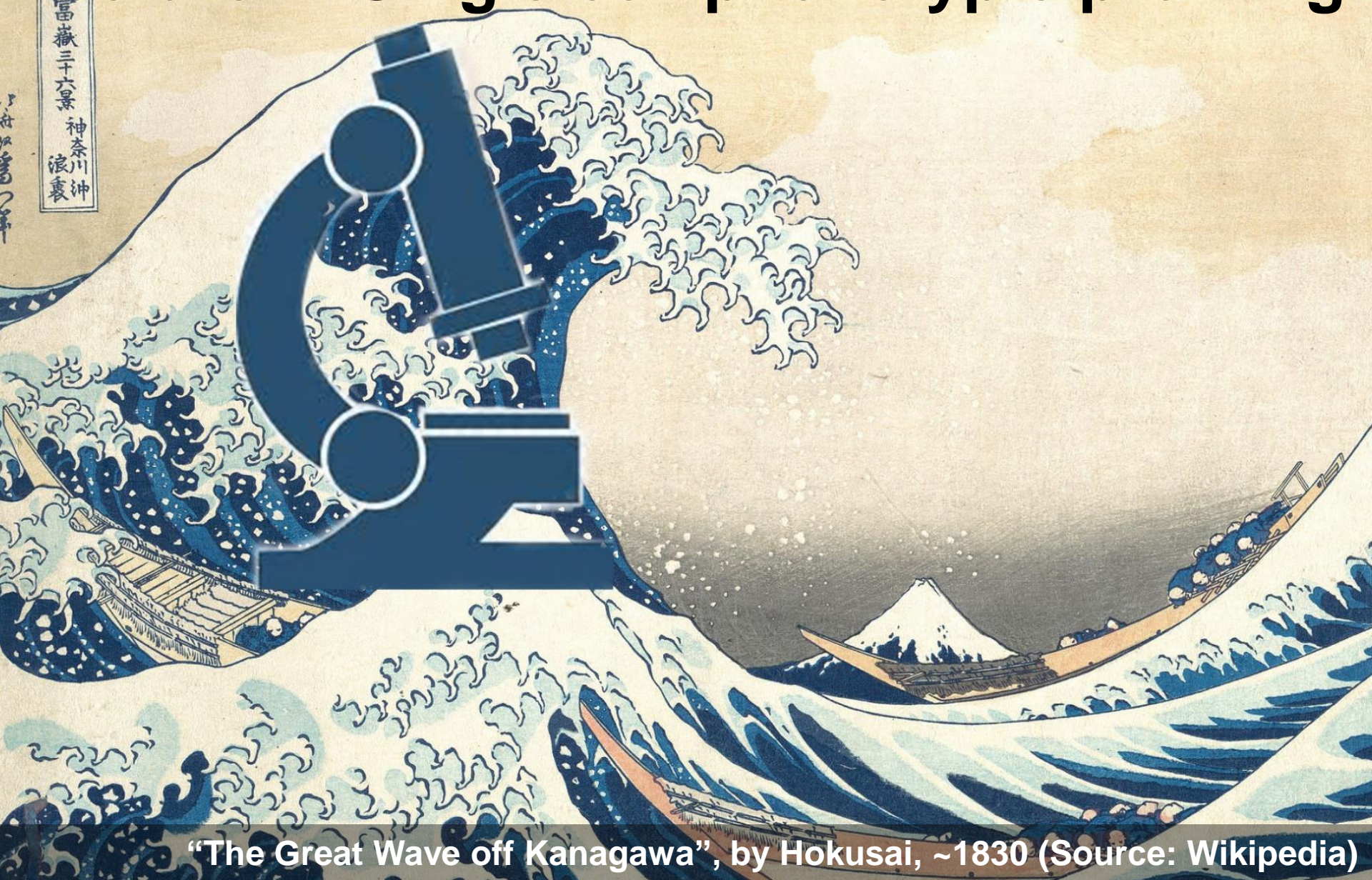


Data science in cell imaging

Lecture 2: Single cell phenotypic profiling



“The Great Wave off Kanagawa”, by Hokusai, ~1830 (Source: Wikipedia)

All slides are open under the cc-by
license

You are free to share and adapt any
content from this presentation provided
that you attribute the work to its author
and respect the rights and licenses
associated with its components

PPTX slides available [here](#)

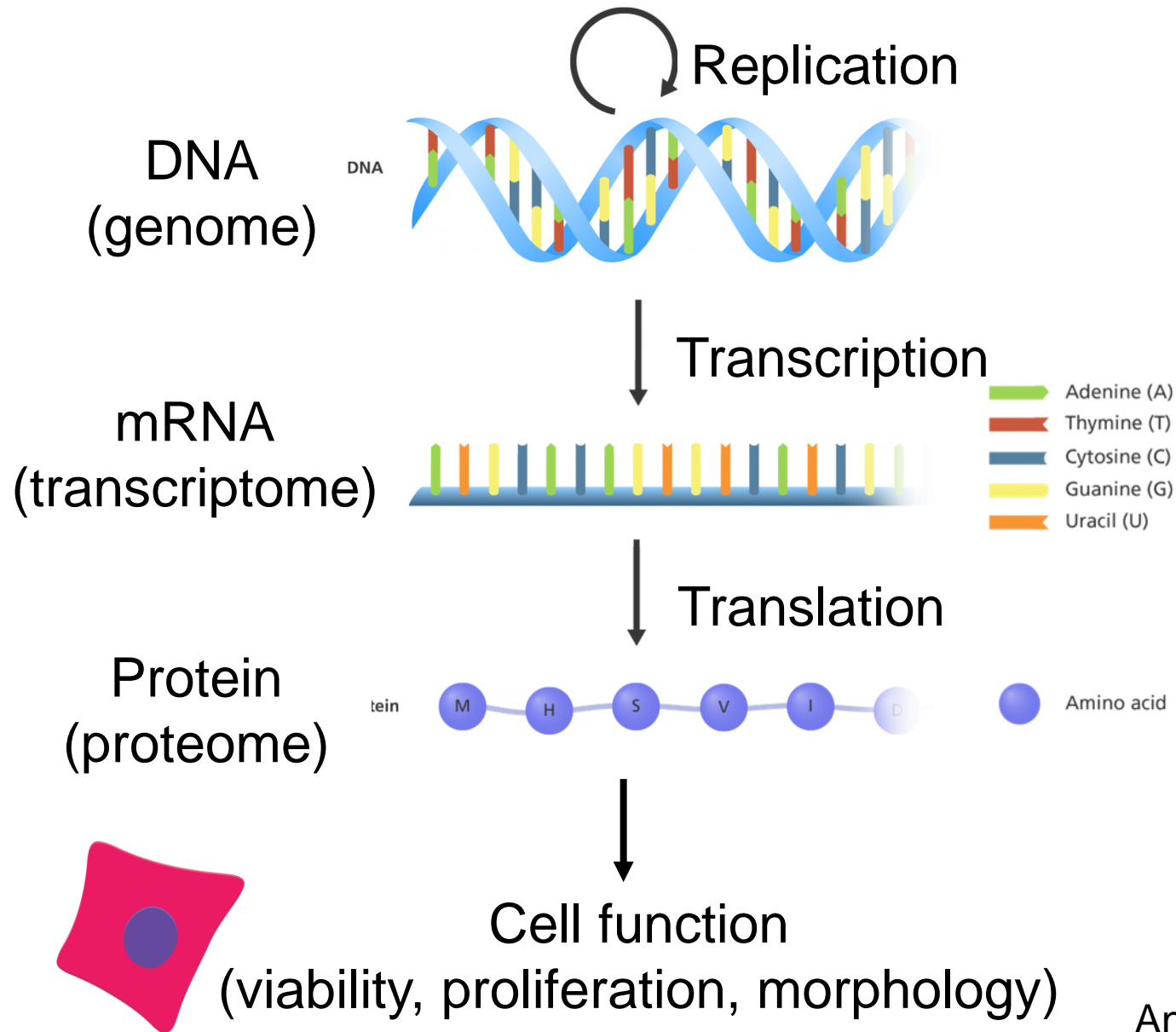


Thank you and caveat toward today's
lecture

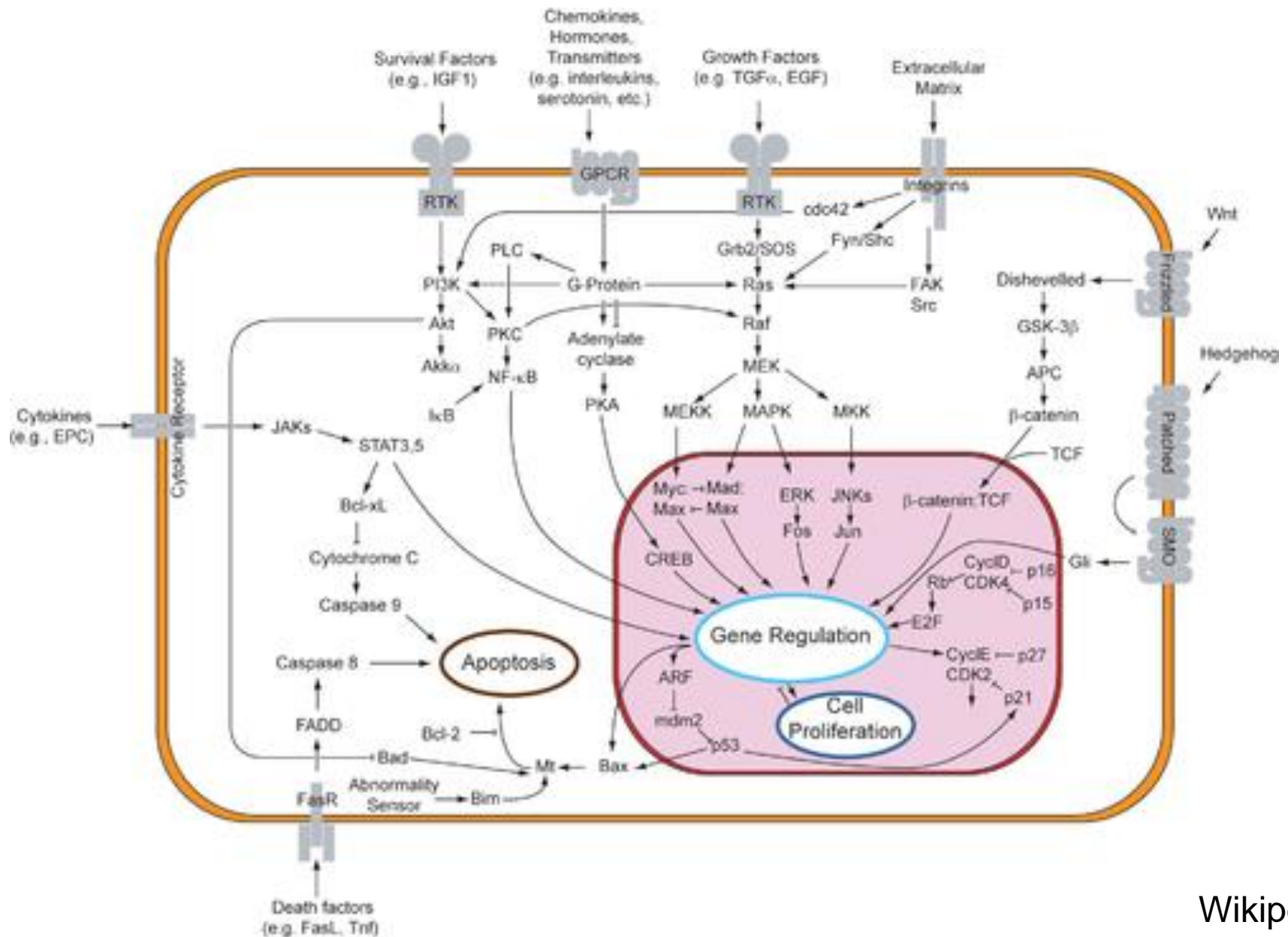
Most of today's slides were adapted from
Anne Carpenter (Broad Institute)!

Heavily biased toward her lab's work
(to be improved next year...)

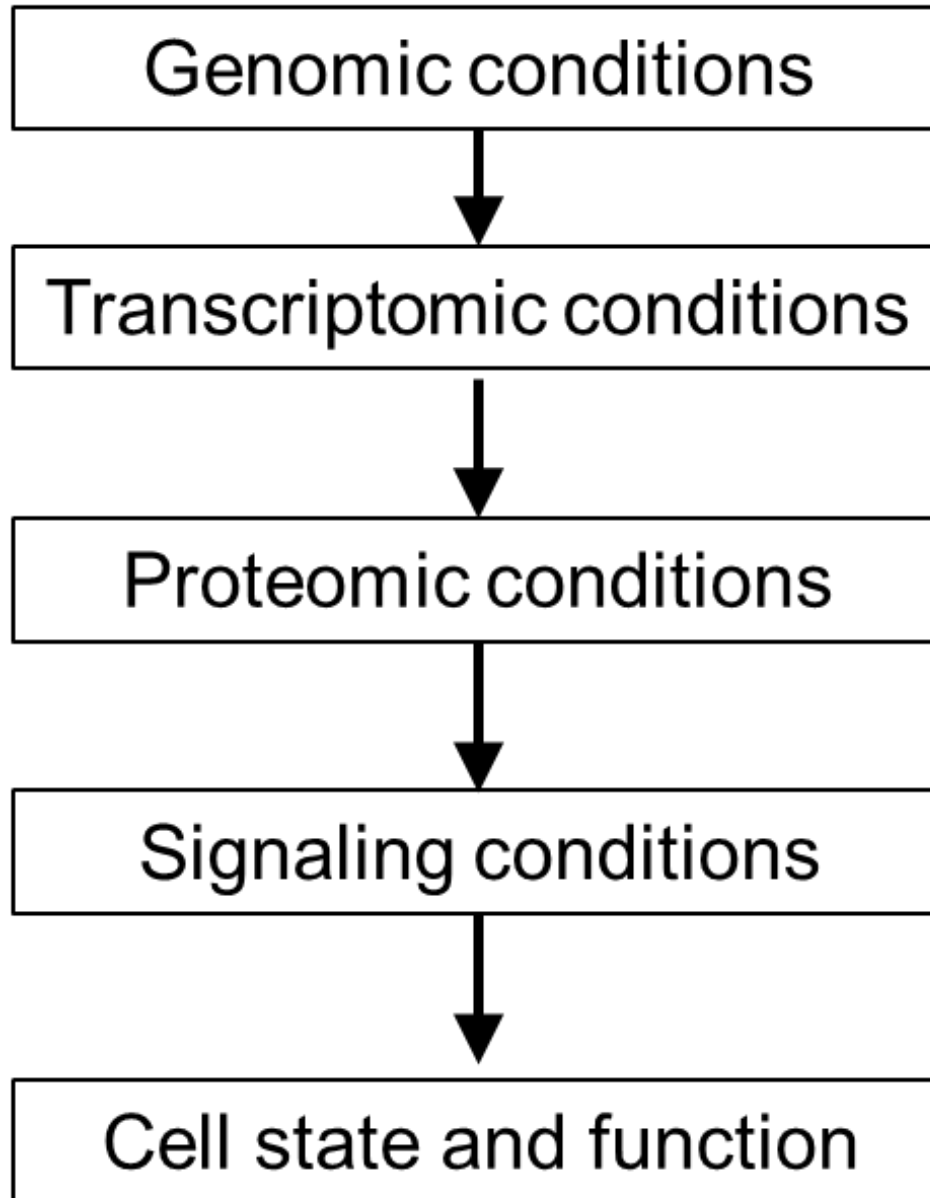
The central dogma of molecular biology



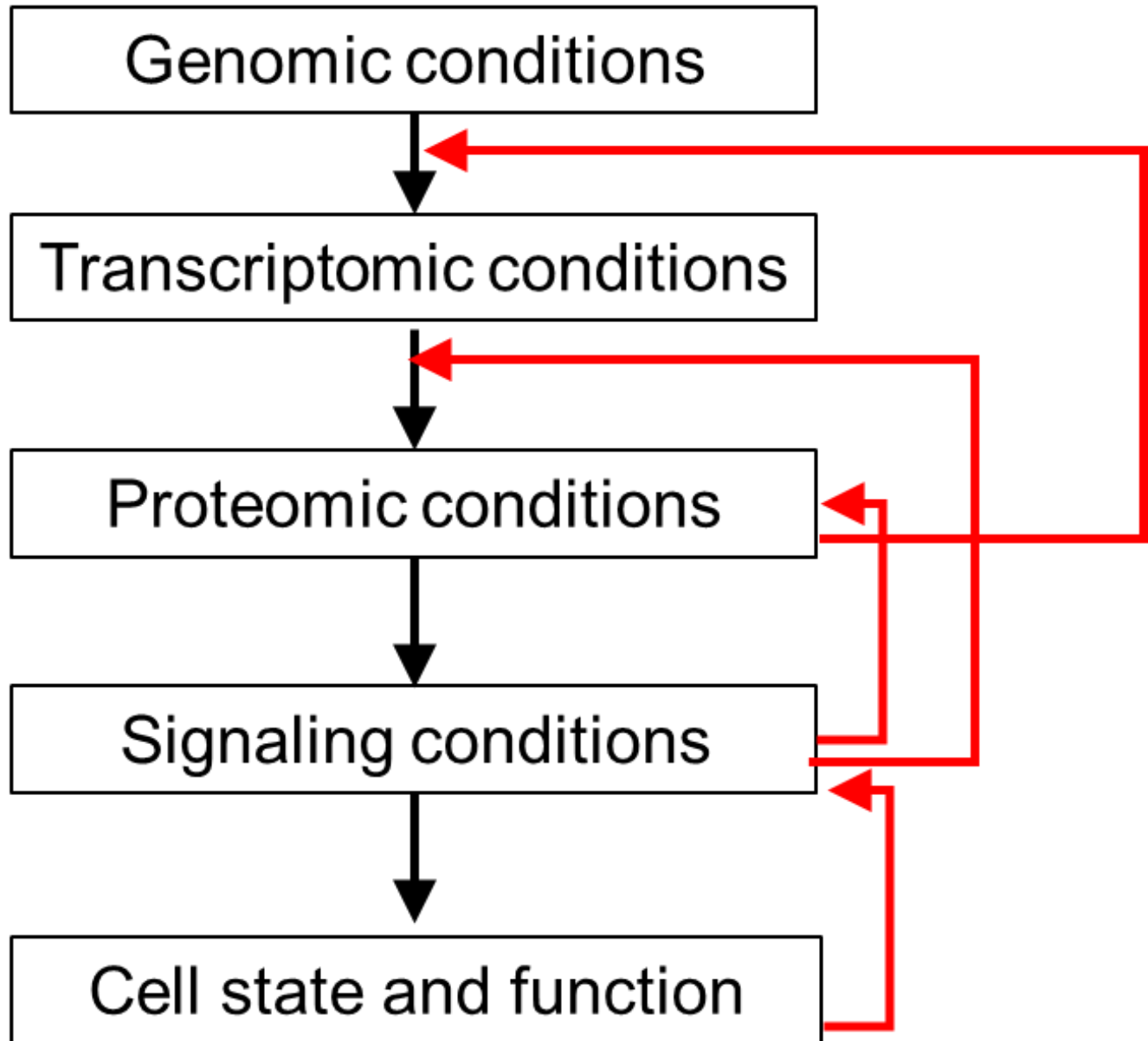
Cell signaling



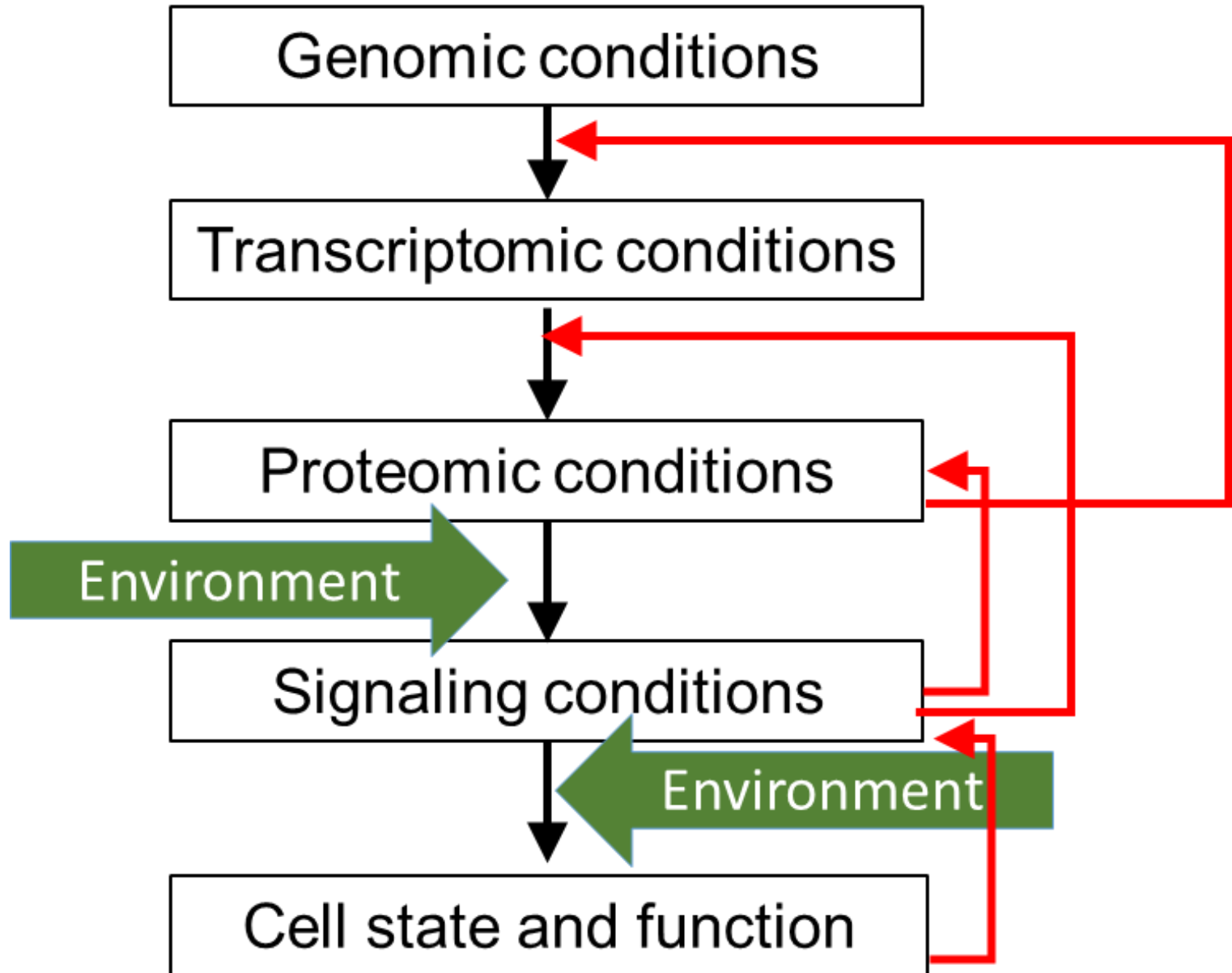
The central dogma of biology



But in reality...



But in reality...



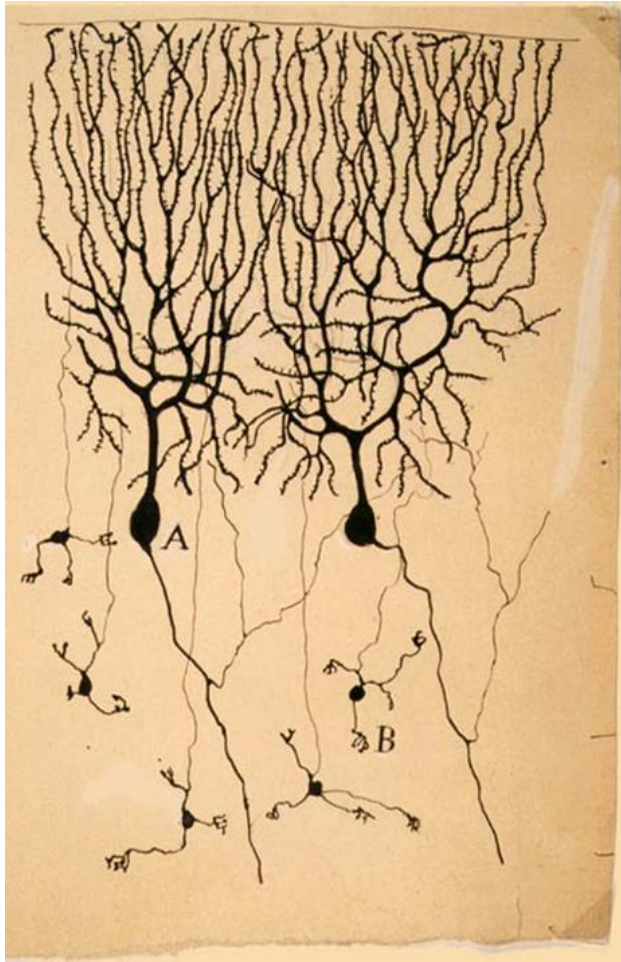
We need functional
readouts!

Why microscopy?

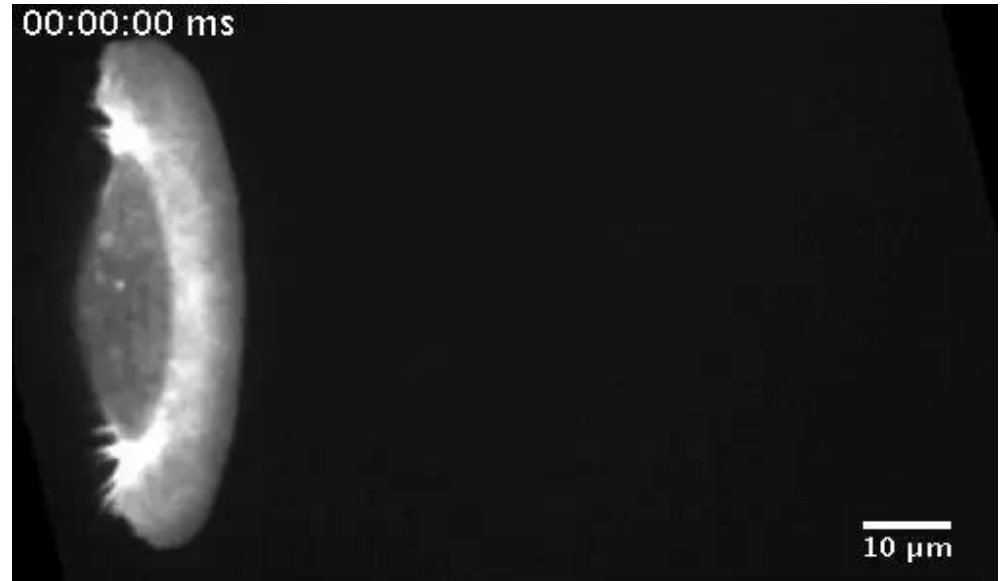
The only technology that allows us to see live cell behavior and to correlate cell function to intracellular (protein quantity and location) and extracellular (environment) factors

Morphology is a marker of
the cell's functional state

Diverse cell morphologies enable diverse functions

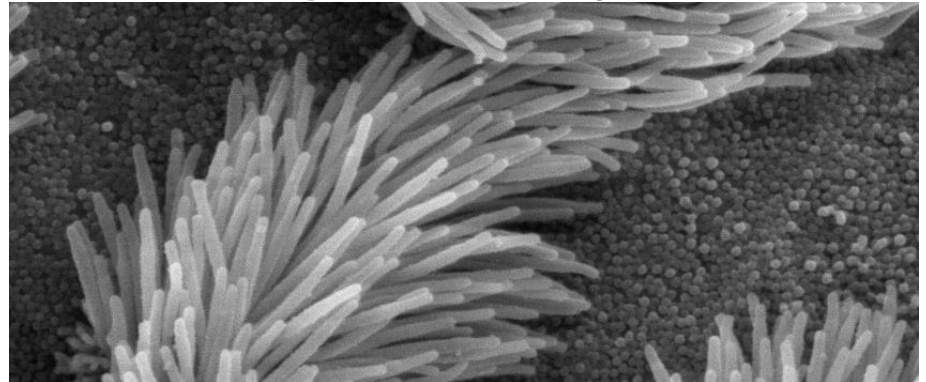


Purkinje cells
Santiago Ramón y Cajal



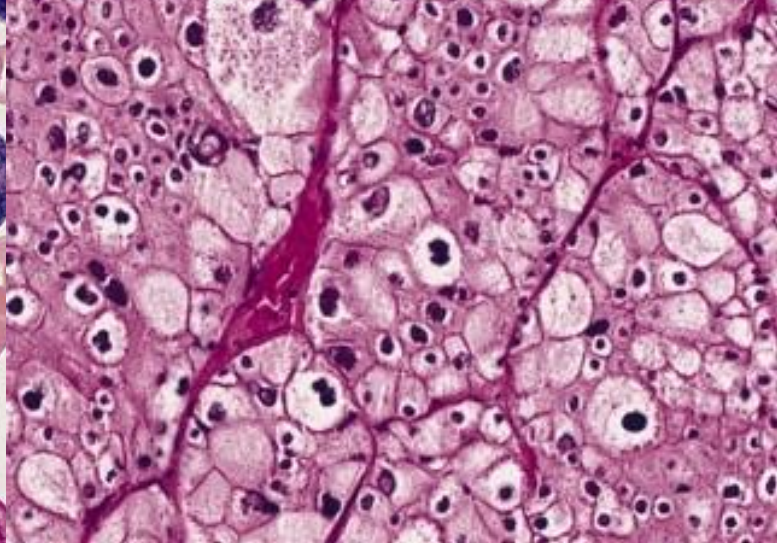
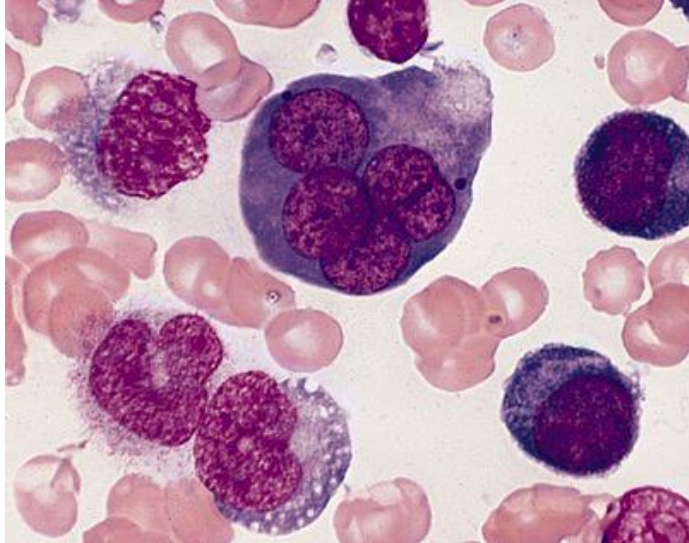
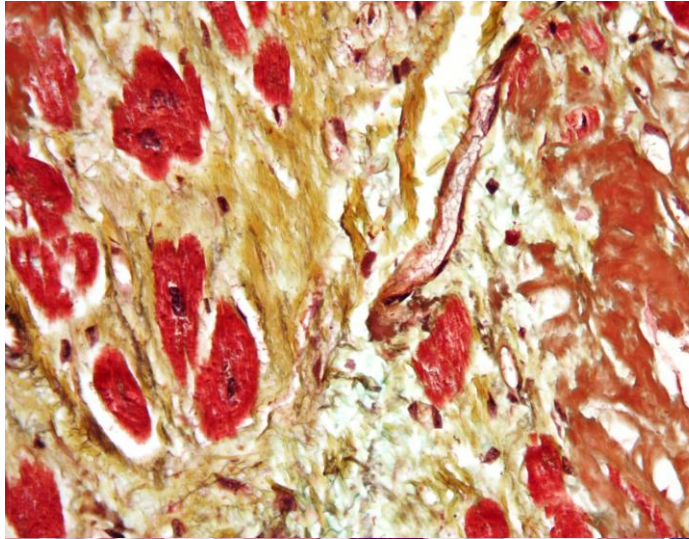
Zebrafish keratocyte; Mueller *et al* (2017)

Cilia in the lungs; Charles Daghljan, Dartmouth



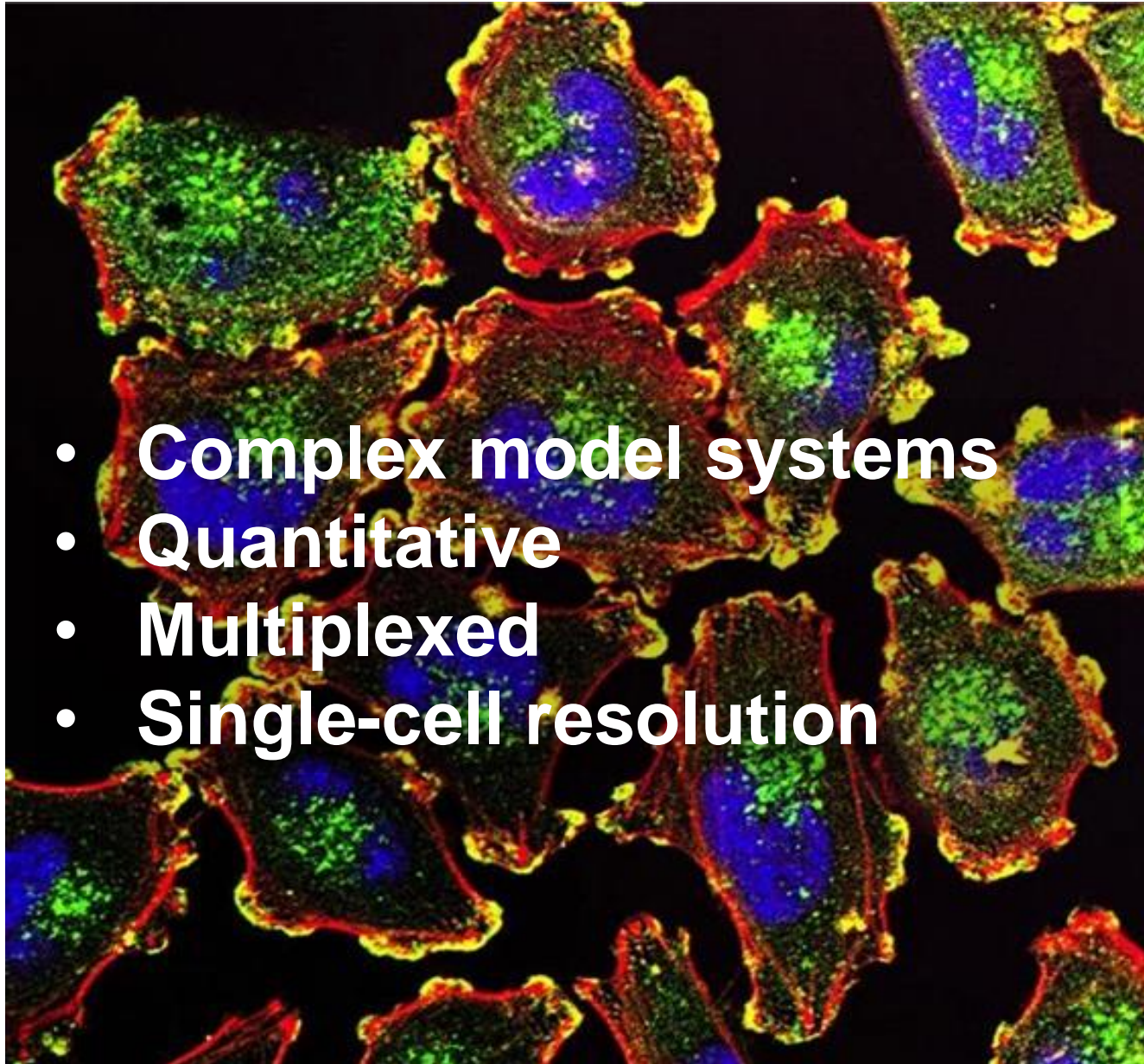
Meghan Driscoll

Visual appearance indicates cell (and disease) state



Anne Carpenter, images from Wikipedia

Images contain a wealth of information

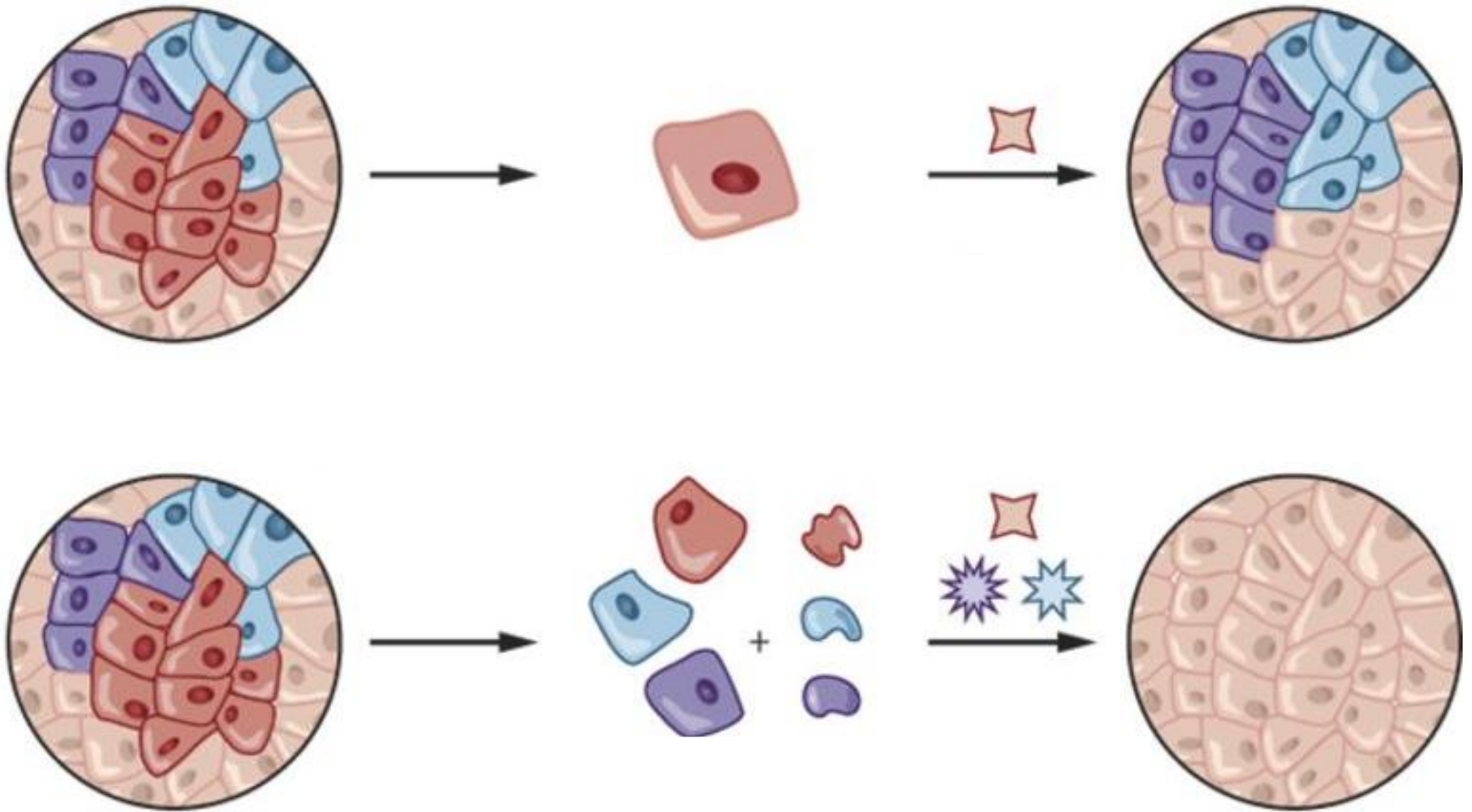


- **Complex model systems**
- **Quantitative**
- **Multiplexed**
- **Single-cell resolution**

Why do we care about single cell resolution?

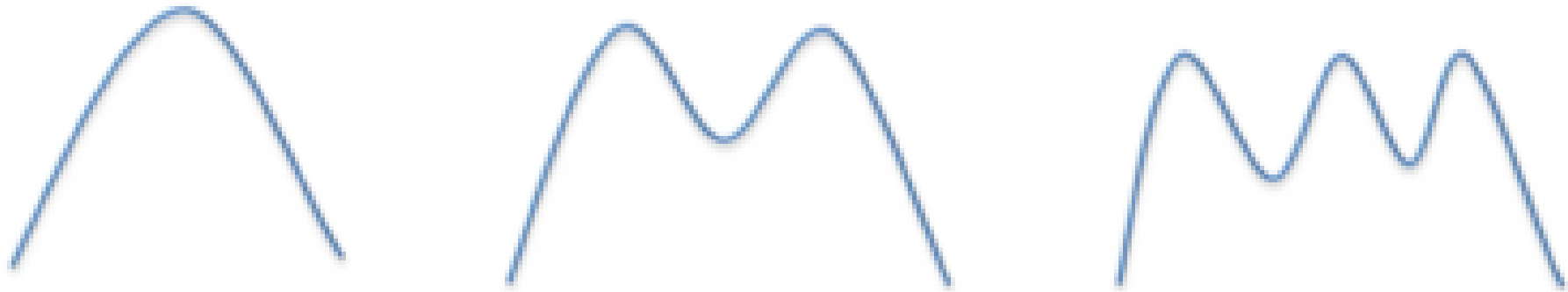


Genetic heterogeneity



Phenotypic heterogeneity

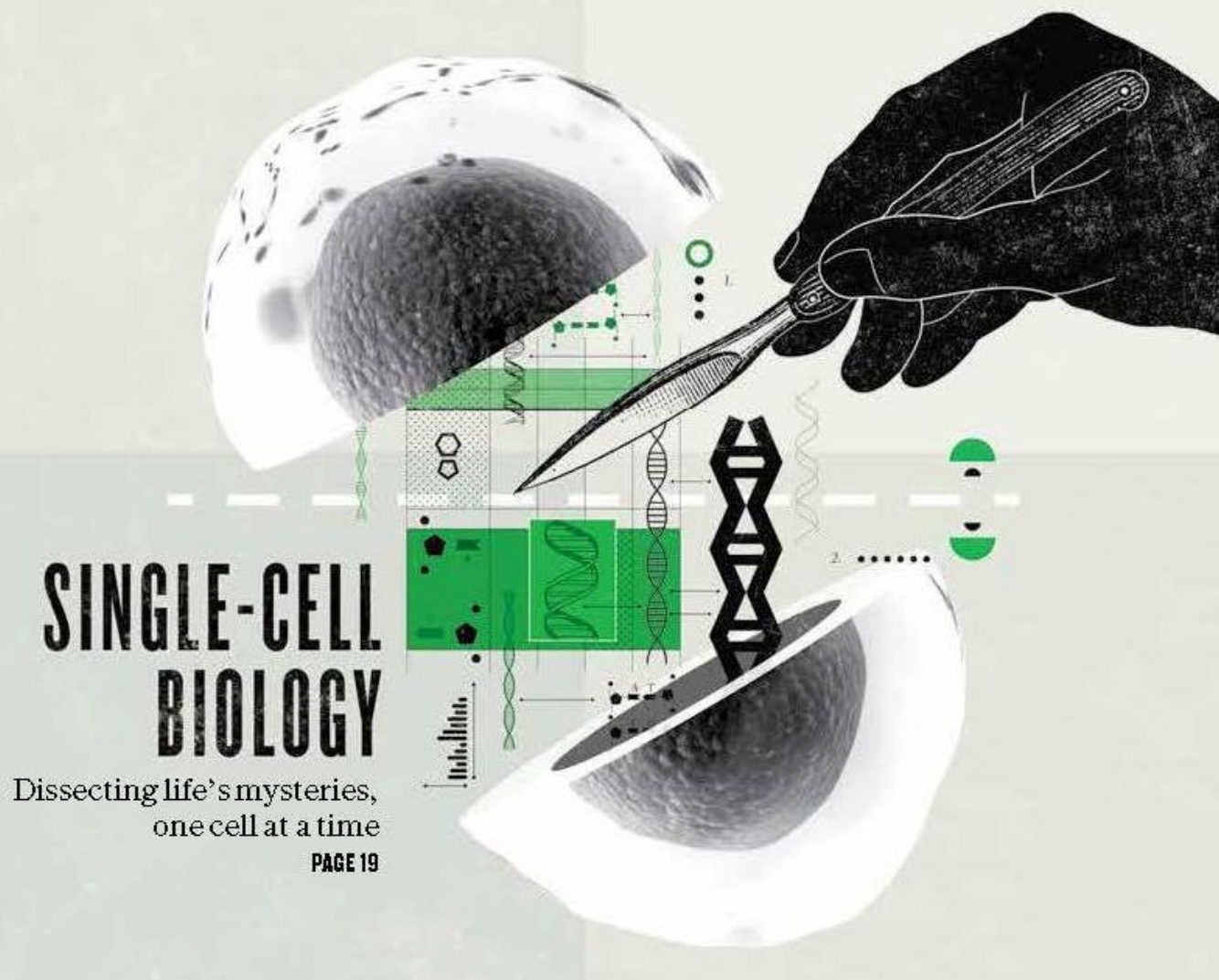
Occurrences



Phenotype

nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

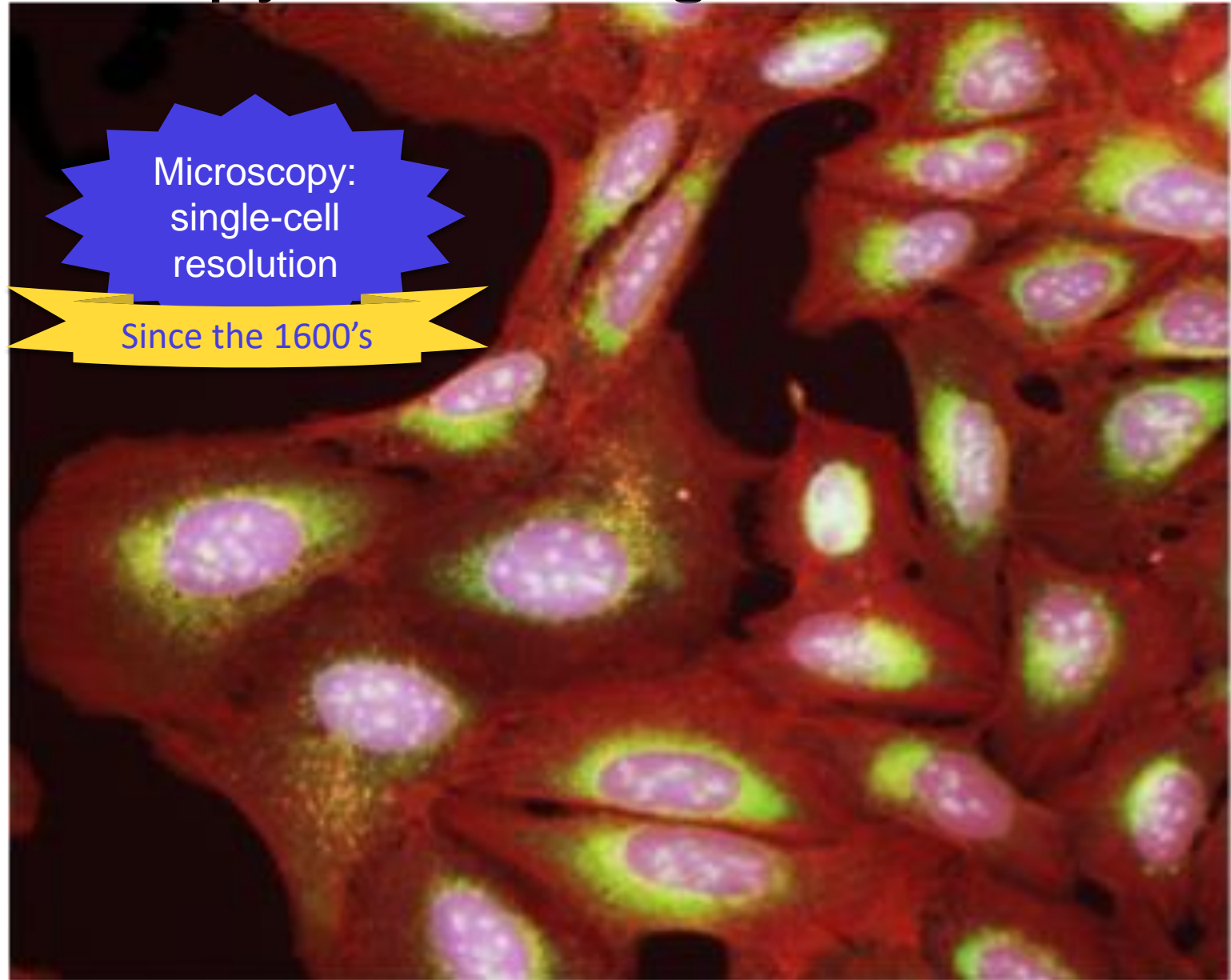


SINGLE-CELL BIOLOGY

Dissecting life's mysteries,
one cell at a time

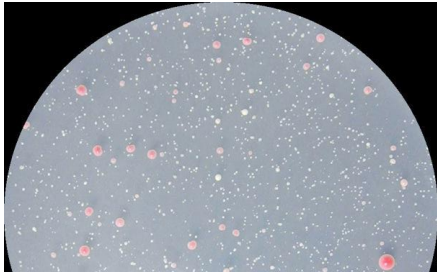
PAGE 19

Microscopy offers single cell resolution

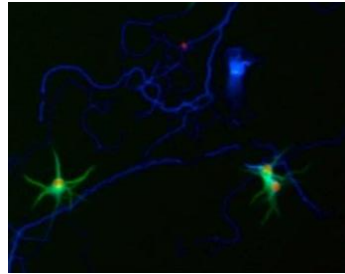


Complex cell models can be quantified by imaging

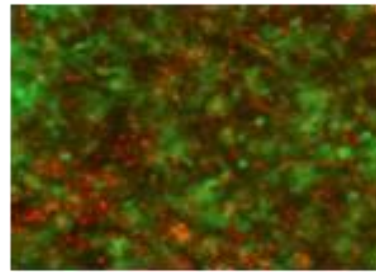
Yeast



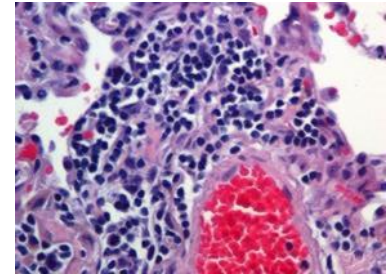
Neurons



Co-cultures



Tissue



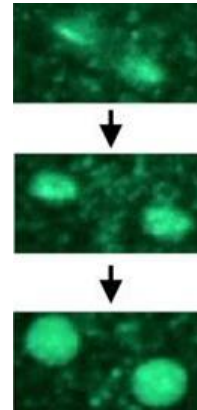
C. elegans



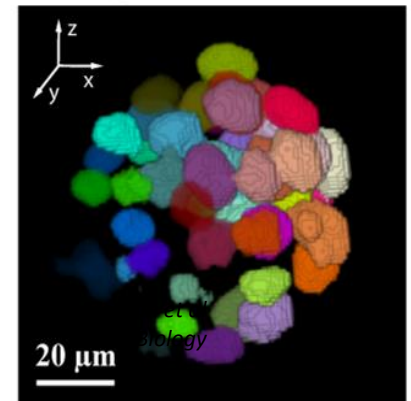
Zebrafish



Time-lapse



3D



Software (e.g., CellProfiler) can quantify cells



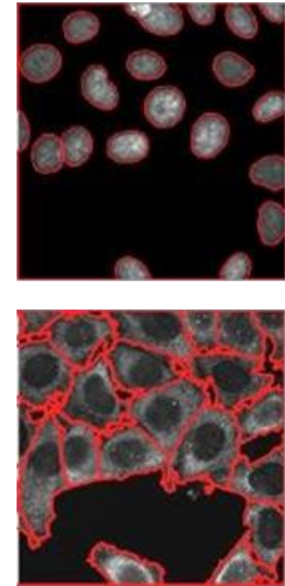
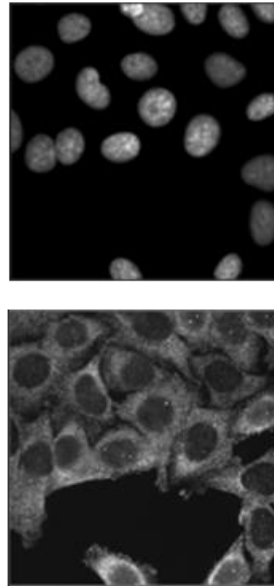
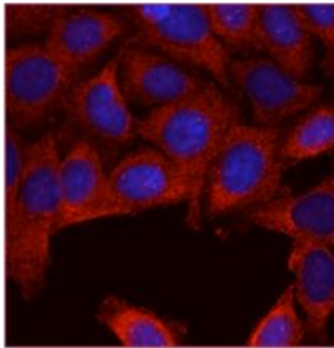
CellProfiler™
cell image analysis software

Split
colors

Correct
illumination

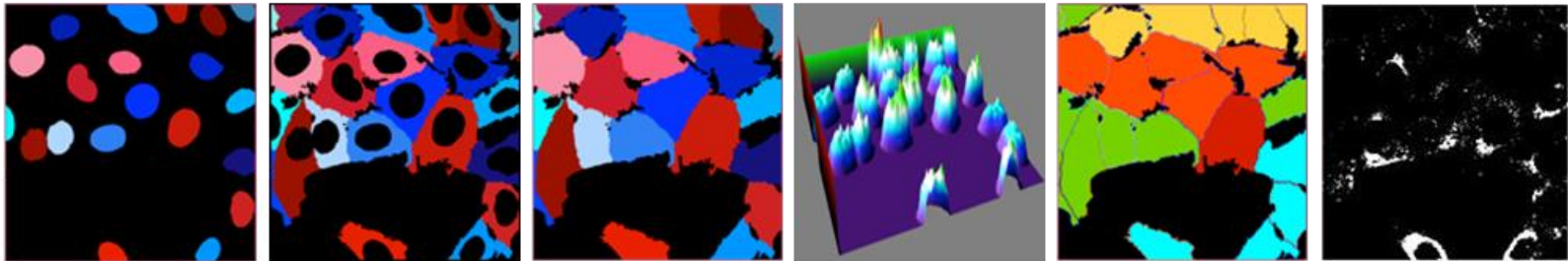
Identify cells/
compartments

Original image

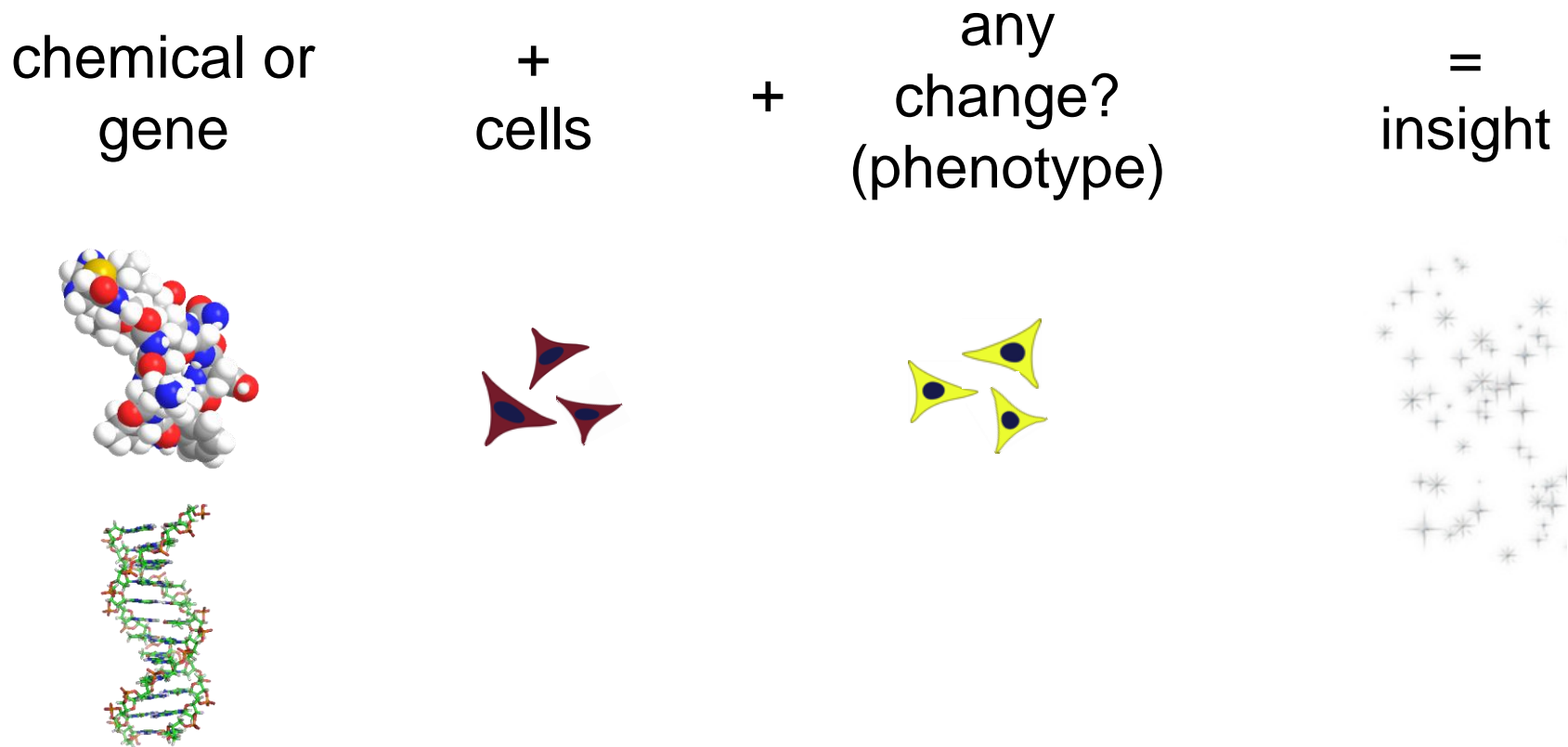


Measure everything

Counts, Shapes, Sizes, Intensities, Textures, Correlations, Relationships

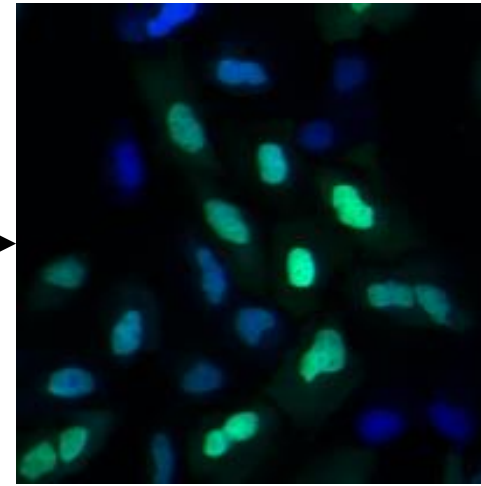
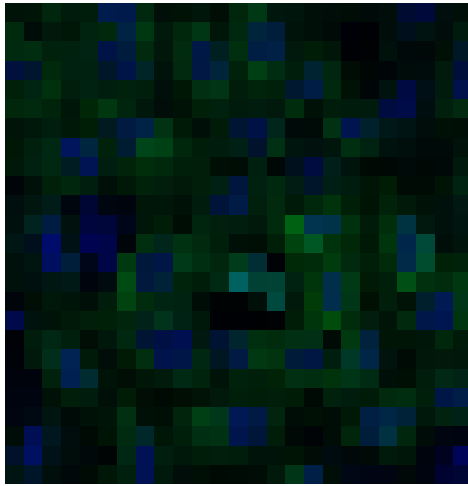


Cell based biomedical research in a nutshell (over simplified)

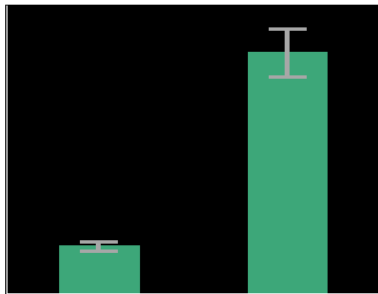


Visual appearance indicate cell state and can be quantified

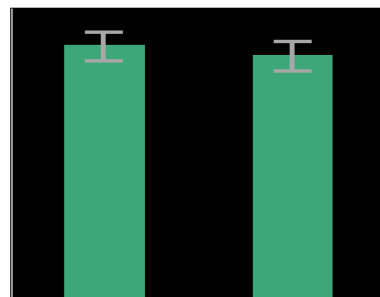
GFP labeled protein which undergoes a translocation from the cytoplasm to the nucleus
in response to perturbations



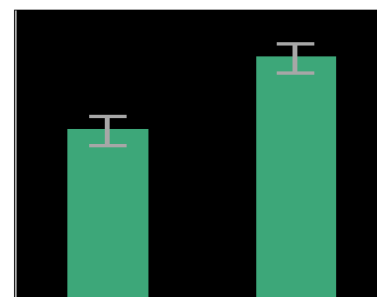
localization



protein levels



morphology



... + hundreds
of other features

Automated image analysis is

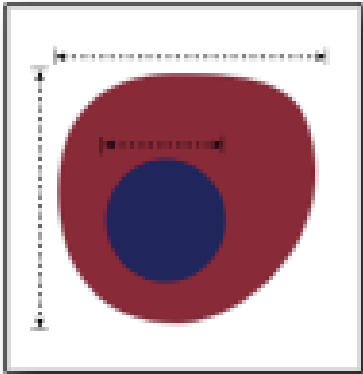
- Objective
- Quantitative, with statistics
- Measure multiple properties at once
- Distinguishes subtle changes, even those undetectable by eye
- Faster, less tedious

Over automation? Advise for scientists who are not tool builders

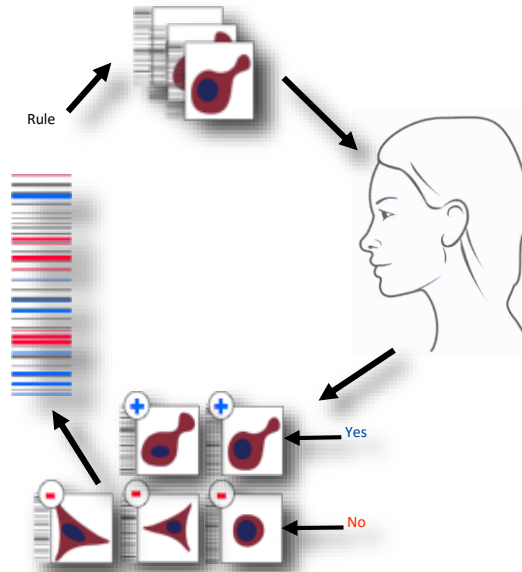
- Your goal is answering a biological question?
- Fully-automated analysis is not always worth the time/effort. Sometimes manual is better!
 - Quantity of data?
 - Accuracy needed?
- Semi-automation can save development time, and improve accuracy
- Application specificity: use specific application knowledge (at the cost of reduced generalization)
- We are not all tool builders!

Three waves of quantitative image analysis

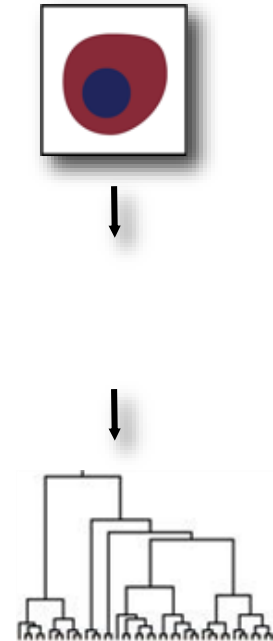
Measure known phenotypes



Train for known phenotypes



Discover new phenotypes



Drug discovery (we wish!)



Tradeoffs in cell biology (& drug discovery)

Physiological relevance vs. experimental complexity

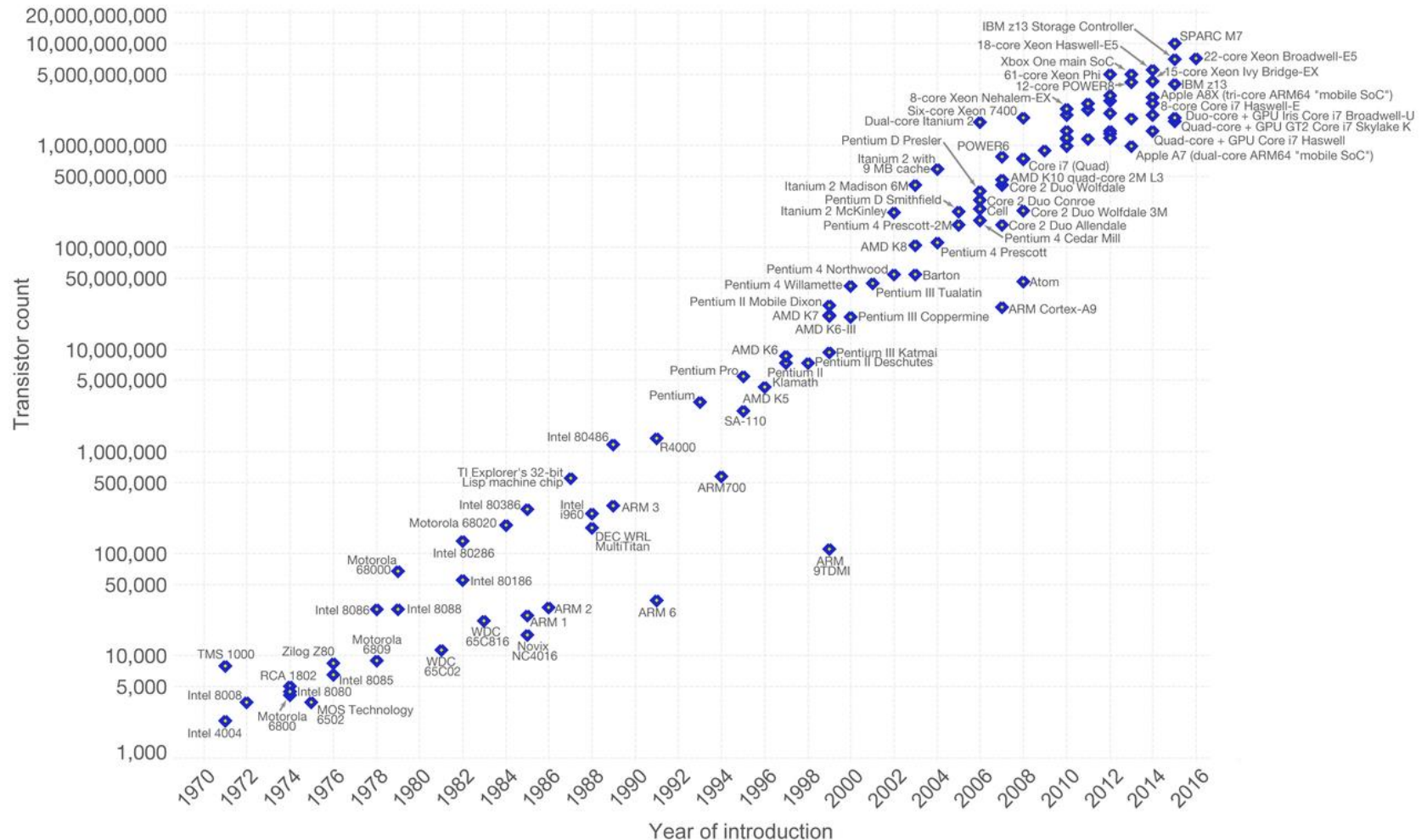
- Physiological relevance:
 - In vivo
 - Resolution (spatial and temporal)
 - 3D
 - In context of microenvironment
- Experimental complexity:
 - Technology
 - Complexity of experiment
 - Costs
 - Amount of data collected
 - Complexity of data analysis

The Moore's law

Moore's Law – The number of transistors on integrated circuit chips (1971-2016)



Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.



Data source: Wikipedia (https://en.wikipedia.org/wiki/Transistor_count)

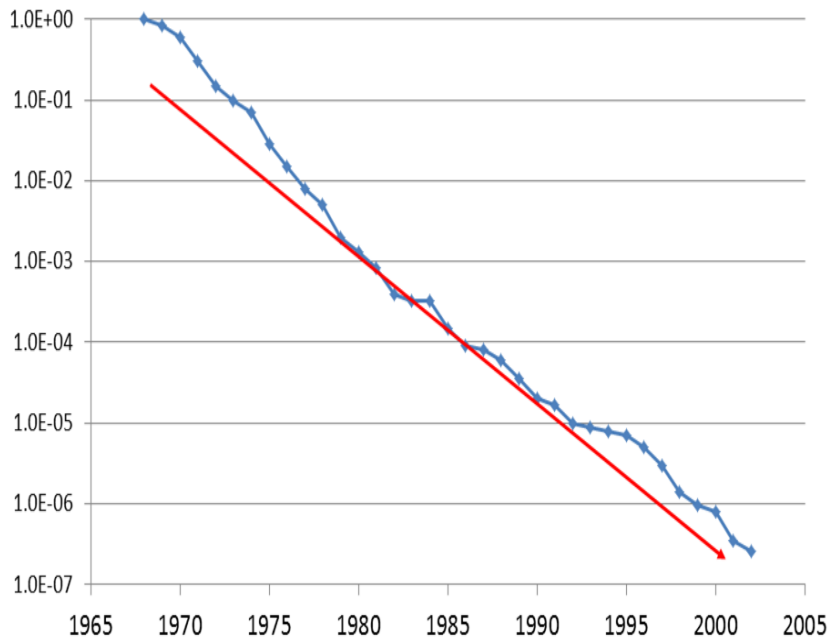
The data visualization is available at [OurWorldinData.org](https://www.ourworldindata.org). There you find more visualizations and research on this topic.

Licensed under CC-BY-SA by the author Max Roser.

Tale of two industries

Moore's Law

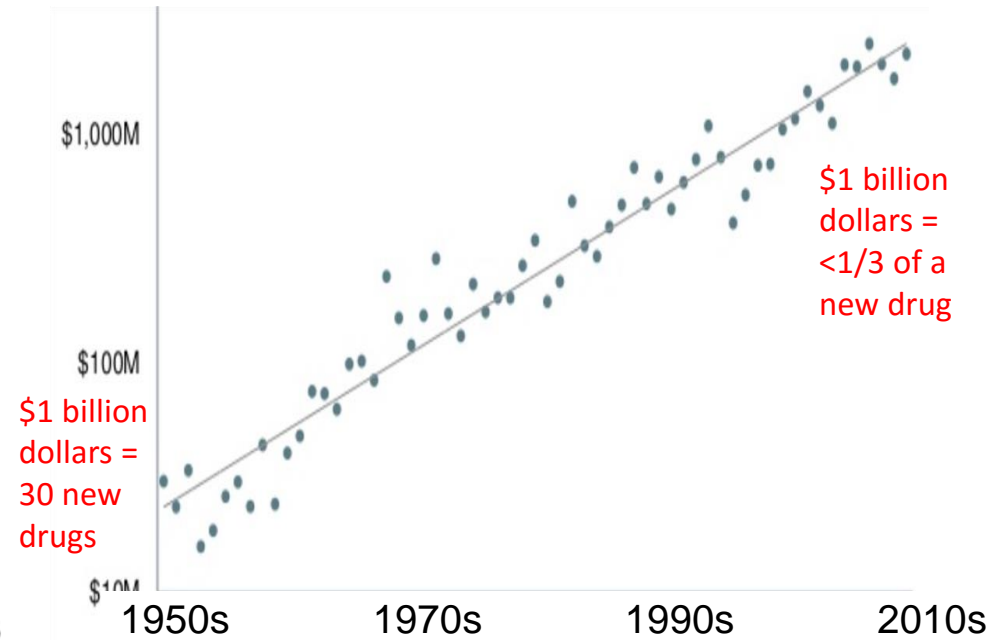
Compute gets cheaper and cheaper



SweptLaser.com

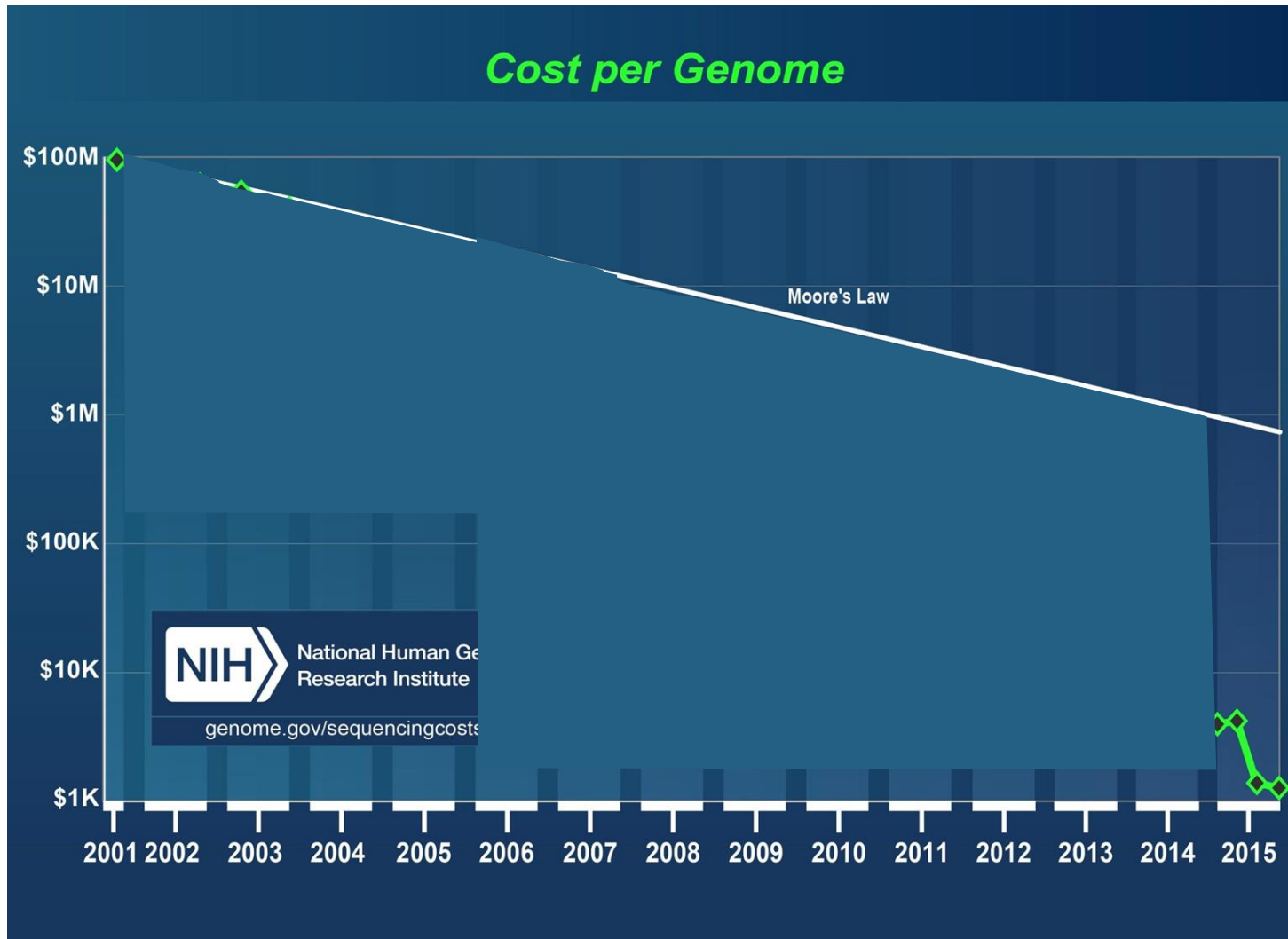
Eroom's Law

Discovering new medicines gets more and more expensive

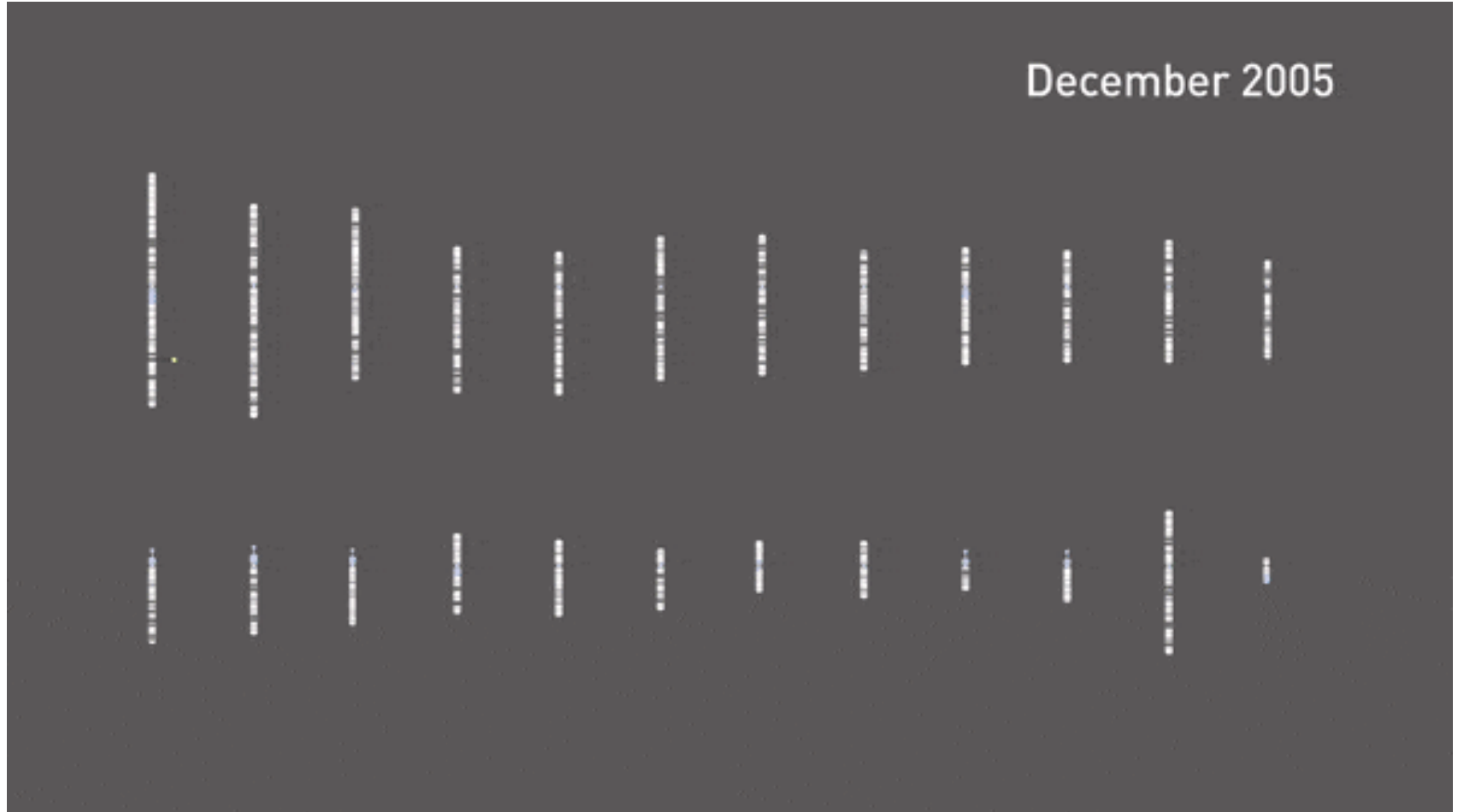


Scannell et al. (2012)

Sequencing costs have dramatically declined



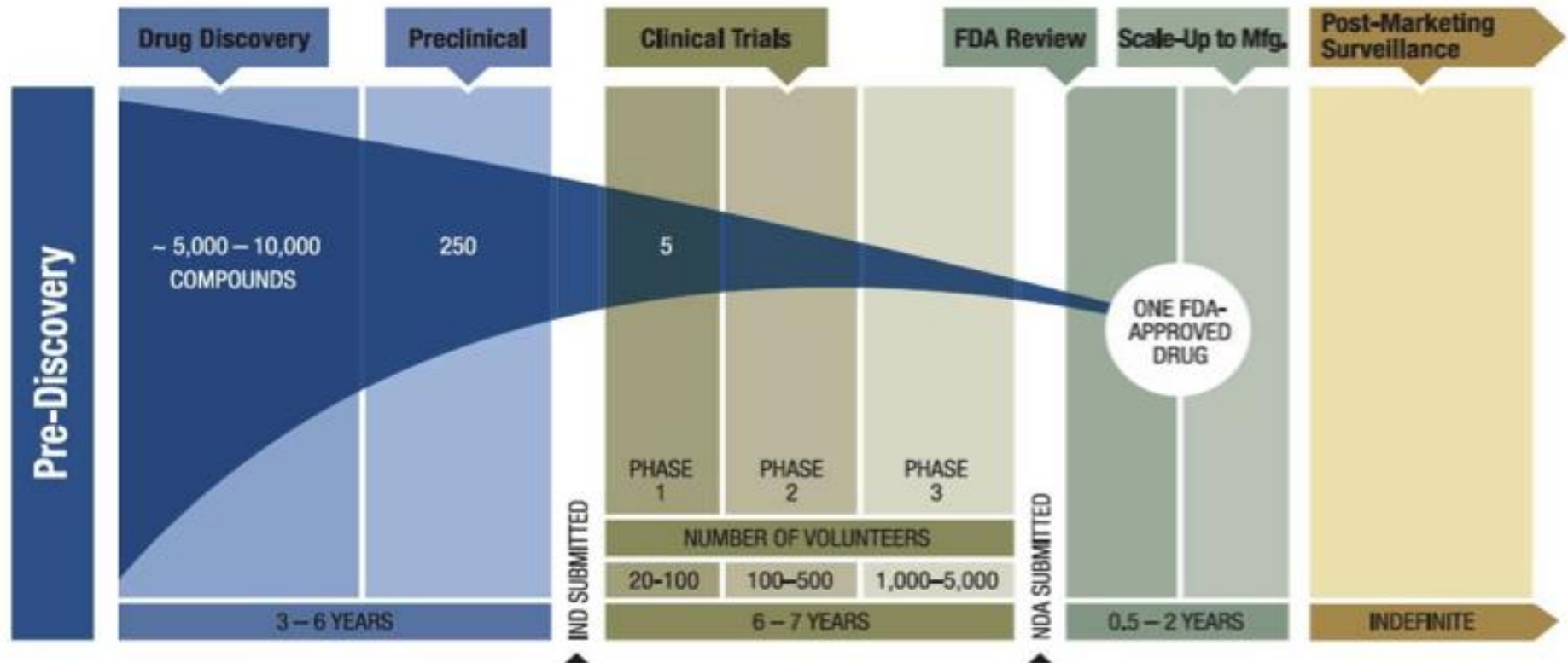
Sequencing is uncovering the cause of many diseases



Genome regions associated with traits or disease through genome-wide association studies

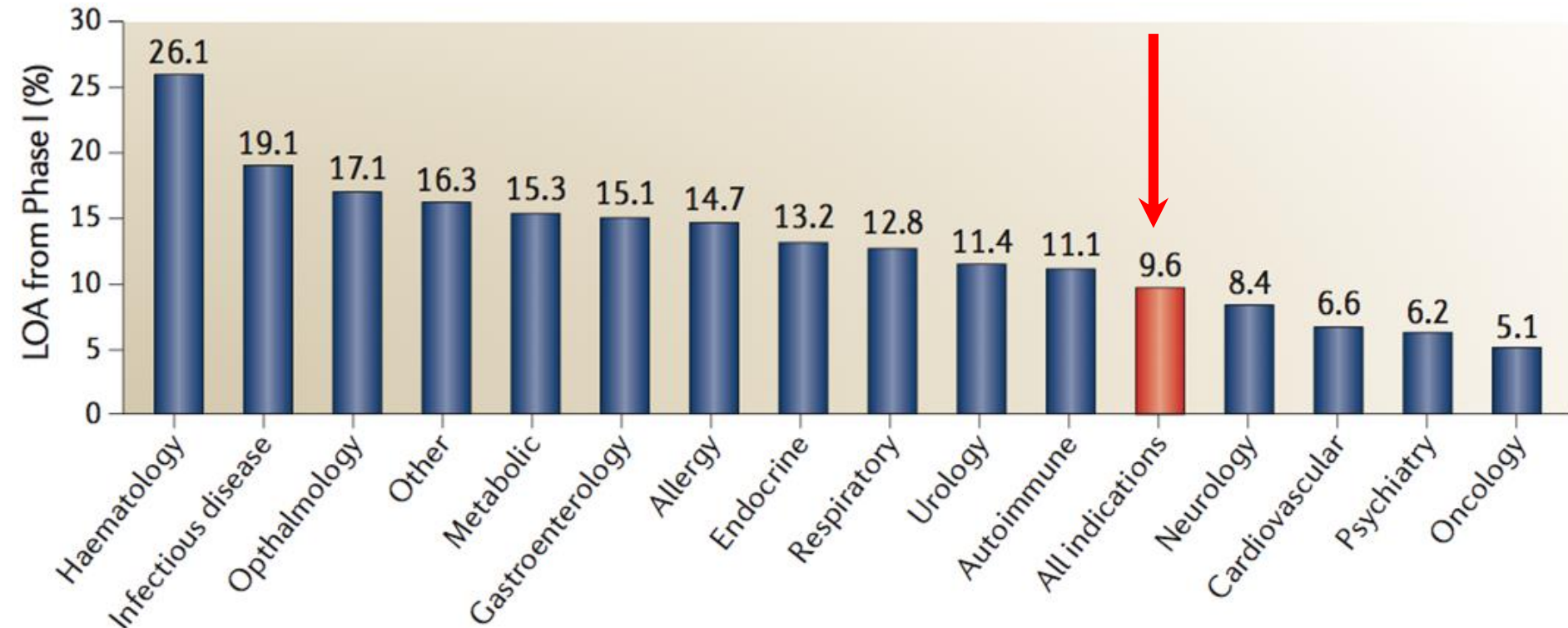
Drug discovery is time consuming

Drug Discovery and Development Timeline



Drug discovery fails most of the time

Success rate of a drug entering Phase I clinical trials
Average = 9.6% make it to market



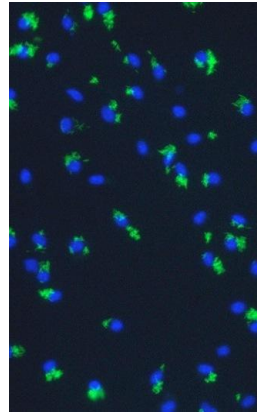
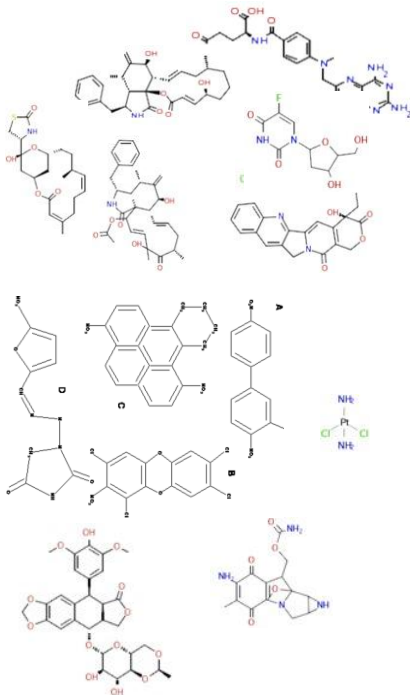
Discovering drugs in high throughput?

millions of
chemicals

+
cells



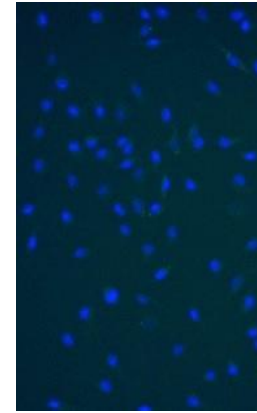
readout
(phenotype) =
drug!



cells

+

tuberculosi
s bacteria



cells are
alive!

+

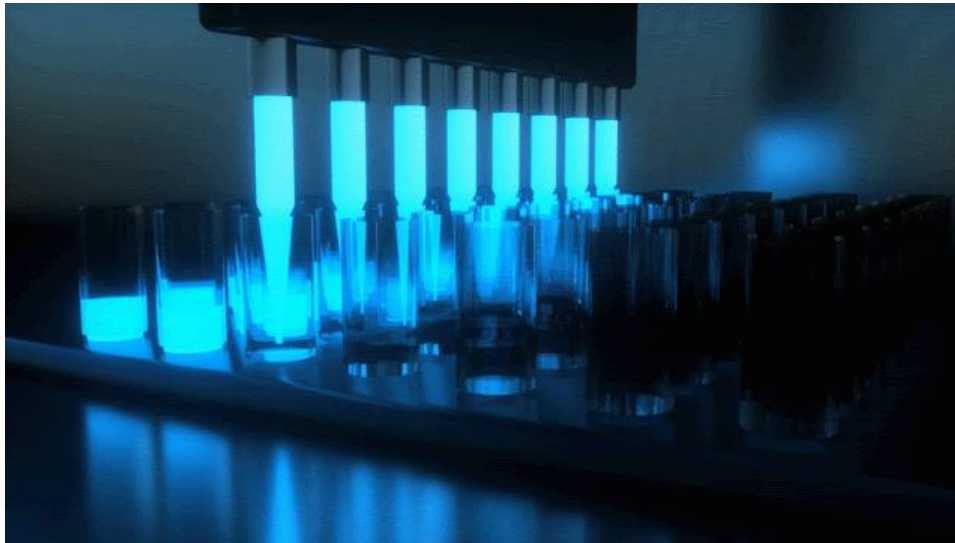
tuberculosi
s is gone!



Discovering drugs in high throughput

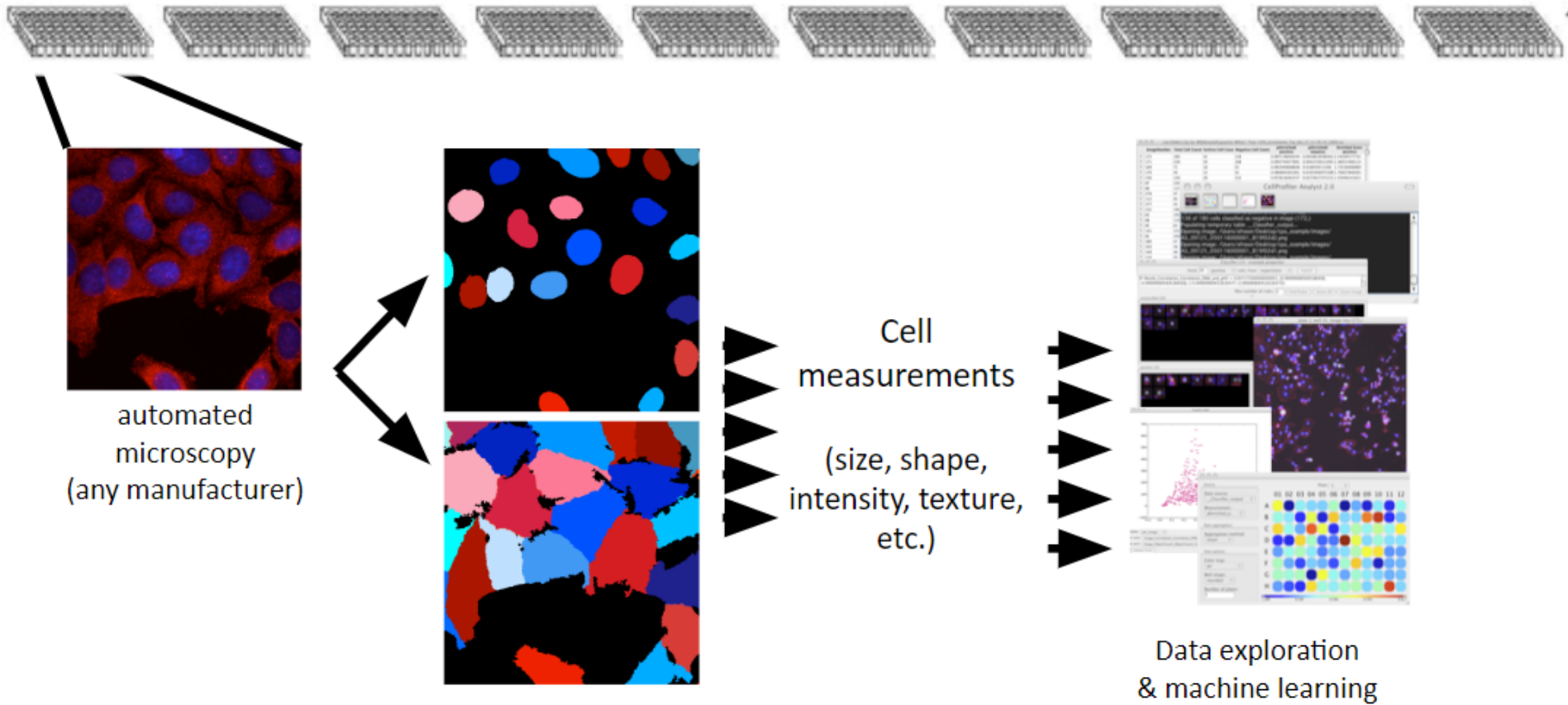


X 384 =



Large scale imaging experiments

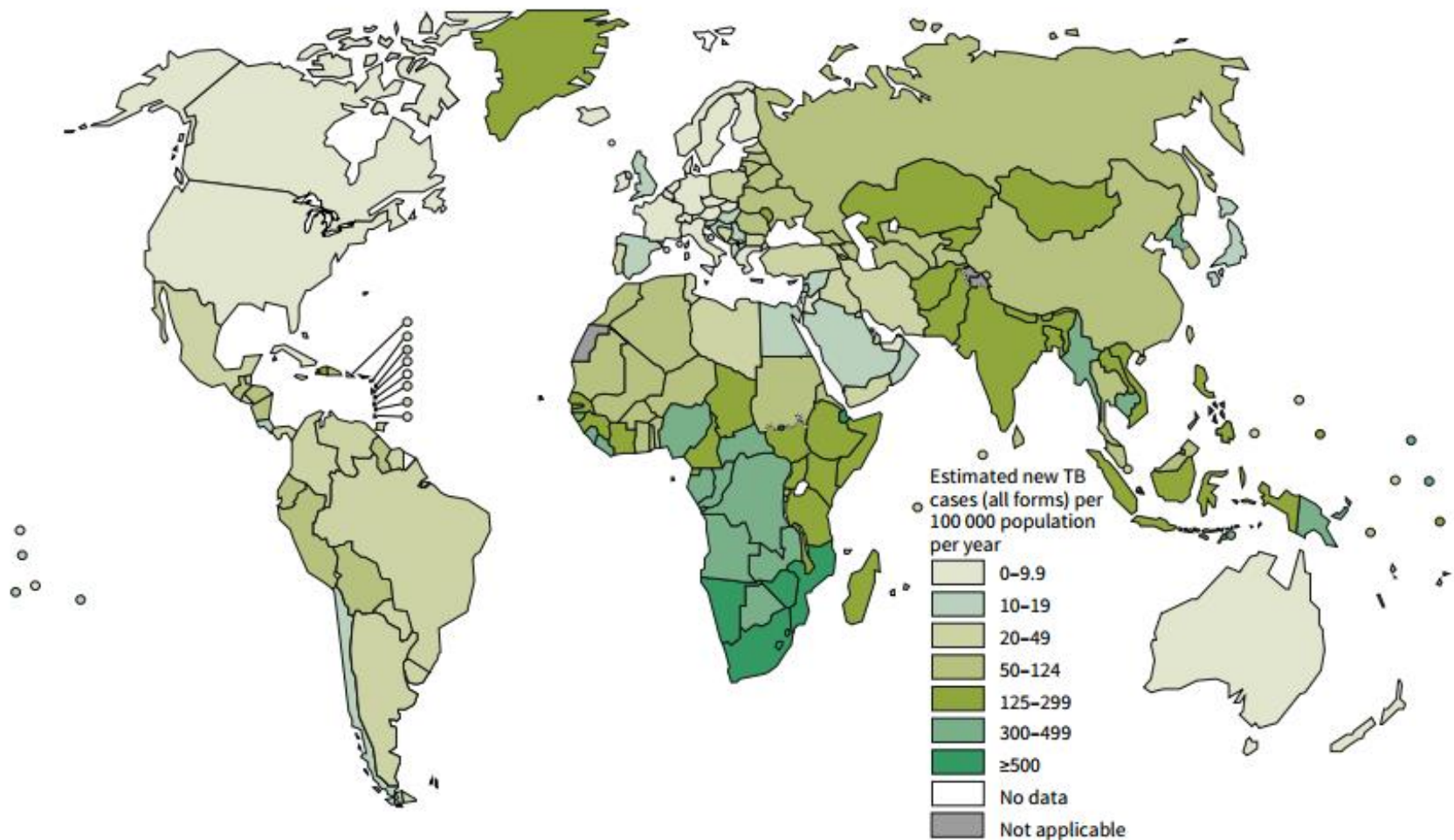
Cells or organisms in multiwell plates, each well treated with a gene or chemical perturbant



Case study: Tuberculosis

Remains a leading cause of mortality globally
1/3 of the world is latently infected

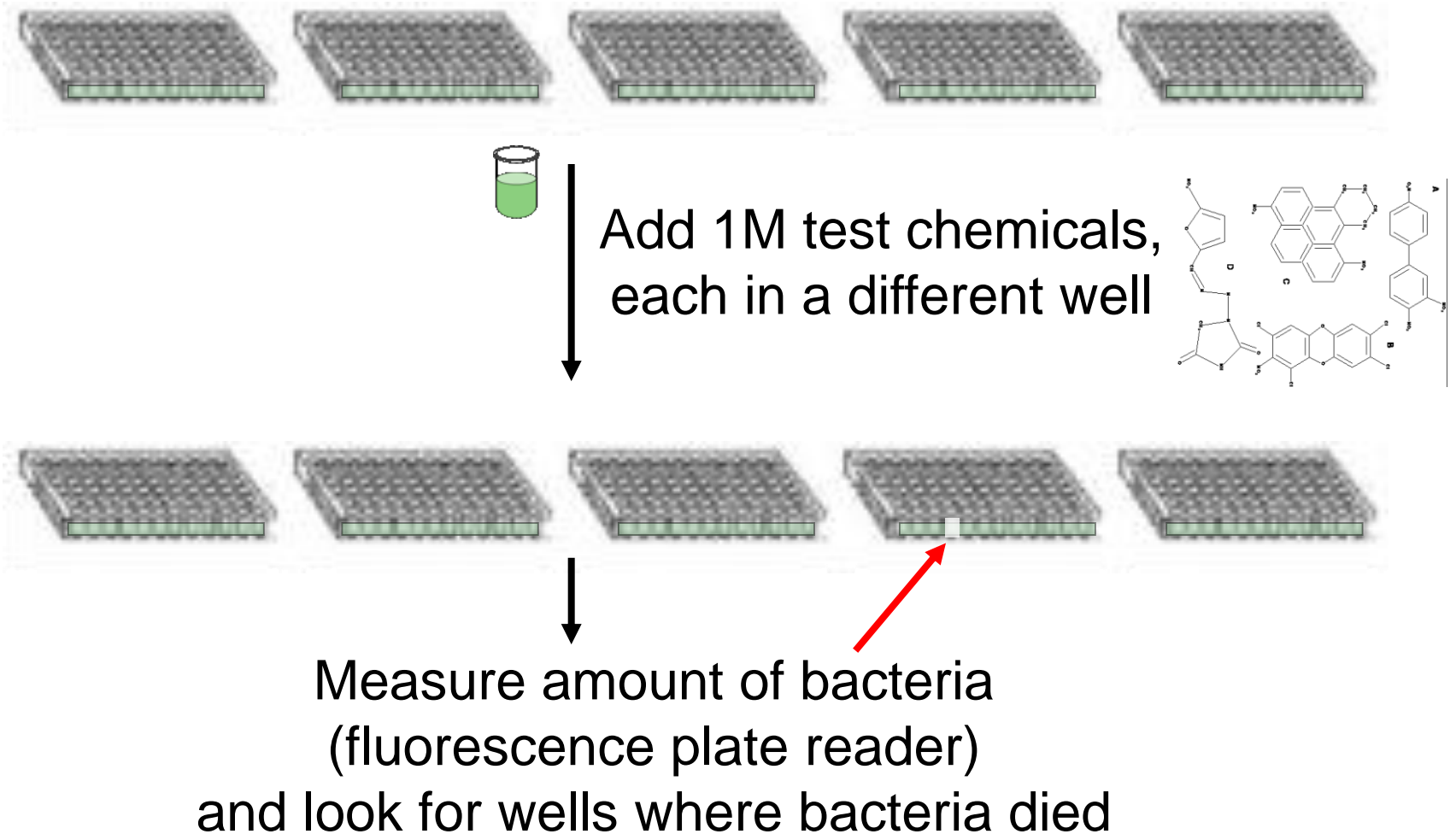
Estimated TB incidence rates, 2013



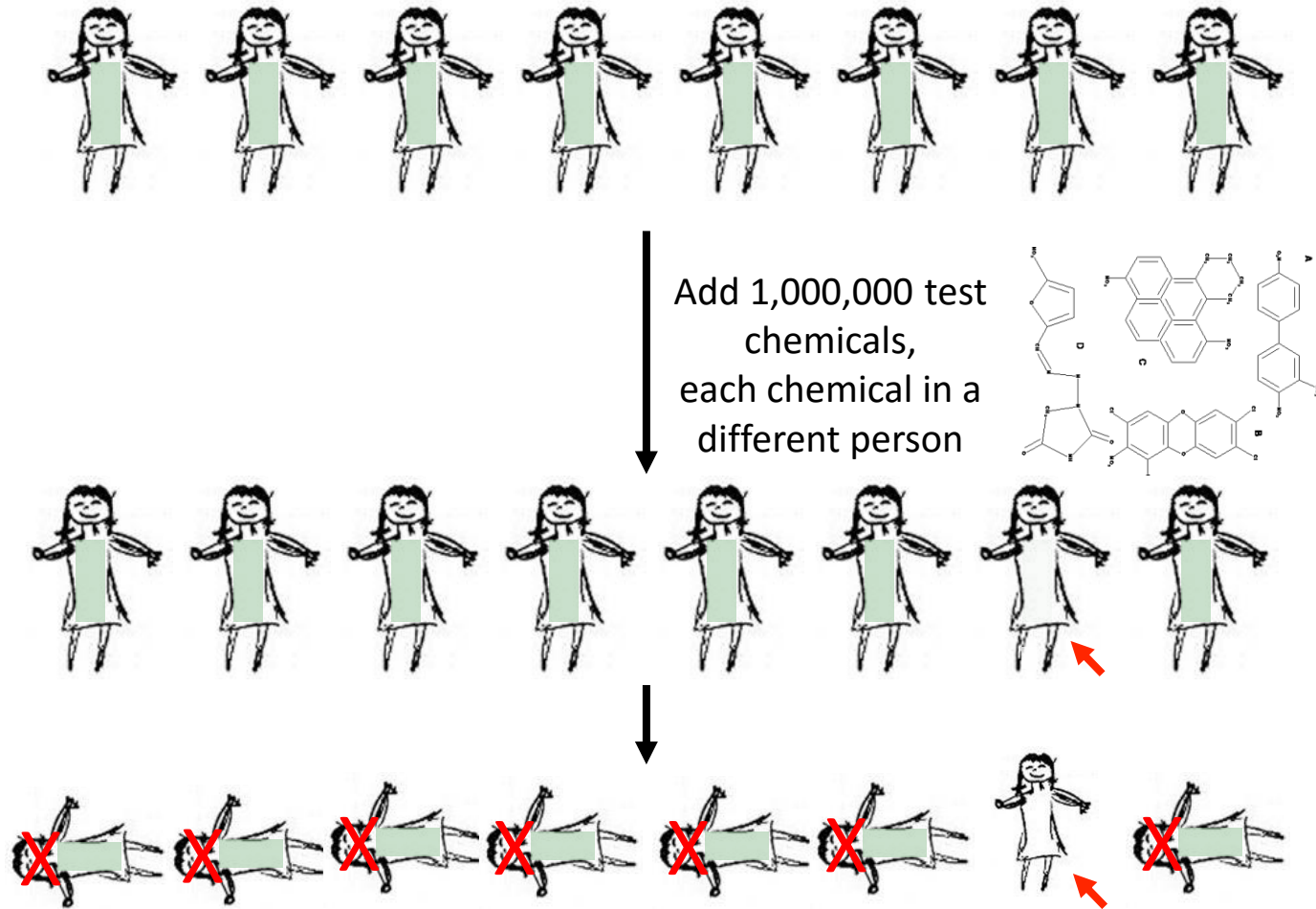
Anne Carpenter, from WHO report, Global Tuberculosis Control 2013

Traditional approach to discover new antibiotics

Try to kill **bacteria** in individual wells of multi-well plates



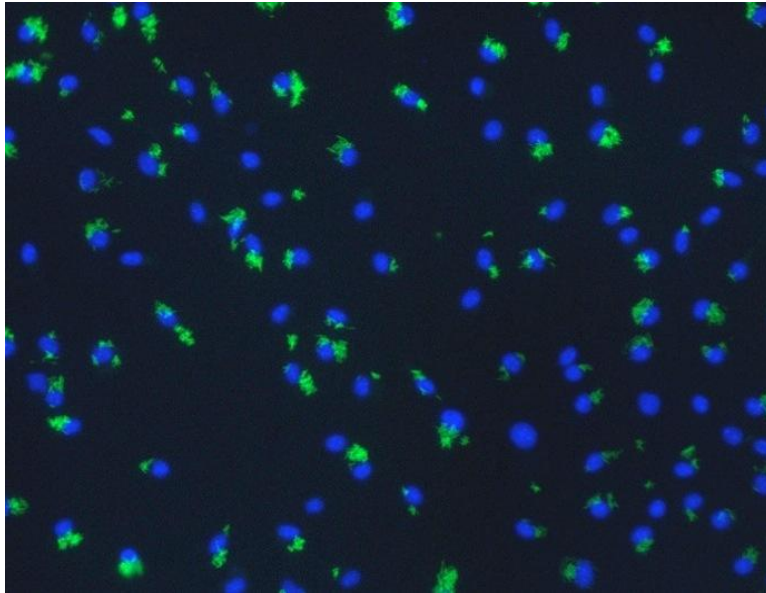
Alternative approach to discover new antibiotics (effective, not ideal)



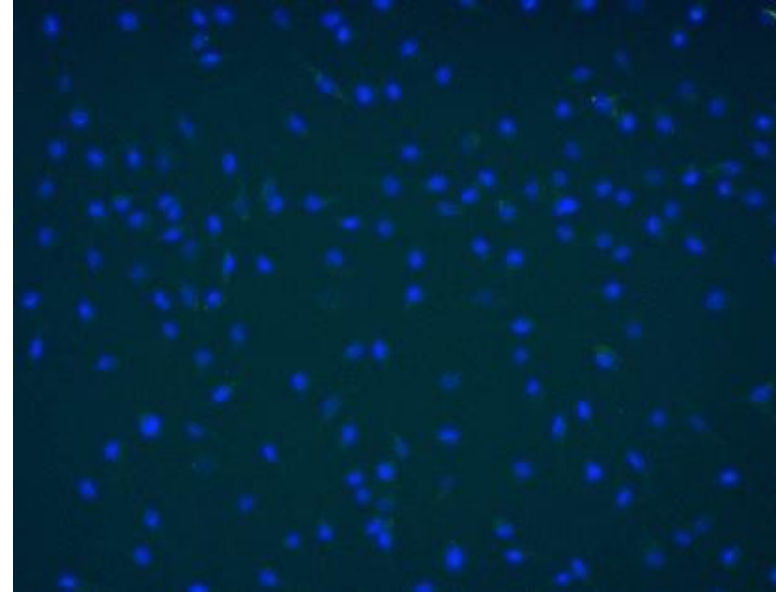
Search for tuberculosis treatments

that prevent infection but do NOT kill the bacterium directly

Without drug



With drug

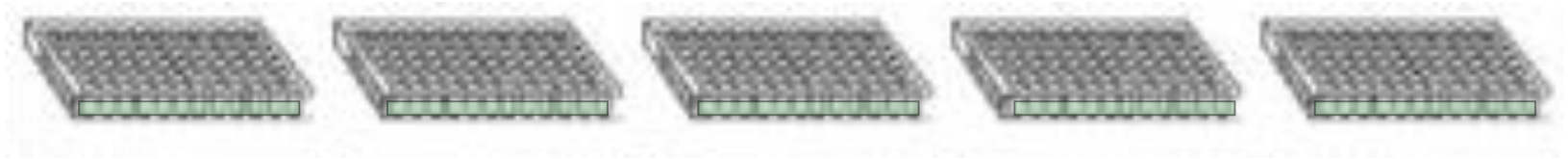


mouse
nuclei

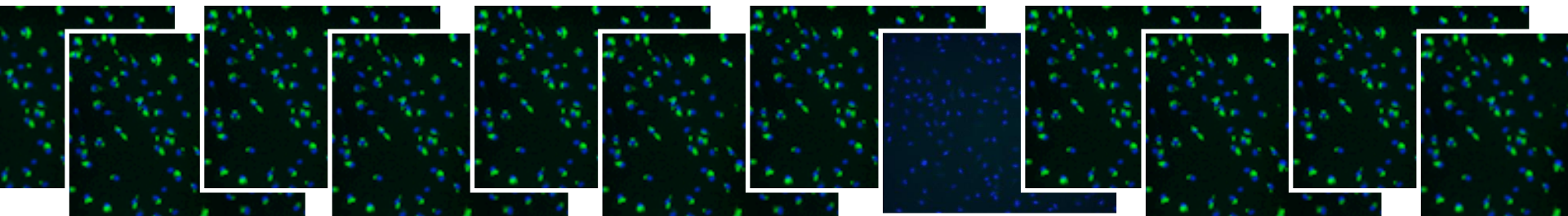
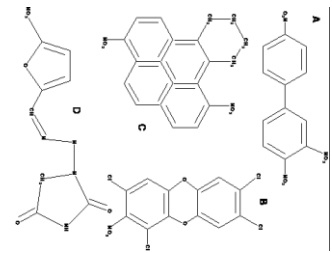
tuberculosis
bacteria

Search for tuberculosis treatments

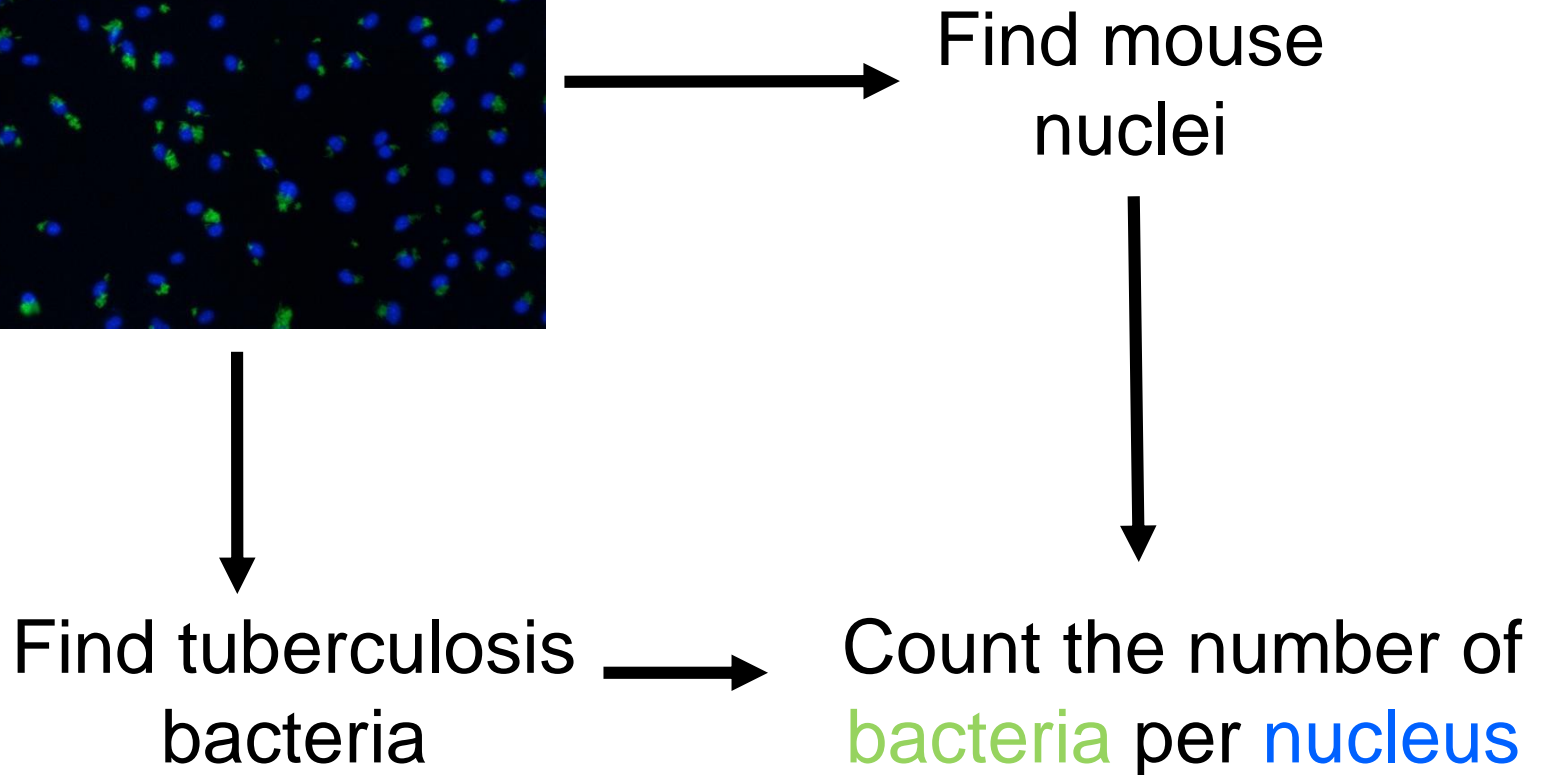
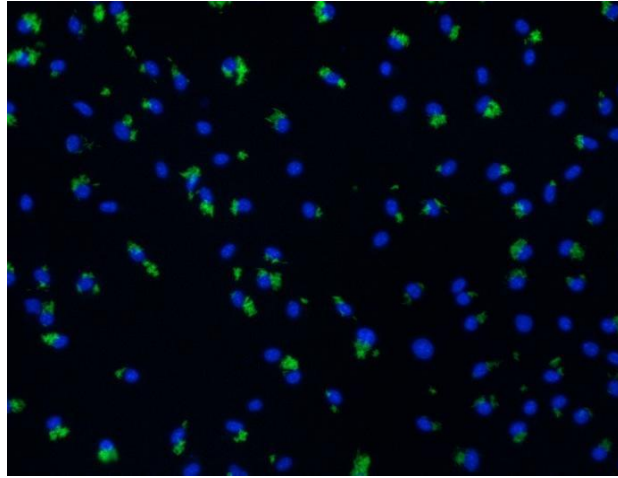
Put **bacteria** and **mouse cells** in individual wells of multi-well plates



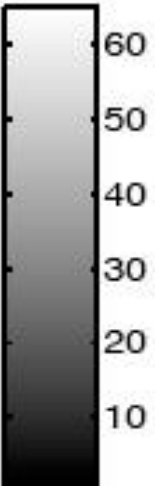
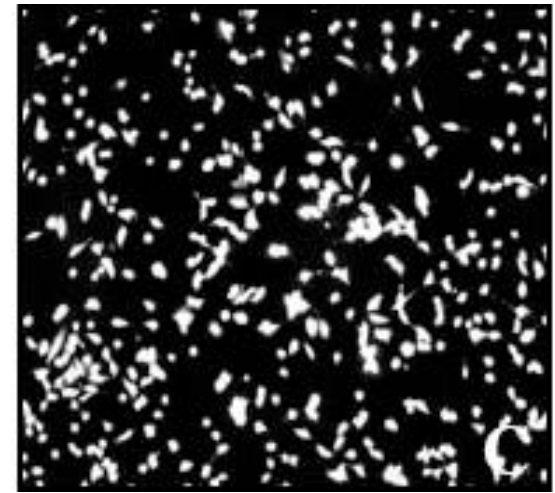
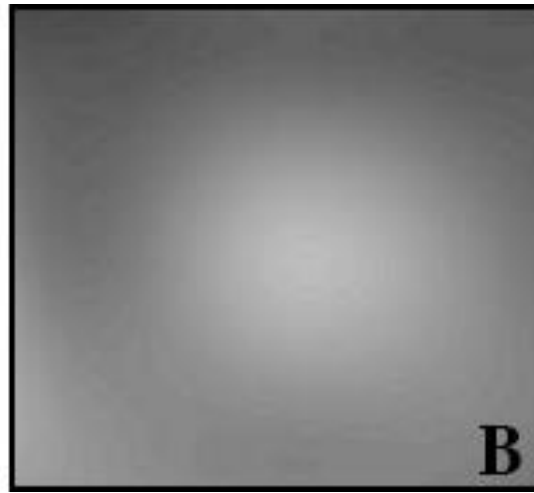
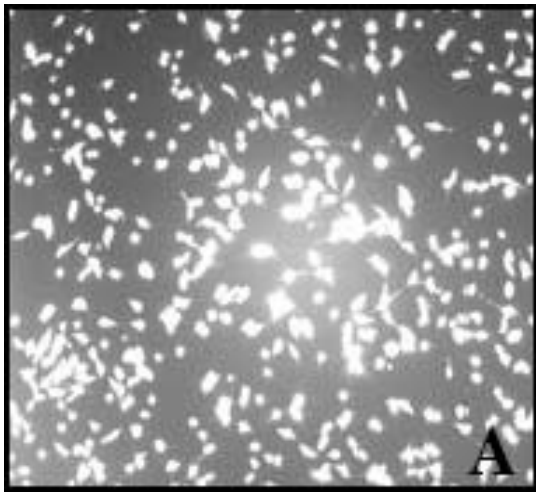
↓
Add 10,000 bioactive chemicals,
each chemical in a different well



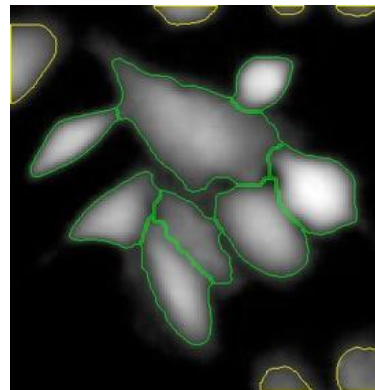
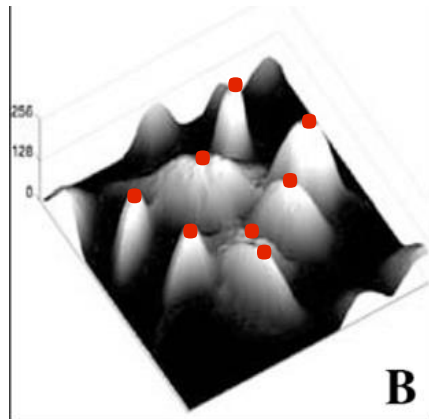
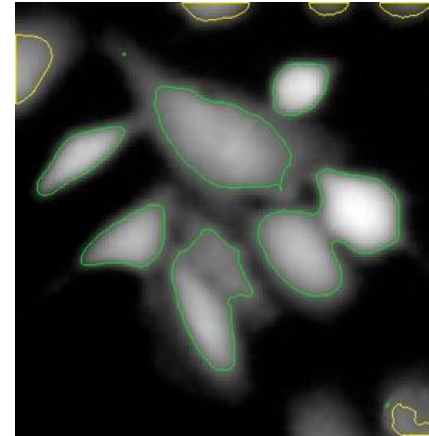
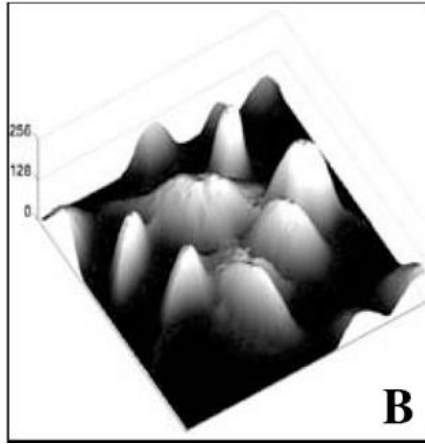
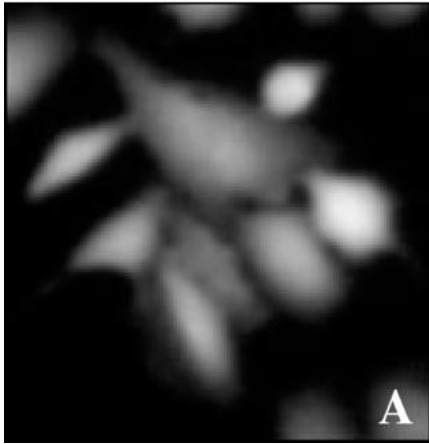
Automated image analysis pipeline



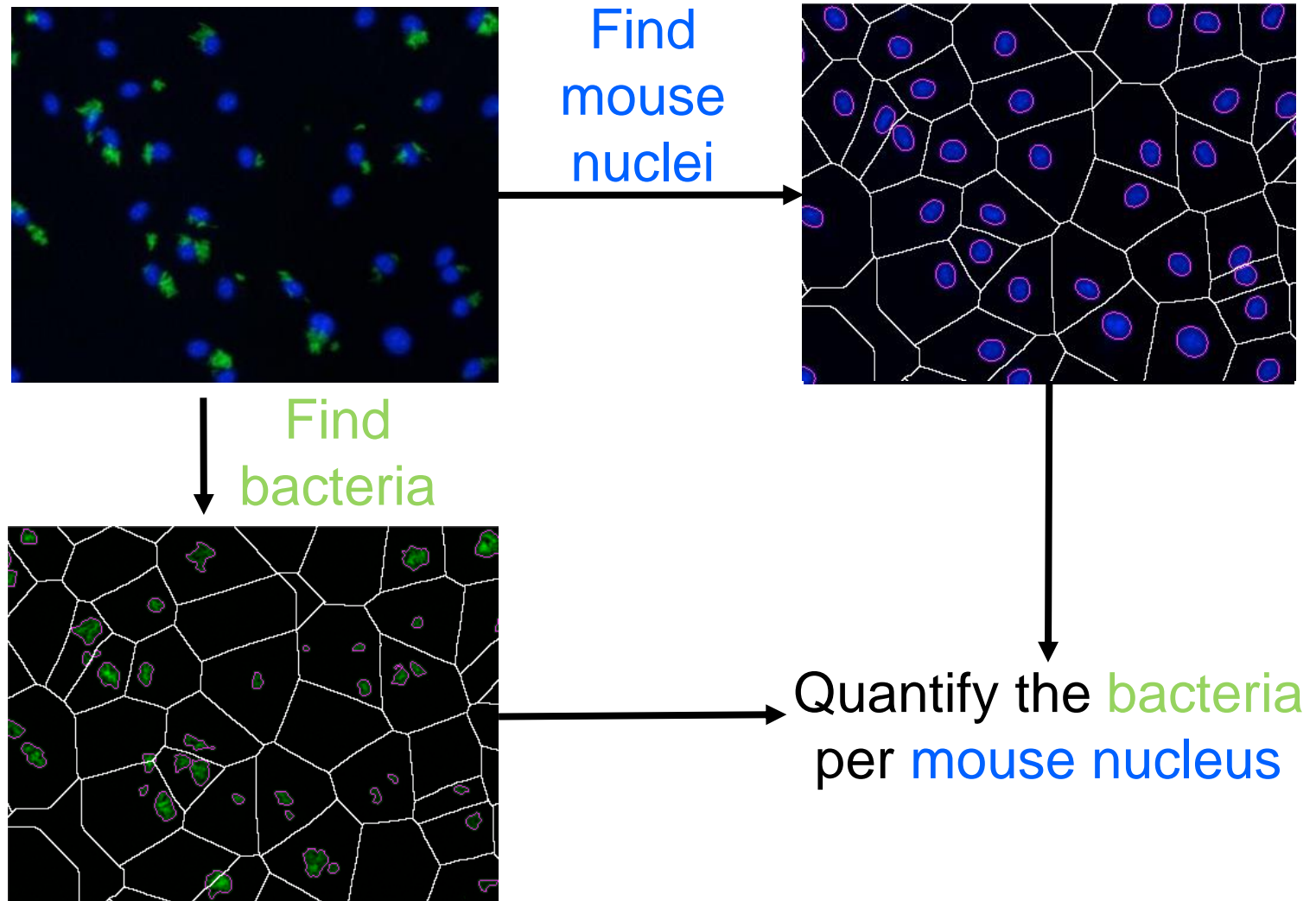
Correct illumination



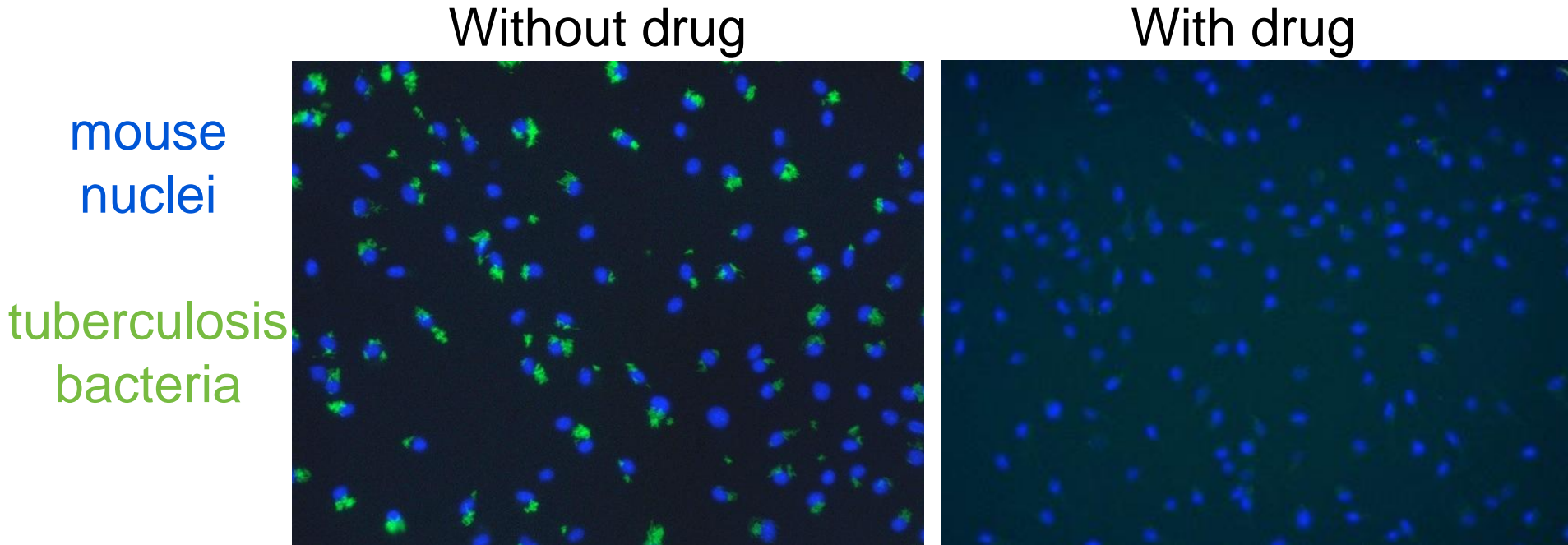
Segment cells/nuclei



Automated image analysis



Search for tuberculosis treatments

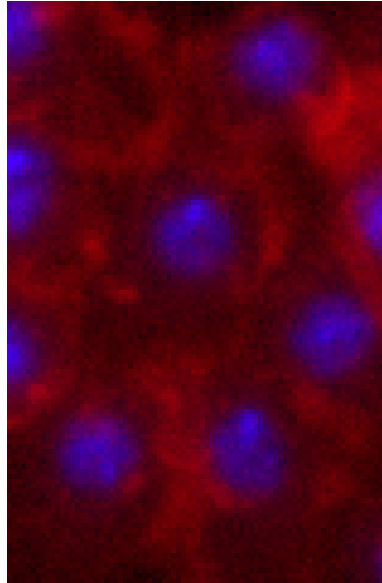


- Identified compounds that prevent bacterial infection/expansion but do NOT kill the bacterium directly
- Gefitinib reduces Mtb growth in the lungs of infected mice

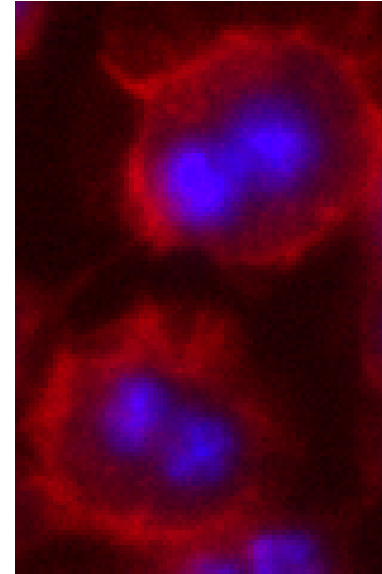
Regulators of cell division

DNA

Actin

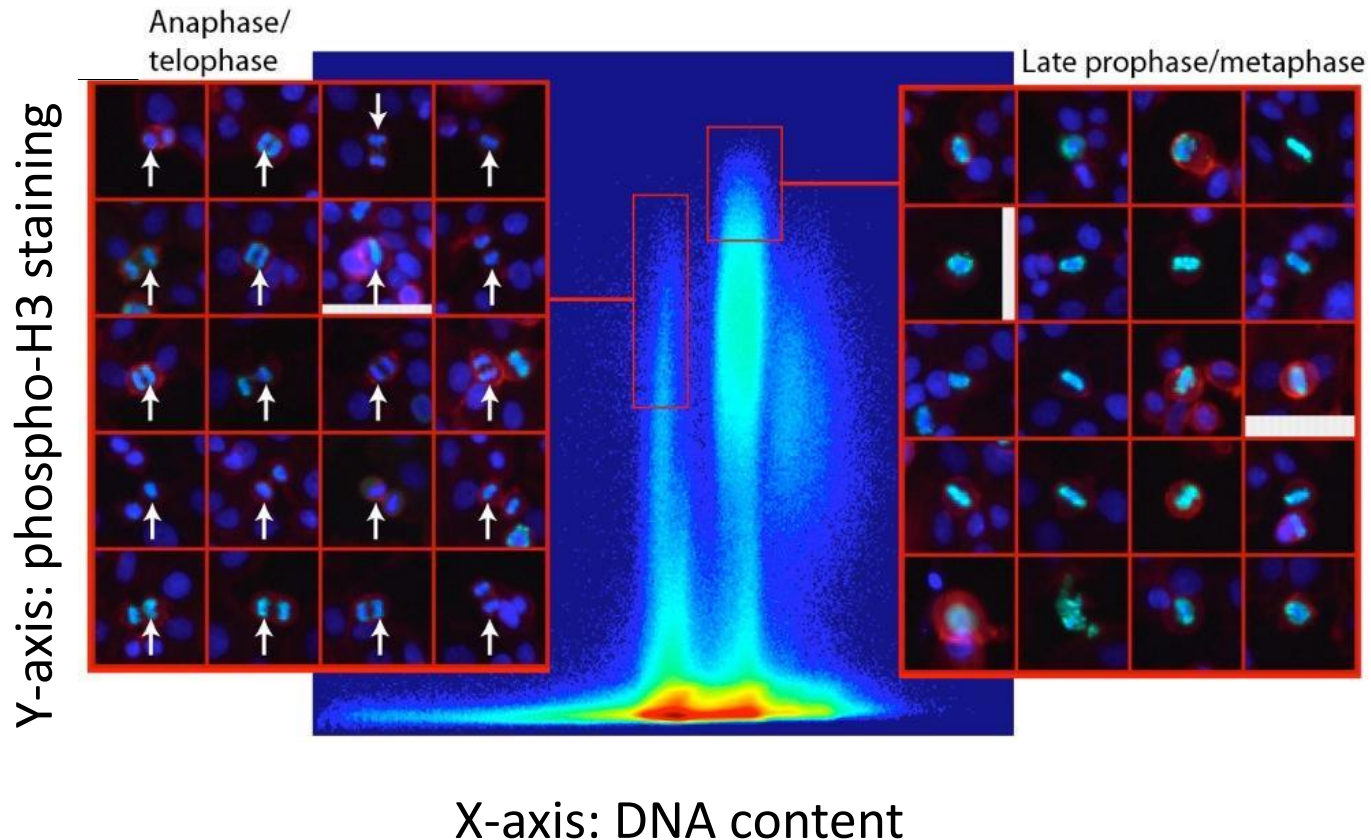


Normal:
one **nucleus**
per **cell**



Abnormal:
two **nuclei**
per **cell**

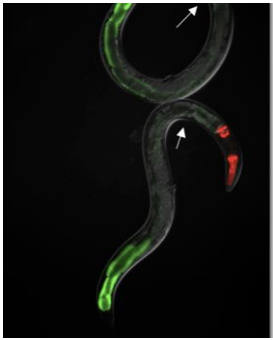
Using > 1 features can identify interesting cells



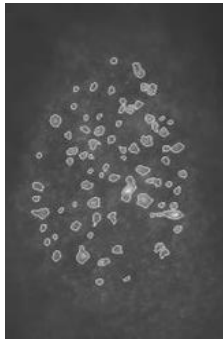
Screen everything!

(Carpenter lab)

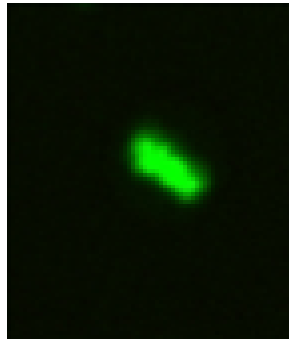
S. aureus
infection



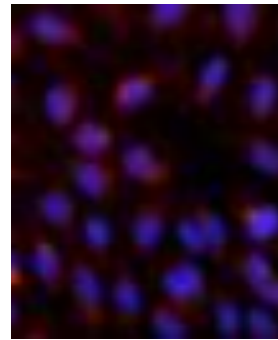
DNA damage



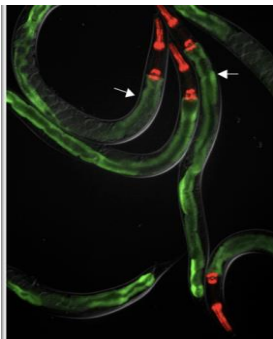
Mitosis



Mitochondrial
abundance



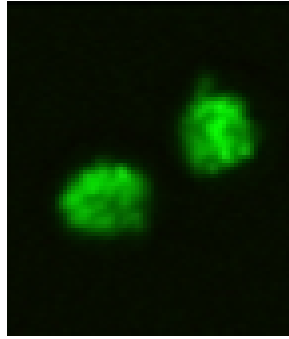
E. faecalis



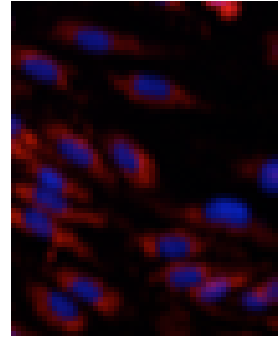
Ausubel/
Irazoqui labs



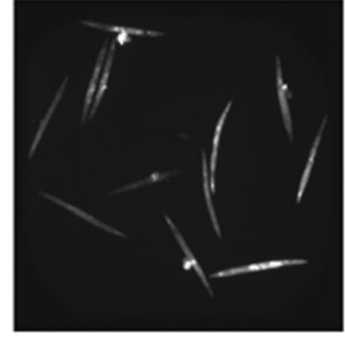
Yaffe lab



Mitchison lab



Mootha lab

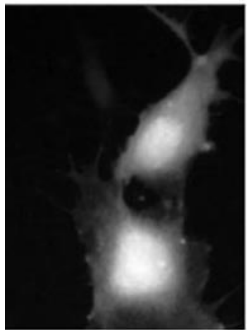
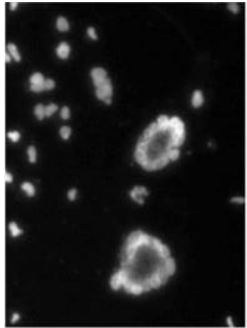


Ausubel lab

Screen everything!

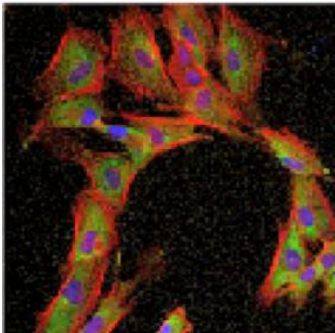
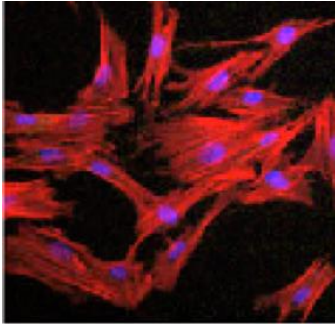
(Carpenter lab)

HIV
neutralization



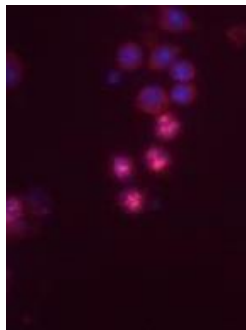
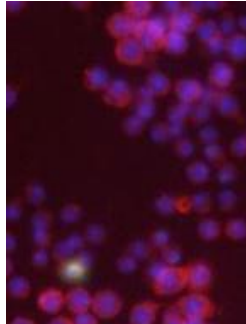
Fenyo lab

mTOR pathway
activation



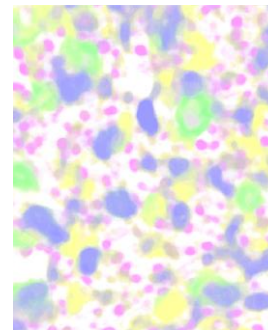
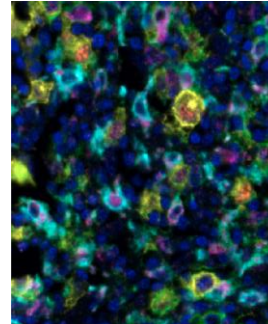
Sabatini lab

Meiosis



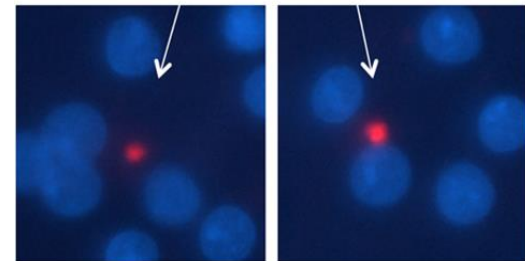
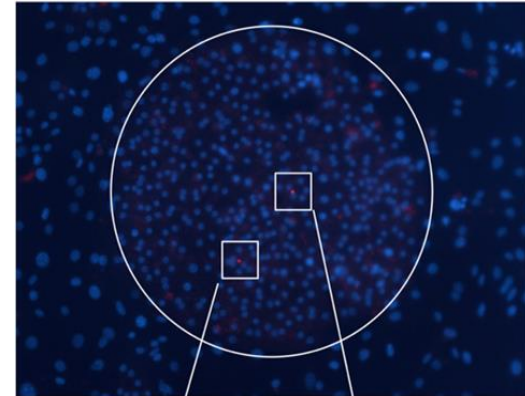
Orr-Weaver
lab

Immune
environment



Shipp/Rodig
labs

Malaria

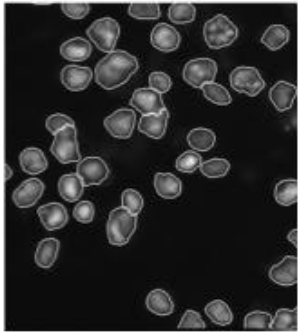


Bhatia lab

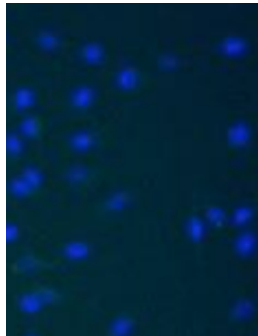
Translational impact

(Carpenter lab, cellprofiler.org/impact)

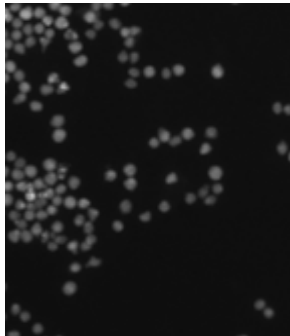
AMKL



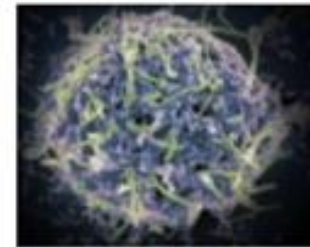
Tuberculosis



Leukemia

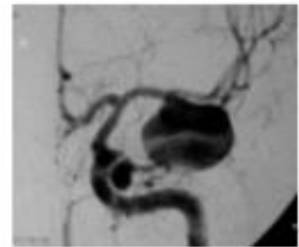


Ebola

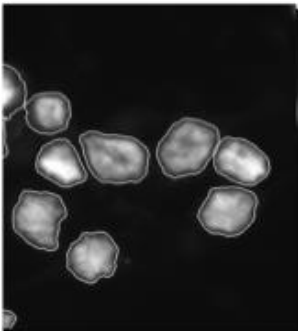


Texas Biomedical
Research Institute

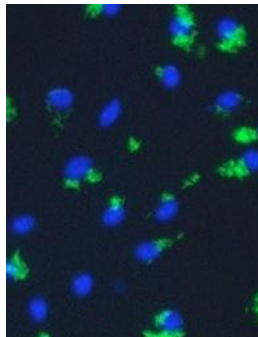
Cerebral
cavernous
malformation



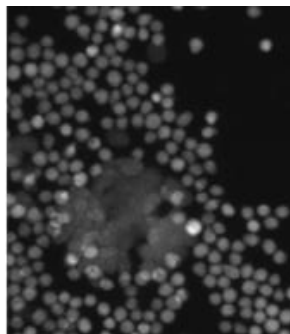
University of
Utah



Crispino lab

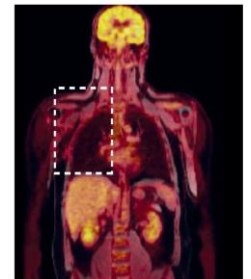
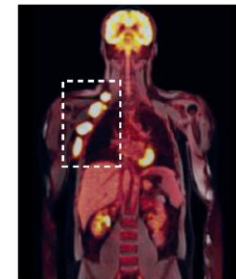


Hung lab



Gilliland/Scadden
/Golub/Schreiber labs

Leukemias &
Lymphomas



Vienna hospitals & medical
institutes, ETH Zurich

A key step is high throughput cell pheontyping is nuclei segmentation

2018 Data Science Bowl: Towards a Universal Nucleus Finder

The screenshot shows the Kaggle website for the 2018 Data Science Bowl competition. The header includes the Kaggle logo, a search bar, and navigation links for Competitions, Datasets, Kernels, Discussion, and Learn. The main banner features the competition title "2018 Data Science Bowl" with the subtitle "Find the nuclei in divergent images to advance medical discovery". It also displays the prize money of "\$100,000" and the number of teams, "3,403 teams". The banner is presented by Booz Allen Hamilton and Kaggle. Below the banner, there are tabs for Overview, Data, Kernels, Discussion, Leaderboard, Rules, and Host. The Overview tab is selected, showing a description of the competition: "Spot Nuclei. Speed Cures." and a brief explanation of the mission: "Imagine speeding up research for almost every disease, from lung cancer and heart disease to rare disorders. The 2018 Data Science Bowl offers our most ambitious mission yet: create an algorithm to automate nucleus detection."

Secure | <https://www.kaggle.com/c/data-science-bowl-2018>

Bookmarks | Defaults | Lit Searches | LabWebsites | Our software | Admin | Personal | Google Maps | ImgP Confluence | MiscCurrent | N

kaggle Search kaggle Q Competitions Datasets Kernels Discussion Learn ...

Featured Prediction Competition

2018 Data Science Bowl
Find the nuclei in divergent images to advance medical discovery

\$100,000
Prize Money

Passion. Curiosity. Purpose.

Presented by
Booz Allen Hamilton · 3,403 teams · 11 days to go (4 days to go until merger deadline) | Allen | Hamilton | kaggle

Overview | Data | Kernels | Discussion | Leaderboard | Rules | Host | **Join Competition**

Overview Edit

Description	Spot Nuclei. Speed Cures.
Evaluation	Imagine speeding up research for almost every disease, from lung cancer and heart disease to rare disorders. The 2018 Data Science Bowl offers our most ambitious mission yet: create an algorithm to automate nucleus detection.
Prizes	

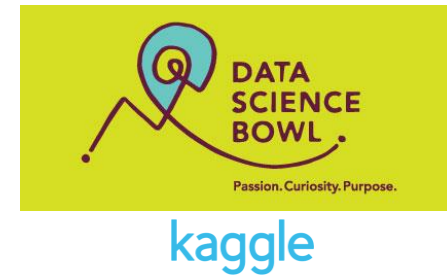
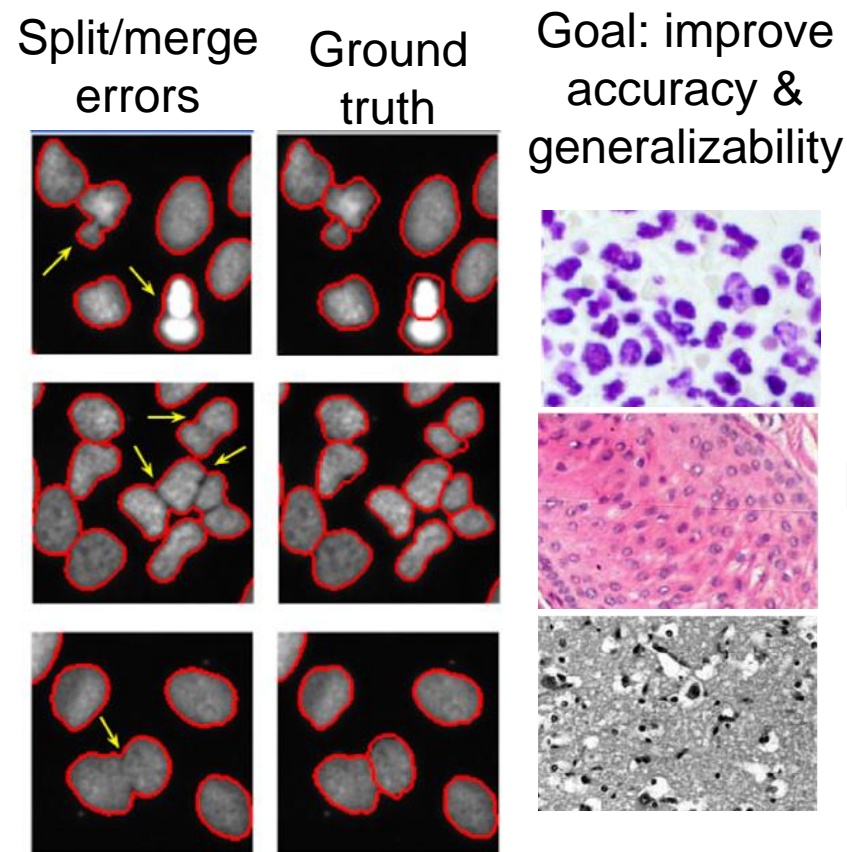
The infographic is titled "Finding the nucleus helps to..." and lists six benefits, each accompanied by a colorful icon. The benefits are: 1. "locate cells in varied conditions to enable faster cures" (target icon), 2. "free biologists to focus on solutions" (wrench and clock icon), 3. "improve throughput for research and insight" (microscope icon), 4. "reduce time-to-market for new drugs— currently 10 years" (syringe icon), 5. "increase # of compounds for experiments" (test tube icon), and 6. "improve health and increase quality of life" (heart icon).

Finding the nucleus helps to...

- locate cells in varied conditions to enable faster cures
- free biologists to focus on solutions
- improve throughput for research and insight
- reduce time-to-market for new drugs— currently 10 years
- increase # of compounds for experiments
- improve health and increase quality of life

A key step is high throughput cell pheontyping is nuclei segmentation

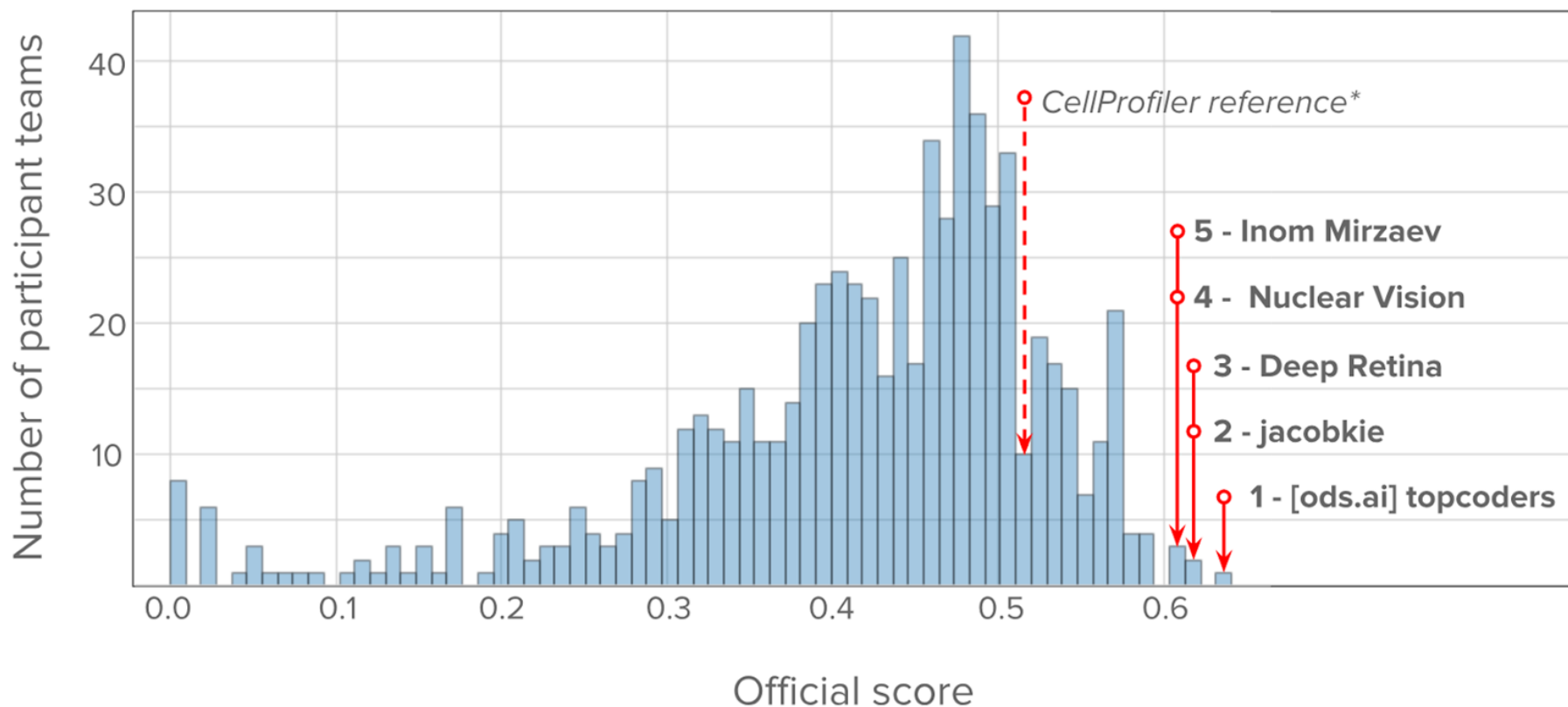
2018 Data Science Bowl: Towards a Universal Nucleus Finder



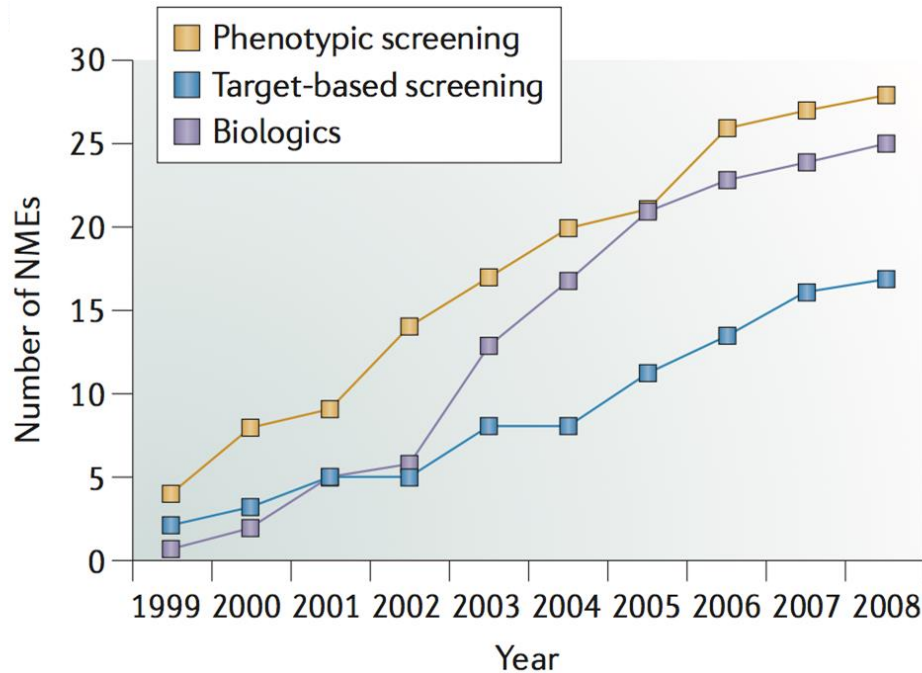
Public dataset: 37,333 outlined nuclei
3,919 teams competed
65,333 submissions
Implemented web app:
www.NucleAIzer.org (Horvath)

Deep learning excels at segmentation

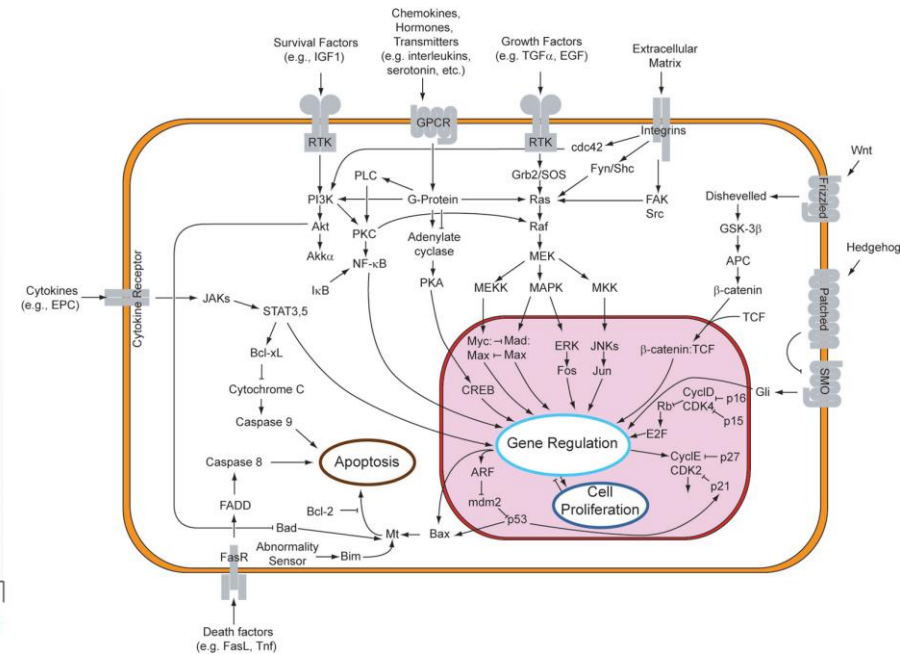
Distribution of scores in second-stage evaluation



Target-based vs. phenotypic drug discovery



Swinney and Anthony (2011),
Zheng et al. (2013),
Lee et al. (2013)

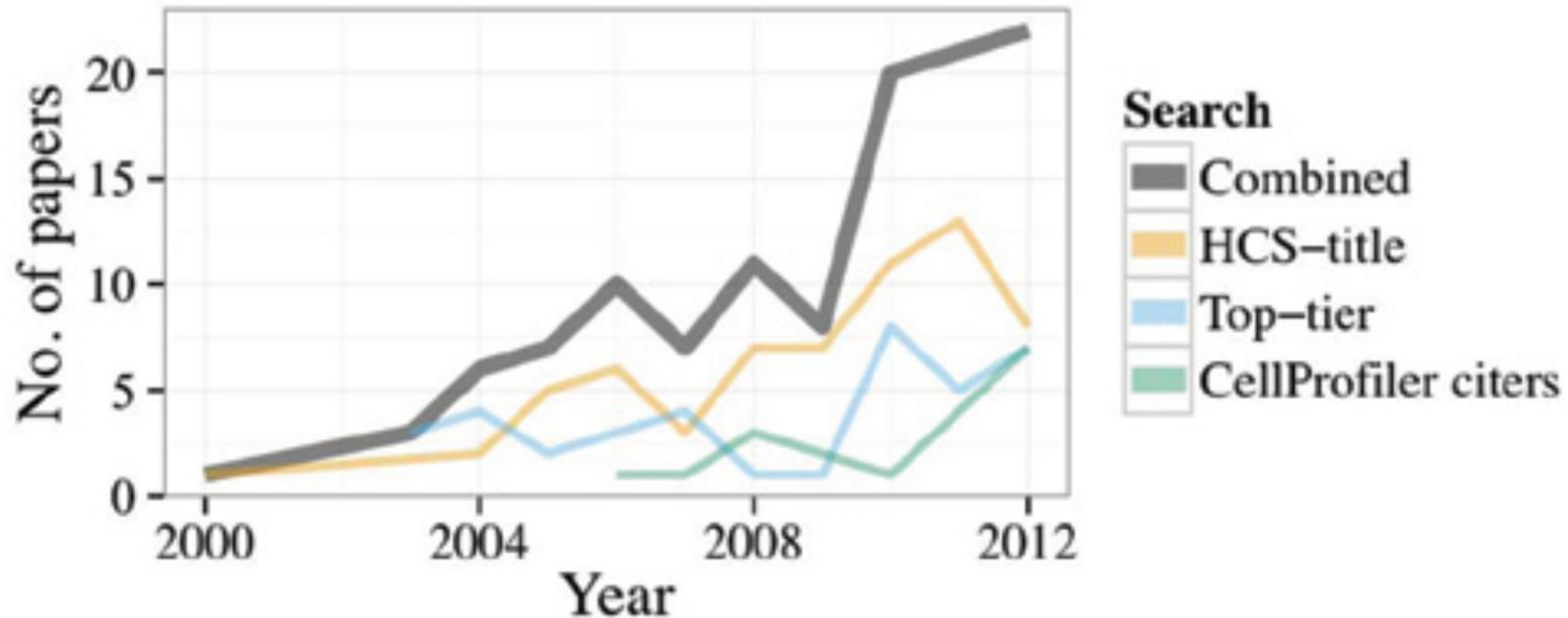


?

Genotype → phenotype

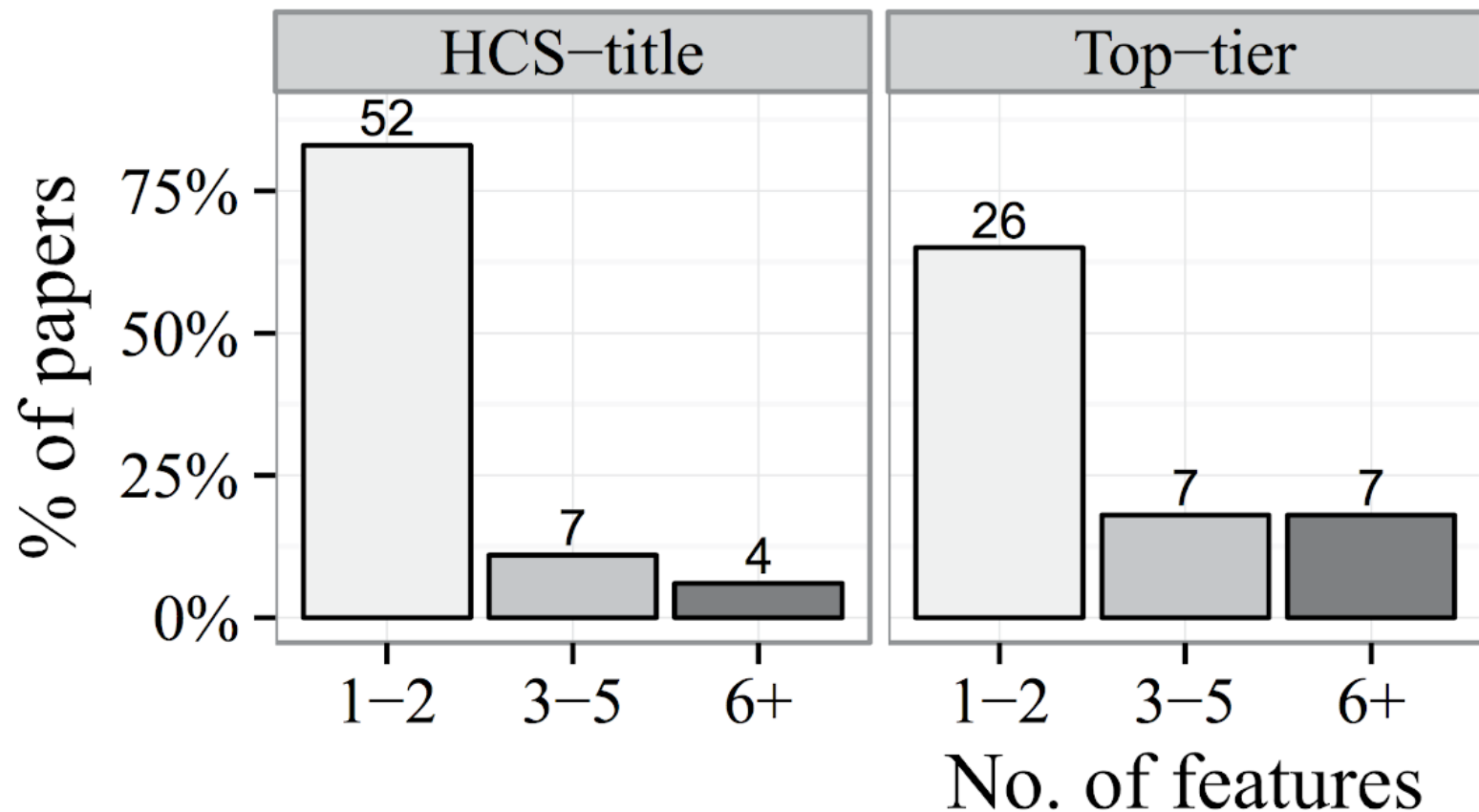
Increased interest in high-content screens

Number of papers in which a high-throughput, image-based experiment was used toward a discovery



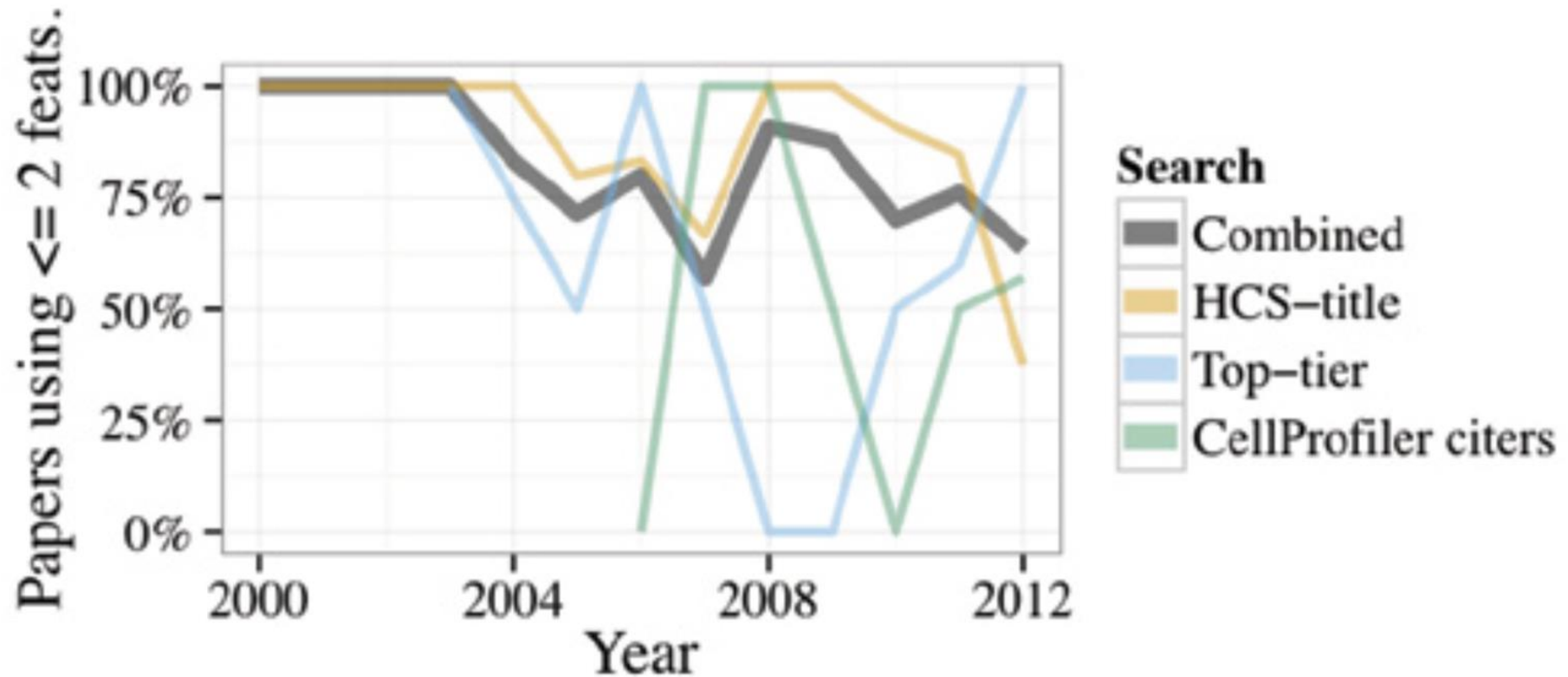
Most high contact screens are not very high in information content...

60-80% of “high-content” studies use only 1 or 2 cellular features

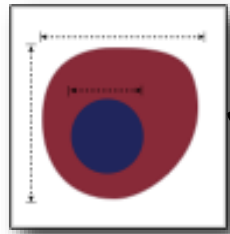


Most high contact screens are not very high in information content...

No dramatic improvement in information content in HCS over the past decade



The data science approach: measure everything, ask questions later...



Cells

Morphological features



“Cytological profile”: collection of measurements describing the appearance of a cell

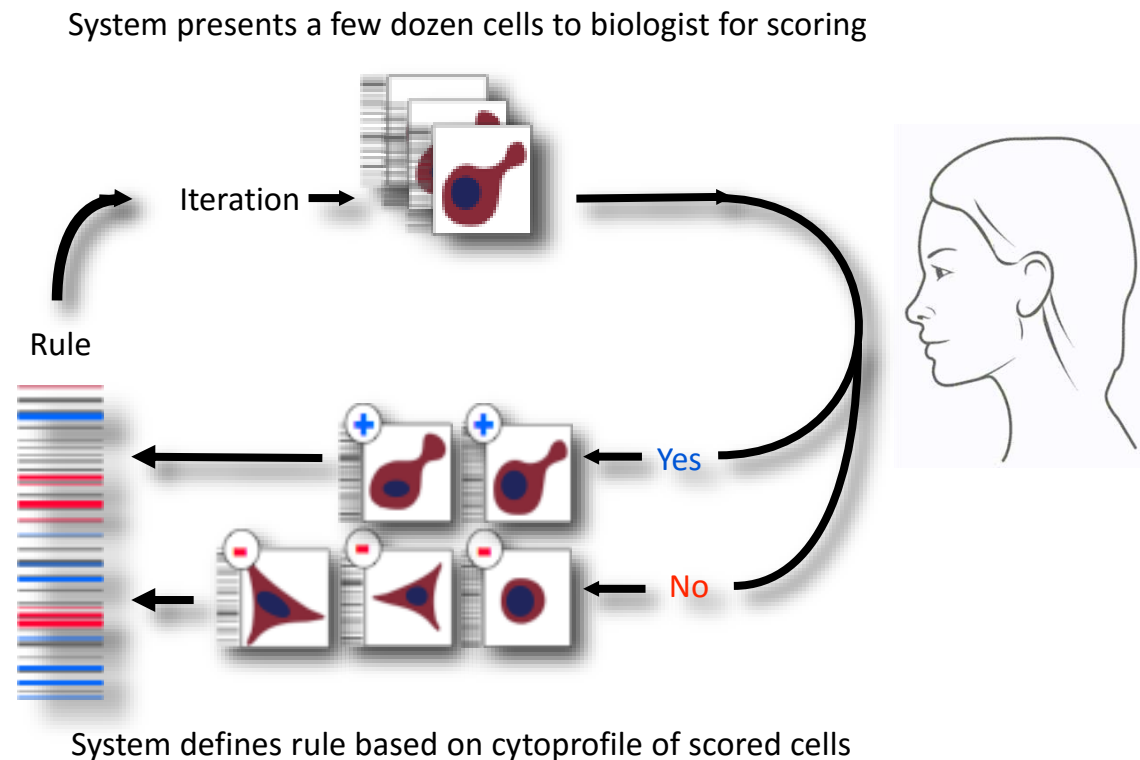
Perlman et al. (2004)

Measuring everything, asking questions later

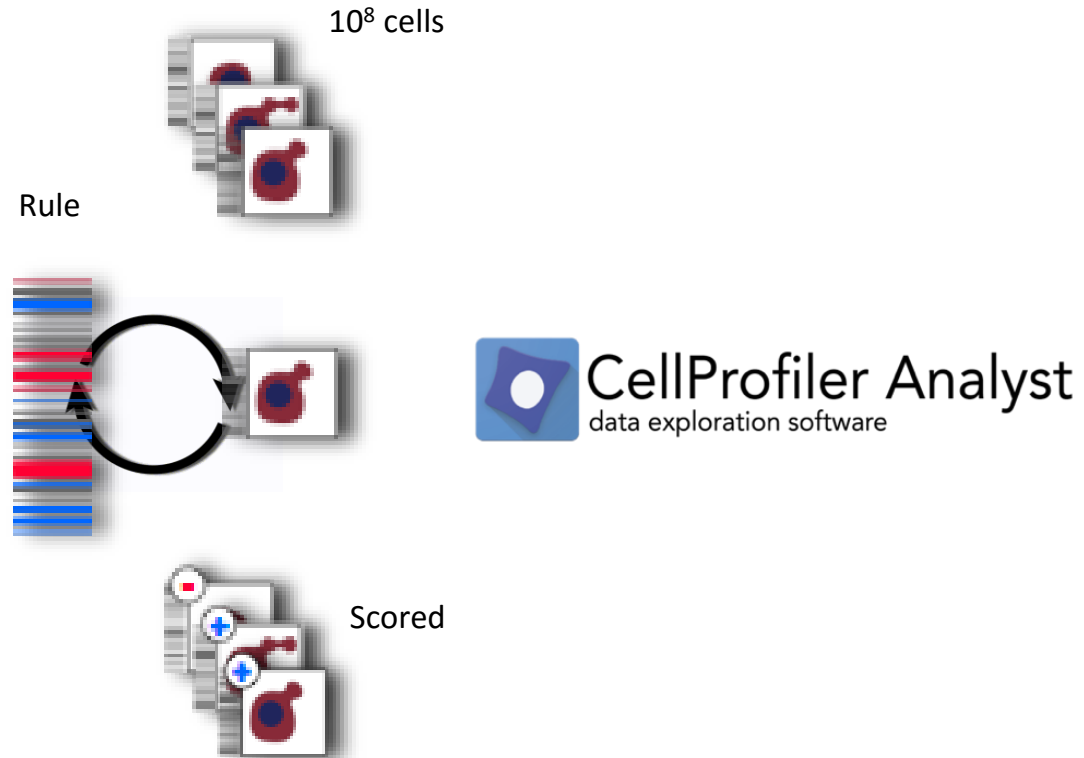
Hundreds of features: size, shape, staining intensity, texture, etc.

- Why?
 - Several features may be necessary to score the phenotype
 - Virtual secondary screens can help characterize hits
 - Later re-screening for new phenotypes
 - The measurements required to score the phenotype of interest may not be known a priori
 - The full spectrum of cellular responses to each treatment (even those not visible by eye) may be useful for data mining/machine learning/clustering...systems biology

Iterative machine learning

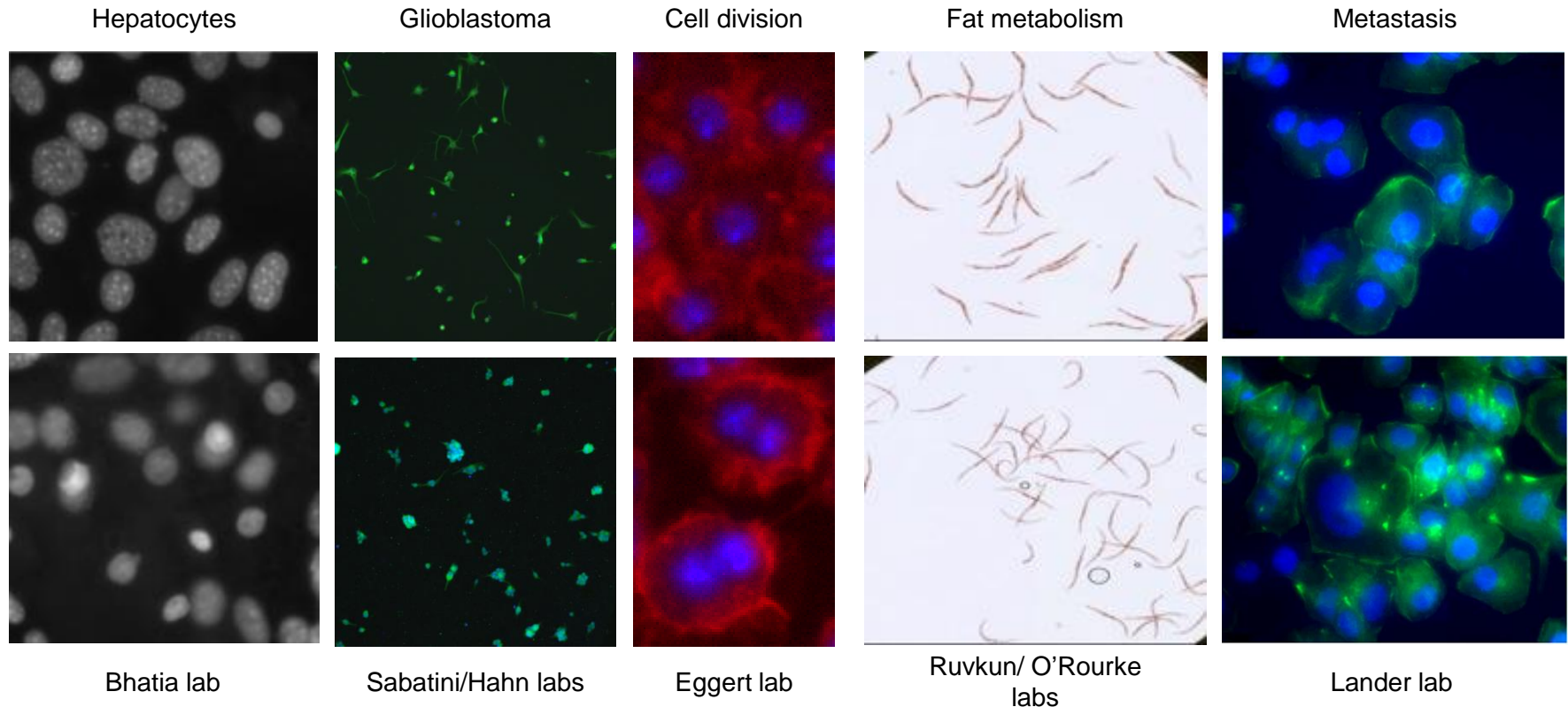


Iterative machine learning



Scored cells are sorted by well:
identify samples with a high proportion of positive cells

Screen more complex phenotypes



CellProfiler Analyst
data exploration software



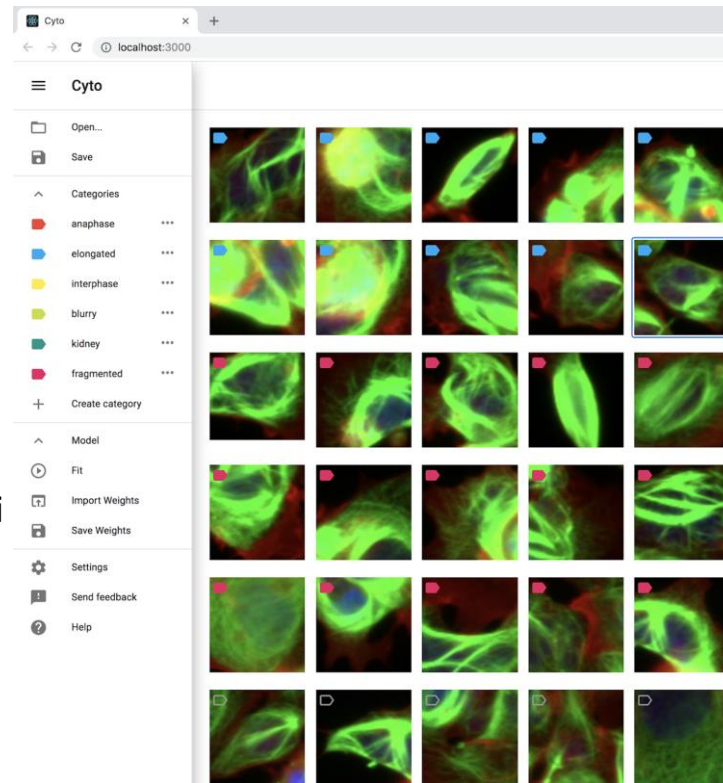
Piximi, <https://github.com/piximi/application>

Piximi: deep learning cell classifier (under construction)

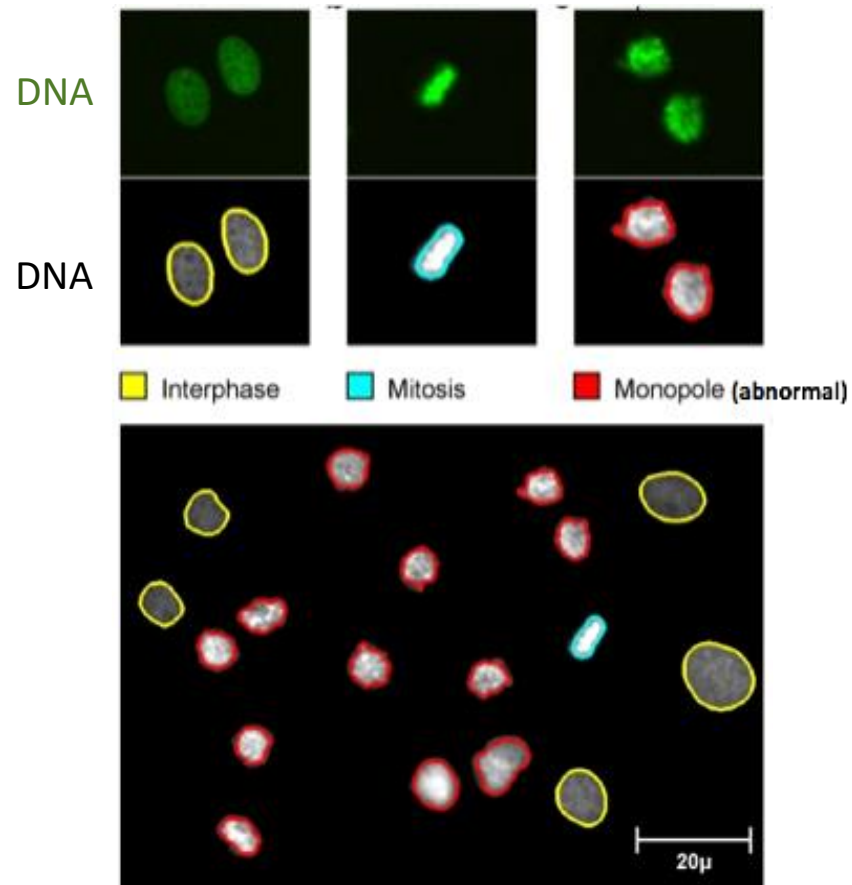
To replace CellProfiler Analyst (Carpenter lab)
+ Advanced Cell Classifier (Horvath lab)

www.piximi.app
(under construction)

- Phenotype classifier for images
- **Web-based**
- Uses **deep learning** algorithms
- Interactive / intuitive user interface
- Source code: <https://github.com/piximi>
- Collaboratively developed!



Regulators of cell division

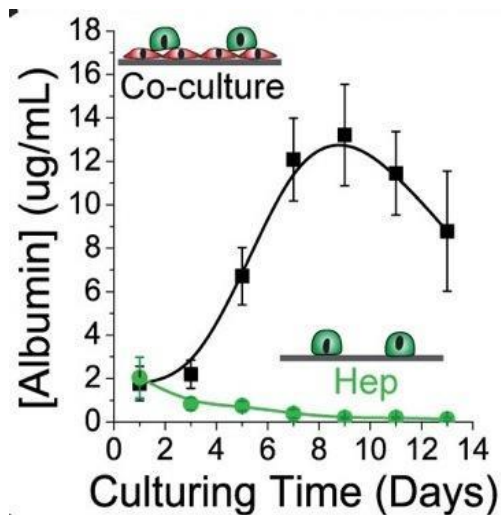


Co-cultured cell systems

Two or more cell types cultured together in order to maintain physiological conditions

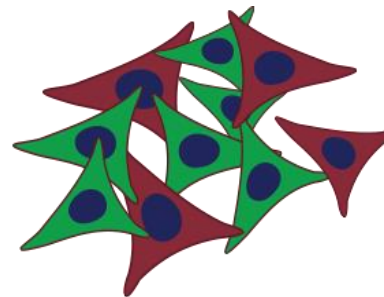
Necessary

- Many primary cell types lose their physiological functions when grown in isolation



Challenging

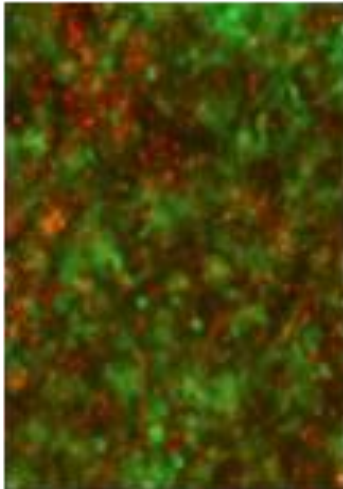
- Culture conditions are difficult to optimize and less robust
- Need to distinguish the cell type of interest from the co-cultured cells, ideally without using additional cellular stains



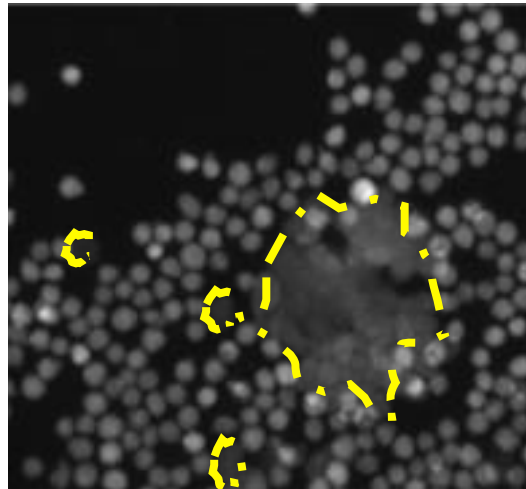
Leukemic & hematopoietic stem cells (HSCs/LSCs)

using mouse primary HSCs or LSCs co-cultured with stromal cells

Co-cultured **LSCs**
and **stroma**



LSC channel only:
live, no DNA stain



Cobblestones

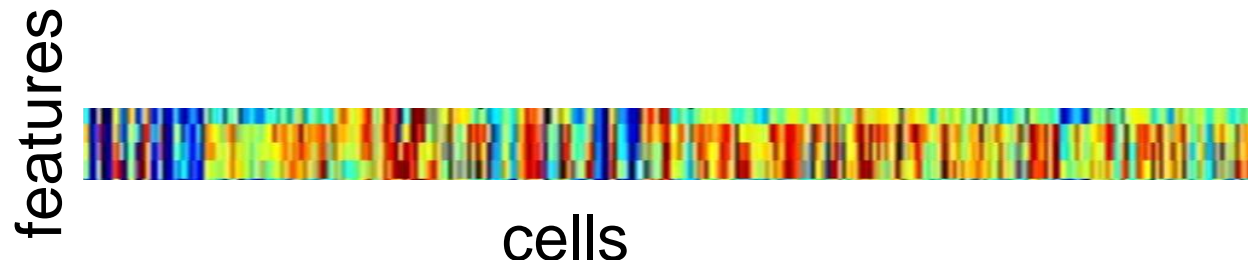
Identified drugs that preferentially reduce
leukemic cell growth; lovastatin extends
lifespan of mice given leukemic bone marrow
cells

Interpretability

- Check what features

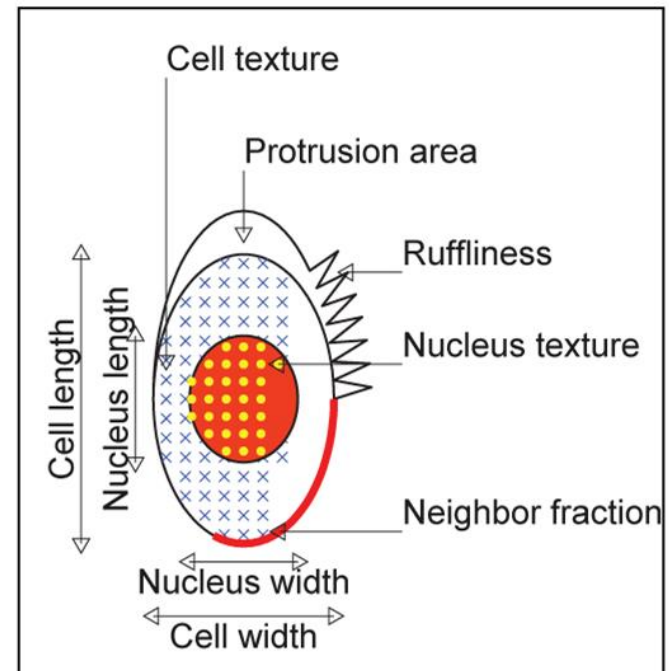
Visualization of high dimensional cellular data

- High dimensional cellular data
- Processing pipeline
- Why do we need visualization?
 - Data interpretation
 - Hypothesis formulation
 - Communication of results
- Visualization methods:
 - Bar charts, scatter plots (3D)
 - Heatmaps

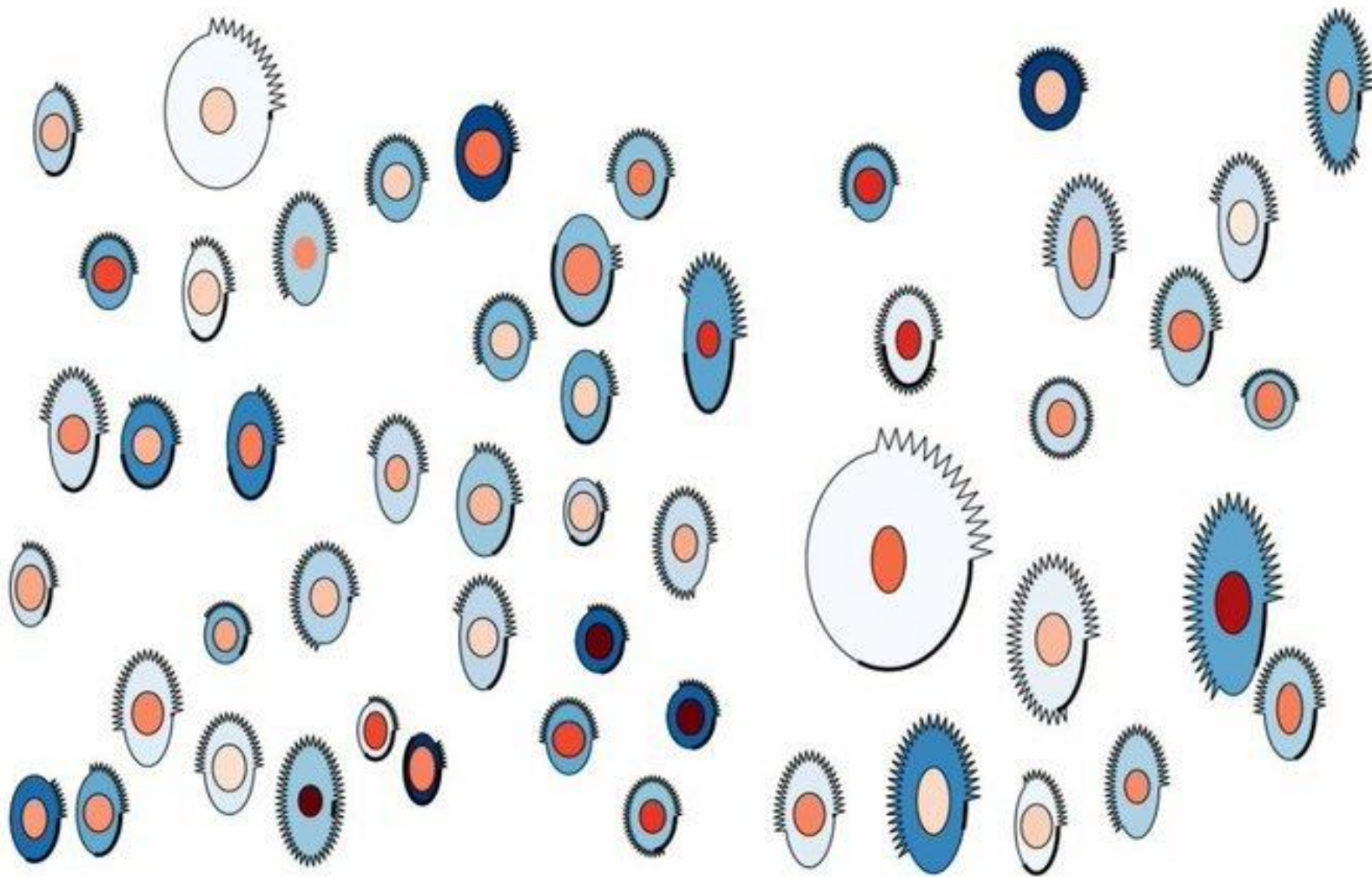


PhenoPlot

- A glyph-based method: “use a collection of visual elements such as size, colour, texture and/or orientation to depict multi-dimensional data”
- Free open source Matlab toolbox (+ GUI)
- 21 variables encoded
- Allows customization

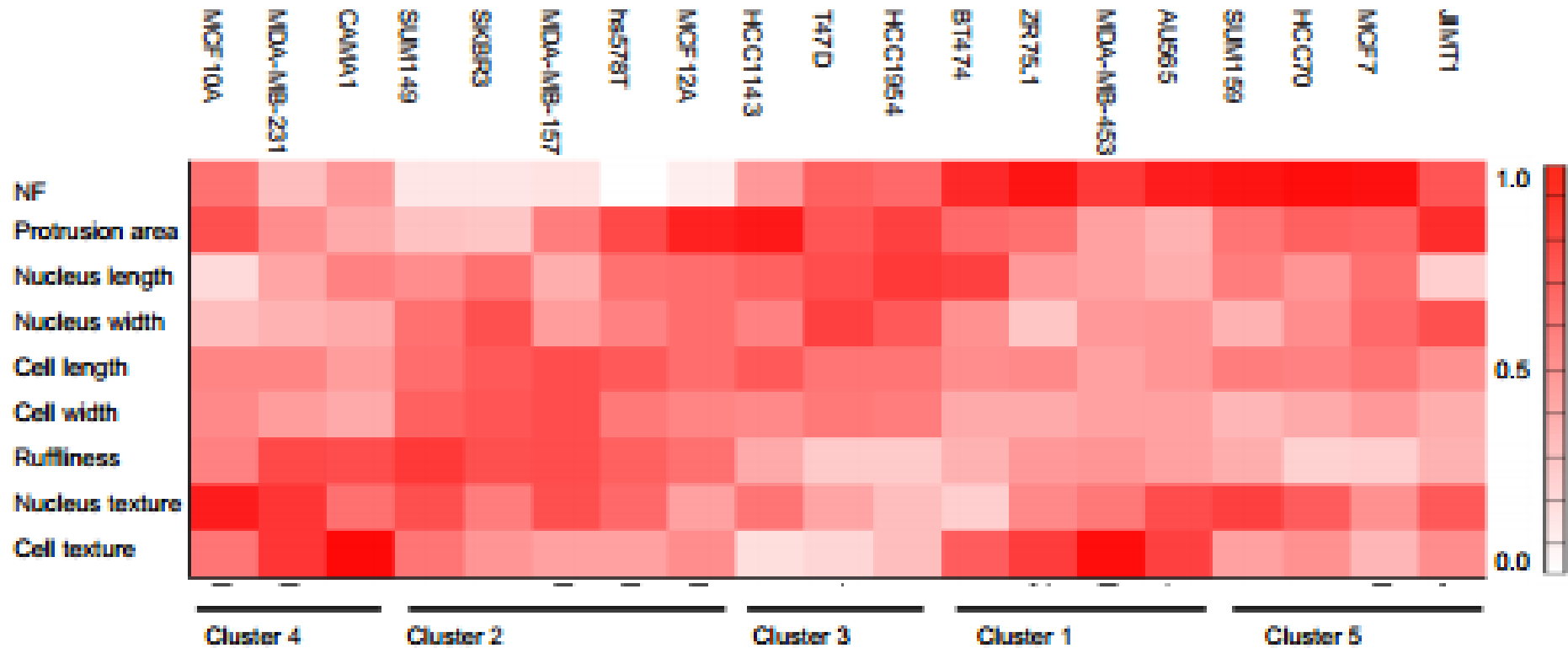


PhenoPlot



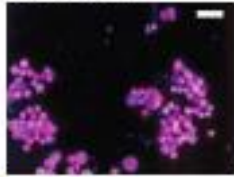
Hierarchical Clustering

N ~ 156,000 cells

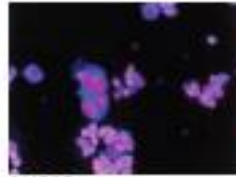


Clusters Visualization

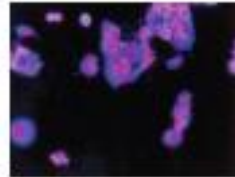
Cluster 1



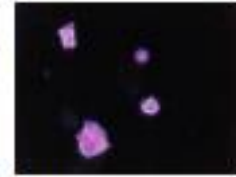
MDA-MB-453



AU565

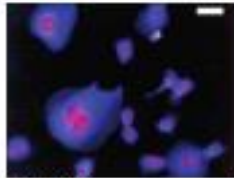


ZR75.1

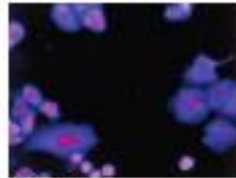


BT474

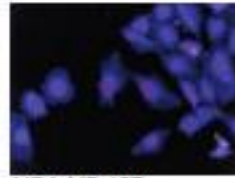
Cluster 2



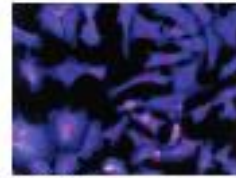
SUM149



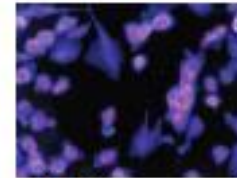
SKBR3



MDA-MB-157

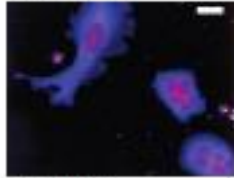


hs578T

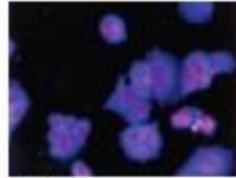


MCF-12A

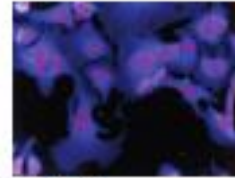
Cluster 3



HCC1954

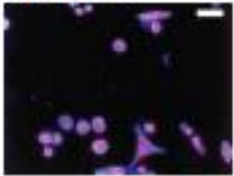


T47D

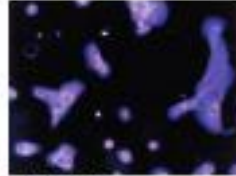


HCC1143

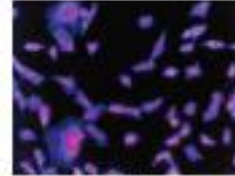
Cluster 4



CAMA1

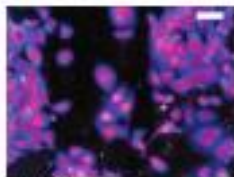


MCF-10A

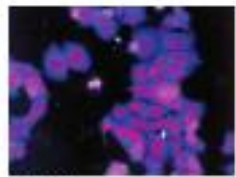


MDA-MB-231

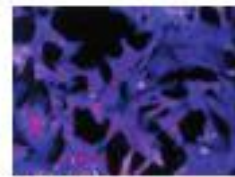
Cluster 5



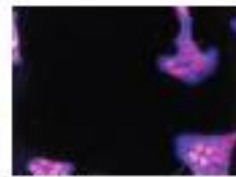
JIMT1



MCF7



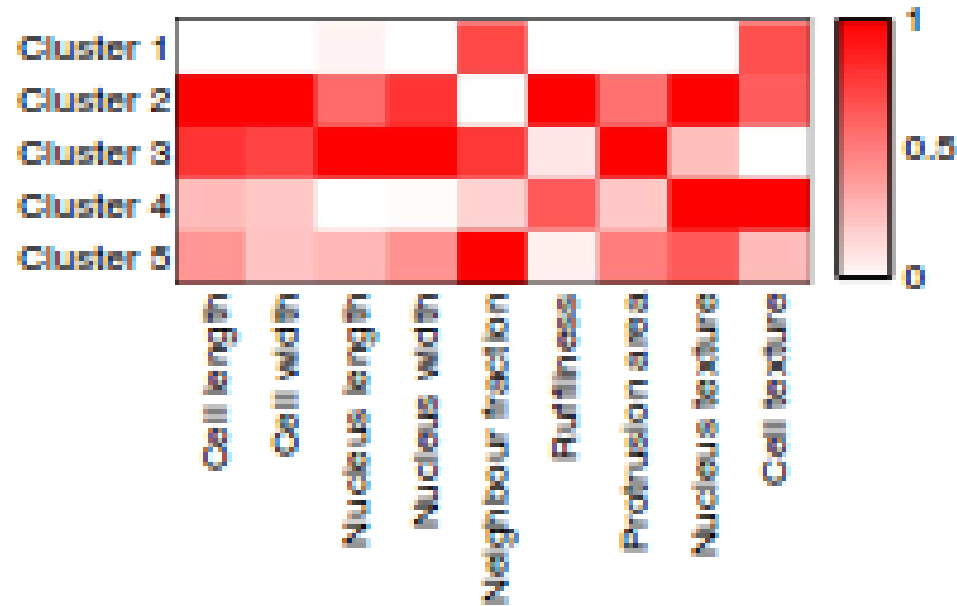
SUM159



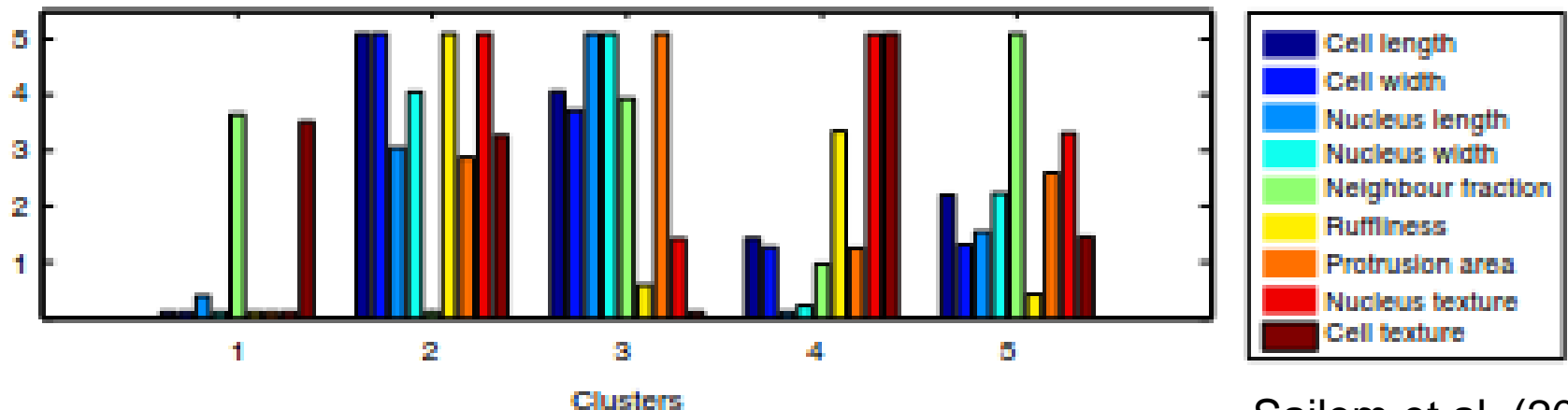
HCC70

Clusters visualization

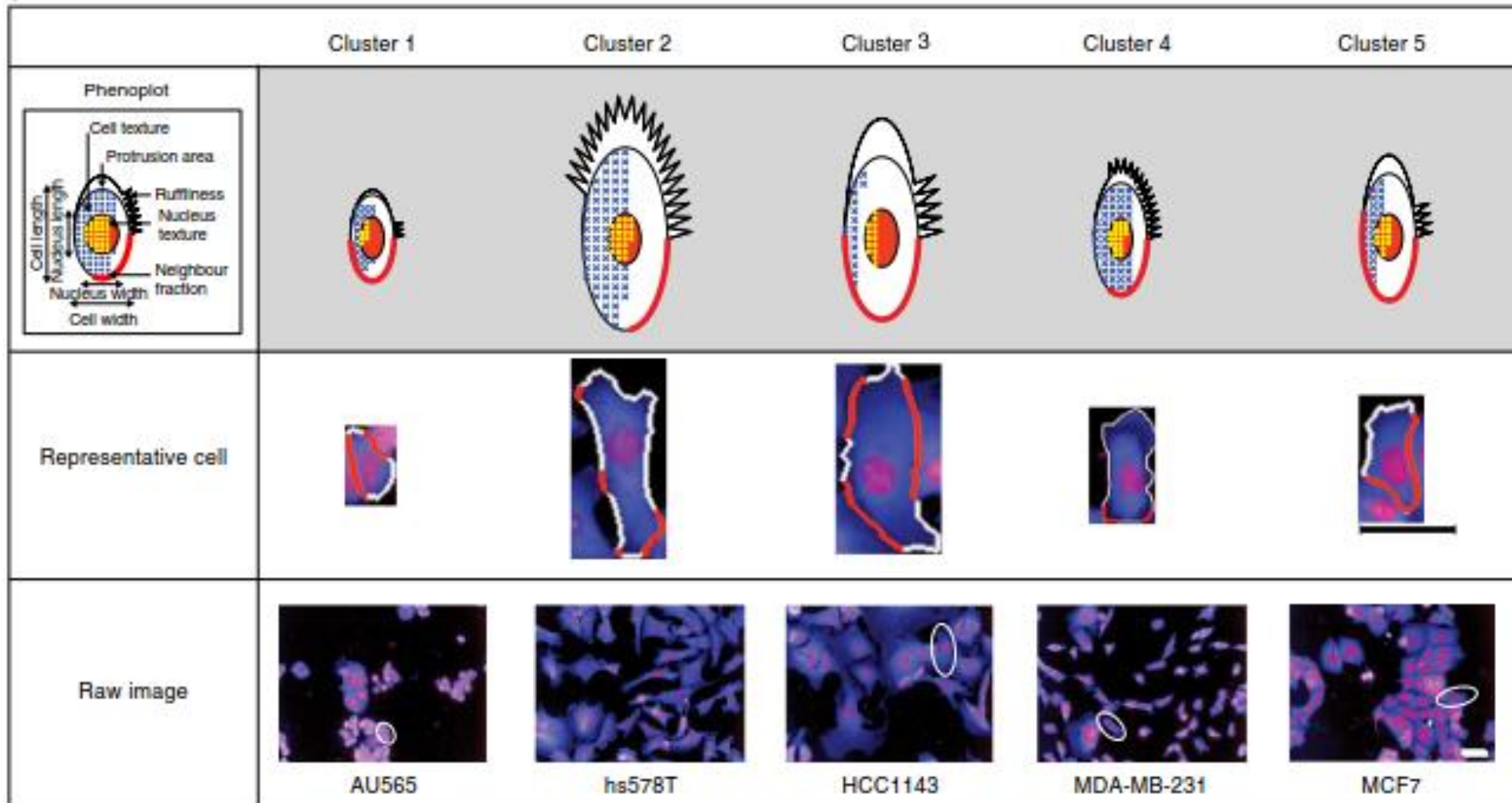
b



c

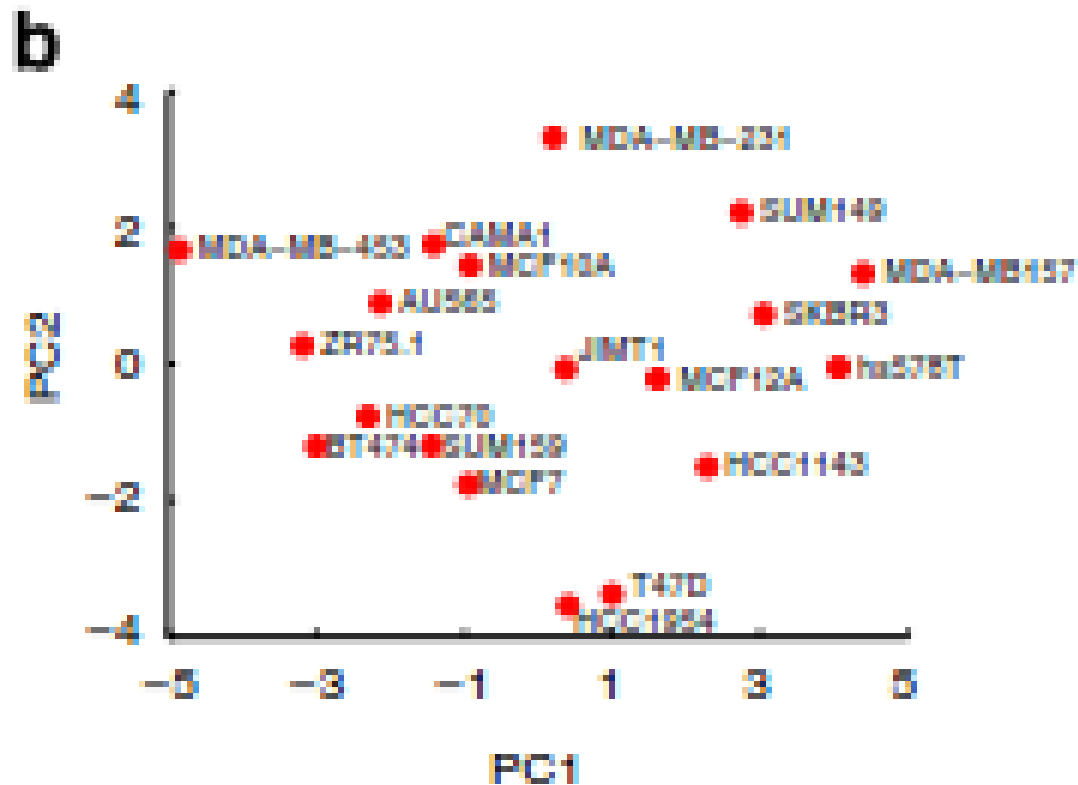


Clusters Visualization

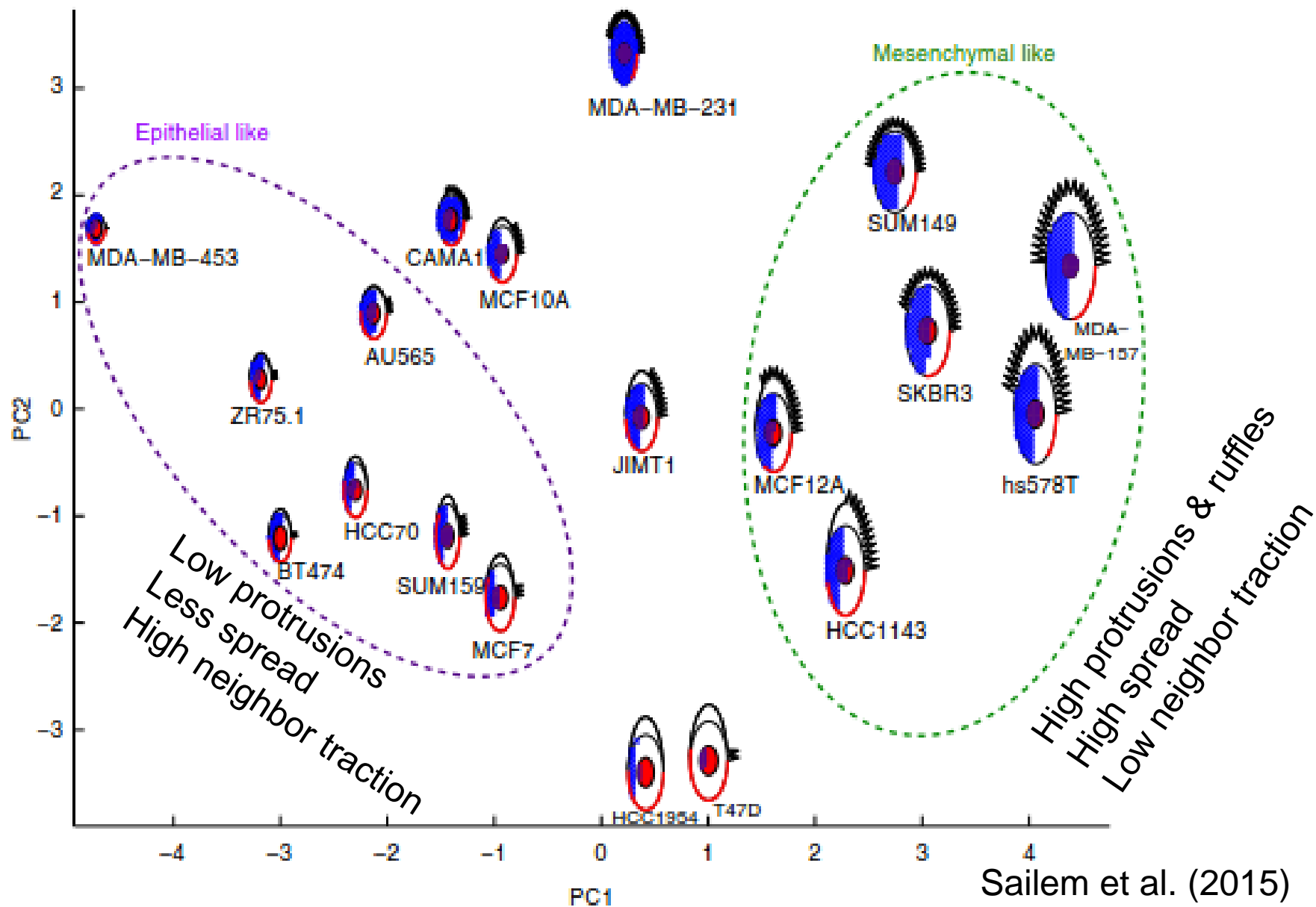


Easier to interpret than “representative” cell

Scatter plot is not informative

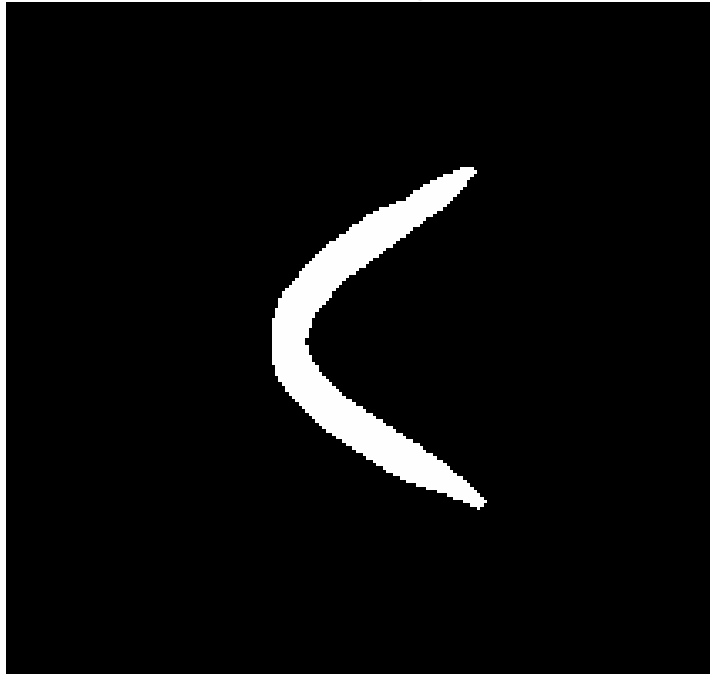


Clusters visualization



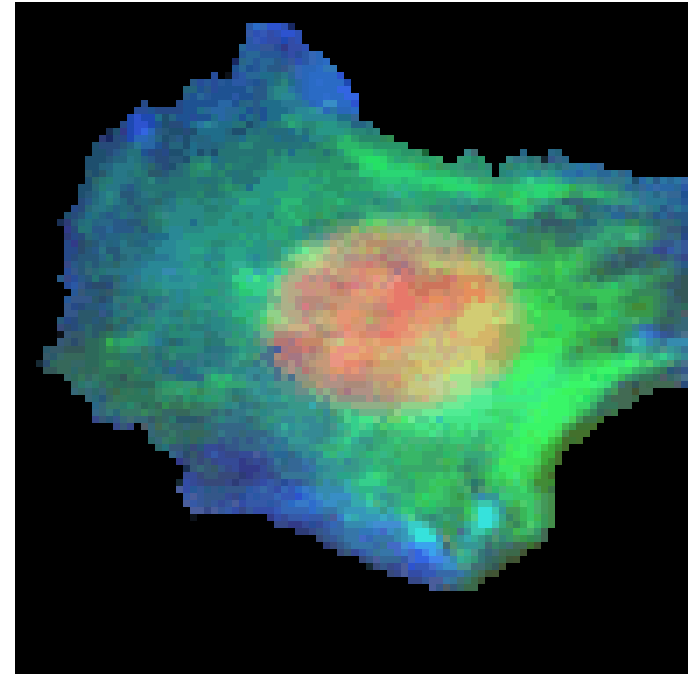
Visual interpretation of complex pheontypes

View examples along a phenotypic continuum (“Eigenworms”)



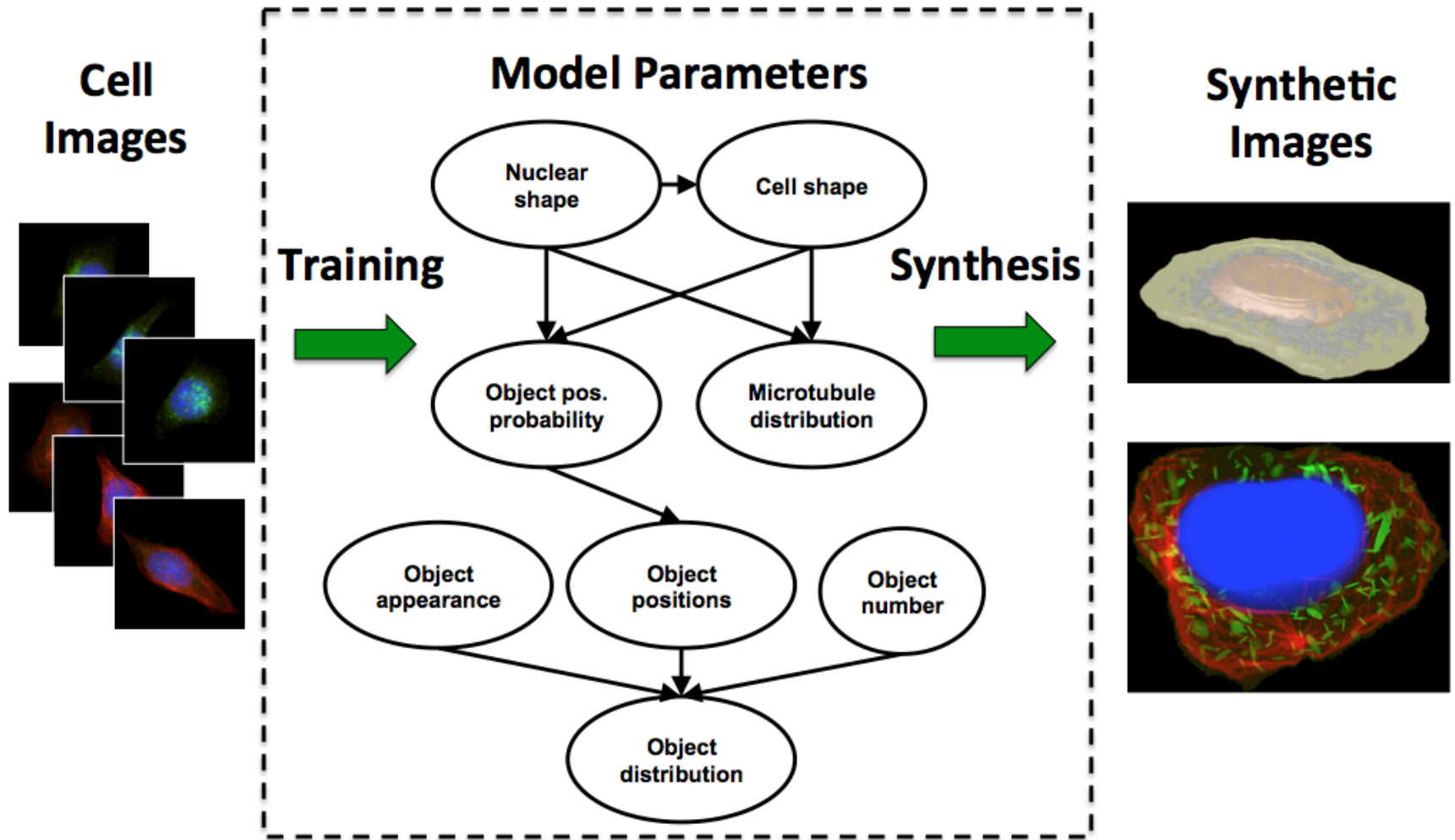
Anne Carpenter,
from Carolina Wahlby

Generative modeling



Goldsborough et al. (2017),
Lafarge et al. (2019)

Generative modeling



Adversarial Autoencoders

Research project types

- Data mining and integration in public repositories
- Collaborative bioimage analysis projects
- Tool building projects

Open resources for your research projects

- Image data resource (IDR), Williams et al. (2017), Image Data Resource: a bioimage data integration and publication platform
<https://idr.openmicroscopy.org/>
- The Allen Institute of Cell Science
- The Human Protein Atlas
- Branda Andrews datasets
<http://sites.utoronto.ca/andrewslab/data.shtml>

Open resources for your research projects

- Bray et al. (2017), A dataset of images and morphological profiles of 30 000 small-molecule treatments using the Cell Painting assay. Data, <https://github.com/gigascience/paper-bray2017>
- Pascual-Vargas et al. (2017), RNAi screens for Rho GTPase regulators of cell shape and YAP/TAZ localisation in triple negative breast cancer Data via IDR
- Pizzagalli et al (2018), Leukocyte Tracking Database, a collection of immune cell tracks from intravital 2-photon microscopy videos (via figshare)