project number, Federal agency, and grants or contracting officer. 25 points

*Components of a Complete Application*

A complete application consists of the following items in this order:
1. Application for Federal Assistance (Standard Form 424, REV 4–88);
2. Budget Information—Non-construction Programs (Standard Form 424A, REV 4–88);
3. Assurances—Non-construction Programs (Standard Form 424B, REV 4–88);
4. Table of Contents;

Dated: July 9, 1997.

**David F. Garrison,**

*Principal Deputy Assistant Secretary for Planning and Evaluation.*

[FR Doc. 97–18528 Filed 7–14–97; 8:45 am]

BILLING CODE 4151–04–M

## DEPARTMENT OF HEALTH AND HUMAN SERVICES

### Office of the Secretary

### Findings of Scientific Misconduct

**AGENCY:** Office of the Secretary, HHS.

**ACTION:** Notice.

**SUMMARY:** Notice is hereby given that the Office of Research Integrity (ORI) has made a final finding of scientific misconduct in the following case:

*Amitav Hajra, University of Michigan*: Based upon a report from the University of Michigan, information obtained by the Office of Research Integrity (ORI) during its oversight review, and Mr. Hajra's own admission, ORI found that Mr. Hajra, former graduate student, University of Michigan, engaged in scientific misconduct by falsifying and fabricating research data in five published research papers, two published review articles, one submitted but unpublished paper, in his doctoral dissertation, and in a submission to the GenBank computer data base. Mr. Hajra's doctoral training and research was supported by two Public Health Service (PHS) grants, and his experiments were conducted at and submitted for publication from the

• Wijmenga, C., Gregory, P.E., Hajra, A., Schröck, E., Ried, T., Eils, R., Liu, P.P., and Collins, F.S. "Core binding factor β-smooth muscle myosin heavy chain chimeric protein involved in acute myeloid leukemia forms unusual nuclear rod-like structures in transformed NIH 3T3 cells." *Proc. Natl. Acad. Sci.* USA 93(4):1630–1635, 1996; and

• Liu, P.P., Wijmenga, C., Hajra, A., Blake, T.B., Kelley, C.A., Adelstein, R.S., Bagg, A., Rector, J., Cotelingham, J., Willman, C.L., and Collins, F.S. "Identification of the chimeric protein product of the CBFB-MYH11 fusion gene in inv(16) leukemia cells." *Genes, Chromosomes, and Cancer* 16:77–87, 1996 (Erratum in *Genes, Chromosomes, and Cancer* 18(1):71, 1997).

Mr. Hajra included fabricated and falsified data in the following review articles:

• Hajra, A., Liu, P.P., and Collins, F.S. "Transforming properties of the leukemic Inv(16) fusion gene CBFB–MYH11." In Molecular Aspects of Myeloid Stem Cell Development in Current Topics in Microbiology and Immunology (L. Wolff and A.S. Perkins, Eds.) 211:289–298, 1996 (Review). Berlin and New York: Springer-Verlag; and

• Liu, P.P., Hajra, A., Wijmenga, C., and Collins, F.S. "Molecular pathogenesis of the chromosome 16 inversion in the M4Eo subtype of acute myeloid leukemia." *Blood* 85:2289–2302, 1995 (Review).

Mr. Hajra submitted a fabricated nucleotide sequence in computer data

6

# 24 years later, this falsified data still being cited...

# After 42 citations... the data is still in GenBank...

# Falsified sequencing data to support a false phylogeny



Biochemical Systematics and Ecology
Volume 96, June 2021, 104263

**Scientific data laundering: Chimeric mitogenomes of a sparrowhawk and a nightjar covered-up by forged phylogenies**

George Sangster [a] ✉, Jolanda A. Luksenburg [b, c] ✉

Show more ∨

Outline | + Add to Mendeley | ⤳ Share | 🗨 Cite

*"The evidence indicates that Liu et al. (2017) published phylogenies that were not based on existing data **but were fabricated to reflect preconceived ideas** about phylogenetic relationships."* – Sangster & Luksenburg (2021)

Liu and colleagues in a paper in *Biochemical Systematics and Ecology* in 2017 is not an authentic sequence of this species but represents a chimera of three different species (a

Sangster, G. and Luksenburg, J.A. (2021) 'Scientific data laundering: Chimeric mitogenomes of a sparrowhawk and a nightjar covered-up by forged phylogenies', *Biochemical Systematics and Ecology*, 96, p. 104263. doi:10.1016/j.bse.2021.104263.

# Unfortunately, the falsified mitogenome is still in GenBank...



**UNVERIFIED: Accipiter gularis mitochondrion sequence**

GenBank: KX585864.1

FASTA   Graphics

Go to: ☑

```
LOCUS       KX585864              17918 bp    DNA     linear   VRT 31-AUG-2021
DEFINITION  UNVERIFIED: Accipiter gularis mitochondrion sequence.
ACCESSION   KX585864
VERSION     KX585864.1
KEYWORDS    UNVERIFIED; UNVERIFIED_ORGANISM.
SOURCE      mitochondrion Accipiter gularis (Japanese sparrowhawk)
  ORGANISM  Accipiter gularis
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Archelosauria; Archosauria; Dinosauria; Saurischia; Theropoda;
            Coelurosauria; Aves; Neognathae; Accipitriformes; Accipitridae;
            Accipitrinae; Accipiter.
REFERENCE   1  (bases 1 to 17918)
  AUTHORS   Liu,G.
  TITLE     The complete mtDNA of Accipiter gularis
  JOURNAL   Unpublished
REFERENCE   2  (bases 1 to 17918)
  AUTHORS   Liu,G.
  TITLE     Direct Submission
  JOURNAL   Submitted (21-JUL-2016) School of life science, Anhui Medical
            University, 81 Meishan Rd, Hefei, Anhui 230032, China
COMMENT     GenBank staff is unable to verify source organism and sequence
            and/or annotation provided by the submitter.
FEATURES             Location/Qualifiers
     source          1..17918
```

Labeled as "Unverified", but the sequence still remains in GenBank and, for example, will come up in a BLAST search...

ATCC

# Intentional falsification is rare… but accidents happen right?

*(>2 million times) …*

# Poor quality genomes result in taxonomic misclassification

*Multiple papers (more than the two listed here) have found widespread misclassification in GenBank*



**~7.8% of genomes misclassified at the species level**

**~4% at the genus level**

**~7% of genomes misclassified at genus *or higher***

# Challenging traceability of most public genomics data



1970     1980     1990     2000     2010     2020

??Transfer    ??Transfer    ??Transfer

Benton

Combs

**2009: Initial Deposit "Draft" Reference Genome ATCC XYZ1**

**3 decades of unknown transfers & research**

**1974: Strain Deposit ATCC XYZ1**

**1995 MTA**

Yarmosh

**2015: Updated Genome "High Quality" Assembly ATCC XYZ1**

**2 decades of lab research**

ATCC

Authenticated source material

**ATCC Genome Portal Assembly**

ATCC

**100% authenticity and traceability**

**2020: Authenticated Genome Reference ATCC XYZ1**

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

| SNPs differences | ATCC Reference |
|---|---|
| Draft genome | >13,000 |
| "High-quality" reference genome | >10,000 |

13

Names Changed To Protect The Innocent.

# A reminder on the growth of public genomics data

1.6B sequences in WGS
232M sequences in GenBank



1.6B WGS sequences

GenBank is doubling in size every 18 months…

232M GenBank sequences



SRA database growth

58,621,760,107,125,660 total bases
25,806,046,433,756,229 open access bases

NCBI's SRA database is over 15 Petabytes in size…

Data curation is a huge challenge

ATCC

# Genomics data quality issues impact many disciplines

**FACTORS**

- Misclassification of sequences

- Chimeric genome assemblies

- Sample contamination

- Sequencing errors

- Mislabeling or data errors

- Data omission

- Data obfuscation

- Intentional misconduct

**Critically Impacted Areas**

- Basic research (hypothesis generation)

- Biodiversity and environmental sciences

- Diagnostics & epidemiology

- Forensics

- Food safety

- Biodefense

- Many other areas…

ATCC

These are not "new" problems.

Many groups have sought solutions.

None, however, have sought to create **Authenticated Genomics Data**

What is
"Authenticated Genomics
Data"?

**Authenticated Genomics Data:**
1.   Traceable to physical materials
2.   Produced with defined quality assurance metrics
3.   Reproducible across multiple tests

# Authenticated genomics data at ATCC

*ATCC is focused on data provenance and closing the reproducibility gap*

**Authenticated Genomics Data**

- Standardized laboratory methods
- Quality Assurance (ISO)
- Traceable to materials in a biorepository
- *Maximum data provenance*
- *Maximum reproducibility*

**Expert Curated Data**

- Standardized metadata
- Standardized biofx methods
- Improved reproducibility
- Less risk, more results
- FAIR data model

**Focused Public Data**

- Improved metadata
- Moderate risk
- Often access-controlled
- Limited scope

GenBank
ENA
SRA
EGA
GEO

**Uncontrolled Public Data**

- Unknown quality
- Missing or non-standard metadata
- Risky to use

ATCC

# Authenticated genomics data at ATCC

*ATCC is focused on data provenance and closing the reproducibility gap*

**RefSeq**
**proGenomes**
**dbGAP**
**BluePrint**
**ICGC**
**TCGA**
**TARGET**
**CCLE**
**GTEx**

## Authenticated Data

- Standardized laboratory methods
- Quality Assurance (ISO)
- Traceable to materials in a biorepository
- *Maximum data provenance*
- *Maximum reproducibility*

## Expert Curated Data

- Standardized metadata
- Standardized biofx methods
- Improved reproducibility
- Less risk, more results
- FAIR data model

## Focused Public Data

- Improved metadata
- Moderate risk
- Often access-controlled
- Limited scope

## Uncontrolled Public Data

- Unknown quality
- Missing or non-standard metadata
- Risky to use

**ATCC**

# Authenticated genomics data at ATCC

*ATCC is focused on data provenance and closing the reproducibility gap*

**OmicSoft Ingenuity (IPA) HGMD**

QIAGEN

**Authenticated Data**
- Standardized laboratory methods
- Quality Assurance (ISO)
- Traceable to materials in a biorepository
- *Maximum data provenance*
- *Maximum reproducibility*

**Expert Curated Data**
- Standardized metadata
- Standardized biofx methods
- Improved reproducibility
- FAIR data model
- *Less risk, more results*

**Focused Public Data**
- Improved metadata
- Moderate risk
- Often access-controlled
- Limited scope

**Uncontrolled Public Data**
- Unknown quality
- Missing or non-standard metadata
- Risky to use

ATCC

# Authenticated genomics data at ATCC

*ATCC is focused on data provenance and closing the reproducibility gap*

**ATCC Genome Portal**
**ATCC Cell Line Land**

**OmicSoft**
**Ingenuity (IPA)**
**HGMD**

**QIAGEN**

**ATCC®**

**Authenticated Data**

- Standardized laboratory methods
- Quality Assurance (ISO)
- Traceable to materials in a biorepository
- *Maximum data provenance*
- *Maximum reproducibility*

**Expert Curated Data**

- Standardized metadata
- Standardized biofx methods
- Improved reproducibility
- FAIR data model
- *Less risk, more results*

**Focused Public Data**

- Improved metadata
- Moderate risk
- Often access-controlled
- Limited scope

**Uncontrolled Public Data**

- Unknown quality
- Missing or non-standard metadata
- Risky to use

**ATCC®**

# The ATCC Genome Portal

Tackling the reproducibility gap in microbial genomics

# ATCC Genome Portal

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

 Download whole-genome sequences and annotations of ATCC materials

 Search for nucleotide sequences or genes within genomes

 View genome assembly metadata and quality metrics

## genomes.atcc.org

**2,522 Authenticated Reference Genomes**

2145 bacteria
221 viruses
155 fungi
1 protist

*New genomes released every month!*

# Authenticated physical material coupled with reference-quality genome sequences

| Isolate | gDNA | NGS and Bioinformation | | |
|---------|------|------------------------|---|---|
| **1** | **2** | **3** | **4** | **5** |
| Growth and QC of isolates | gDNA extraction and QC | **Hybrid Sequencing**<br>Illumina +<br>Oxford Nanopore | QC, Assembly & Annotation | Monthly release of new Authenticated Genomes to the ATCC Genome Portal |

# Sequencing QC – Read trimming/filtering



Only keep high-quality base calls

Only keep long, good quality reads

Quality controlled data

Raw data

# Hybrid genome assembly



**Illumina-only genome assembly**
**150 bp reads**

Long reads mapped to a tangled region creates a resolved bridge
Successively applying bridges resolves the structure of the genome

**Completed hybrid assembly**

Image reproduced from https://github.com/rrwick/Unicycler

# Advantage of hybrid assemblies

| Illumina-only assembly | Hybrid assembly |
|---|---|

*Neisseria meningitidis* (ATCC® 53417™)



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



| Illumina-only assembly | Hybrid assembly |
|---|---|

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



*Streptococcus gordonii* (ATCC® 35105™)

# Quality of ATCC Genome Portal assemblies



**All bacterial and fungal genomes are sequenced on both Illumina and Oxford Nanopore**

miseq

nextseq

No. of ONT Reads

No. of Illumina Reads

Fewer contigs

More Complete

CheckM % Completeness

N50/Genome Size

**All genomes are required to be at least 95% complete (CheckM/Busco)**

# Comparison of ATCC vs. RefSeq bacterial assemblies

>98% of our assemblies are more complete and of higher quality than RefSeq



Yarmosh DA et. al. *Comparative Analysis and Data Provenance for 1,113 Bacterial Genome Assemblies.* *mSphere* **2022**, e00077-22. https://doi.org/10.1128/msphere.00077-22.

# ATCC Cell Line Land

A partnership with QIAGEN Digital Insights

# ATCC cell biology collection

ATCC has **3,000+ authenticated** mammalian cell lines, genetic engineered cell lines, primary cells, stem cells, iPSCs, hTERT-immortalized cells, and tumor organoids representing various species, cell types, tissues origins, and diseases.

| 70+ Species | 100+ Cell types | 100+ Tissue types | 400+ Diseases types |
| --- | --- | --- | --- |

**Cell line models for over 400 disease types**

**2,111 cell lines for human disease models (long tail)**

For human cells :
- 84% have karyotyping information
- 89% include at least some clinical data
    - 1,749 known age
    - 751 female, 878 male
    - 870 with ethnicity
    - Additional metadata and biomarker data available as well

# ATCC Cell Line Land

| KEY FEATURES | 1. Repository of *authenticated 'omics data traceable to physical materials* <br> 2. Data production, curation, and analysis uniformly standardized <br> 3. Enables the highest level of **scientific reproducibility** <br> 4. End-to-end **data provenance** |
|---|---|

**ATCC Cell Biology Collection**

ATCC Biomanufacturing

Comparative Transcriptomics Projects

Customer Sponsored Projects

**Standardized cell culture, RNAseq, and bioinformatics**

Cell Growth

RNA Extraction

**RNAseq**

Bioinformatics & Curation

**ATCC CELL LINE LAND**

**Strict quality acceptance criteria at each step**

ATCC

# ATCC Cell Line Land

*A partnership with QIAGEN Digital Insights*

Q4 '2022 — **Kidney cell lines**
+250 reference data sets

Q1 '2023 — **Blood Lymph Spleen**
+250 reference data sets

Q2 '2023 — **Lung Liver**
+250 reference data sets

Q3 '2023 — **Liver Brain**
+250 reference data sets

Q4 '2023 — **TBD**
+250 reference data sets

- Current road-map for data production is subject to change
- Based on customer feedback

- **1,000+ traceable, authenticated RNAseq datasets per year**

ATCC

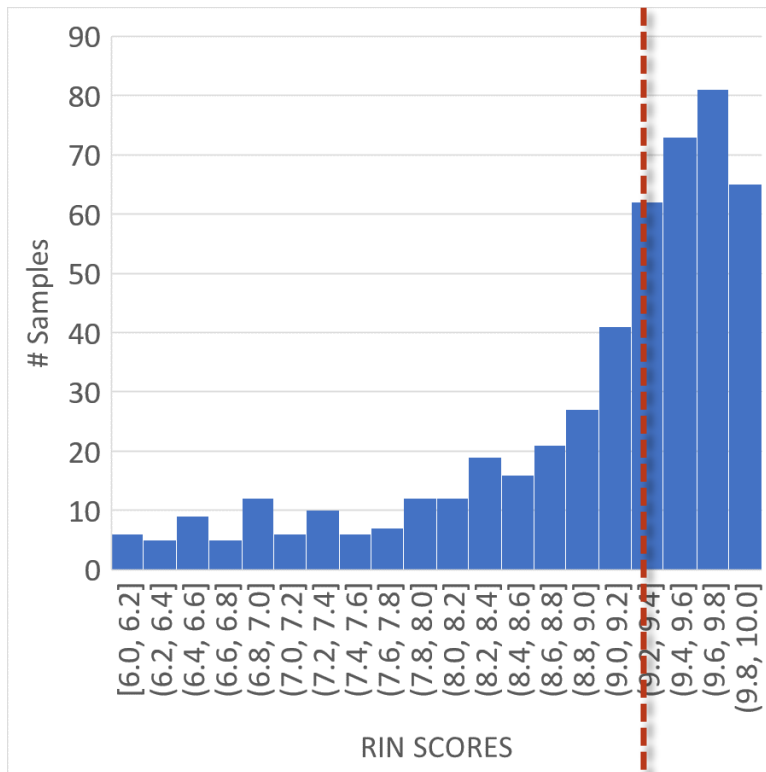# ATCC Cell Line Land
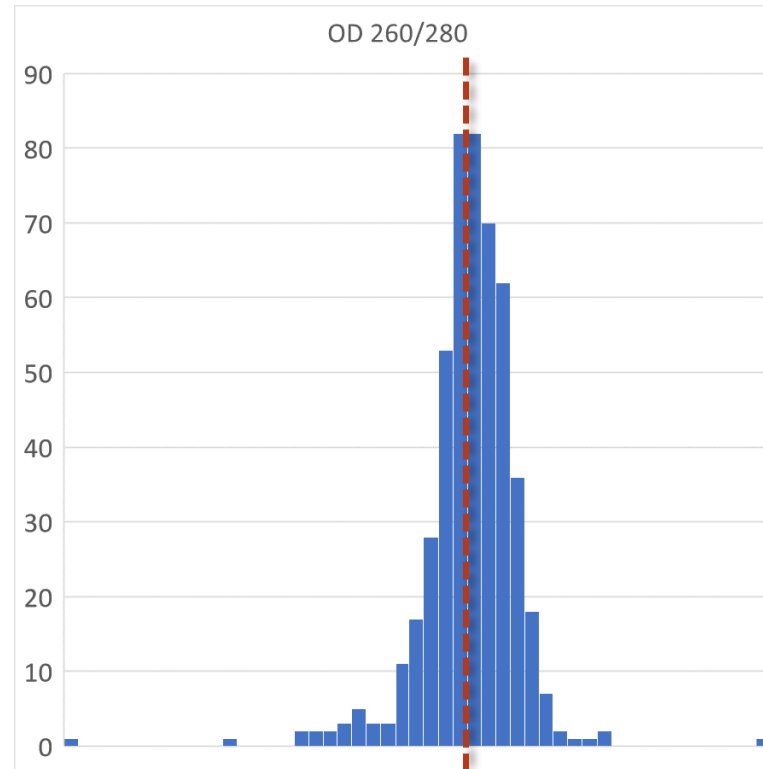
*A partnership with QIAGEN Digital Insights*

- Quality-controlled data from ATCC cell lines

- Over 1000 new datasets added each year, released quarterly

- Careful metadata curation with controlled vocabulary

- Reprocessed and normalized RNAseq expression

- Metadata include standard culture conditions, extraction protocols, sample preparation, and library preparation

- Data grows based on what you, as a researcher, need most:
  - Our team takes your requests to prioritize the cell lines you want added to our ATCC Cell Line Land collection, as well as the type of experimental data you want curated
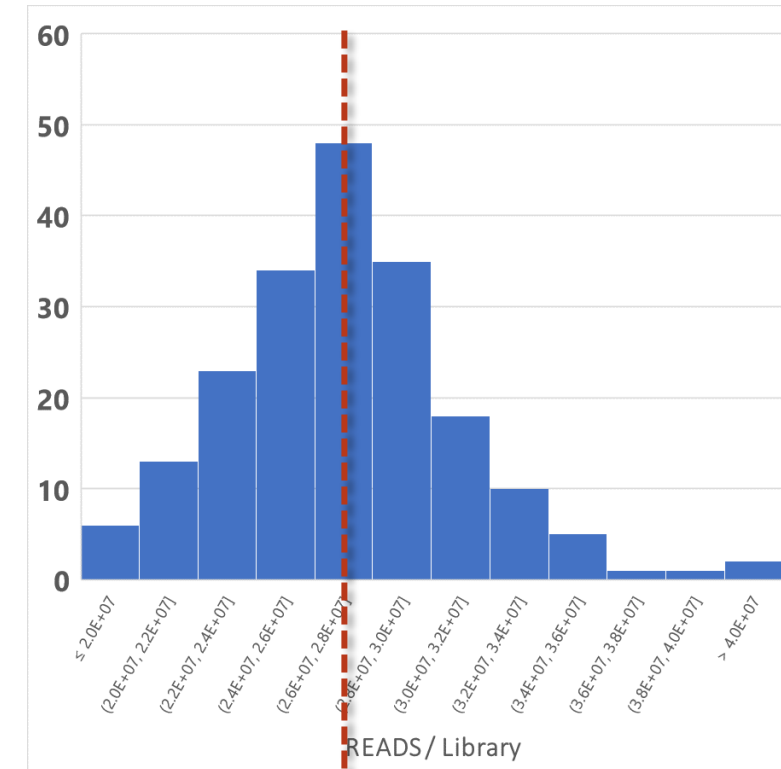
# ATCC Cell Line Land

*A partnership with QIAGEN Digital Insights*



**9.3 median RIN score**

**2.071 median OD260/280**
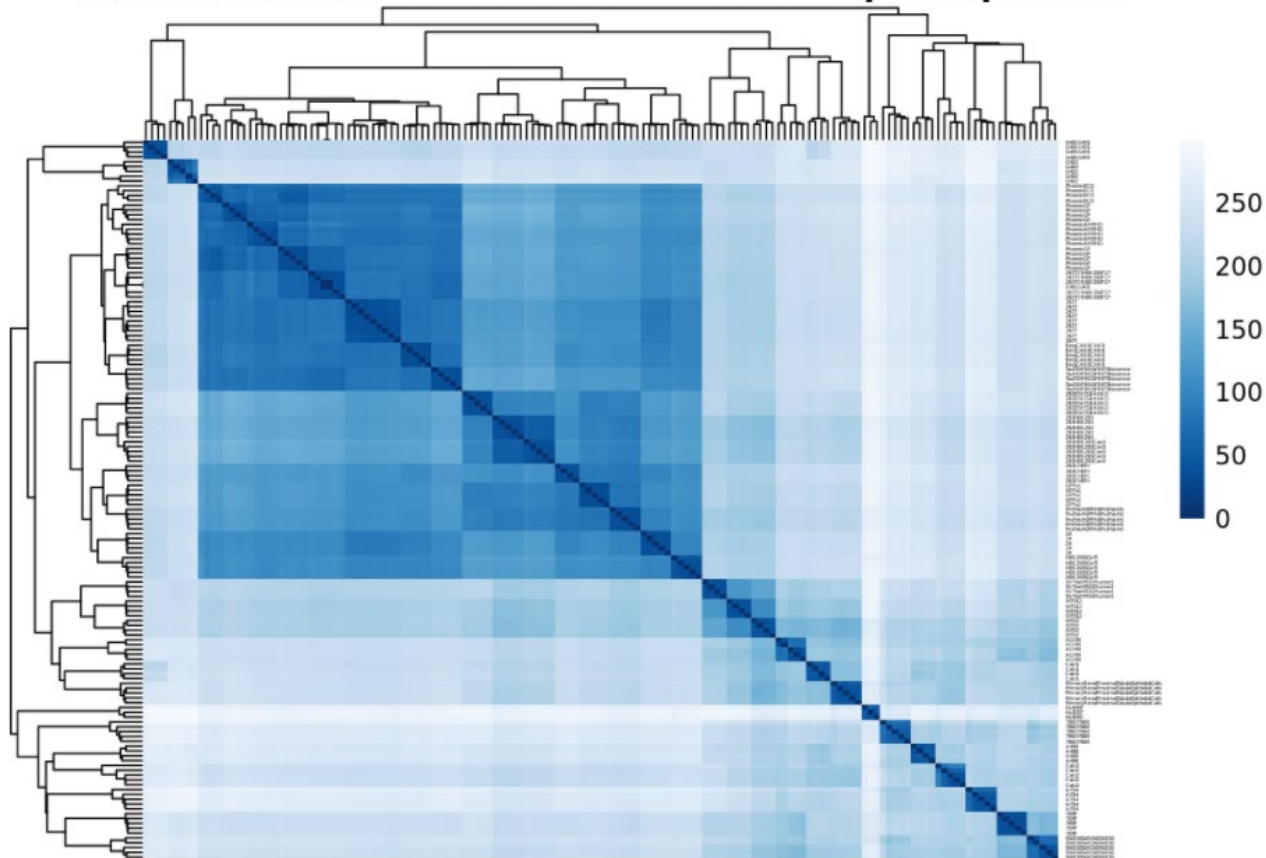
**26.9M reads per library (median)**
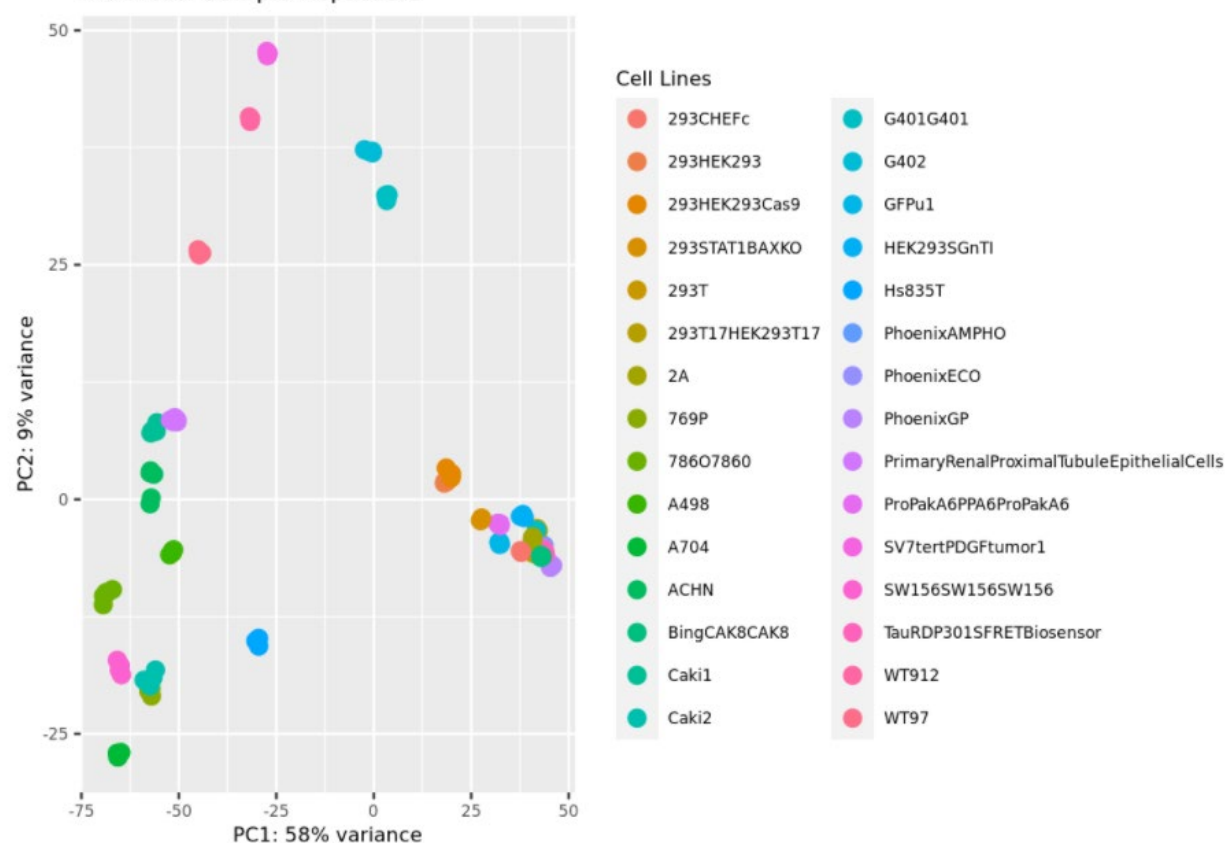
# ATCC Cell Line Land – Example (kidney cell lines)

*A partnership with QIAGEN Digital Insights*
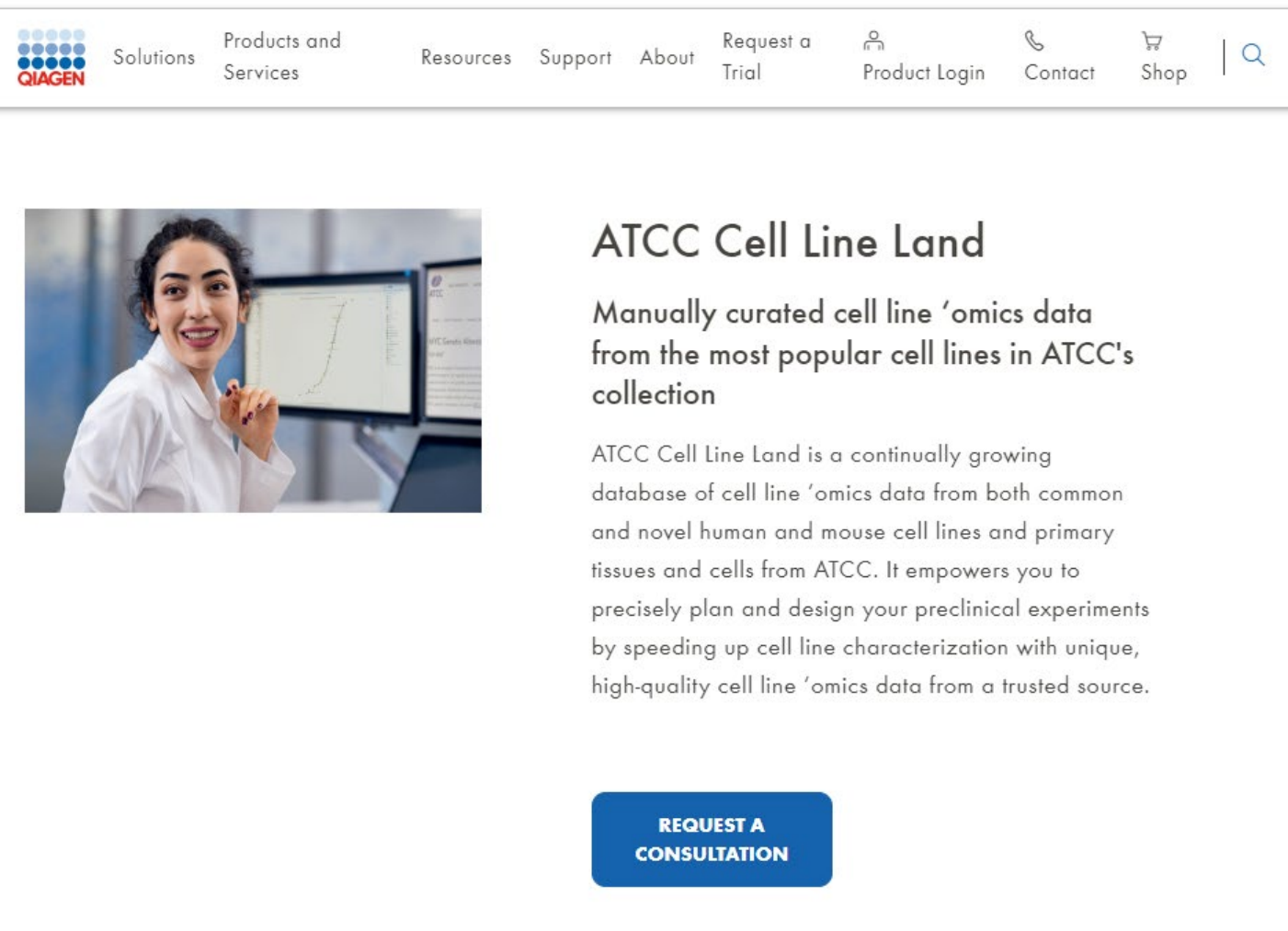


Full data for over 60 kidney cell lines will be presented at the American Society of Cell Biology (December 2022)

# ATCC Cell Line Land – Available through QIAGEN

*A partnership with QIAGEN Digital Insights*



Currently includes **Authenticated RNAseq Data** for over 200 ATCC cell lines.

**https://digitalinsights.qiagen.com/atcc-cell-line-land/**

# The ATCC Genomics Team

🐦 *@ATCCgenomics*

**Jonathan Jacobs, PhD**

Senior Director, Bioinformatics

BioNexus Principal Investigator

✉ jjacobs@atcc.org   🐦 *@bioinformer*

| Genomics Team | Bioinformatics Team |
|---|---|
| **Briana Benton, PMP** | **John Bagnoli** |
| Ana Fernandes | Scott Nguyen, PhD |
| Ajeey Singh, PhD | David Yarmosh, MSc |
| Stephen King, MSc | Nikhita Putheveetil, MSc |
| James Duncan, MSc | P. Ford Combs, PhD |
| Samuel Greenfield, MSc | Amy Reese, MSc |
| Corina Tabron, MSc | |
| Noah Wax, MSc | |
| Rula Khairi, MSc | |
| Robert Marlow | |
| Jade Kirkland | |

PARTNERS

**Marco Riojas** (ATCC / BEI)

**James Crill** (Syracuse University)

**Kylene Kehn-Hall** (Virginia Tech)

**Vishnu Chaturvedi** (Wadsworth)

**QIAGEN Digital Insights**

**One Codex**

# Thank you!

| ATCC Genome Portal | https://genomes.atcc.org |
| --- | --- |
| ATCC Cell Line Land | https://digitalinsights.qiagen.com/atcc-cell-line-land/ |

ATCC