# Metrology Meets Image Data Sharing
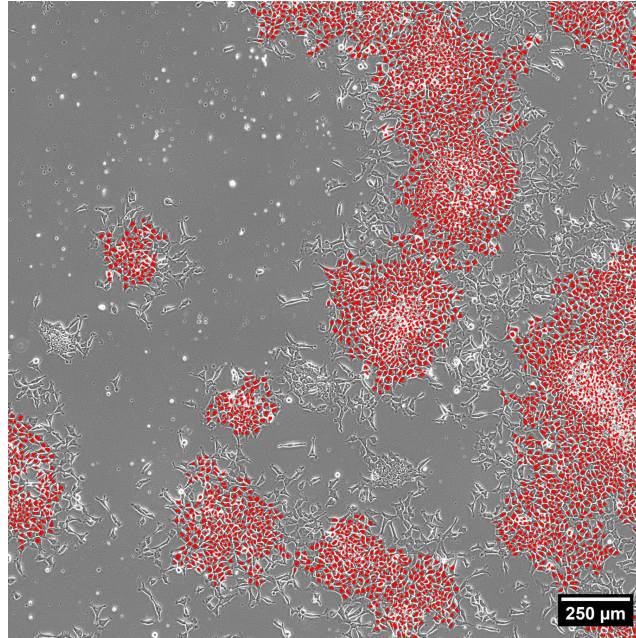
Anne Plant

National Institute of Standards and Technology
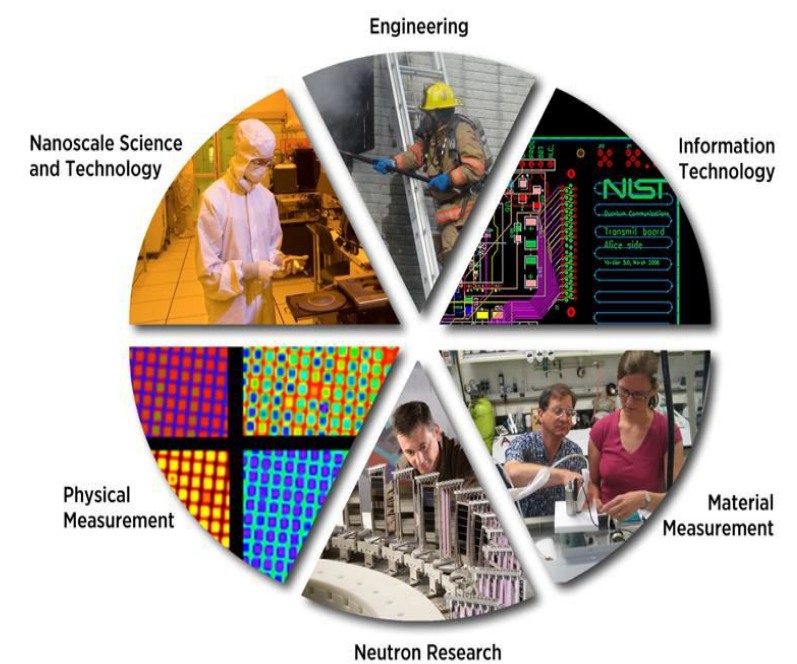
USA

anne.plant@nist.gov

- US National Measurement Laboratory.

- Non-regulatory agency partnering with academia, regulatory agencies and industry to support measurement science needs.

- Focus on measurement "infrastructure"- protocol development tools, reference material, reference data, applied statistics, and new measurement technologies.

- Participate in Standards Development Organizations (i.e. ISO TC276 Biotechnology, ASTM F04 TEMPS, CLSI, USP), and interlaboratory comparison studies.

**Trust in measurements is critical for economic, scientific and manufacturing progress**

# Data sharing > Confidence in measurements > Knowledge transfer

**_Sharing Data Allows:_**

Greater confidence in results

Combination and re-analysis of rich datasets

Development and testing of theoretical models

Comparing results from different experimental systems to test generalizability

Discovering what experimental parameters are responsible for differences in experimental results
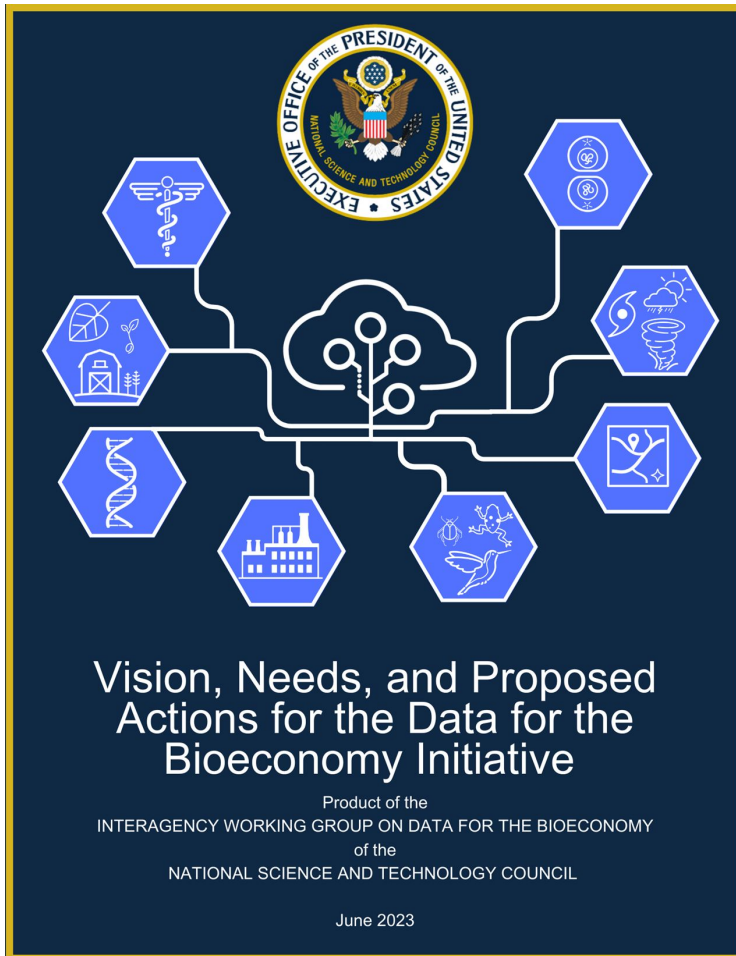
**_For data to be maximally reusable requires data qualification:_**

Sufficiently granular experimental / analytical details

Reporting of uncertainties or other indications of quality

Benchmark data from reference materials and datasets

# USG interagency WG on Data for the Bioeconomy: Recommendations



Vision, Needs, and Proposed Actions for the Data for the Bioeconomy Initiative

Product of the
INTERAGENCY WORKING GROUP ON DATA FOR THE BIOECONOMY
of the
NATIONAL SCIENCE AND TECHNOLOGY COUNCIL

June 2023

- **Dedicated long-term funding mechanisms for data and computational resources and infrastructure.**

- **Standards.**

- **Biodata Catalog.** Data/ metadata/ PIDs.

- **Security.**

- **Workforce.**

- **Strategically Targeted Areas for Rapid Transformation (STARTs)**

- **Coordination of intergovernmental investments, efforts, and resources.**

# New Opportunities for Optical Microscopy

**Easier to generate image data:**

- Ability to acquire large datasets in an efficient manner

Automated microscope technology is commonplace

- Faster/more efficient processing of big datasets, including efficient development of AI/ML pipelines to speed up development of training algorithms

GPU workstations are affordable and training data generation can be automated
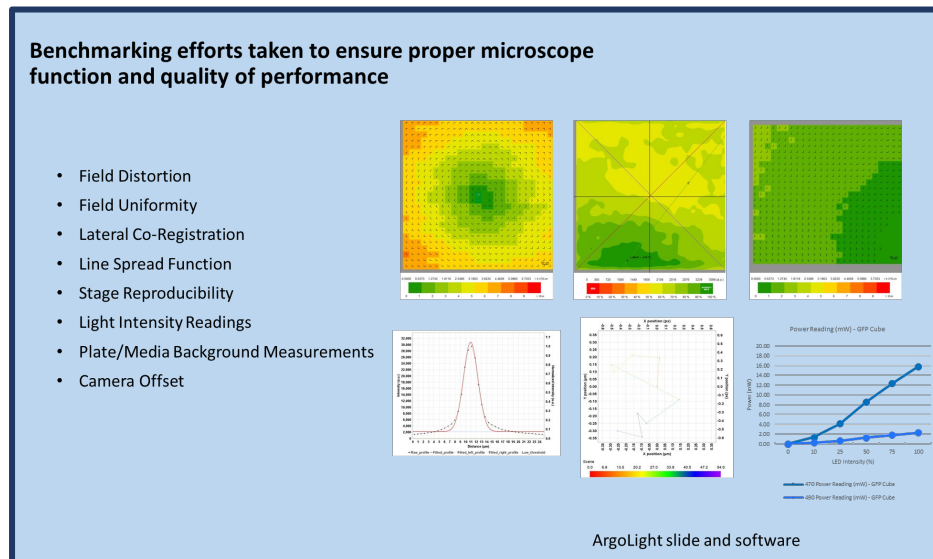
- High quality images which are reproducible

Benchmarking and reference materials

**Benchmarking efforts taken to ensure proper microscope function and quality of performance**

- Field Distortion
- Field Uniformity
- Lateral Co-Registration
- Line Spread Function
- Stage Reproducibility
- Light Intensity Readings
- Plate/Media Background Measurements
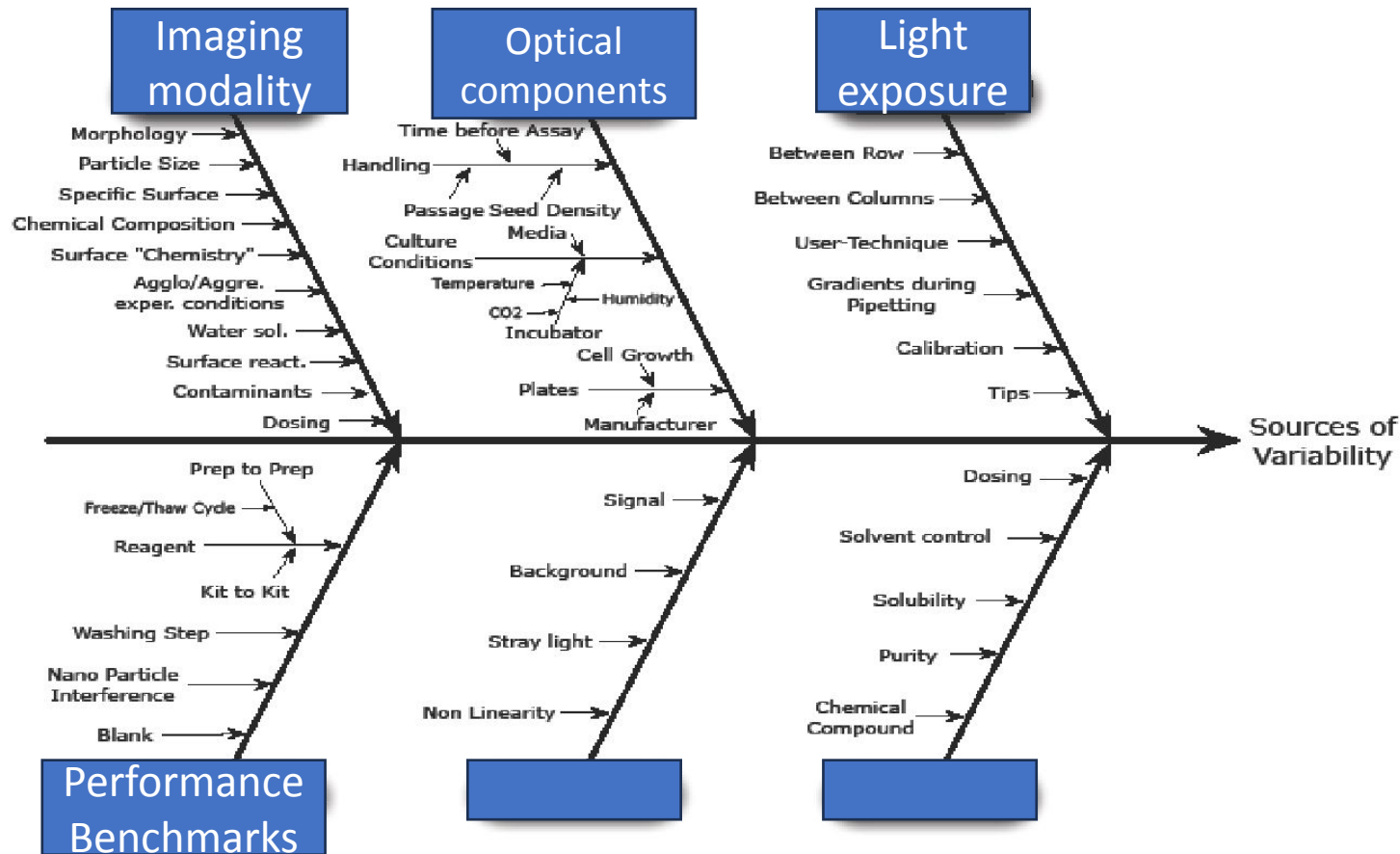- Camera Offset

ArgoLight slide and software

Using computational processing, automated methods **can** provide large scale, **quantitative** metrics for evaluation ....

**BUT** only if **benchmarking** tools are used **and sufficient metadata** are available.

NIST
National Institute of
Standards and Technology
U.S. Department of Commerce

# What metadata to report?

# What metadata to report?



**Sources of variability:**

<mark>Microscopy instrumentation</mark>………

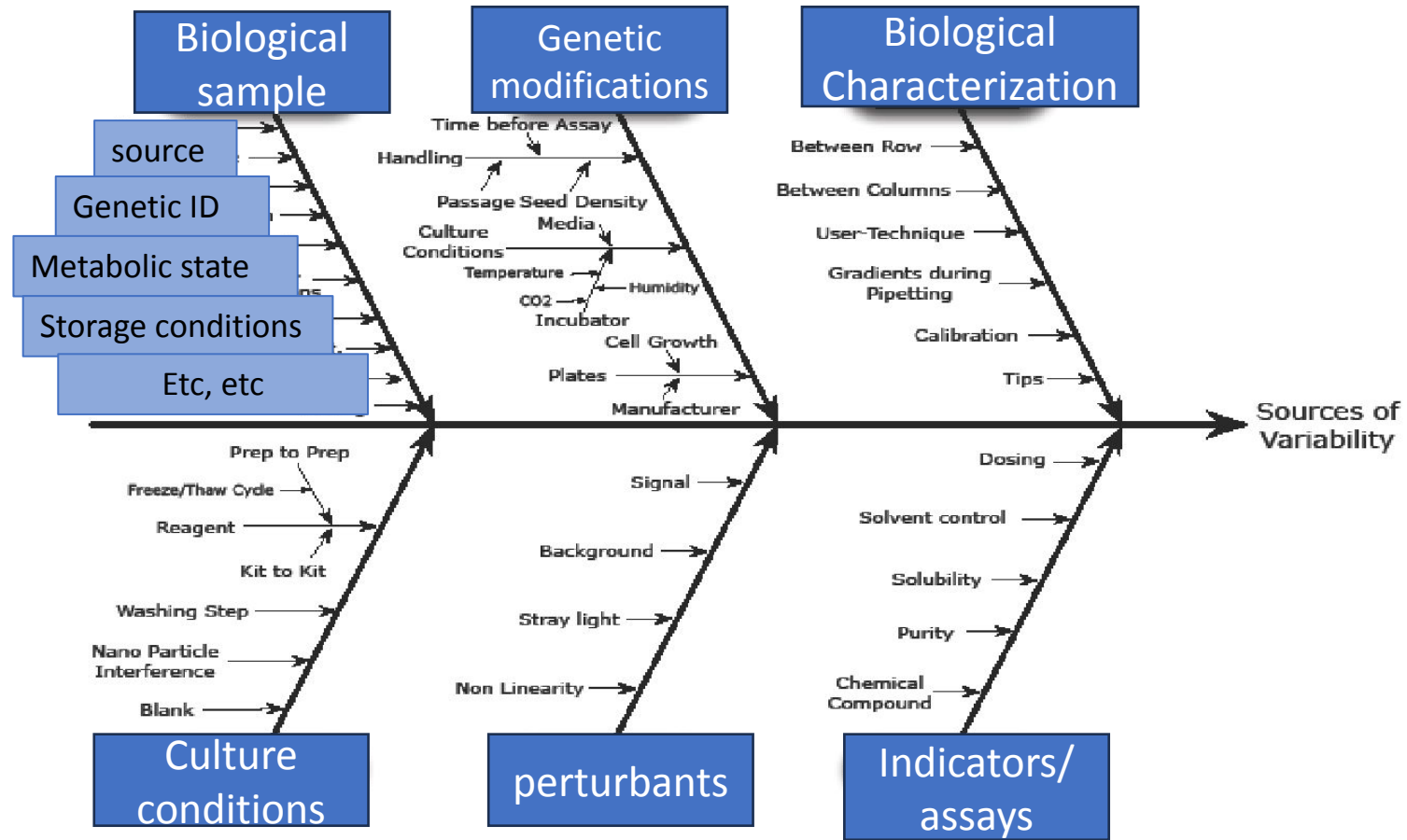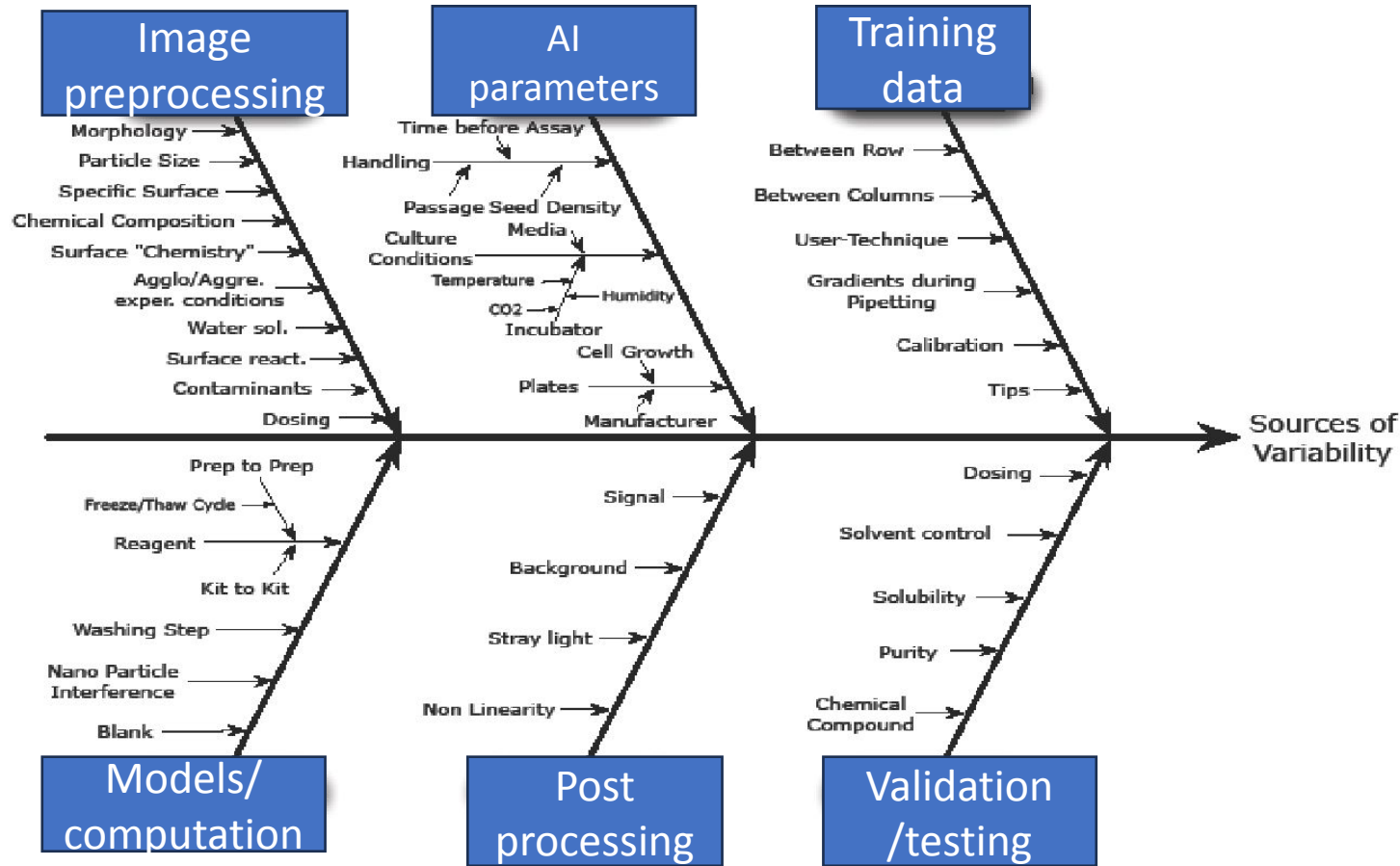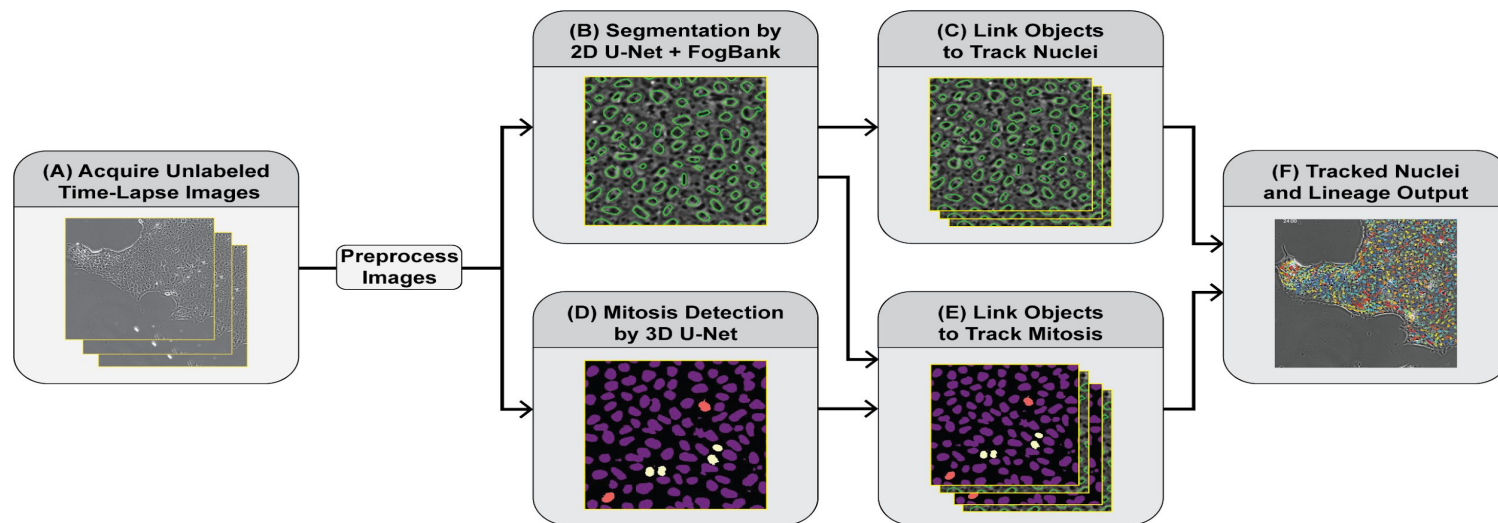Analytical pipeline………………………

**And how those variables are mitigated:**

Benchmarking microscope performance
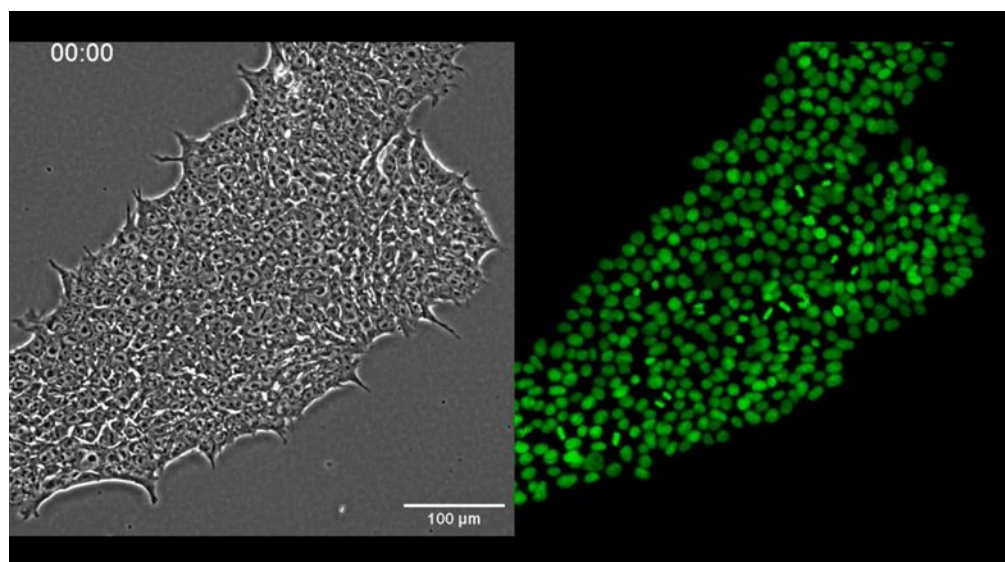
Testing analytical and computational parameters

QC of materials and assays; control of experimental details

# What metadata to report?



**Sources of variability:**

Microscopy instrumentation………

Analytical pipeline……………………

**And how those variables are mitigated:**

Benchmarking microscope performance

Testing analytical and computational parameters

QC of materials and assays; control of experimental details

# What metadata to report?



**Sources of variability:**

Microscopy instrumentation………

Analytical pipeline…………………………

**And how those variables are mitigated:**

Benchmarking microscope performance

Testing analytical and computational parameters

QC of materials and assays; control of experimental details

# High-volume, label-free imaging for quantifying single-cell dynamics in iPSC colonies



Models trained with H2B-EGFB WTC11 (Allen Institute for Cell Science / Coriell)

- Evaluate pipeline performance (imaging + analysis):
- Is classical automated segmentation equivalent to ground truth?
- Effect of some model and processing parameters.
- Reproducibility with replicate data.
- Generalizability wrt cell lines, microscopes.
- Sensitivity to cell density/cell number.

# Motivation for single cell imaging: Correlated dynamics of transcription factor expression



Induces differentiation ?

Poised for differentiation ?

Supports self-renewal ?

$W(\{x\}, t)$

SOX2

OCT4

|        | OCT4          | SOX2          | NANOG         |
|--------|---------------|---------------|---------------|
| OCT4   | $\sigma^2_O$  | $\sigma_{OS}$ | $\sigma_{ON}$ |
| SOX2   | $\sigma_{OS}$ | $\sigma^2_S$  | $\sigma_{SN}$ |
| Nanog  | $\sigma_{ON}$ | $\sigma_{SN}$ | $\sigma^2_N$  |

Each cell will report on the rate of fluctuations in expression of different genes, and their covariances are a measure of their causative relationship.

Theory is based on the Boltzmann H theorem; ties the steady state population distribution to a low relative free energy state

$$\frac{dH(t)}{dt} = -\frac{1}{2}\sum_{i,j}^{N} \int dx^N W(\{x\}, t) \frac{\partial}{\partial x_i} \ln R \cdot \boldsymbol{D}_{ij}(\{x\}) \cdot \frac{\partial}{\partial x_j} \ln R$$

Where $R = \frac{W_1(\{x\}, t)}{W_2(\{x\}, t)}$

The results of this analysis:
- **Kinetic and thermodynamic metrics indicate the relative stability of each microstate in the landscape.**
- **Prediction of time for cells to move between microstates.**
- **Identifies and quantifies causal relationships between genes.**
- **Allows identification of the most important network contributors.**
- **Allows quantification of the relative thermodynamic cost associated with maintaining the homeostatic network.**

**Hubbard et al J Phys Chem 2013; Hubbard et al PLoS 2020**

NIST
National Institute of Standards and Technology
U.S. Department of Commerce

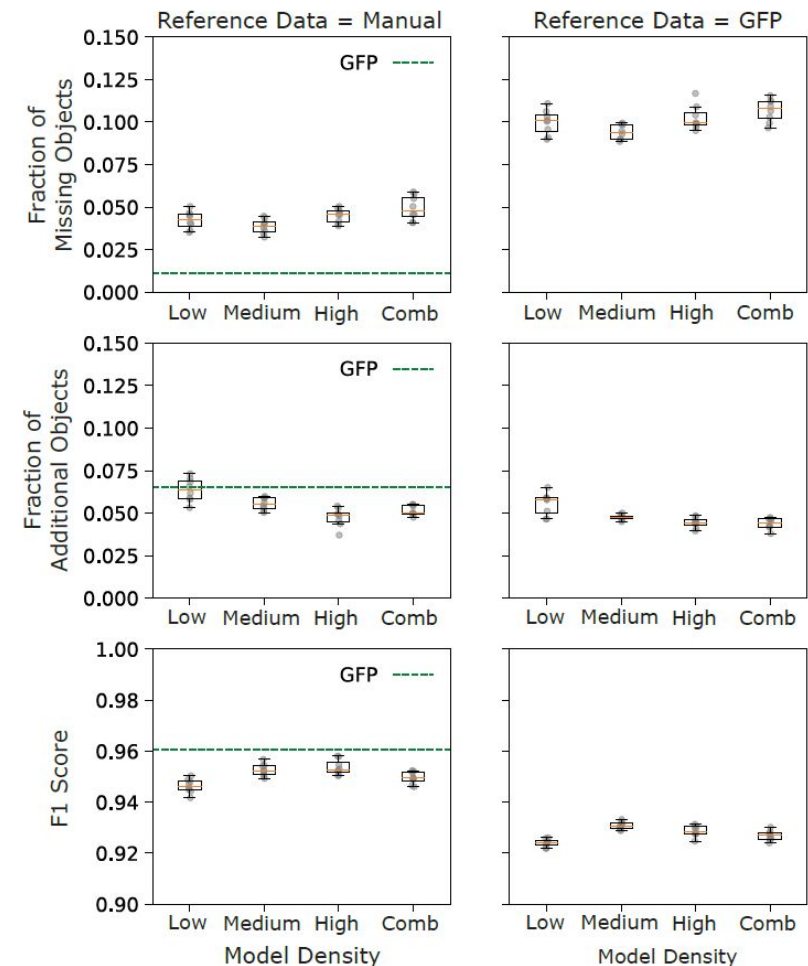# How good is auto-segmentation? Is it ground truth?



Many regions show high concordance between the manually annotated, GFP labeled, and inferred nuclei.

The GFP fluorescence-based automated image analysis tends to merge nuclear objects compared to the manual annotations and the AI-based analysis of the phase contrast images.

Some areas are more ambiguous than others
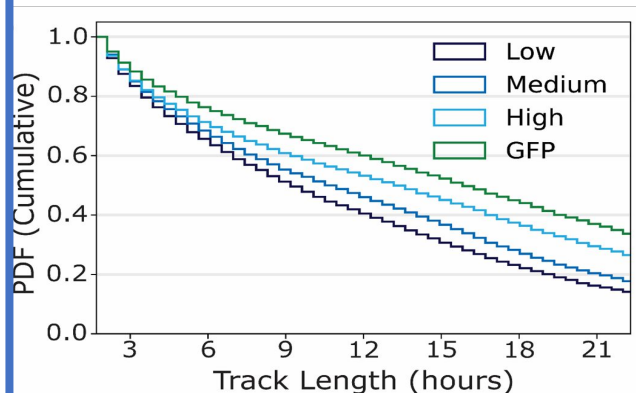  Scale bars = 10mm.

**Comparison of inferenced objects with auto-segmented or manually segmented objects.**

# (What evaluation should be expected/reported?)
# Other model parameters explored:

## Accuracy of tracking cells with the three different models

Compared to tracking using GFP fluorescence, 79%, 52%, and 42% of tracks were correctly inferenced for the high, med and low-cell-density models.



The total number of pixels used for training the 3 models was kept constant.

The number of cell objects in the training sets :
**38,875** for the low-cell-density model
**80,997** for the medium-cell-density model
**214,944** for the high-cell-density model

These results suggest that the U-Nets are sensitive to number of objects used in training.

## Influence of post-processing on inferenced data



Binary mask      Post-fogbank

The Fogbank algorithm is applied after the 2D U-Net to separate two or more nuclei that share a boundary and are considered one object.



The 'erode_size' parameter is varied from 1 to 5 and for each 'erode_size' value, the 'min_size' parameter (size filter) is varied from 5 to 15.

The highest F1 scores for segmentation accuracy can be obtained with 'erode-size' in the range of 1 to 4 and 'min_size' in the range of 8 to 10.

## Testing threshold level for concordance between different models

(scale bar = 25μm).



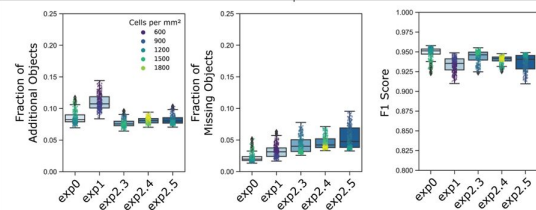The color scale indicates the number of times a trained U-Net inferred that a pixel was classified as a nucleus.
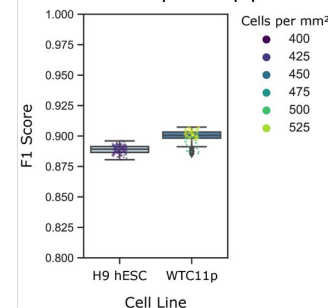
The 'Fraction of Missing Objects' increases with threshold value (oversegmenting), the 'Fraction of Additional Objects' decreases with threshold value, and the 'F1 Score' is highest for intermediate threshold values.

We inferenced with 3 instances and thresholded by 2, which was practical vis a vis computing time, and produced reasonable results.

## Reproducibly and Generalizability

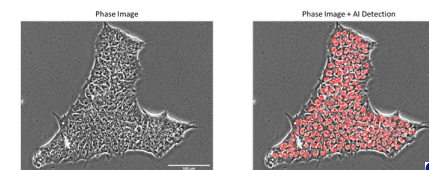Reproducibility of pipeline: within-day and between-day variability



Generalizability of the pipeline



We also compared 2 microscopes (differed slightly in transmitted light source and condenser, and in fluorescence excitation source.

Differences were negligible (F1-scores of 0.95 and 0.94)



H9 Embryonic Stem Cells

NIST
National Institute of Standards and Technology
U.S. Department of Commerce

# Availability of data: [doi:10.18434/mds2-2960](doi:10.18434/mds2-2960)

- We post only the experimental data used in Figures.
- The image process operations being done on these images : select best focal plane, normalize phase images, and stitch multiple FOV.

**Excel file**

- Lists each final image dataset designation and the figure(s) where those data appear.
- Note phase + fluorescence (and radiant exposure) or only phase
- Note if that dataset is used for training.
- Initial/compressed data file size (87/34GB).
- Cell count, detected mitoses count.

# Outstanding questions

- How much to report re models and processing parameters.

- How much data to make available?

- How to best report the evaluation of the image analysis model?

- **How to improve the capture and reuse of experimental metadata terms?**

**Future????    New technologies for metadata capture and retrieval?**

- Use of LLMs to identify relevant metadata terms/ definitions.
- RAGs or other methods for retrieving terms from a database.
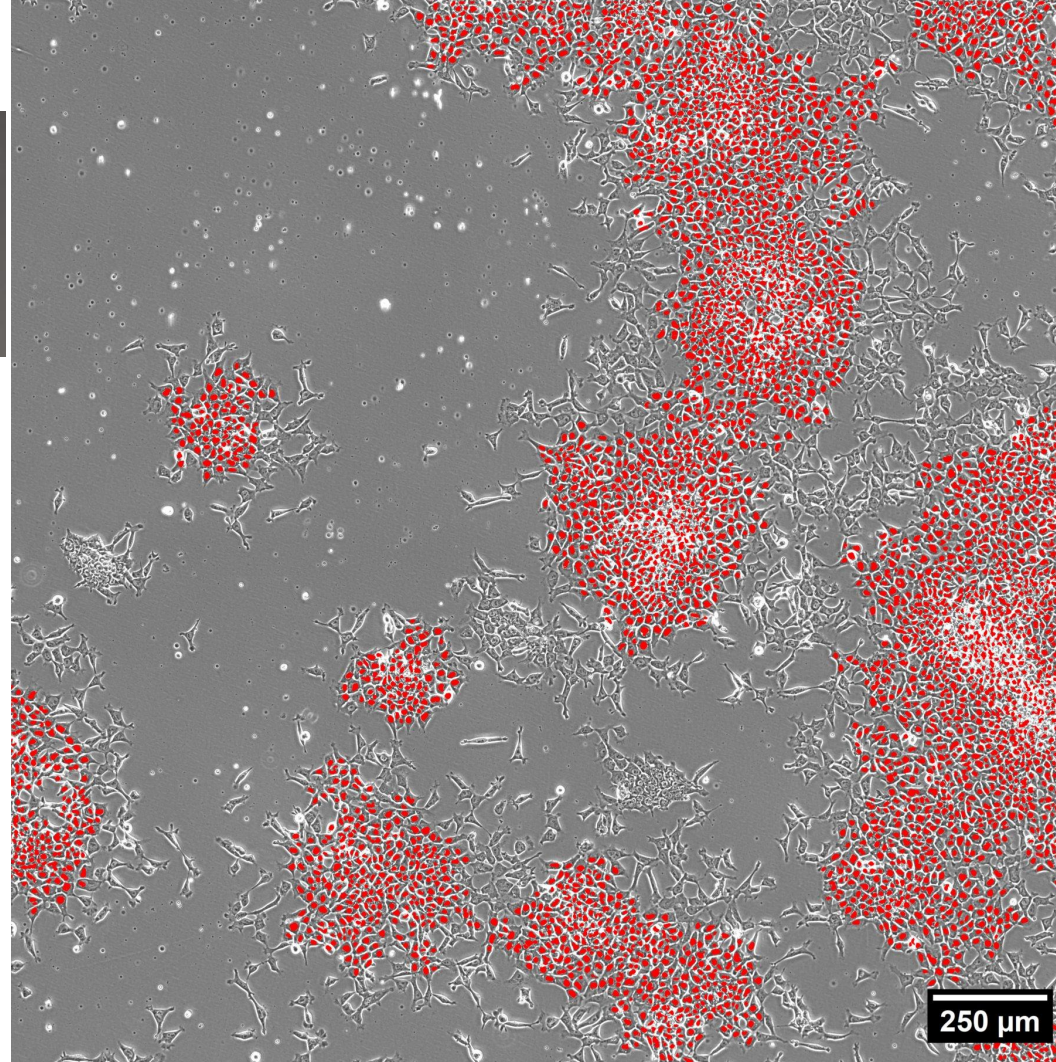
# Thank you

Anthony Asmar

Michael Halter

Zack Benson

Joe

Adele Peskin



250 µm

PLOS ONE |
https://doi.org/10.1371/journal.pone.0298446
February 20, 2024

Postdoctoral opportunities available
anne.plant@nist.gov

**NIST**
National Institute of
Standards and Technology
U.S. Department of Commerce