# Integrating Information from Diverse Microscope Images: Learning and Using Generative Models of Cell Organization

**Robert F. Murphy**

Ray & Stephanie Lane Professor of Computational Biology and
Professor of Biological Sciences, Biomedical Engineering and Machine Learning
External Senior Fellow, Freiburg Institute for Advanced Studies
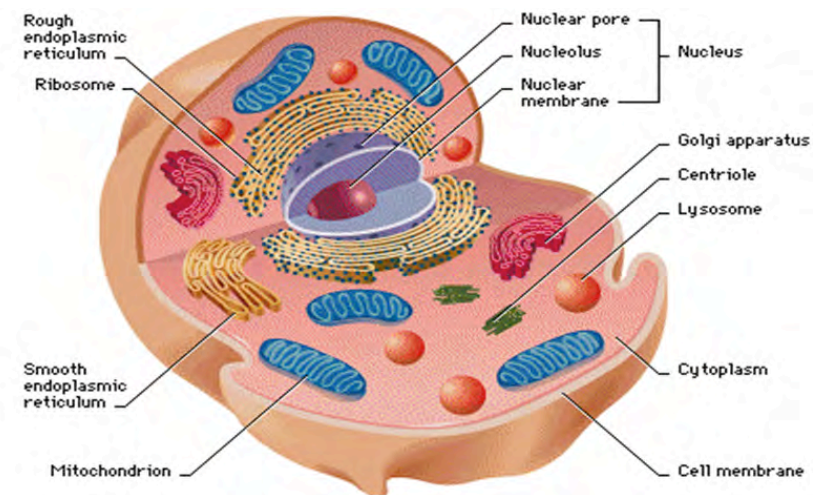Honorary Professor, Faculty of Biology, University of Freiburg, Germany

**Carnegie Mellon University**
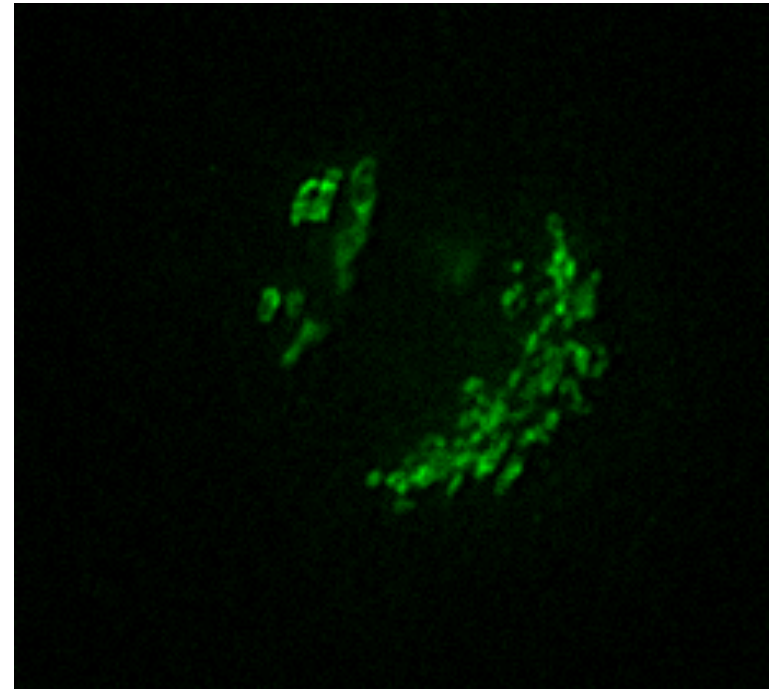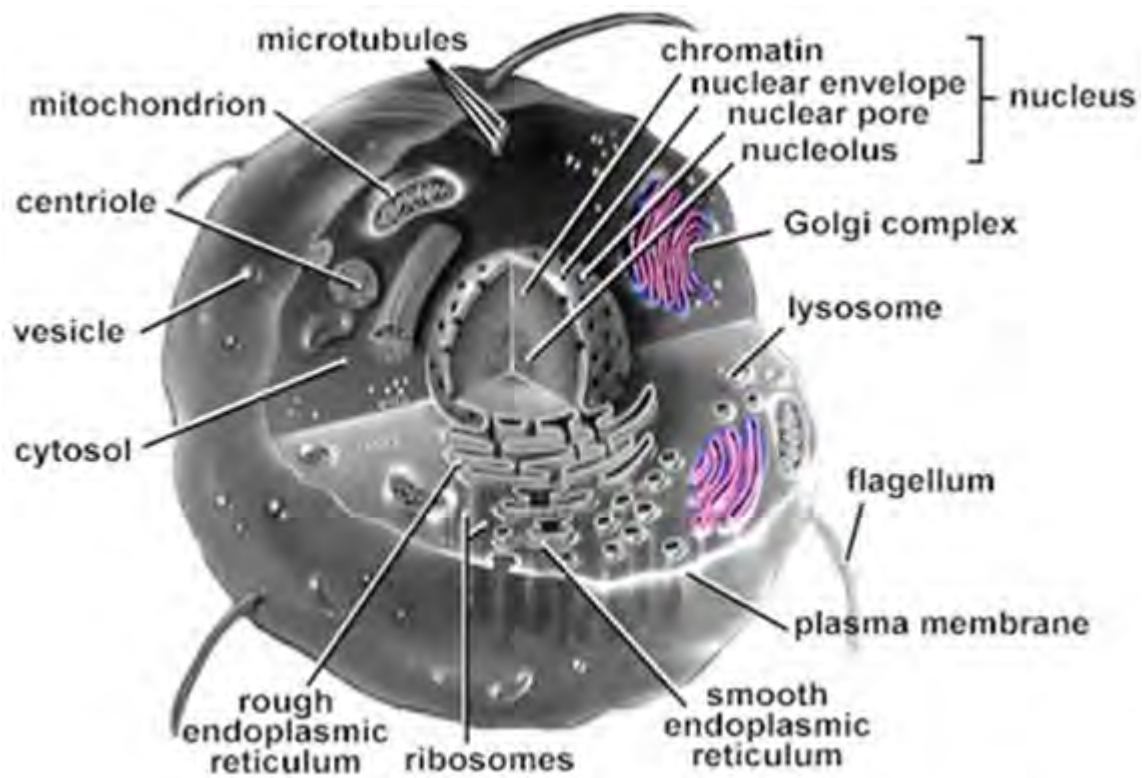Computational Biology Department

**MMBioS**

An NIH Biomedical Technology Research Center

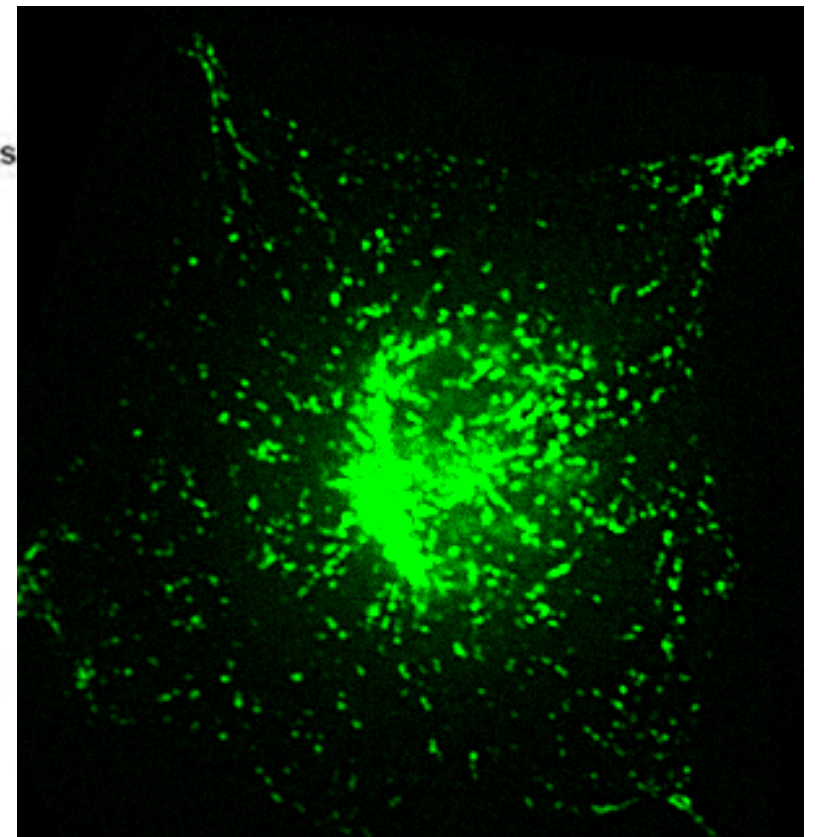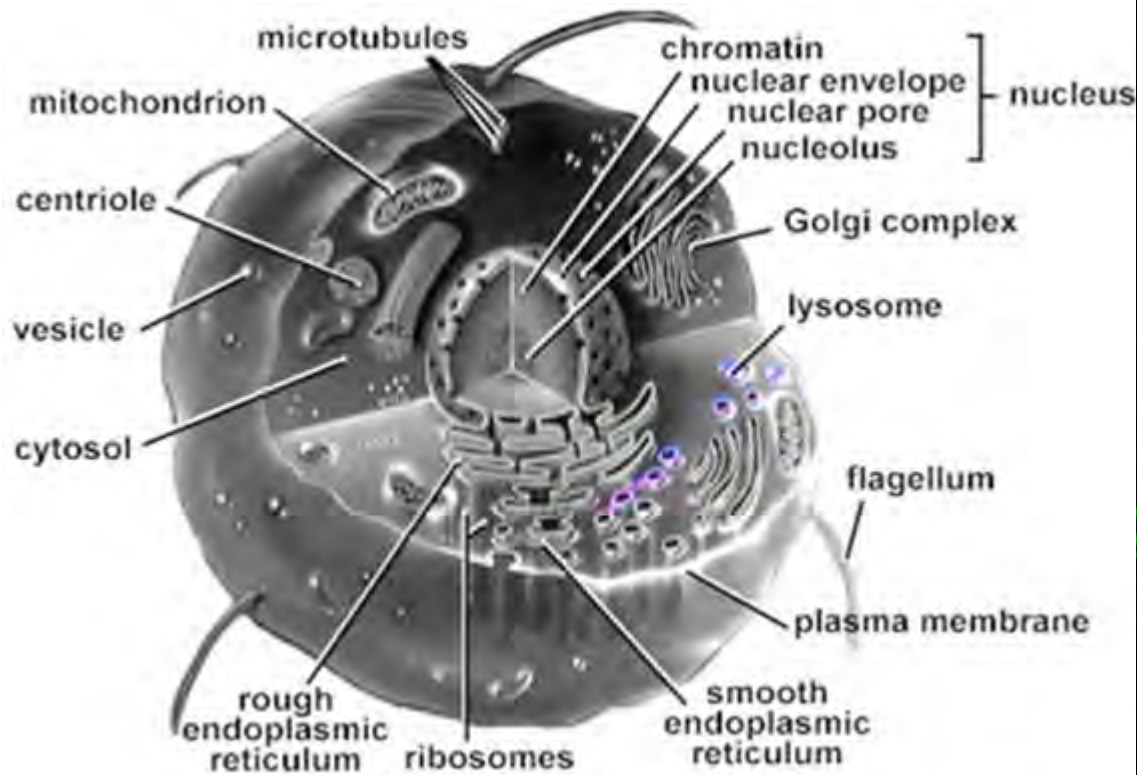# Classic problem in cell and developmental systems biology

- How do we learn and represent
  - sizes and shapes of different cell types
  - number, sizes, shapes, positions of organelles
  - the distribution of proteins across organelles
  - how organelles depend upon each other
  - how any of these vary
    - from cell to cell
    - from cell type to cell type
    - during development
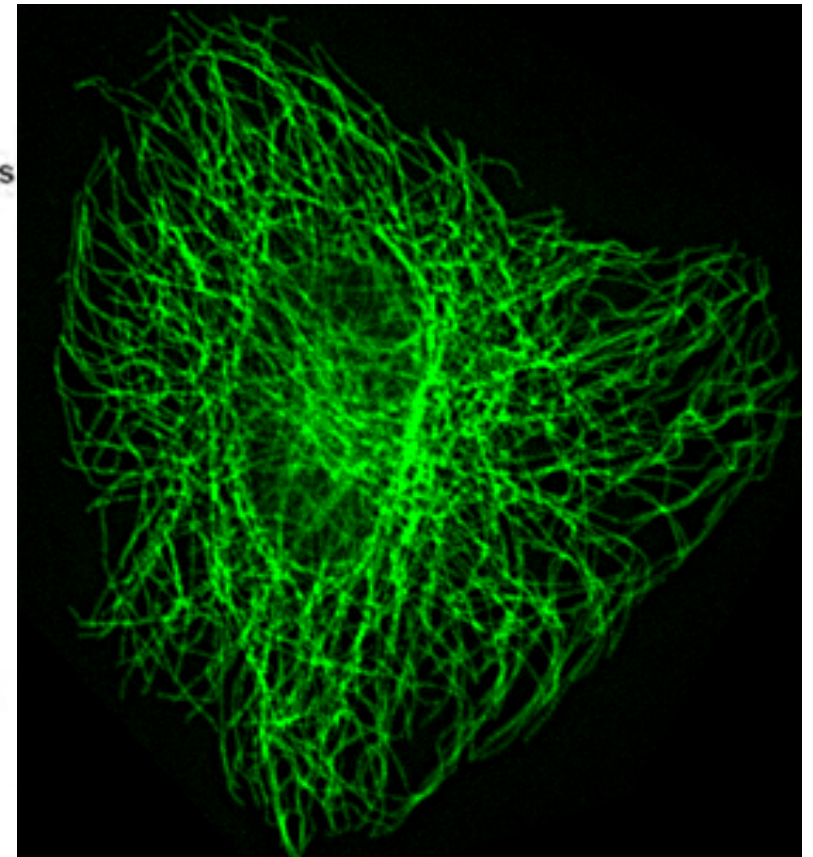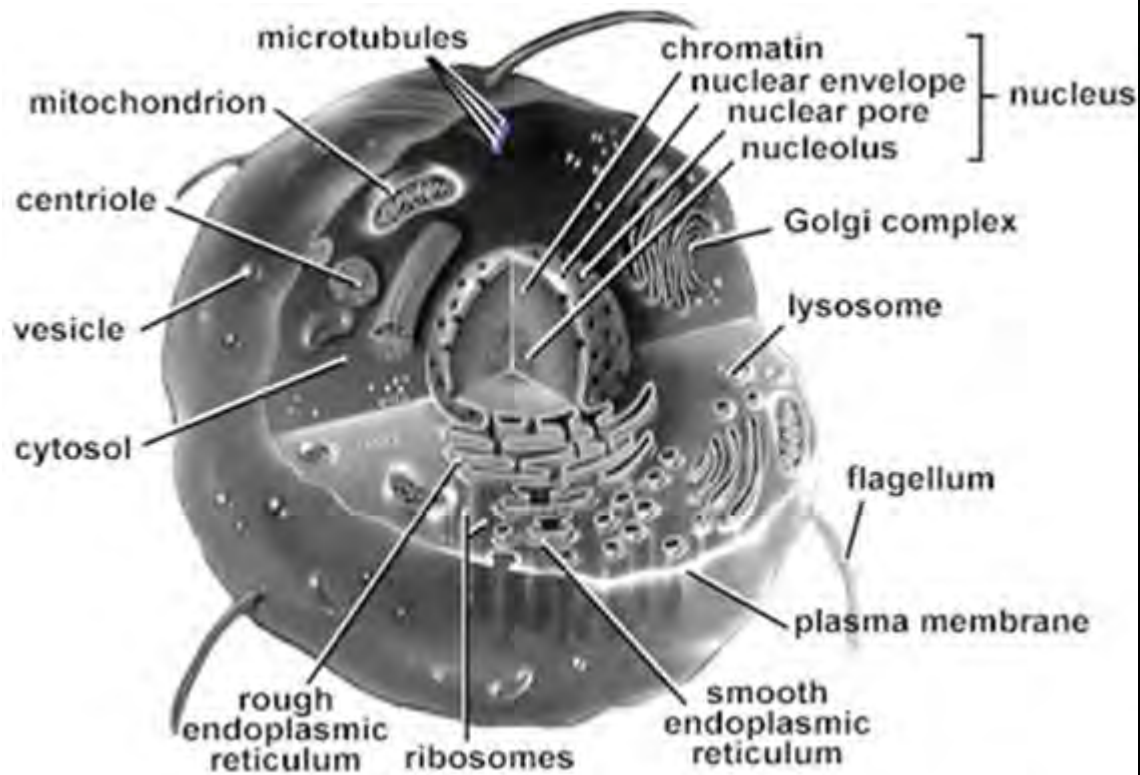    - in presence of perturbagens
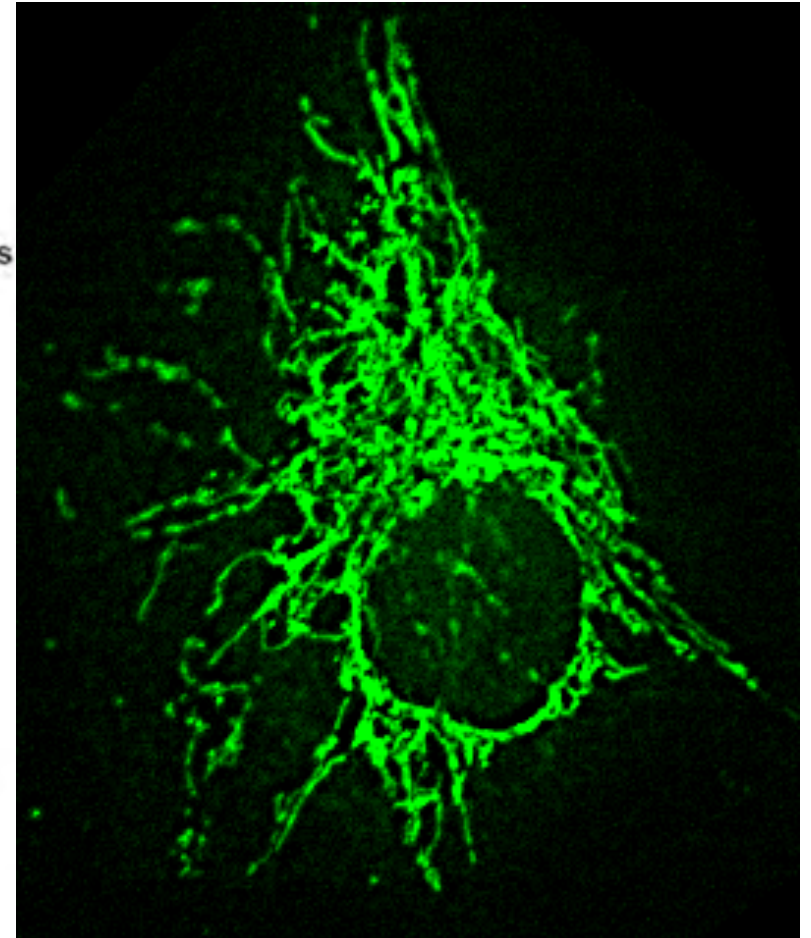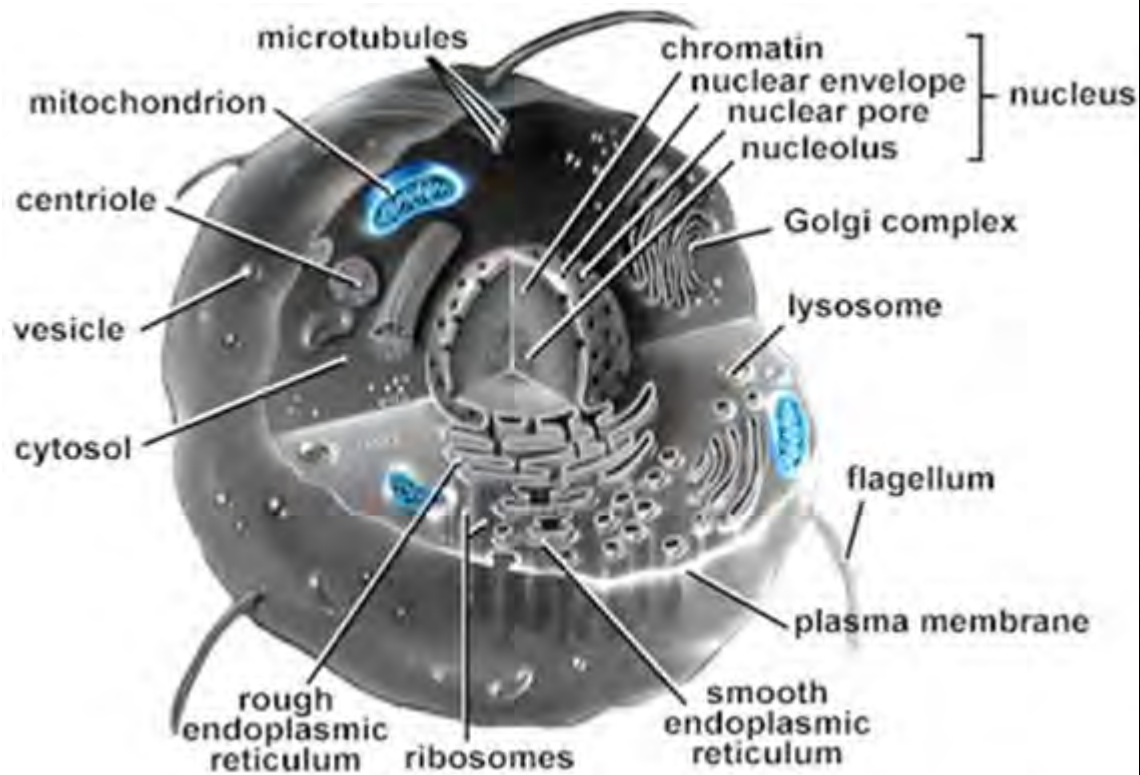
# Subcellular Location

# Subcellular Location

# Subcellular Location

# Subcellular Location

# Classic approach

- Do biochemical or imaging experiments, capture relationships in words
  - "secretory vesicles bind to microtubules"
- Two problems
  - Difficult to establish these relationships from images
  - Does not adequately describe them
- Can we do better via machine learning?

# Cellular Pattern Recognition



- Describe cell patterns using numerical features

- Do classification, etc. to assign terms

- First described in Boland, Markey & Murphy (1998) and Boland & Murphy (200

- Later popularized in packages such as CellProfiler, WND-CHARM, Ilastik, CellCognition, etc.

# Drawback

- Image features are typically not transferable across images from different sources (widefield vs. confocal vs. superresolution, differences in magnification or camera pixel size, pixel bit depth, etc.)

# Traditional High Content Screening/ Analysis

# Different HCS systems or sites



Traditional HCS design focused on finding hits *within* a given screen, not on comparing results *between* screens or learning *generalizable* effects

# Another drawback

- Term assignment/classification approaches are incomplete and do not make full use of information in images

- "Is this an apple or an orange?" is a *discriminative* question; can be answered with 1 or 2 features

- "What does an apple look like?" requires a *generative model*

# Generative models?

- Human cognition

 examples

Learn

"mental model"

Write

Generated examples

- Image-based models

 Training images

**MODEL** Statistical generative model

Generated image

Zhao & Murphy, Cytometry 2007

# Open source project: CellOrganizer

Cell Images

**Statistical Model**

Synthetic Images

Nuclear shape → Cell shape

**Training**

**Synthesis**

Object pos. probability

Microtubule distribution

Object appearance

Object positions

Object number

Object distribution

http://CellOrganizer.org

**CellOrganizer**
Images ⟷ Models

Carnegie Mellon University — Center for Bioimage Informatics

Home   News   People   Publications   Downloads   Documentation

The **CellOrganizer** project provides tools for

- learning generative models of cell organization directly from images
- storing and retrieving those models
- synthesizing cell images (or other representations) from one or more models

Model learning captures variation among cells in a collection of images. Images used for model learning and instances synthesized from models can be two- or three-dimensional static images or movies.

**CellOrganizer** can learn models of

- cell shape
- nuclear shape
- vesicular organelle size, shape and position
- microtubule distribution
- average protein distributions

These models can be *conditional* upon each other. For example, for a given synthesized cell instance, organelle position is dependent upon the cell and nuclear shape of that instance.

Cell types for which generative models for at least some organelles have been built include human HeLa cells, mouse NIH 3T3 cells, Arabidopsis protoplasts and mouse T lymphocytes.

2D HeLa
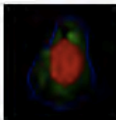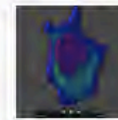(endosomes)

3D HeLa
(mitochondria)

3D protoplast
(chloroplasts)

3D HeLa
(microtubules)

3D HeLa movie

# Generative vs. discriminative HCS

not comparable

**Differences** ← **Differences**

Features    Features          Features    Features

Control images from HCS1    Hit images from HCS1          Control images from HCS2    Hit images from HCS2

Params    Params          Params    Params

comparable

**Differences** ← **Differences**
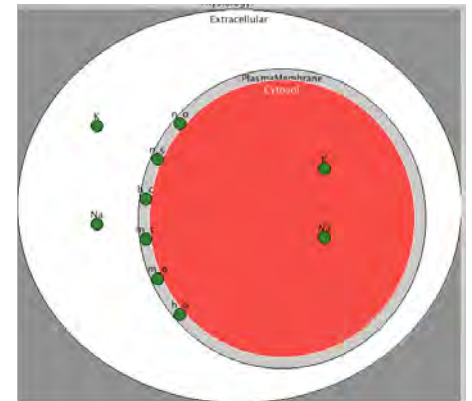
# Compartmental models for cell simulations

- Use the assignments to put each protein "in" its compartment and
  - use a cartoon compartmental model
  - use real image to determine compartment volume/surface area
  - use PDEs for each pixel of a real image
- These geometries are not very realistic

# CellOrganizer modeling goals

- Cell models should be

    - **Automated**: learned directly from images,

    - **Generative**: able to synthesize new examples,

    - **Statistically accurate**: reflect variation from cell to cell,

    - **Compact**: can be communicated with significantly fewer bits than the training data.
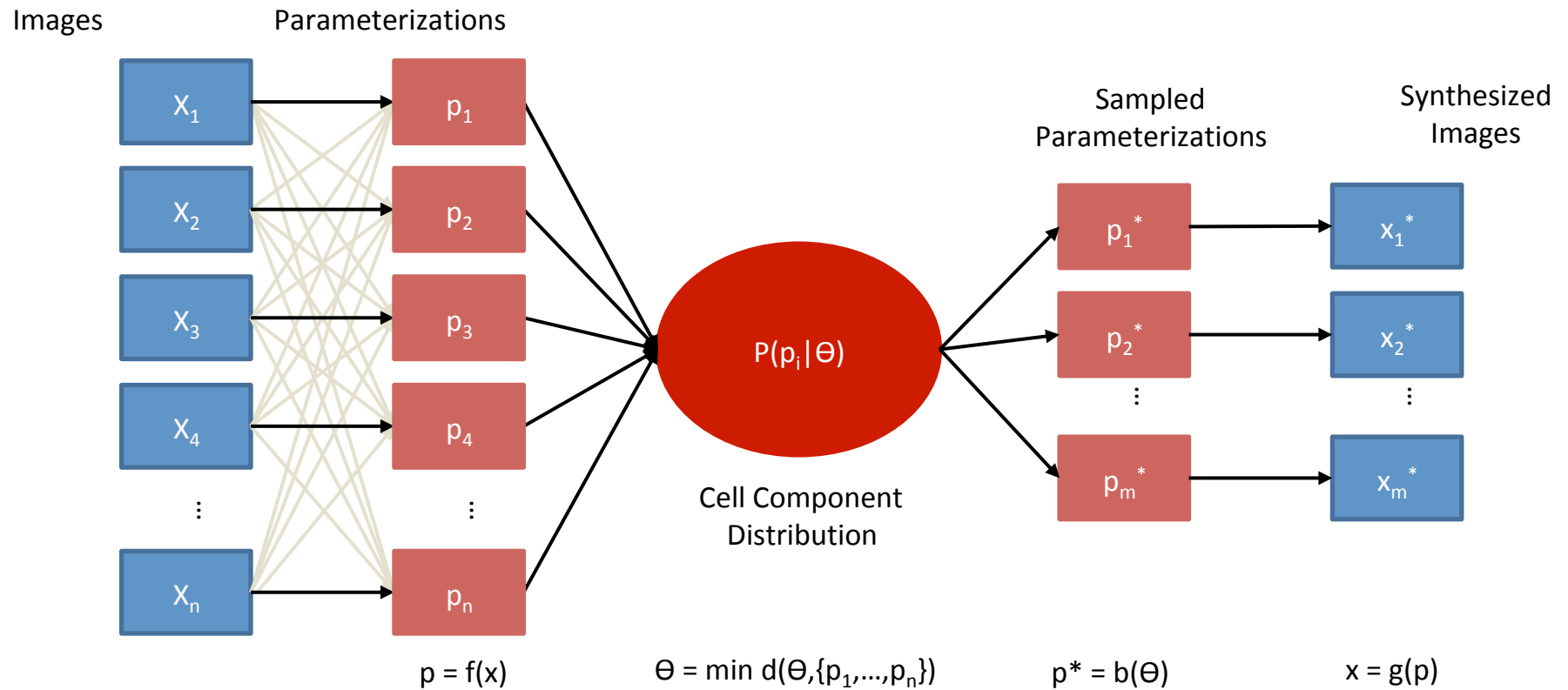
# Classical inverse problem

- Learn underlying reality observed via imaging
- Extensive work on image **reconstruction** to create a (higher resolution?) model of a conserved structure (e.g., nuclear pore, ribosome) by *removing noise and variation*
- Our goal is learning *statistical,* **generative** model of reality sampled via imaging by removing noise but *keeping* variation

# Parametric models

- Computer vision problems such as this have traditionally been tackled by hand-constructing models and learning their parameters from images

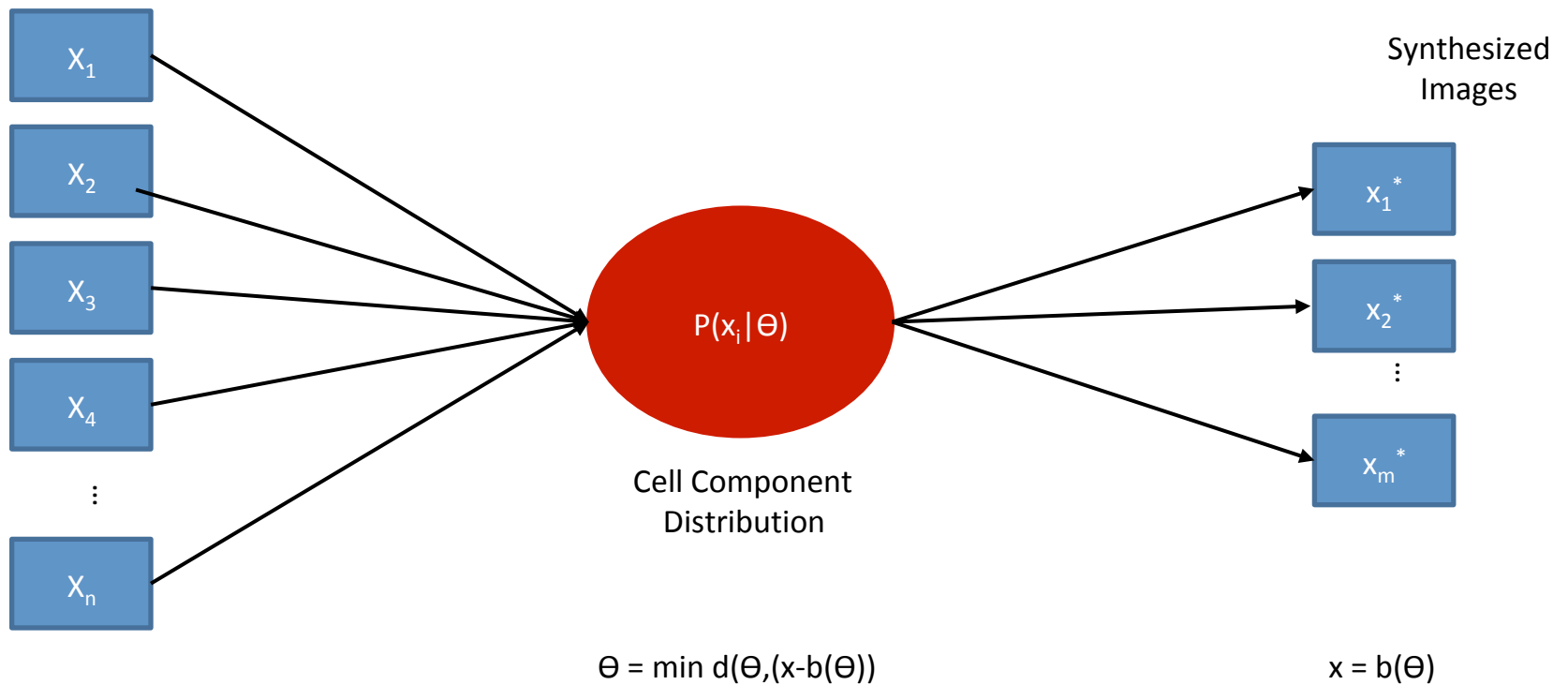# Parametric modeling (e.g., CellOrganizer)

# "Deep" learning

- If large numbers of training examples are available, "deep learning" methods can learn directly from images without need for custom design

# Deep learning models (e.g., autoencoders)

Images

$X_1$

$X_2$

$X_3$

$X_4$

$\vdots$

$X_n$

$P(x_i|\Theta)$

Cell Component
Distribution

Synthesized
Images

$x_1^*$

$x_2^*$

$\vdots$

$x_m^*$

$\Theta = \min d(\Theta,(x-b(\Theta))$

$x = b(\Theta)$

# But...

- Large numbers of "labeled" training images are often not available
- Deep learning models only understand pixels, not structures/objects
  - Not easily compared/combined across diverse images
  - Cells are not made of probabilistic "blobs" of macromolecules
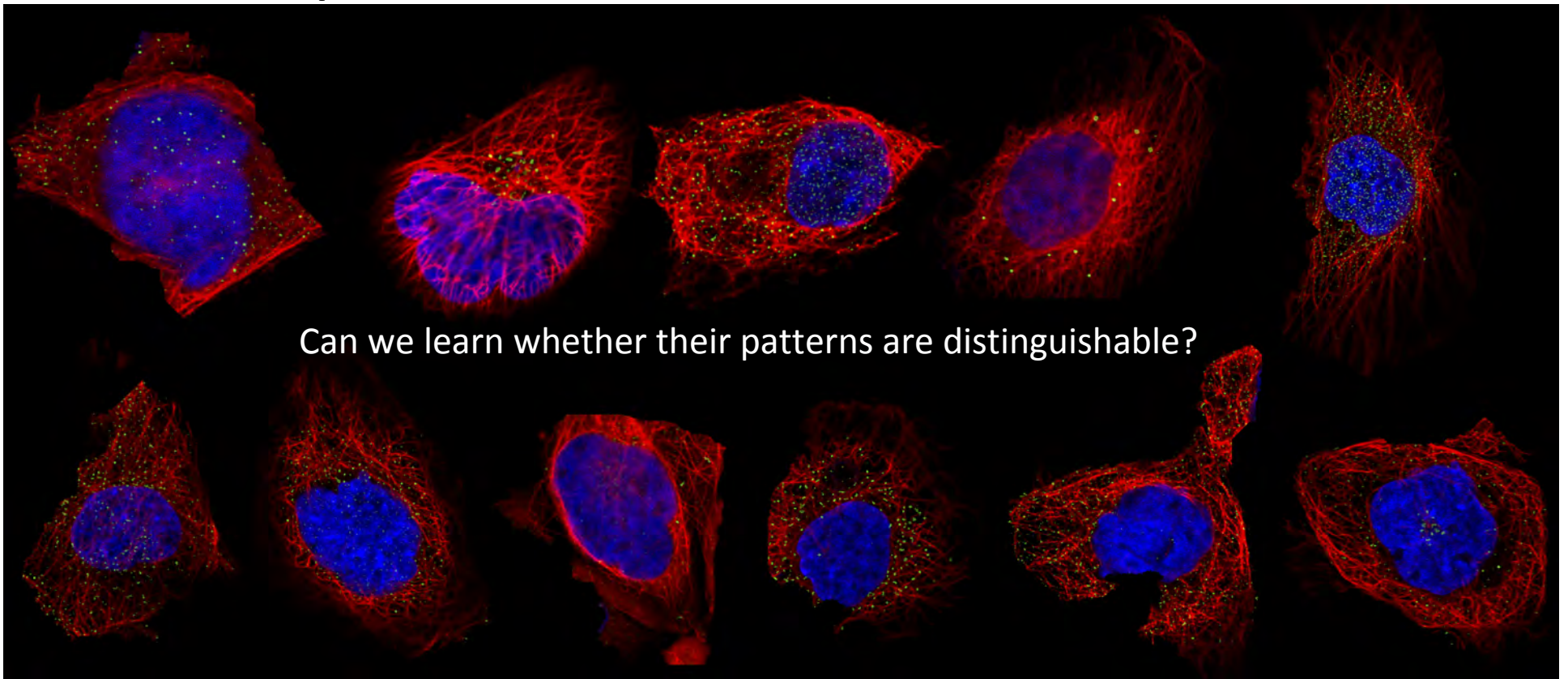  - Many organelles have discrete boundaries/structures

# Challenge

- Fluorescent microscopy provides very useful information about cell organization and processes
- But the number of molecules that can be imaged at the same time in live cells is smaller than the number involved in many processes
- How do we combine information from different images to provide coherent picture?

# Solution?

- Merging information through generative models built upon a common reference

- Two examples:
  - Distinguishing different punctate structures from separate images
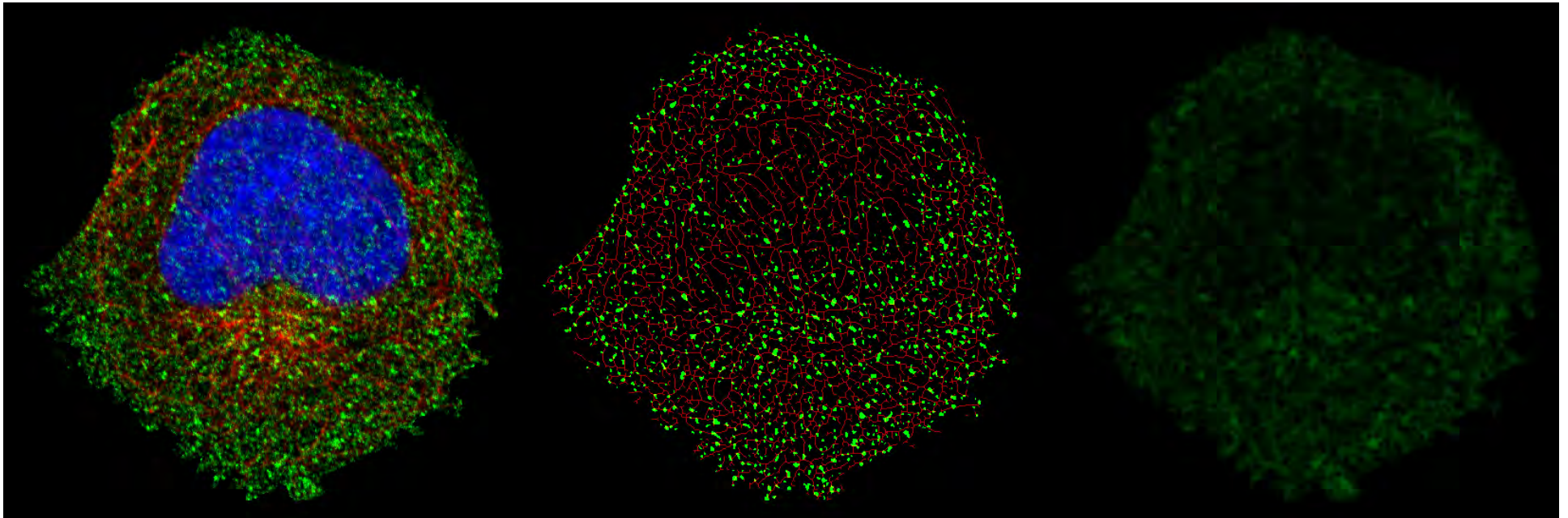  - Learning potential spatial causal relationships involving in cell signaling

# MODELING PUNCTATE ORGANELLE DISTRIBUTIONS

# Images of 11 different "vesicle" proteins from Human Protein Atlas

Can we learn whether their patterns are distinguishable?

# Segmentation of punctate organelles

- Use high pass filter



Original image

segmented puncta
and microtubules

remaining
fluorescence

# Point process models

- Capture relationship between position of an organelle and positions of organelles of different types (**"inhomogeneous Poisson process"**)

$$f\left(X^{(n)}|n\right) = \frac{1}{Z_\theta}\prod_{i=1}^{n} b_\theta\left(X_i\right)$$

- Positions of $n$ organelles depend upon $b_\vartheta$ functions

# Factors for point process models

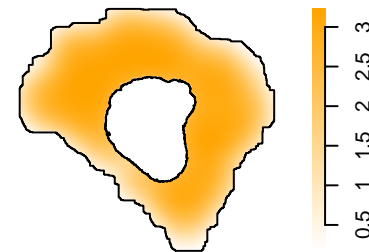- The functions depend upon specified factors, variables for which values are known at all positions in the cell



**Distance to nuclear boundary**

**Distance to cell boundary**
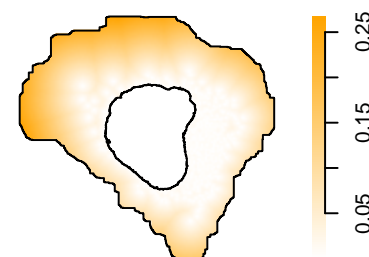
**Kernel density of microtubules**

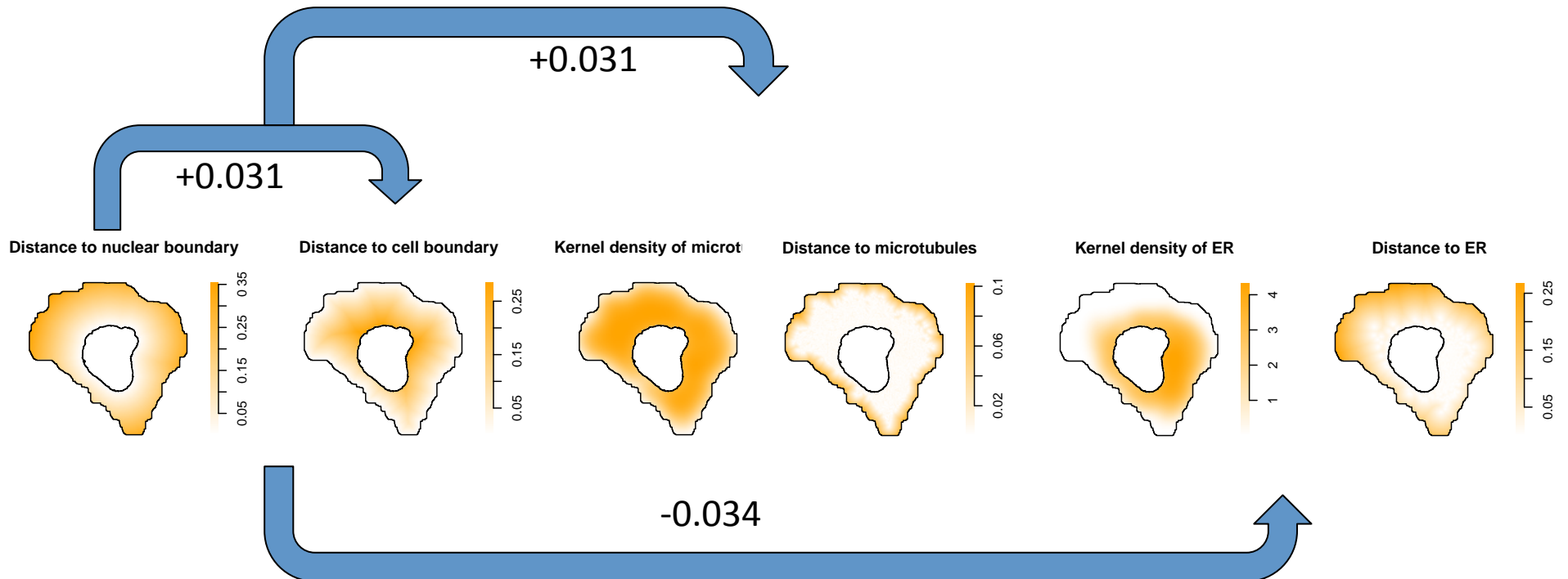**Distance to microtubules**

**Kernel density of ER**

**Distance to ER**

# Learning dependencies on factors

- An important question is to learn on *which* factors a particular pattern depends
- Can do this by cross-validation: for each combination of factors
  - Estimate parameters from training data
  - Estimate likelihood of test data being generated by that model
  - Average likelihoods

# Contributions of different factors

# How different are the 11 punctate patterns?

- Can also assess by cross-validation (only 2 images available in HPA!)

- Train 11 models using 1 image of each protein

- Assign remaining test image of each protein to the model that it has the highest likelihood of it having been produced by
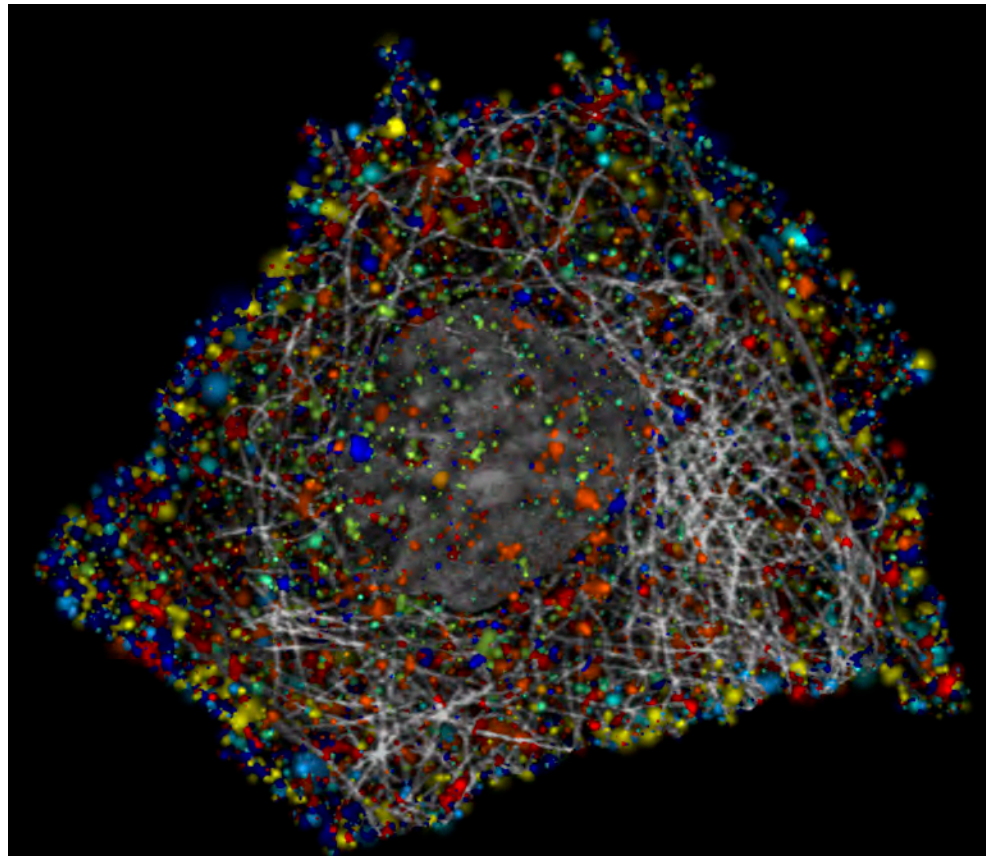
# 11 distinct punctate patterns using relationship to microtubules

| U-251 MG | COPI | COPII | Caveolae | Coated Pits | Early Endosomes | Late Endosomes | Lysosomes | Peroxisomes | RNP bodies | Recycling Endosomes | Retromer |
|---|---|---|---|---|---|---|---|---|---|---|---|
| COPI | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| COPII | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Caveolae | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Coated Pits | 0 | 0 | 0 | 0.67 | 0 | 0 | 0 | 0 | 0 | 0 | 0.33 |
| Early Endosomes | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| Late Endosomes | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| Lysosomes | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| Peroxisomes | 0.08 | 0 | 0 | 0 | 0 | 0 | 0 | 0.77 | 0 | 0.08 | 0.08 |
| RNP bodies | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| Recycling Endosomes | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| Retromer | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

Overall accuracy:

| A-431 | 0.73 |
|---|---|
| U-2 OS | 0.90 |
| U-251 MG | 0.86 |

# Example synthetic cell image with 11 punctate organelles

# MODELING SUBCELLULAR DISTRIBUTION CHANGES DURING CELL SIGNALING
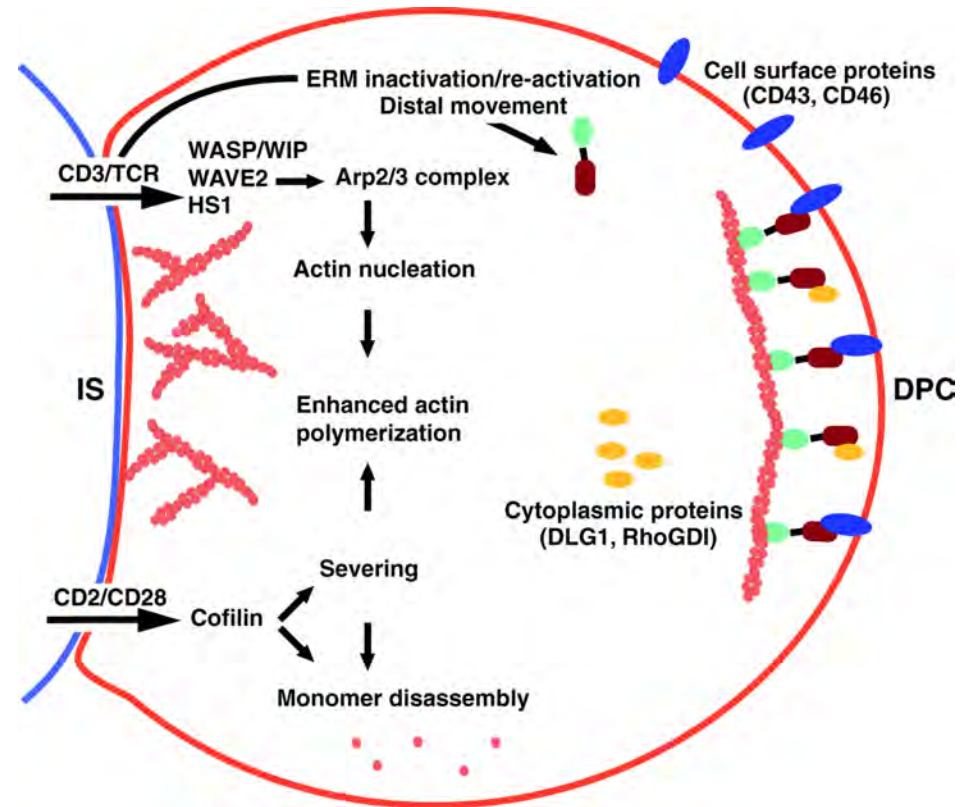
RESEARCH RESOURCE

TECHNIQUES

# Computational spatiotemporal analysis identifies WAVE2 and cofilin as joint regulators of costimulation-mediated T cell actin dynamics

Kole T. Roybal,[1,2*†] Taráz E. Buck,[3†] Xiongtao Ruan,[3†] Baek Hwan Cho,[3‡] Danielle J. Clark,[1] Rachel Ambler,[1] Helen M. Tunbridge,[1] Jianwei Zhang,[3§] Paul Verkade,[4] Christoph Wülfing,[1,2,5¶‖] Robert F. Murphy[3,6,7¶‖]

# Background

- T cells bind to APC cells triggering stimulation
- Actin and its regulators are recruited to the interacting region.

Question: how do the proteins involved in Actin dynamics regulate each other?



Huang & Burkhardt, 2007

# Data



- We start from DIC and fluorescence movies of GFP-tagged proteins at different time points before and after immunological synapse formation (~100 cells per protein)

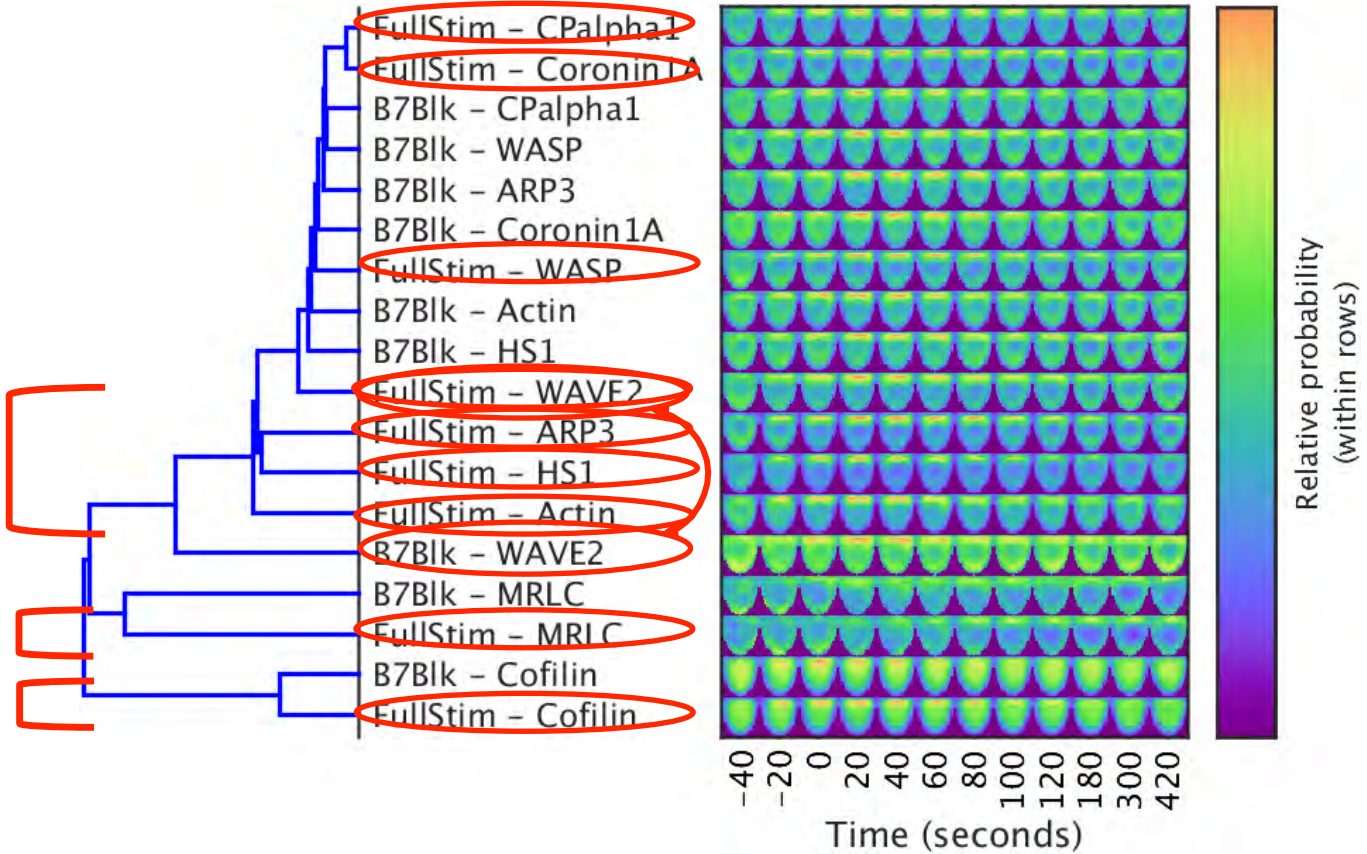# Image processing pipeline

# Summary of the Automatic analysis

- More than 17,000 cell pairs were analyzed.

- Two conditions: Full stimulus, B7 blockade

- Ten proteins: ARP3, Actin, cofilin, Coronin1A, CPalpha1, HS1, MRLC, WASP, WAVE2, LAT

# Applications of the model

- Use voxel concentrations as features to compare different proteins across all time points across different conditions

- Enrichment analysis: see how proteins accumulate in specific locations over time

- Visualize the spatiotemporal dynamics of selected proteins

# Clustering reveals differences
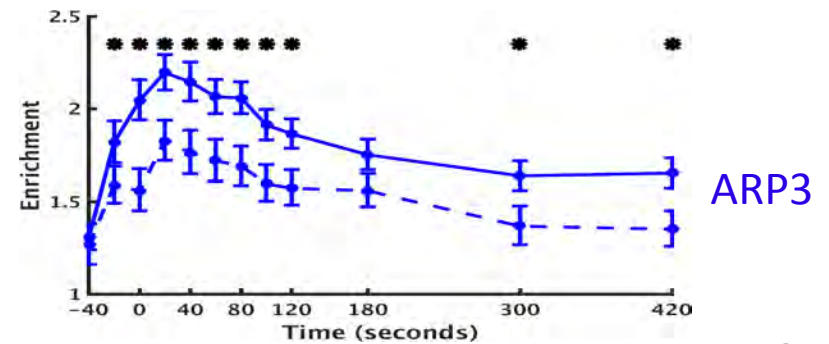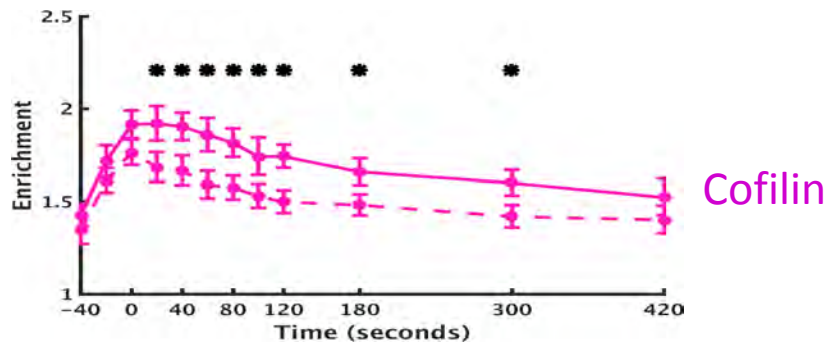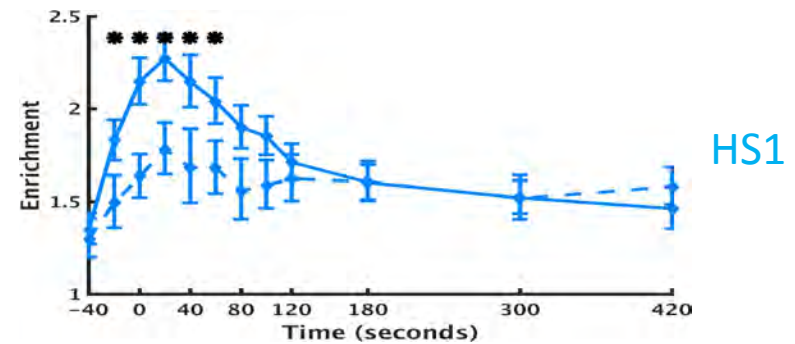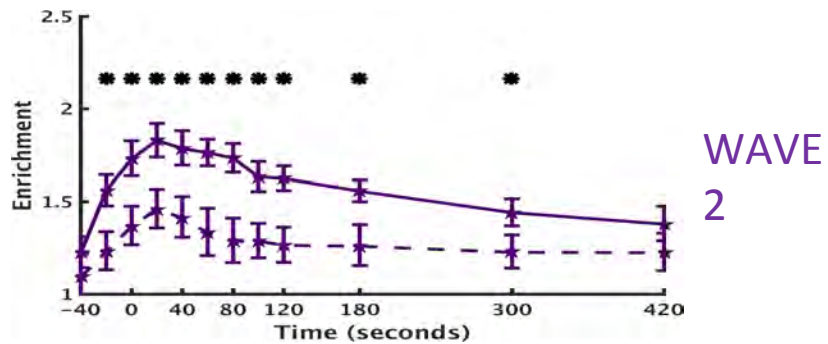
# Enrichment analysis

- Actin is recruited to the synapse region, we would like to measure kinetics of recruitment of other proteins

- Idea: define an enrichment region where it contains the top 90% fluorescence in the average model map for all models.

- The enrichment is defined as the ratio of the mean intensity in the region to the mean intensity out of the region.
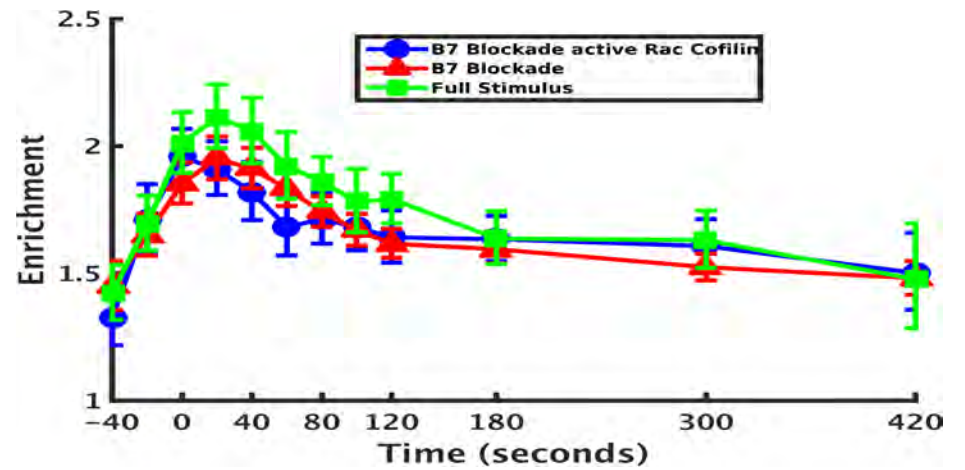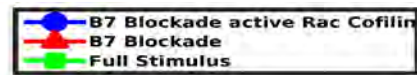
# Enrichment of all proteins

# Several proteins have significant enrichment

# Validation of candidate regulators

- We identified Wave2 and cofilin as candidate regulators in costimulation-mediated Actin dynamics.

- Question: could selective activation of these two regulators promote actin dynamics and synapse formation even under the B7 blockade condition?
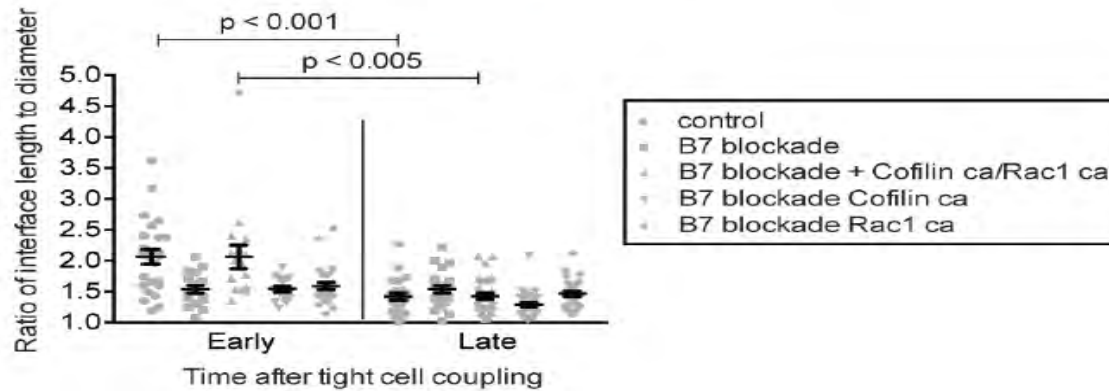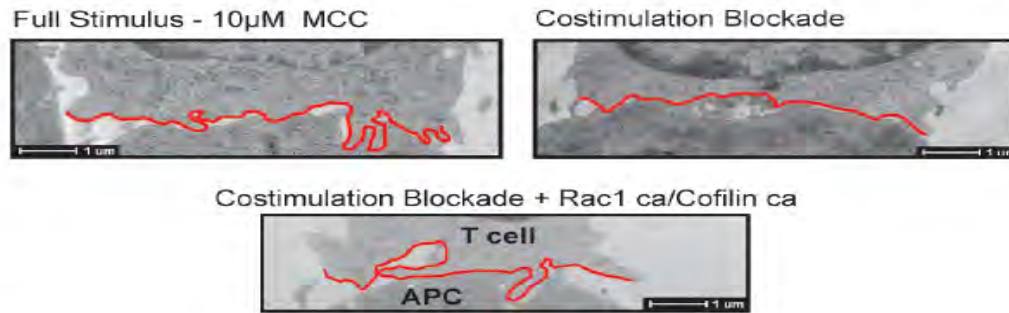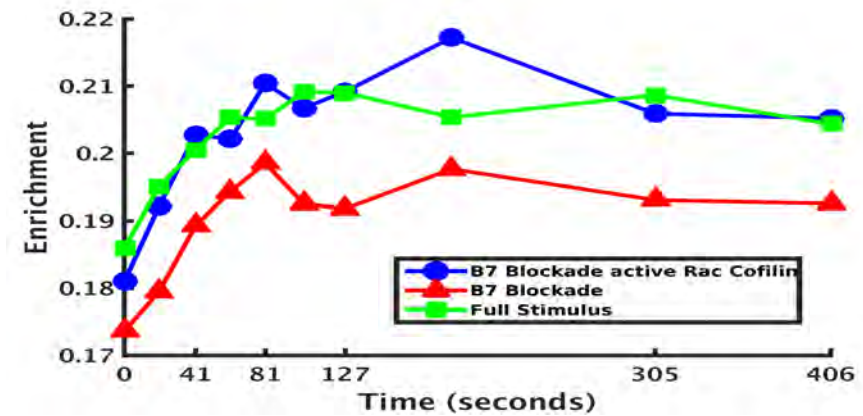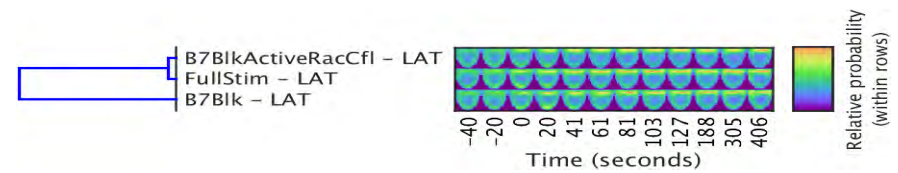
# Reconstruction



**WAVE2**                                    **Actin**

# Reconstruction



Full Stimulus - 10µM MCC

Costimulation Blockade

Costimulation Blockade + Rac1 ca/Cofilin ca

T cell

APC

p < 0.001

p < 0.005

control
B7 blockade
B7 blockade + Cofilin ca/Rac1 ca
B7 blockade Cofilin ca
B7 blockade Rac1 ca

Ratio of interface length to diameter

Early

Late

Time after tight cell coupling

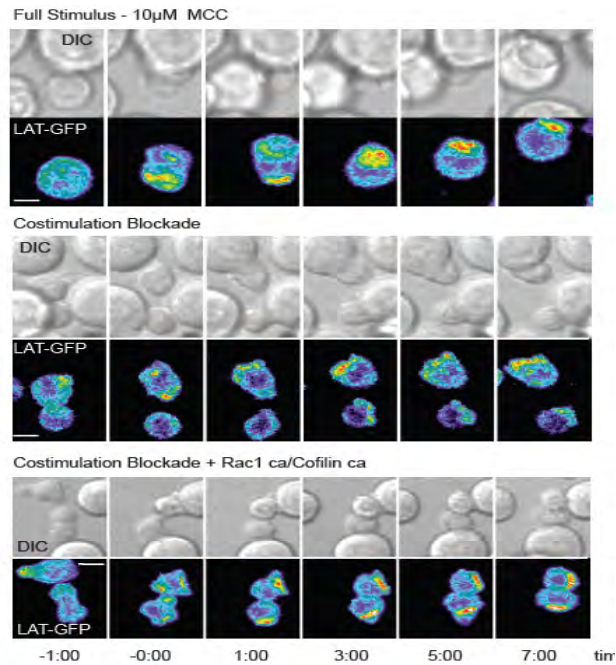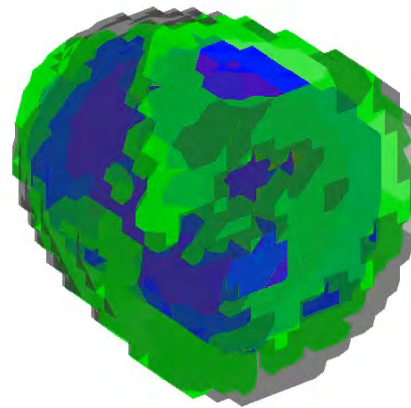# Active Rac and cofilin restore defective LAT location

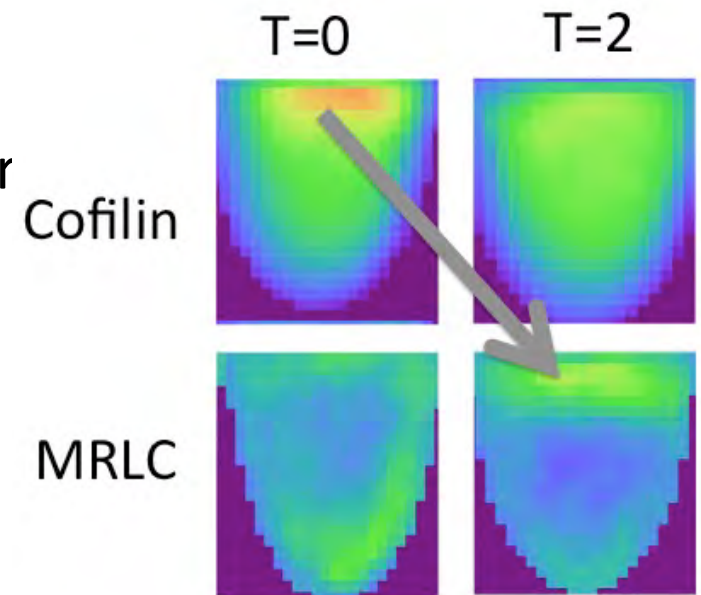# Spatiotemporal distribution of proteins



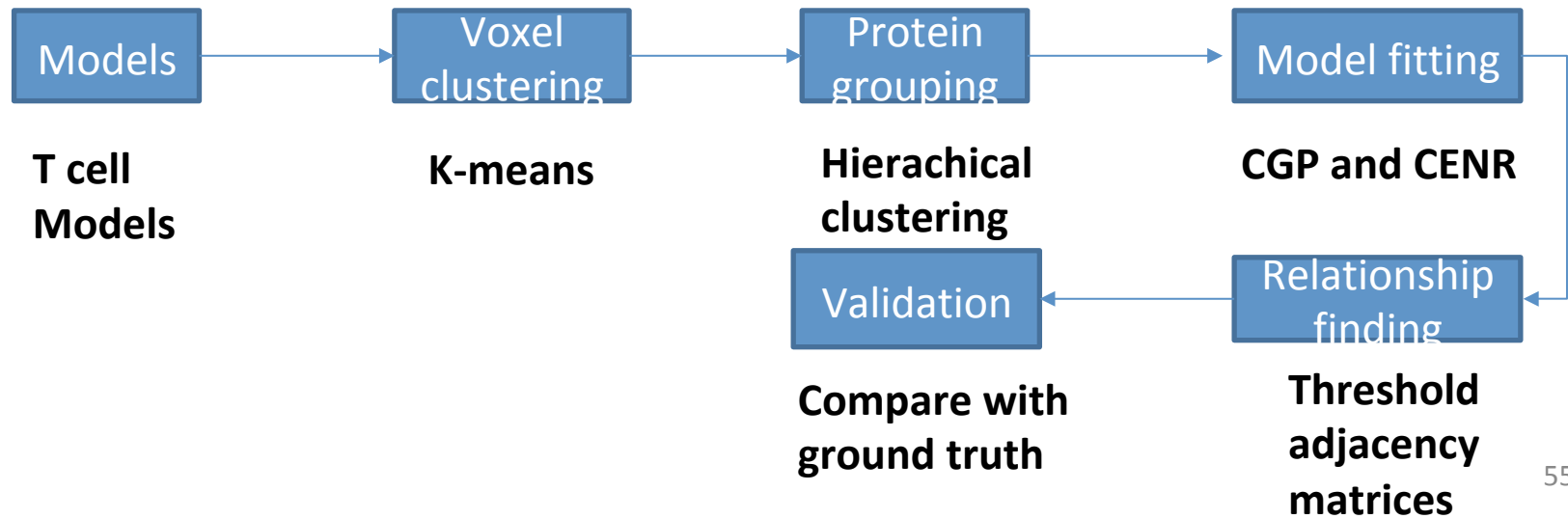**<span style="color:red">cofilin</span> <span style="color:green">MRLC</span> <span style="color:blue">WAVE2</span>**

# Using Spatiotemporal Maps to Learn Putative Regulatory Relationships

- Given spatiotemporal maps for multiple proteins, we sought to determine whether a change in one protein in one region of the cell precedes a change in another protein in another region

- For this we applied methods for learning causal graph process models

- Nodes represent a specific protein at a specific location, edges represent a possible predictive relationship between nodes
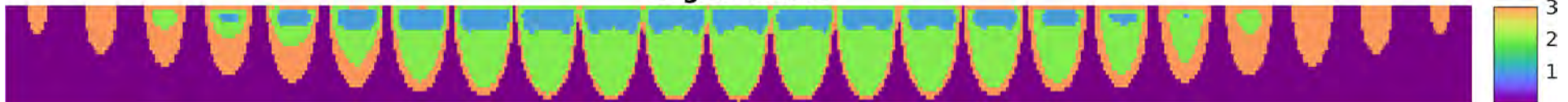
# Approach

```
┌──────────┐     ┌──────────┐     ┌──────────┐     ┌──────────┐
│  Models  │────▶│  Voxel   │────▶│ Protein  │────▶│  Model   │──┐
│          │     │clustering│     │ grouping │     │ fitting  │  │
└──────────┘     └──────────┘     └──────────┘     └──────────┘  │
                                                                  │
  T cell          K-means         Hierachical      CGP and CENR   │
  Models                          clustering                      │
                 ┌──────────┐                     ┌──────────┐    │
                 │Validation│◀────────────────────│Relationship│◀─┘
                 │          │                     │ finding  │
                 └──────────┘                     └──────────┘

                  Compare with                     Threshold
                  ground truth                     adjacency
                                                   matrices
```

55

# Reducing graph size

- The spatiotemporal maps have 6628 voxels per cell, and there is one map for each of 12 proteins
- The graph model requires an edge between every pair of nodes: too many edges, need to reduce the number
- Solution: represent each voxel by a vector of the intensities at all time points for all proteins
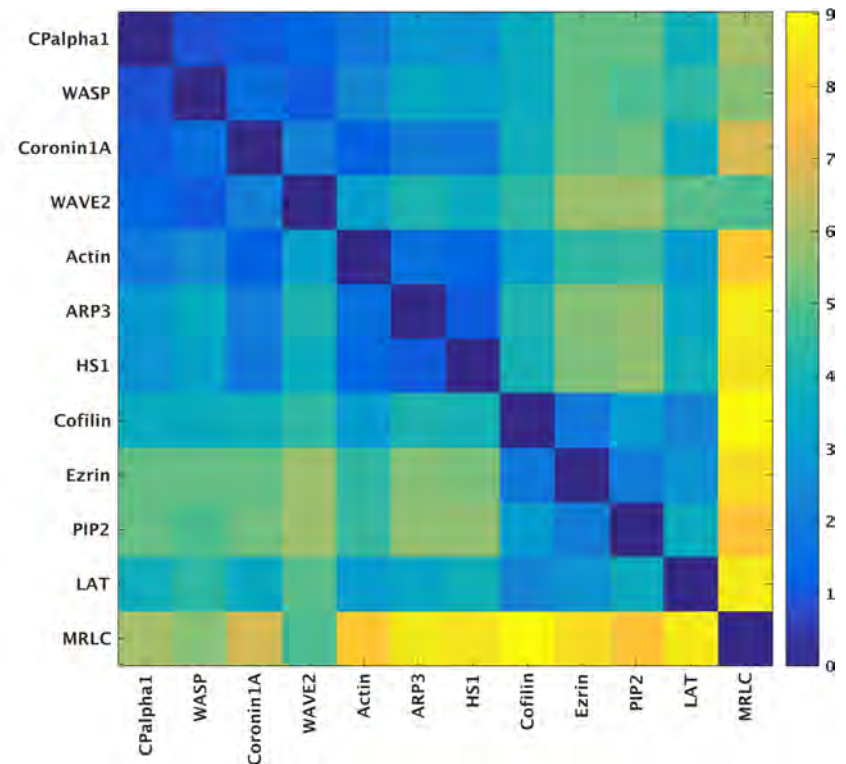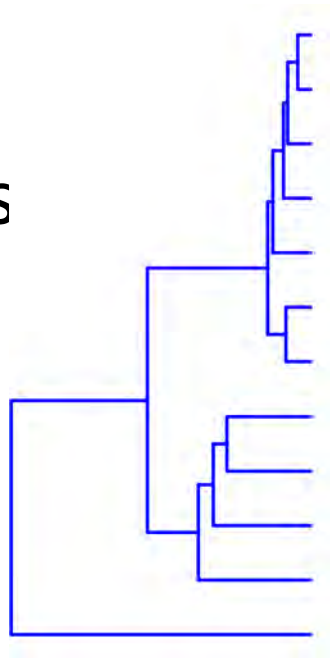- Use K-means clustering of voxels to form regions

**Region labels**



- Represent each region by the average intensity of all voxels in it

# Removing confounding by correlations

- Highly correlated proteins/regions allow self prediction
- Collapse them into one representative

# Causal graph process model

Let $x[t]$ be concentration of all proteins in all regions

Goal is to find single model to predict all times

$$x[t] = w[t] + \sum_{i=1} P_i(\mathbf{A})x[t-i]$$

$$= w[t] + (c_{10}\mathbf{I} + c_{11}\mathbf{A})x[t-1]$$

$$+ (c_{20}\mathbf{I} + c_{21}\mathbf{A} + c_{22}\mathbf{A}^2)x[t-2] + \cdots$$

$$+ (c_{M0}\mathbf{I} + c_{M1}\mathbf{A} + \cdots + c_{MM}\mathbf{A}^M)x[t-M]$$

Mei & Moura, 2015

We also used a second method we developed called CENR

Given $x[t]$ and M, find $\mathbf{A}$ and $c$ to minimize $w[t]$

# Adjacency matrix of CGP method

# Summarizing relationships

- We threshold the strength of the relationships, (elements in the adjacency matrices), and identify the time when each is most strongly observed.

# Evaluation

- Make list of known or suspected regulatory relationships from literature
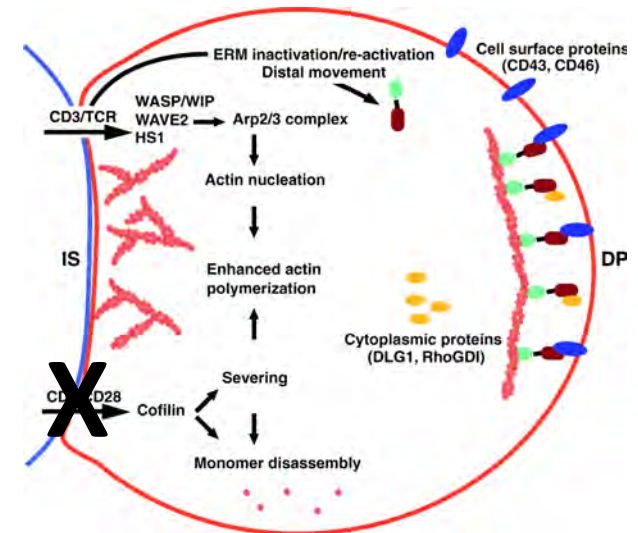
# Evaluation

- Measure how well learned models capture these known relationships using a Receiver-Operator Curve (quantitate by "Area Under the Curve")

| | CGP | CENR |
|---|---|---|
| AUC | 0.709 | 0.644 |

- Note some "False positives" may be real positives that are not yet known

# Evaluation

- Data used so far was from costimulation conditions (stimulation through both TCR and CD28)

- Additional maps available for conditions where CD28 costimulation is blocked ("B7 blocked")

# Evaluation

- Using the model learned from costimulation condition, make predictions from early time points for blocked condition at later time points

| | CGP | CENR |
|---|---|---|
| Prediction error from cross-validation on training images (full stimulation) | 8.6% | 6.9% |
| Prediction error on testing images (costimulation blocked) | 12.0% | 8.5% |

# HIERARCHICAL ASSEMBLY MODELS

# Spatial models



- Most of you probably have built models: LEGO's, K'nex, etc.

- You get different colored parts and a set of instructions

- The instructions are *hierarchical*

**115**

4x    1x

**116**

9

2x

9

9

**117**
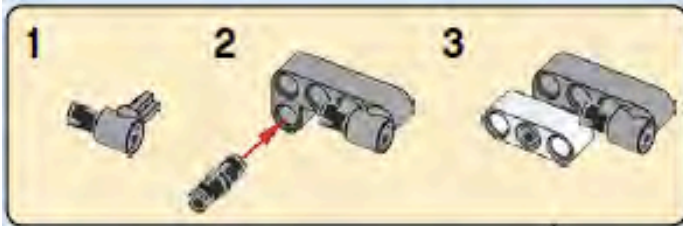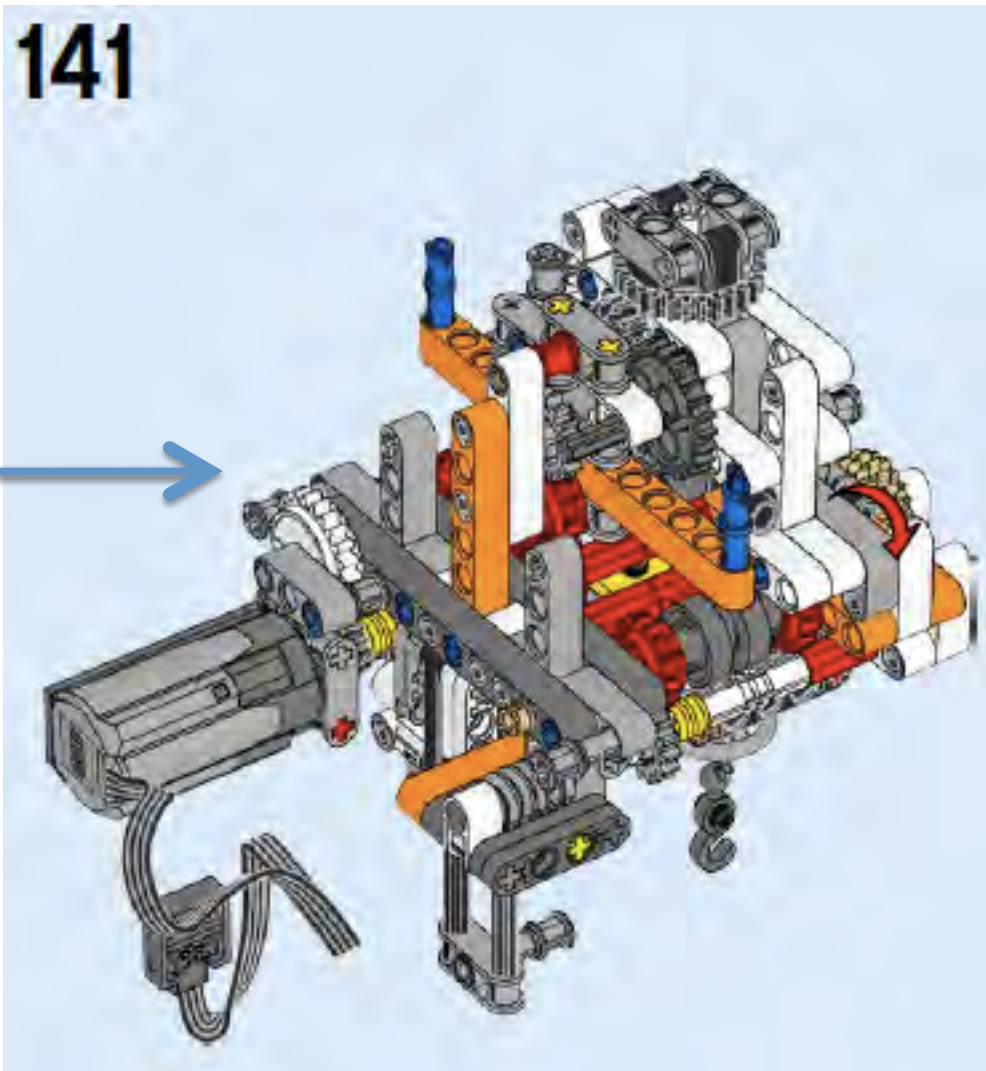
4x    2x    2x
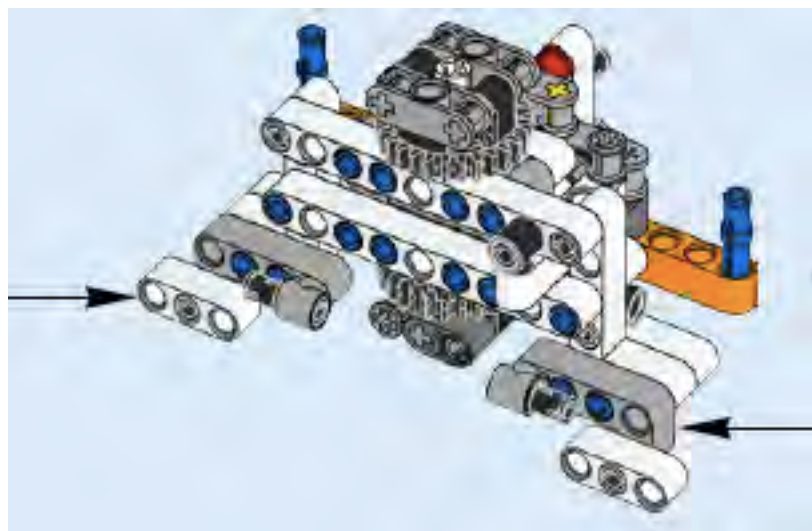
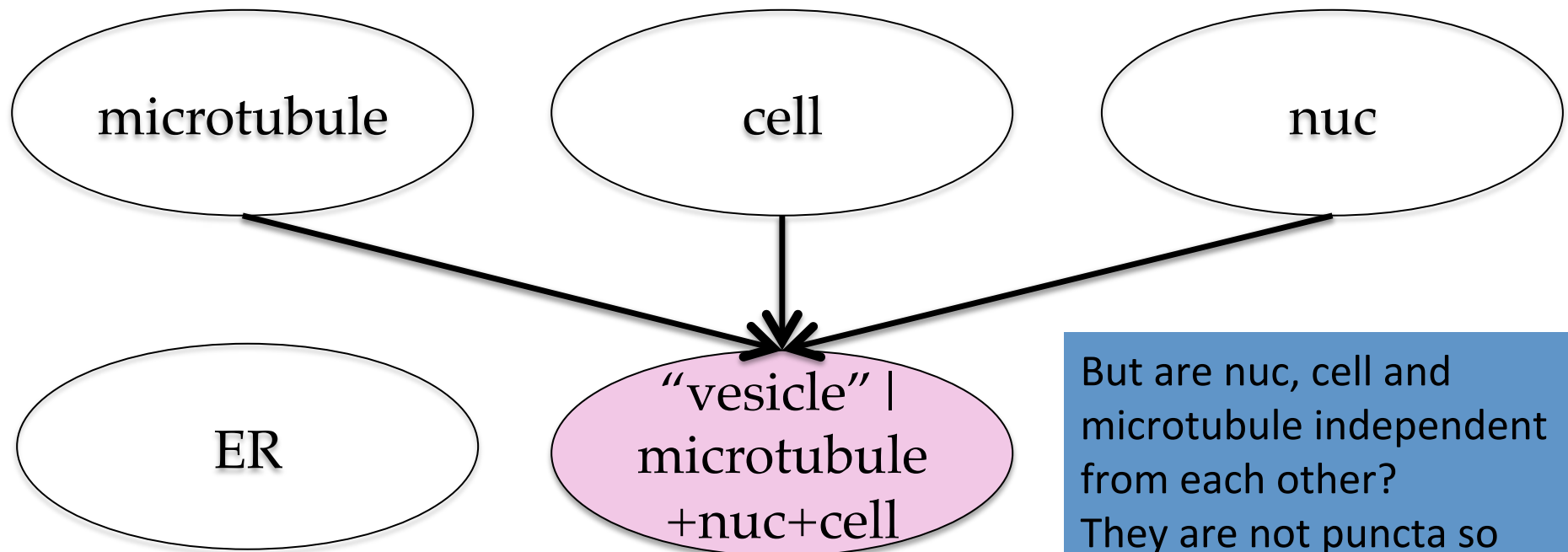1    2                                    2x

1:1    9

**139**

# Bayesian network / graphical model

- One way to think of this hiearchical assembly process is as a graphical model

- Nodes correspond to parts or previous assemblies

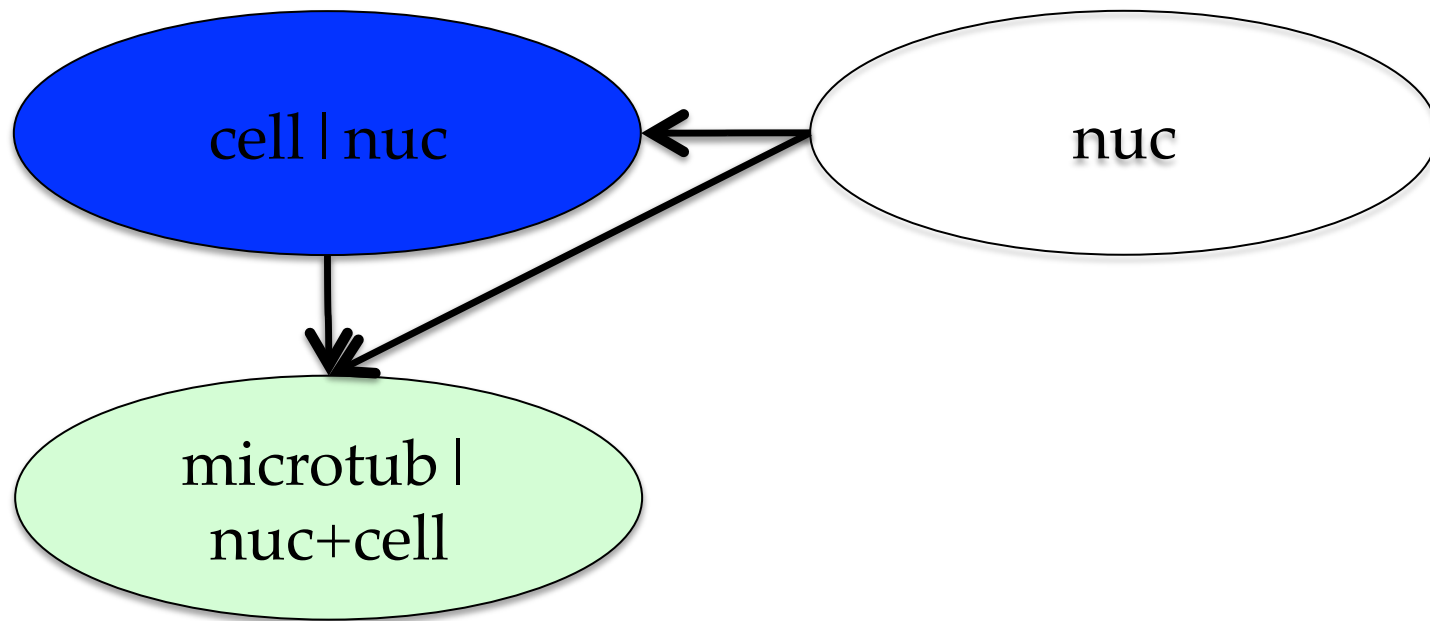- Edges correspond to *dependencies* – parts required to produce/localize an assembly

# Merged Bayes net

cell | nuc

nuc

microtub | nuc+cell

"vesicle" | microtub+nuc +cell

# Merged Bayes net

# From other modeling…

Liu-Huang et al (2017) submitted

cell | nucleus

nucleus

ER

microtub | nuc+cell

Protein $p$ | microtub+ nuc+cell

$p$

Node function=random graph
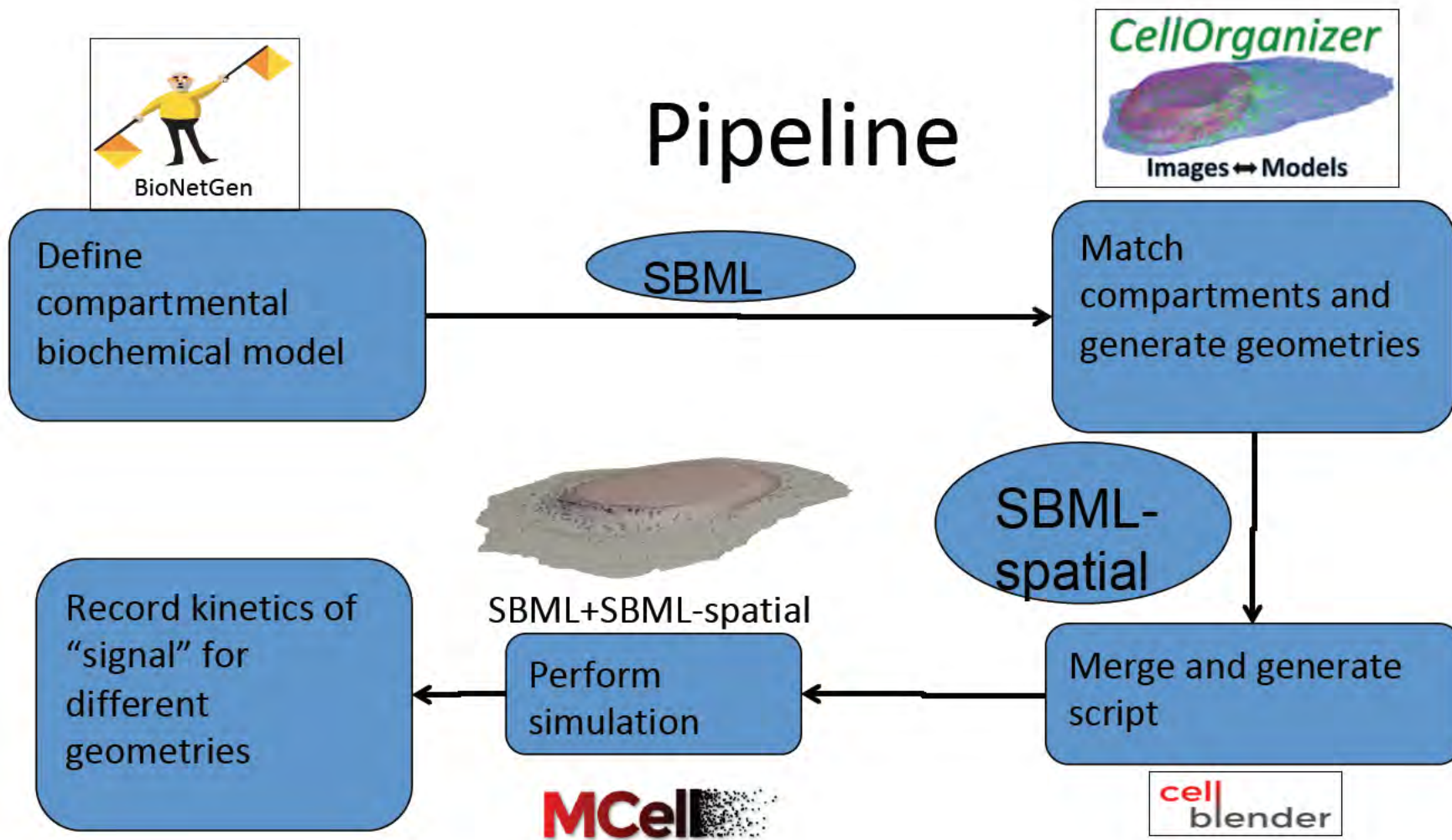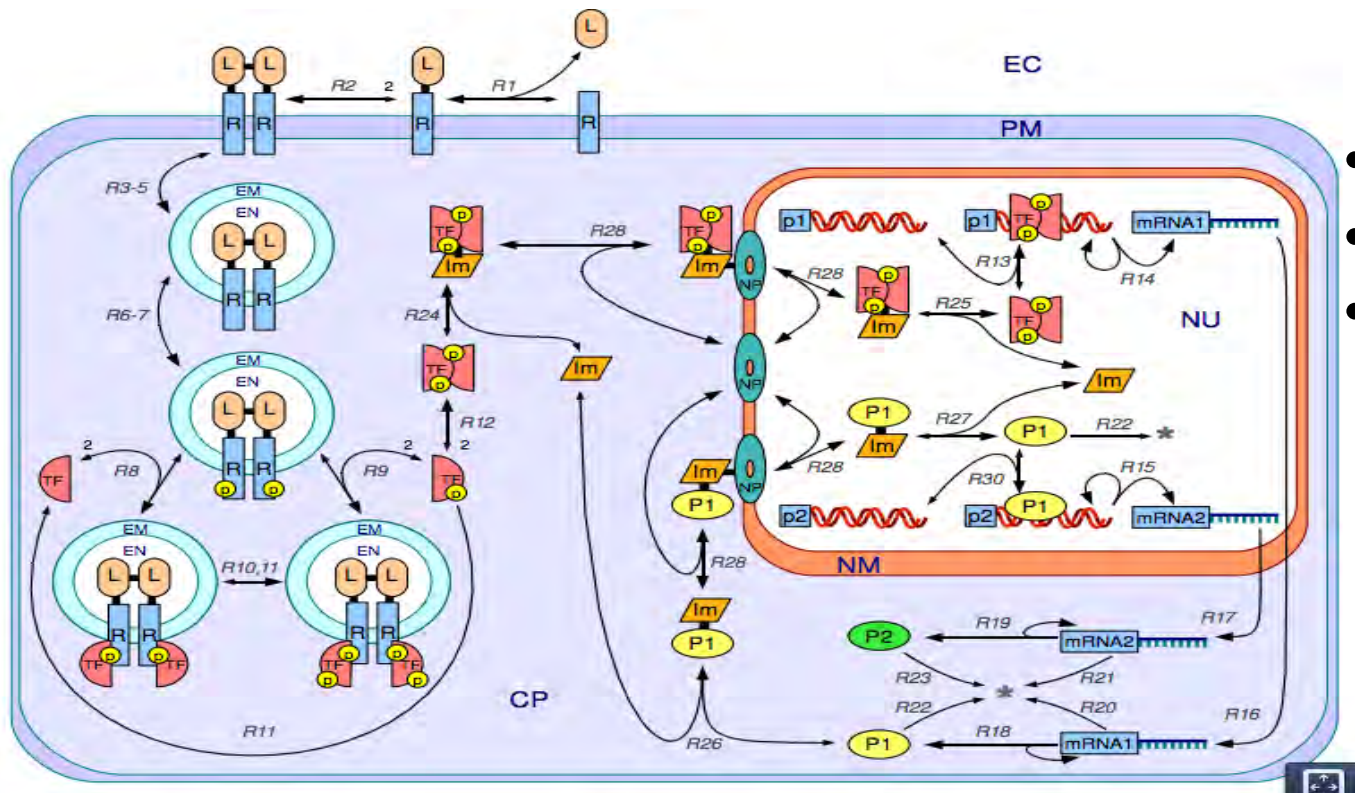
# HIGH-THROUGHPUT CELL SIMULATIONS

# High-throughput spatially realistic simulations

- Study the effects of spatial variance caused by
  - Cell cycle
  - Diseases
  - Drugs
  - Inherent cell variance
- Model large systems with high spatial realism
- Validate generative model accuracies

# Pipeline

BioNetGen

Define compartmental biochemical model

SBML

CellOrganizer
Images ↔ Models

Match compartments and generate geometries

SBML-spatial

Record kinetics of "signal" for different geometries

SBML+SBML-spatial

Perform simulation

MCell

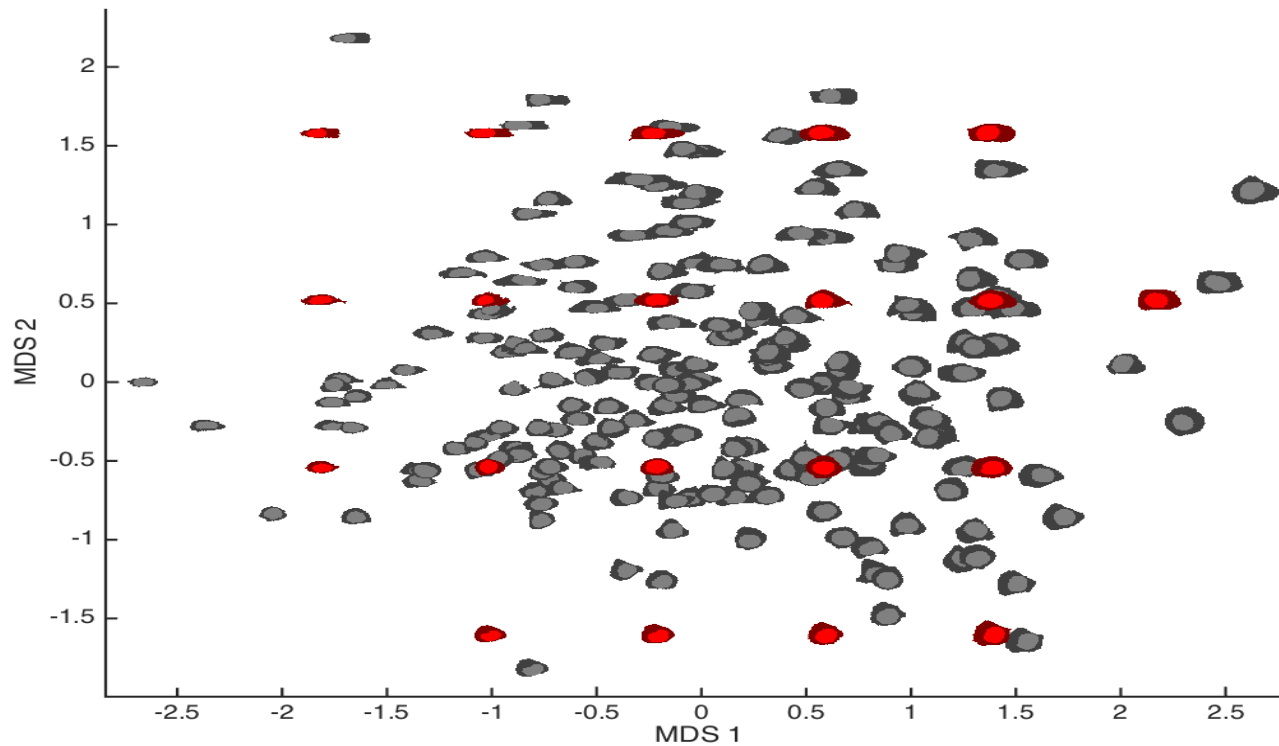Merge and generate script

cell blender

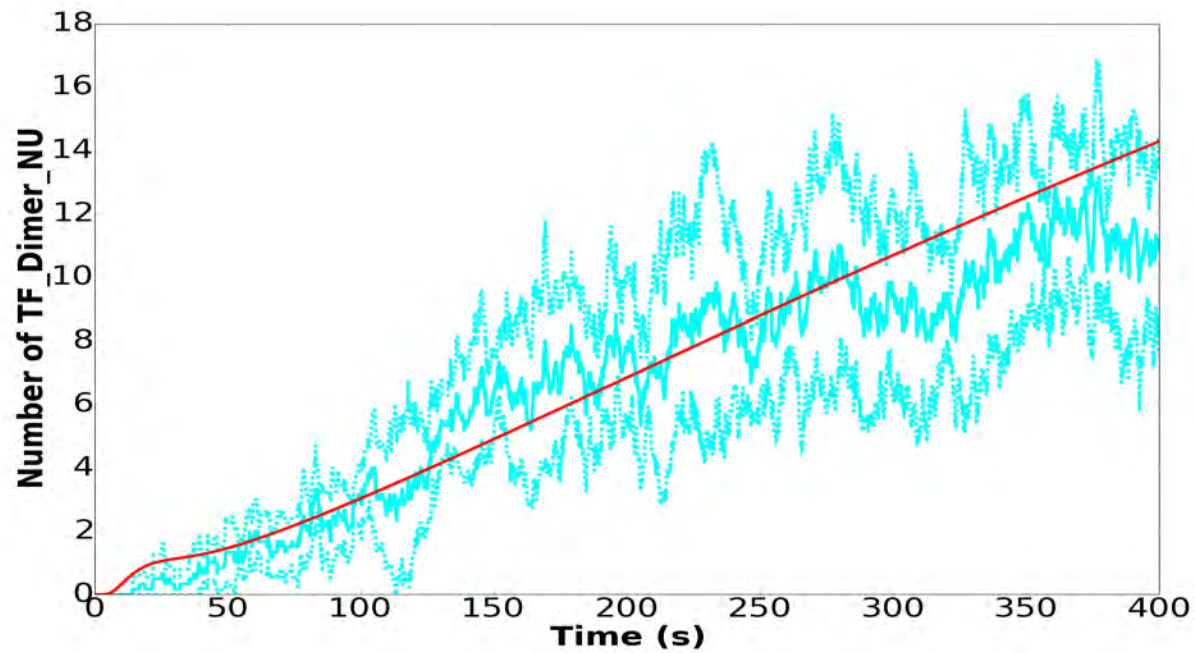# "Simple" biochemical model



- 354 reactions
- 78 species
- 7 "compartments"

# Selected geometries from shape space

# Different geometries lead to variation in signaling



Mean +/- S.D.

# Conclusions

- Tools becoming available to construct models of cell components directly from images
    - Better comparison across instruments/cell types
    - Provide input geometries for cell simulation
    - Provide simulated images for testing algorithms
    - Learn putative spatiotemporal causal relationships
- Need to combine images and data from various other experiments to create overall spatiotemporal models

# CellOrganizer Team

## Project Leaders



**Robert F. Murphy**

**Gustavo Rohde**

**Gregory R. Johnson**

## Current Team Members



Xiongtao Ruan

Kelvin Liu-Huang

Ivan Cao-Berg

## Collaborators

Jörn Dengjel

Christoph Wülfing

## Past Contributors

| | |
|---|---|
| Ting Zhao | Jieyue Li |
| Tao Peng | Baek Hwan Cho |
| Wei Wang | Taraz Buck |
| Aabid Sharif | Devin Sullivan |
| Joshua Kangas | Ying Li |
| Jianwei Zhang | Tim Majarian |

SBML

MMBioS