

Computational Biology: A Measurement Perspective

Alden Dima

Information Technology Laboratory

alden.dima@nist.gov

Problem

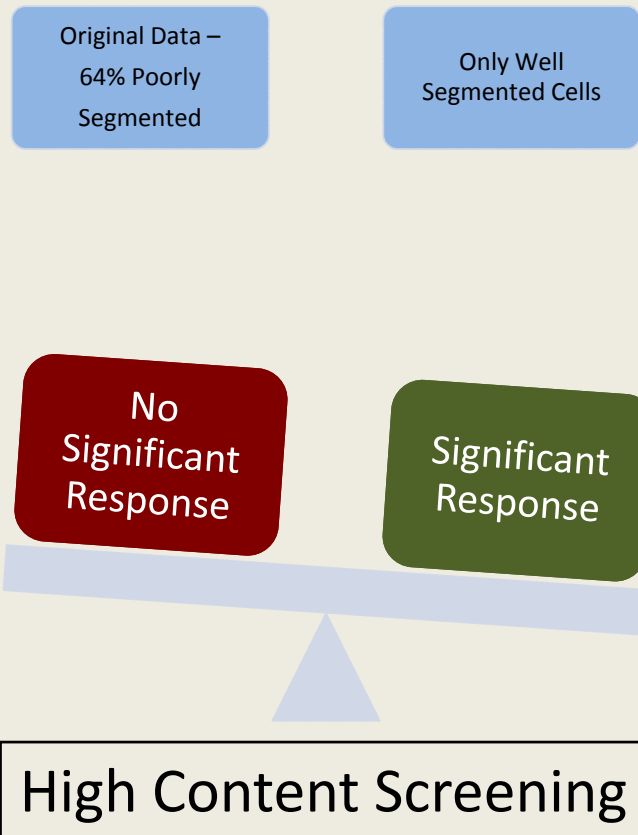
- High-throughput technologies are generating large amounts of complex data that is difficult to process and convert into knowledge
- Issues exist throughout data lifecycle:
 - Acquisition
 - Analysis
 - Archiving
 - Interchange

Imaging Technologies

- Imaging technologies are increasingly being used both as diagnostic tools and as research tools in the biosciences
 - Novel methods are needed for automated analysis and comparisons
 - Correlation among studies is difficult at best
 - Off-the-shelf methods are not well characterized and can contribute significantly measurement uncertainty
 - Need to combine features from images with other biological or medical sources

Example from Literature

- Relies heavily on cell imaging
- Gigabytes of images collected
- Algorithms treated as “black boxes”



- Many algorithms published but rarely compared
- Poor segmentation significantly effects conclusions

Source:

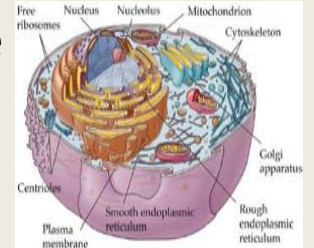
Hill, LaPan, Li, Haney (Wyeth Research)

Impact of image segmentation on high content screening data quality for SK-BR-3 cells

BioMed Central Bioinformatics 2007

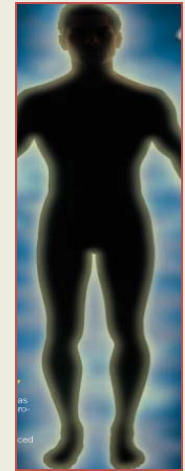
Computational Biology: *Single Cell Analyses*

Intracellular molecular reactions and interactions that control the response and fate of cells and organisms cannot be unambiguously compared and combined due to a lack of standards and validated protocols



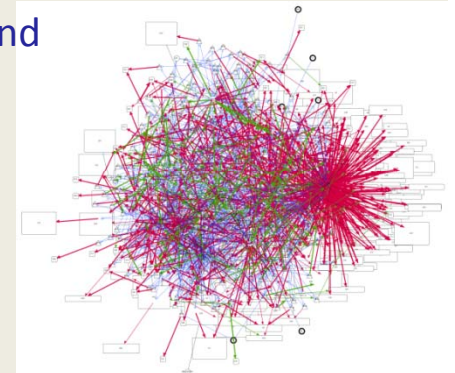
NIST Role

- Provide the measurement tools and standards that enable quantifiable and reproducible measurements of cells and their interactions through:
 - Standard data/metadata format for image capture, storage, retrieval, analysis
 - Software to enable high throughput cell image analysis and interoperability
 - Standards and validation required to ensure reproducible image analysis
 - Multi-site experiments to test software and validation protocols



Technical Approach

- Create and evaluate an integrated data collection, organization and analysis infrastructure for cellular imaging
- Experimentalists and computational scientists - focus on the physical standards and protocols for data collection, image processing and analysis, storage of data and metadata, and evaluation of results

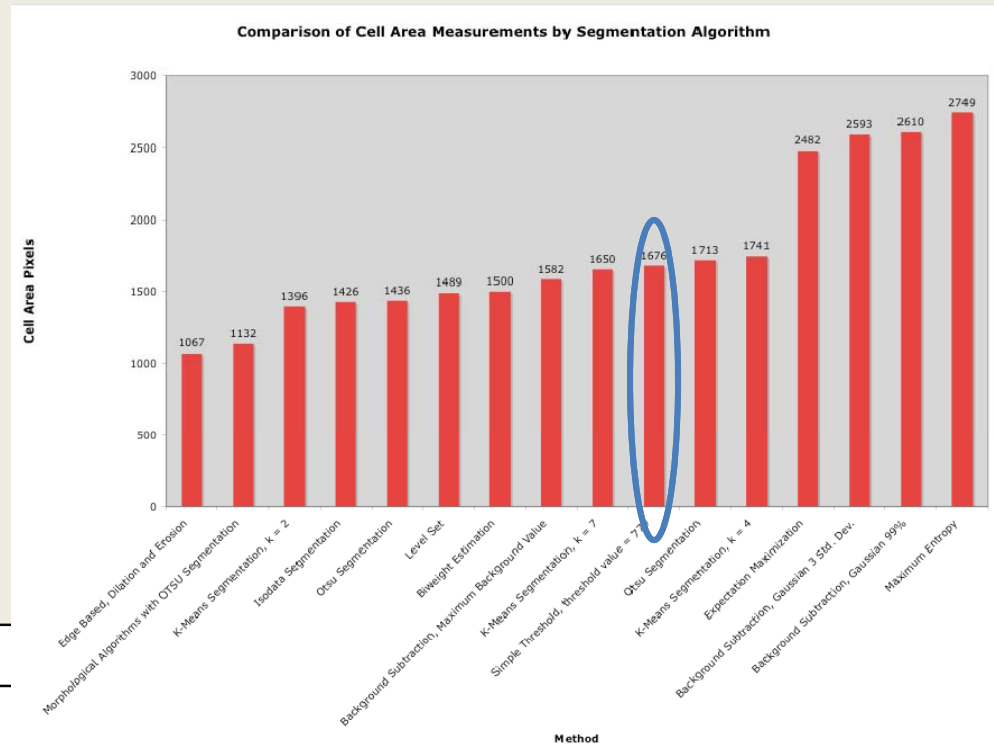
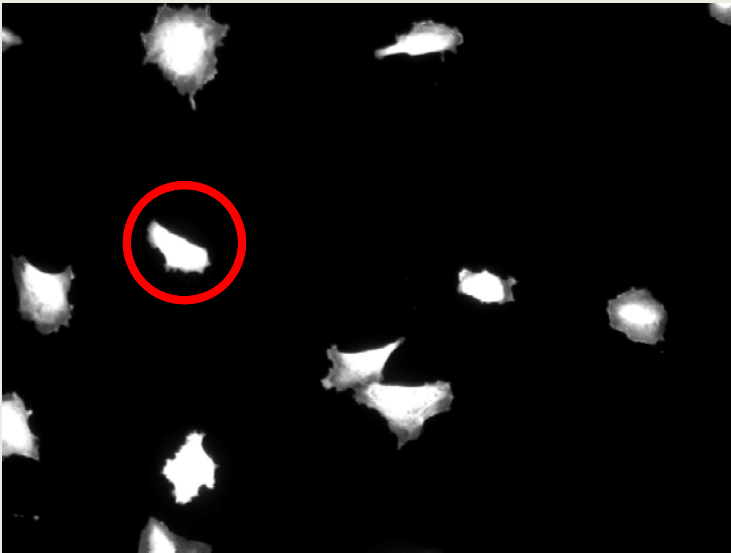


Segmentation Evaluation

- Determine which segmentation technique and associated parameters can be used to reliably determine the morphology of cells for the purposes of comparing cell lines as part of a new standard procedure under development

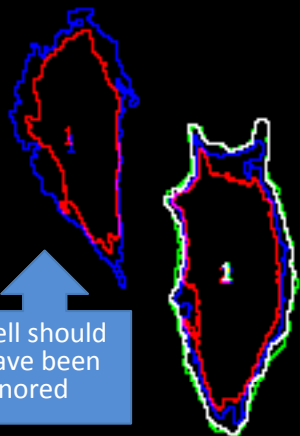
Variability Across Methods

Different segmentation techniques can change results

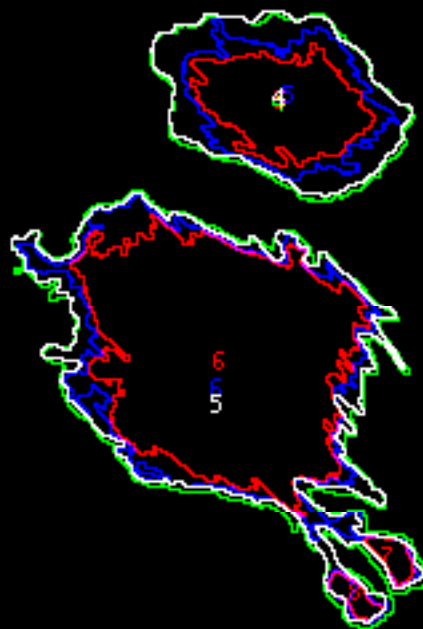
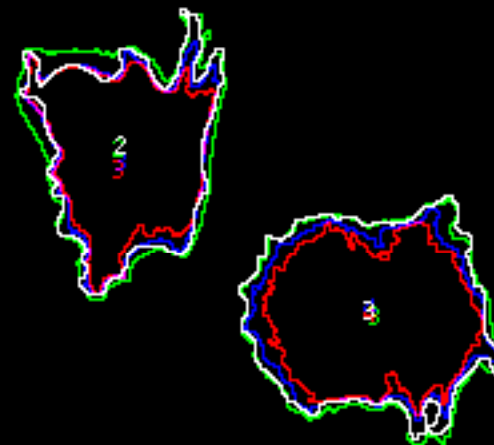


Method ID	Person	Reference Cell Area Value (pixels)	Tool	Method
1	John Elliot	1676	ImageJ	Simple Threshold, threshold value = 779
2	Asim Wagan	1436	Matlab	Otsu Segmentation
3	Asim Wagan	1426	Matlab	Isodata Segmentation
4	Asim Wagan	1741	Matlab	K-Means Segmentation, k = 4
5	Asim Wagan	2482	Matlab	Expectation Maximization
6	Marcin Kociolek	2593	C++	Background Subtraction, Gaussian 3 Std. Dev.
7	Marcin Kociolek	2610	C++	Background Subtraction, Gaussian 99%
8	Marcin Kociolek	1582	C++	Background Subtraction, Maximum Background Value
9	Xiao Lan Li	1489	Matlab	Level Set
10	Xiao Lan Li	1067	Matlab	Edge Based, Dilation and Erosion
11	Rui Fang	1132	Matlab	Morphological Algorithms with OTSU Segmentation
12	Jim Filliben	1396	Dataplot	K-Means Segmentation, k = 2
13	Adele Peskin	1500	C++	Biweight Estimation
14	Javier Bernal	1650	FORTRAN	K-Means Segmentation, k = 7
15	Javier Bernal	2749	FORTRAN	Maximum Entropy
16	Javier Bernal	1713	FORTRAN	Otsu Segmentation

Cell should
have been
ignored



Red = k-means (k = 2)
Blue = Otsu
Green = Maximum Entropy
White = Ground Truth

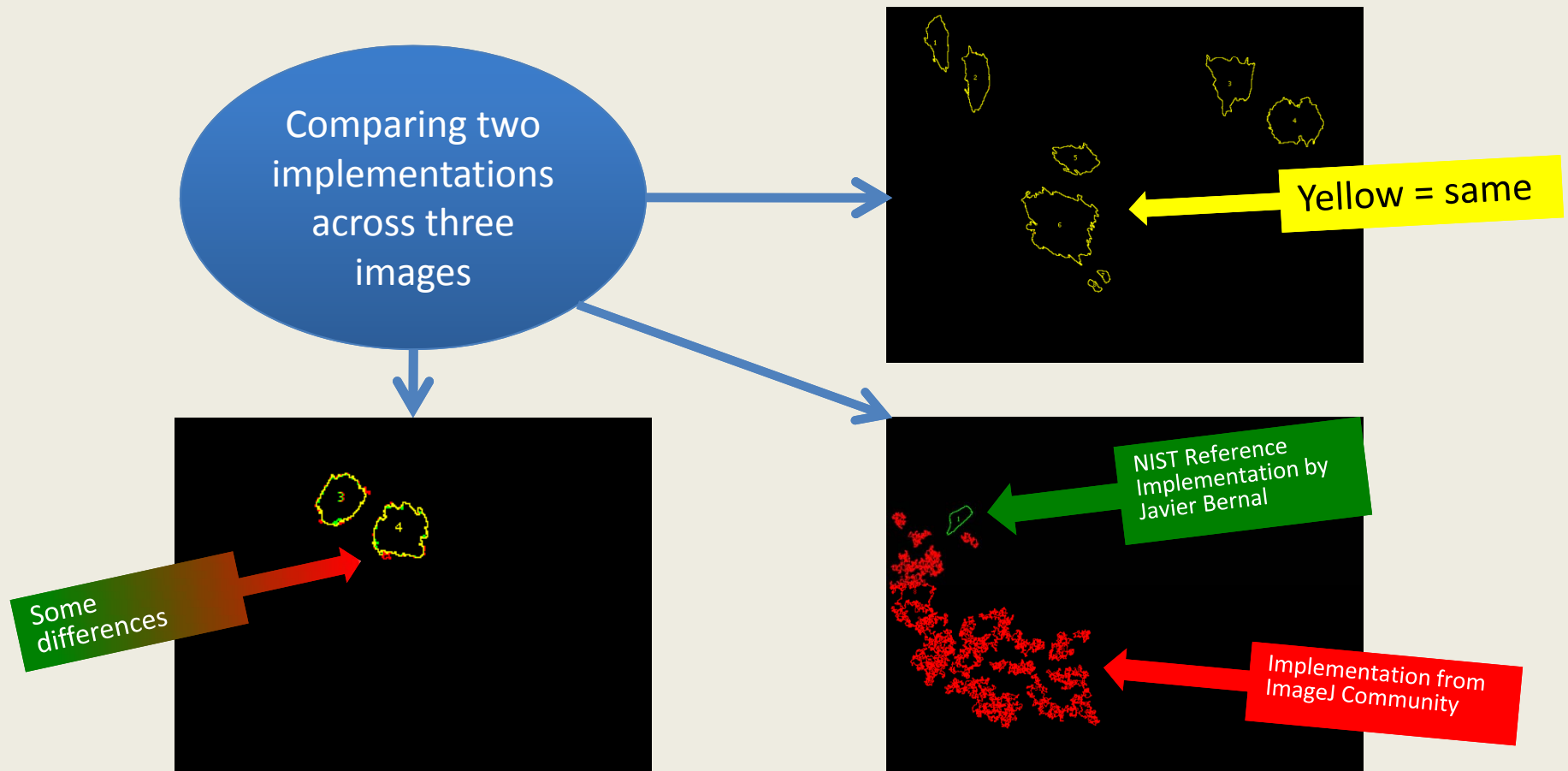


Preliminary evaluation
shows that results can
vary by more than $\pm 40\%$

A10 Cell Line
Three different segmentation
techniques

Variability Across Implementations

Different implementations of the same technique can change results as well

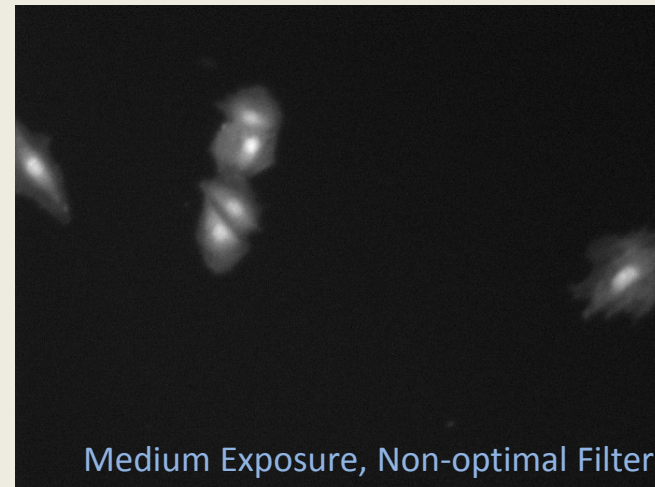
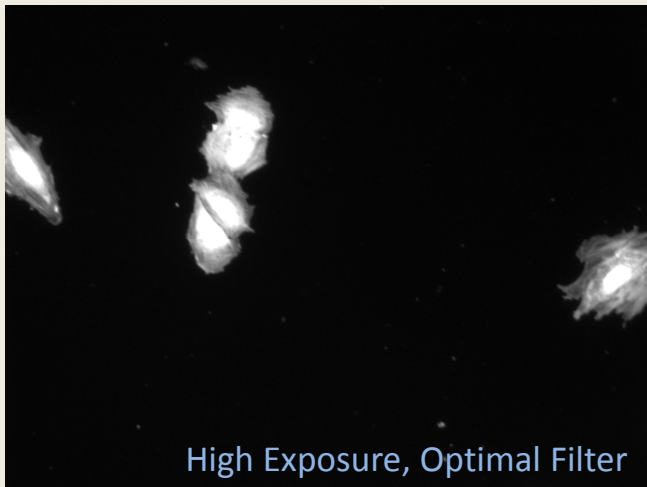
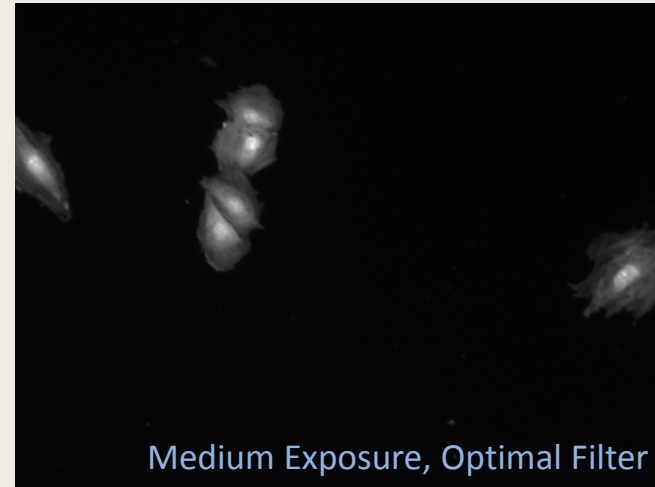
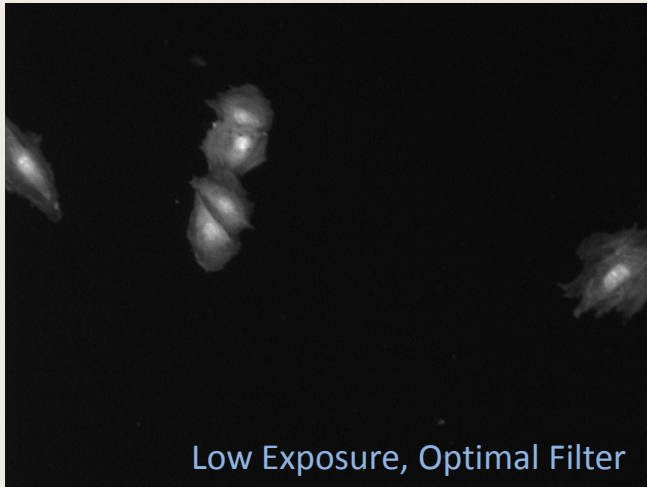


“Software as measurement”

Experimental Design

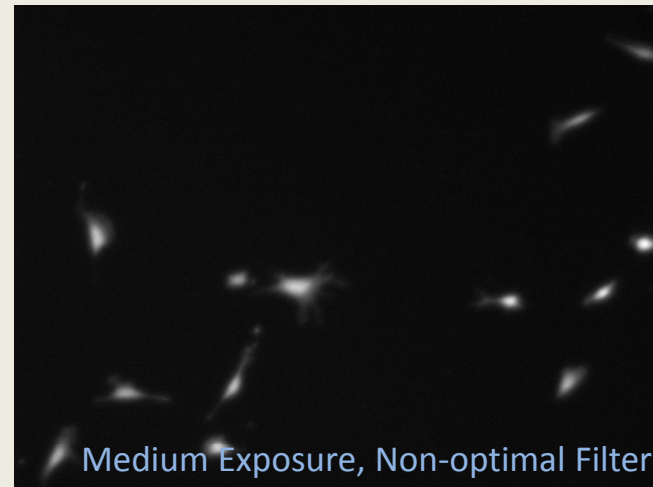
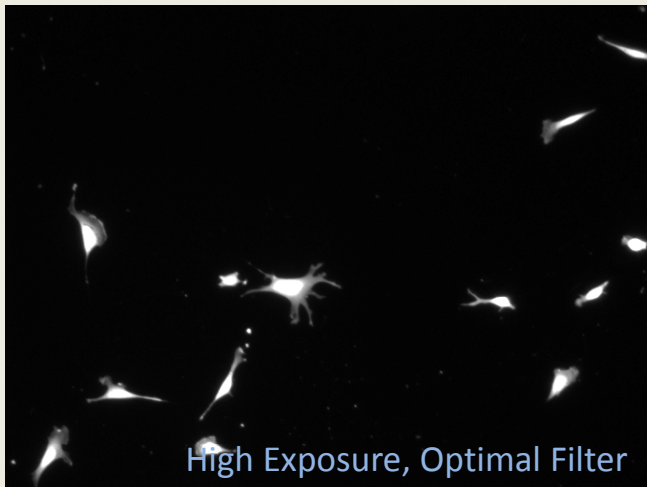
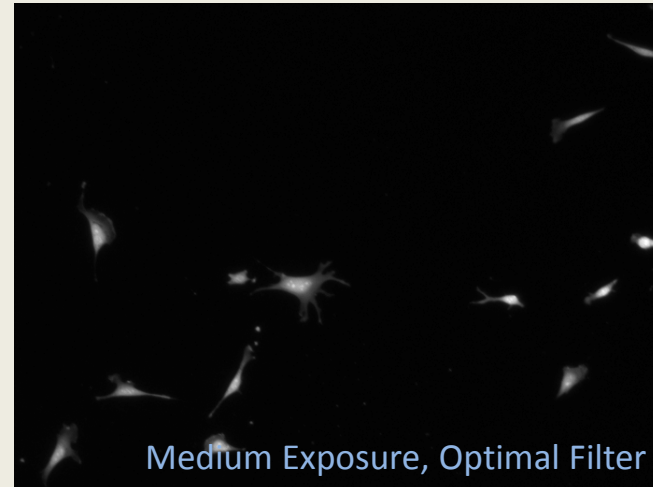
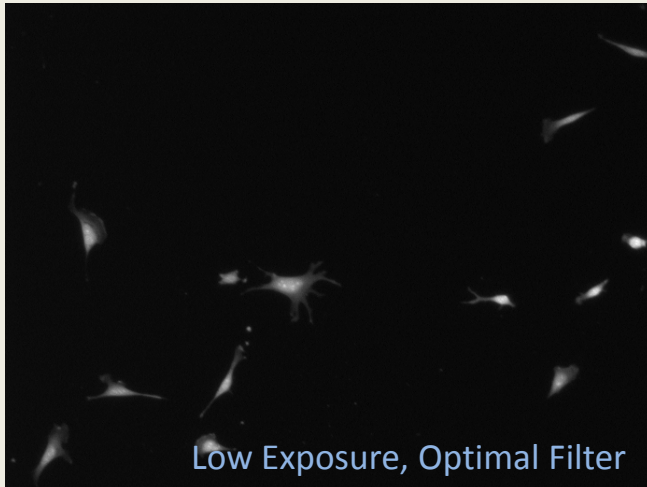
- Two cell lines – A10 & 3T3
- Cell preparation follows ASTM document
- Different image exposure and filter levels
- No flat-field correction
- Multiple sampling to capture noise
- Ground truth: expert manual segmentation
- Result: approximately 8000 images – 80 of which are used for initial evaluation work

Red Channel Images - A10



16-bit Gray-scale Images

Red Channel Images - 3T3



16-bit Gray-scale Images

Planned Progression of Evaluations

Basic Algorithms

```
graph TD; A[Basic Algorithms] --> B[Edge-based Methods]; B --> C[Advanced Techniques];
```

Edge-based Methods

Advanced Techniques

Algorithms Evaluated

Basic Methods

Otsu

Maximum
Entropy

K-Means Clustering

NIST

ImageJ

NIST

ImageJ

K=2

Threshold
Clusters
(K=3,4,5)

Combined
with Otsu
(K=3,4,5)

Combined
with Max.
Entropy
(K=3,4,5)

NIST

ImageJ

NIST

ImageJ

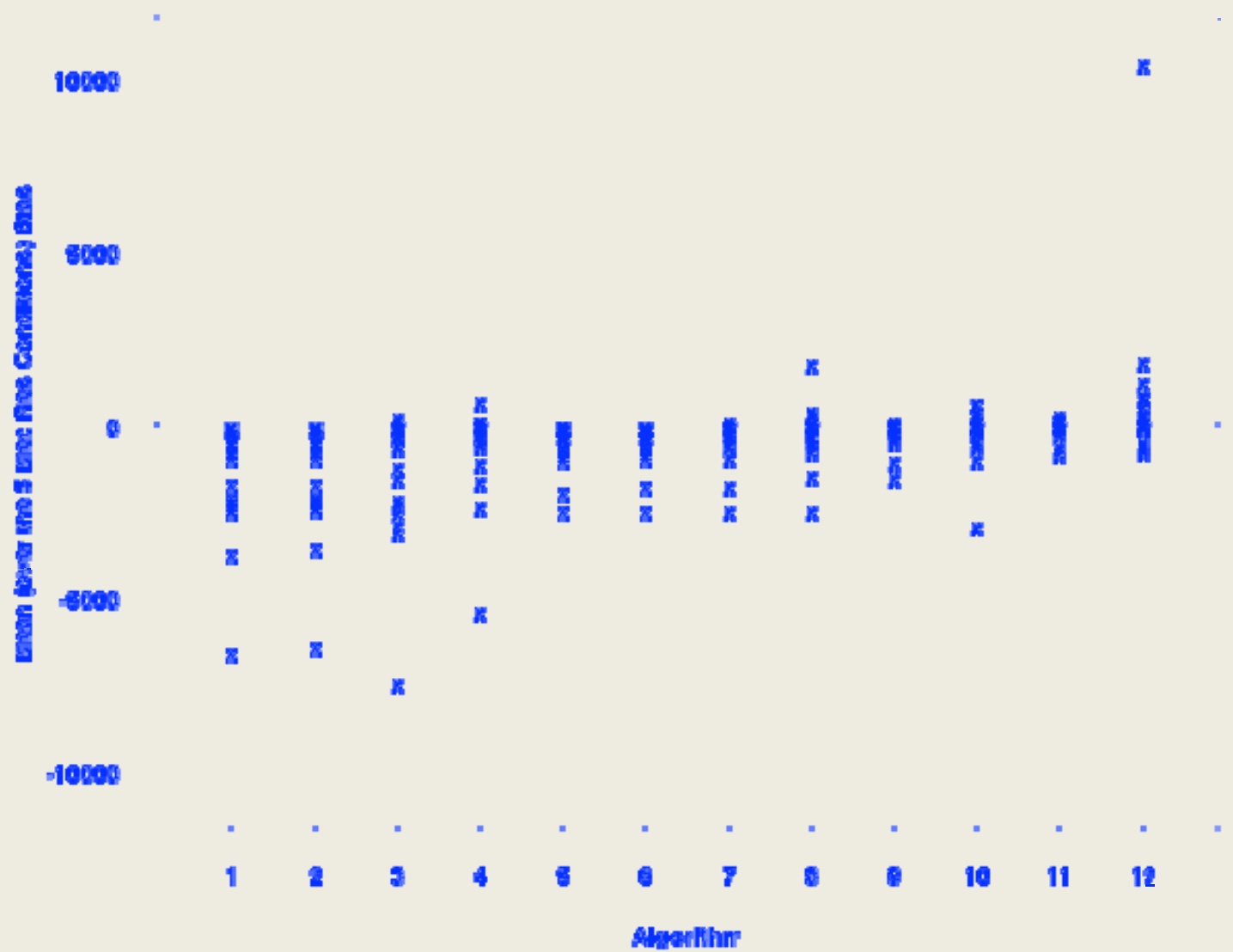
ImageJ

Bio-Image Sensitivity Analysis (Complete Project) (Brady/Dima/Plant)
 Q1. What is the best algorithm?
 Cell Entry = Mean Bias (Over 5 Medium-Res Image Conditions)
 Algorithm

High Res Image ID

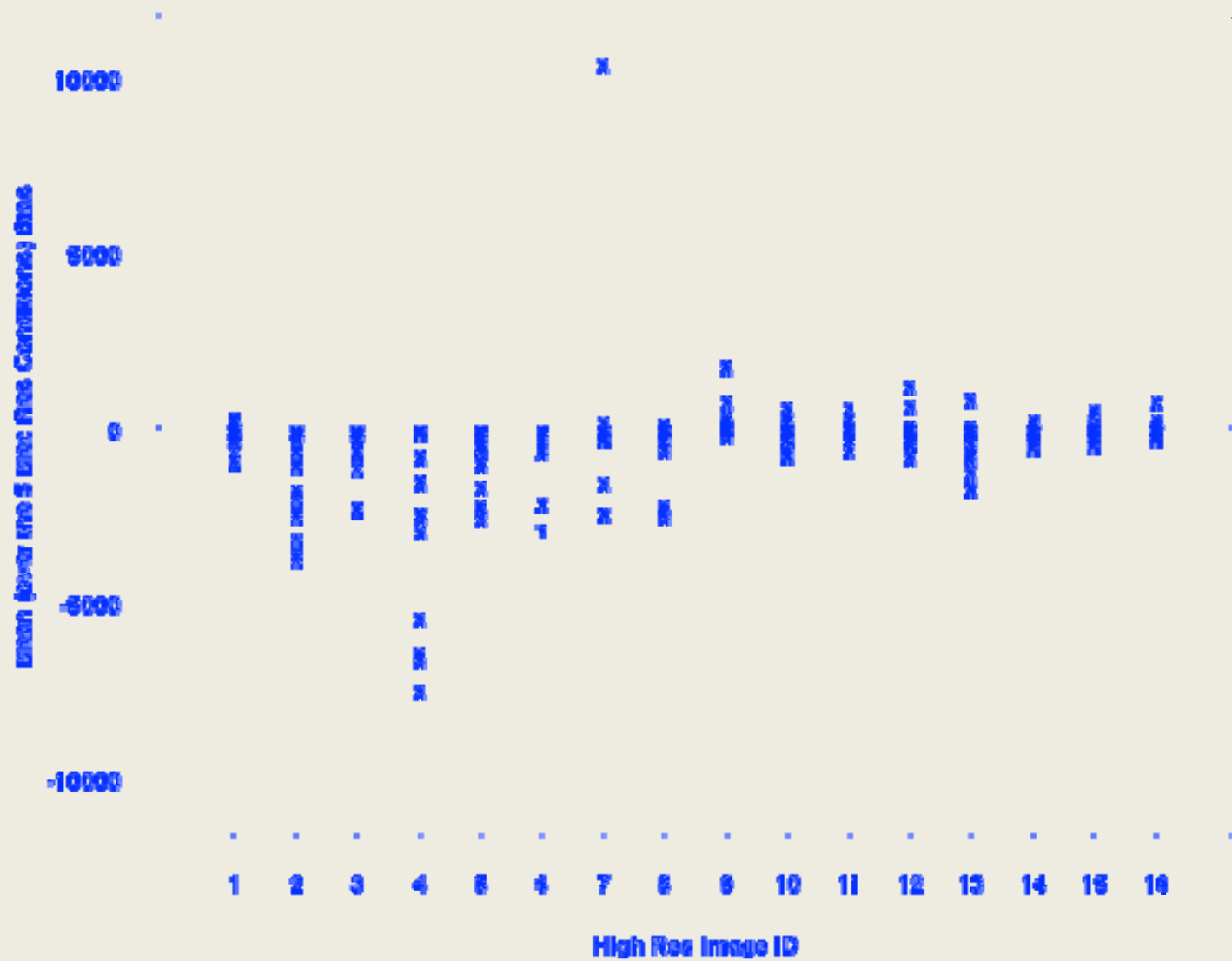
	K2D	Orca	Merlot	K2D	K2H	K2H	K2D+O	K2D+2H	K2H+O	K2H+2H	K2D+O	K2D+2H	Mean	Median	Rank
1	-997.	-999.	-376.	-194.	-244.	-223.	-223.	85	84	493	937	234	-189.	-183.	6
2	-3678.	-3518.	-3083.	-2271.	-1800.	-1728.	-1728.	-2398.	-1017.	-995.	-995.	-354.	-1080.	-1874.	13
3	-2358.	-2178.	-2221.	-1658.	-898.	-818.	-818.	-578.	-282.	-205.	-84.	-682.	-631.	-639.	12
4	-6831.	-6384.	-7483.	-3303.	-2427.	-2378.	-2378.	-1423.	-1481.	-2878.	-788.	-727.	-3348.	-2463.	16
5	-2208.	-2144.	-2448.	-1813.	-928.	-882.	-882.	-418.	-448.	-419.	-348.	-682.	-1102.	-888.	14
6	-2090.	-2040.	-2831.	-463.	-878.	-388.	-388.	-388.	-323.	-285.	-134.	-488.	-635.	-688.	11
7	-2381.	-2343.	-1488.	-811.	-288.	-233.	-233.	-282.	-131.	-188.	-8.	10808	234	-241.	8
8	-2387.	-2302.	-2143.	-882.	-348.	-340.	-340.	48	-188.	240	-84.	177	-681.	-340.	10
9	-108.	-58.	292	183	41	-8.	-8.	1788	140	184	228	1646	372	182	5
10	-721.	-708.	-531.	-218.	-281.	-242.	-242.	431	21	688	232	248	-108.	-231.	7
11	-248.	-231.	-308.	181	-228.	-234.	77	318	-82.	688	128	294	-8.	27	1
12	-748.	-737.	-308.	788	-378.	-377.	-377.	-280.	24	-188.	187	1287	-81.	-244.	9
13	-1872.	-1698.	-1181.	-184.	-843.	-787.	188	-717.	-882.	-814.	-84.	818	-828.	-682.	15
14	-688.	-473.	-331.	-118.	-147.	-141.	-141.	-70.	18	20	137	314	-118.	-128.	4
15	-443.	-430.	-383.	-83.	-168.	-167.	-167.	88	28	478	148	898	-97.	-108.	3
16	-283.	-258.	-37.	143	-114.	-112.	-112.	-87.	18	323	128	648	41	-87.	2
Mean	-1713.	-1683.	-1847.	-688.	-887.	-883.	-485.	-241.	-241.	-188.	-27.	898			
Median	-1288.	-1284.	-883.	-217.	-238.	-231.	-238.	-174.	-77.	-88.	84	280			
Rank	12	11	10	5	6	8	8	4	3	2	1	7			488888

Efo-Image Sensitivity Analysis (Complete Project) (Brady/Dima/Plant)
 G. Best & Worst Algorithms?



	1	2	3	4	5	6	7	8	9	10	11	12
Mean	-1713.	-1999.	-1547.	-999.	-397.	-293.	-493.	-241.	-241.	-199.	-27.	999
Median	-1299.	-1234.	-669.	-207.	-299.	-291.	-239.	-174.	-77.	-69.	64	299
Rank	12	11	10	8	9	8	8	4	3	2	1	7

Bio-image Sensitivity Analysis (Complete Project) (Brady/Dina/Plant)
Q. Easiest and Hardest Images?

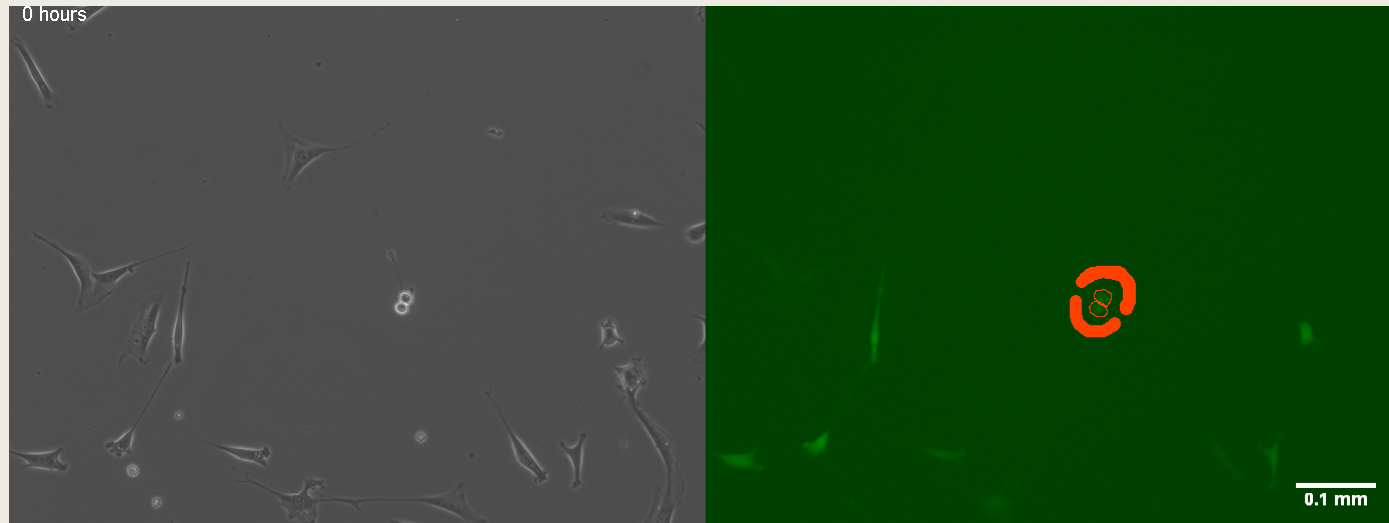


	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
Mean	-199	-1970	-931	-3348	-1102	-913	234	-981	372	-109	-6	-61	-552	-118	-27	41
Median	-183	-1814	-628	-2403	-658	-696	-241	-340	182	-231	27	-244	-488	-128	-106	-67
Rank	6	16	12	18	14	11	8	10	8	7	1	9	13	4	3	2

Live Cell Tracking

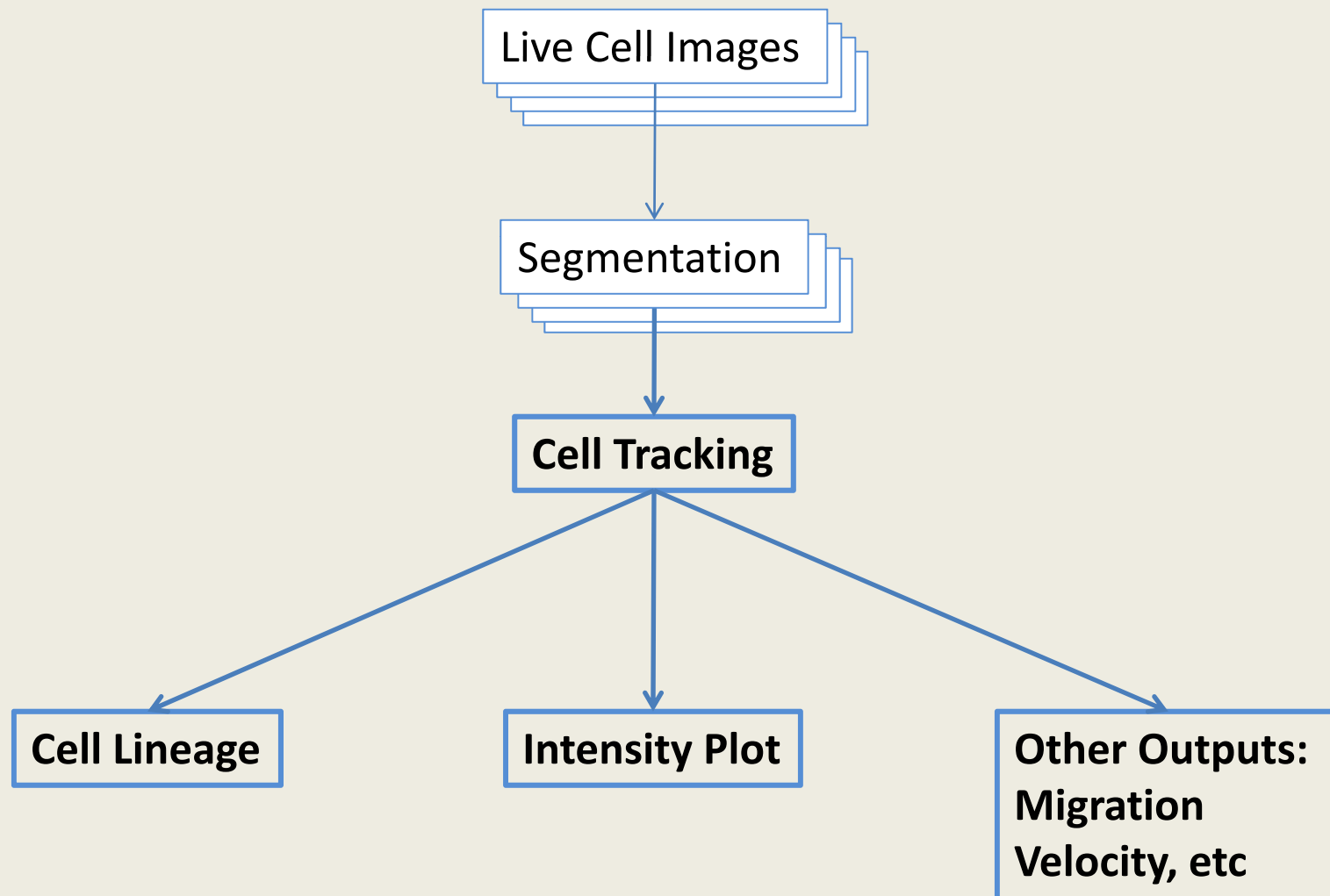
- There are few automated image analysis options available to the cell biologists to quantify live cell image data
- Segment and track cells in an image sequence to quantify the total fluorescence intensity of individual cells over time

Manually Tracking NIH 3T3 Cells

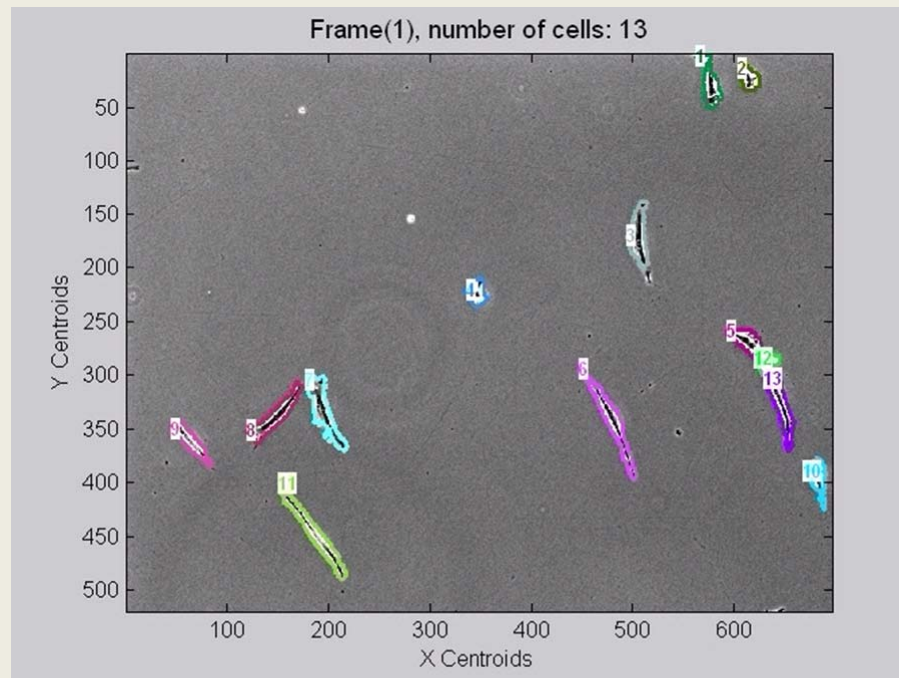


Phase contrast on left, GFP on right

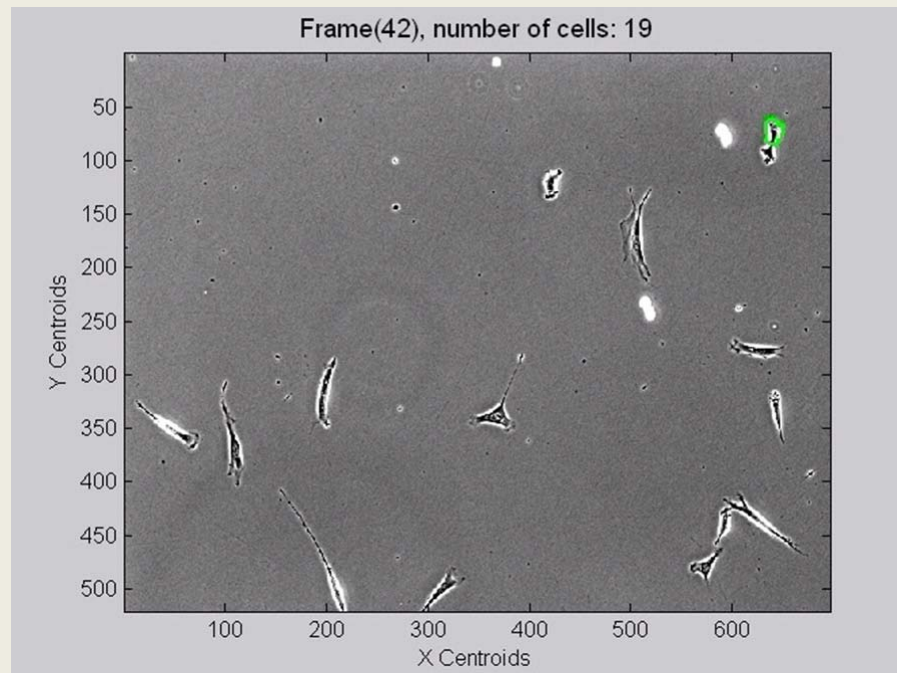
CompBio Cell Tracker



Automatic Tracking of NIH 3T3 Cells



Automatic Tracking of NIH 3T3 Cells



Thank You!

alden.dima@nist.gov