# Grand Challenges of Measurement Science with Big Data

## Presenter: Peter Bajcsy

**Software and Systems Division**

**Information Technology Laboratory**

**National Institute of Standards and Technology**

# Goal & Objective

- **Goal:** to articulate "Grand Challenges" in data-intensive research

- **Objective:** to identify differences between measurement science for Little Data and Big Data

# Concrete Application Domain

- **Analysis of cell biology microscopy images**

- **NIST Project: Computational Science in Biological Metrology**
  - **CS-BIO-MET: nist.gov/itl/ssd/is/computational-science-in-biometrology.cfm**

# Basic Biological Questions

- What are the quantitative dynamic characteristics of cell changes as a function of their surrounding environment, cell signaling, phenotype and genotype?
  - Parameter estimation
  - Bayesian inference and learning methods
  - Development of mathematical models
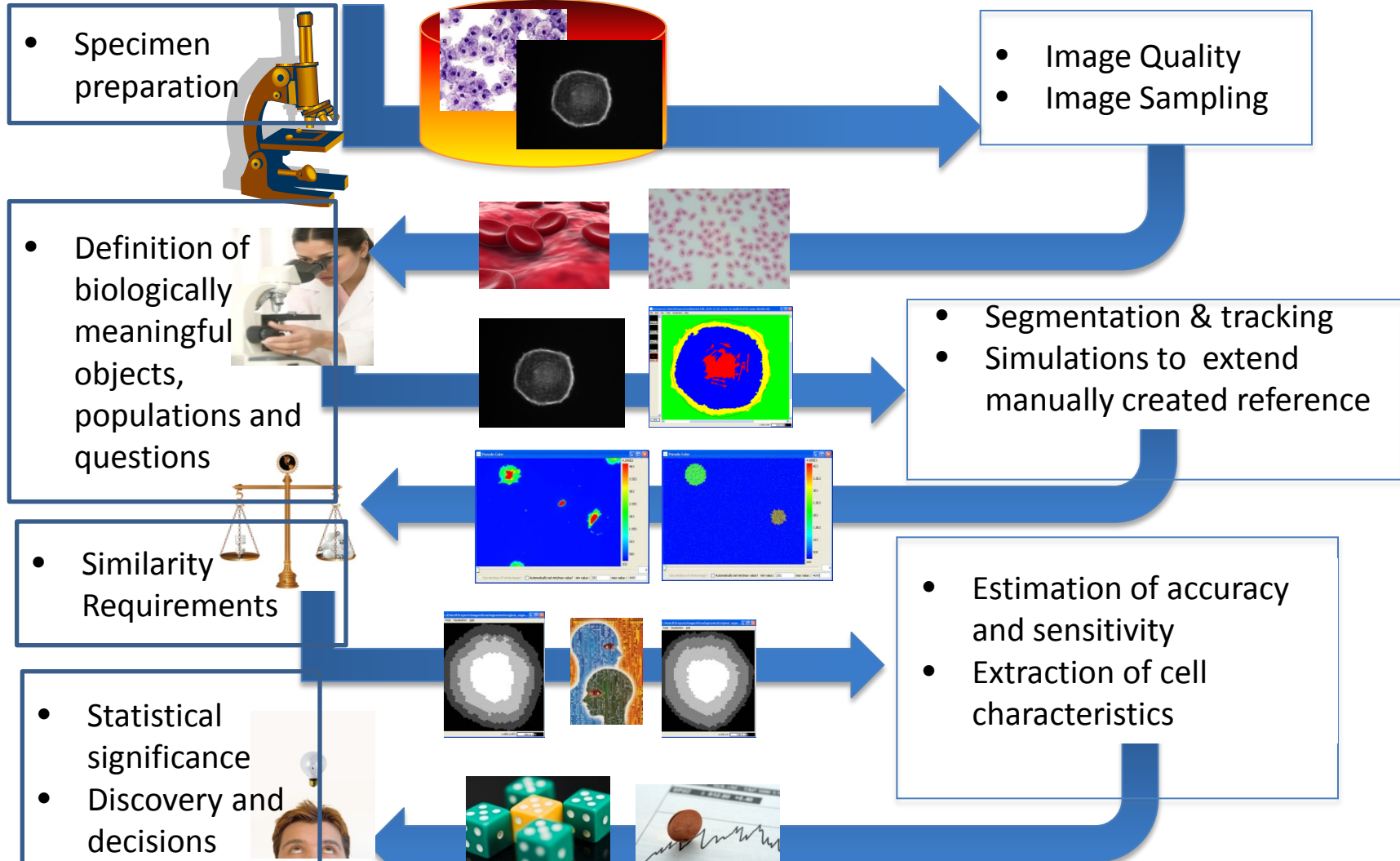
# Mission Oriented Results

- Big Data → **Limited human input** → Need for automation and measurements of automation accuracy

- Automation Outcome
    - → Quantitative cell measurements
    - → Cost reduction
    - → New scientific discovery
    - → Avoid missed opportunities

# Computational Science in Biological Metrology

**Biological Metrology**

**Computational Science**



- Specimen preparation

- Image Quality
- Image Sampling

- Definition of biologically meaningful objects, populations and questions

- Segmentation & tracking
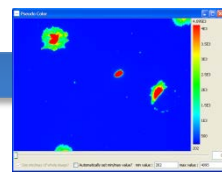- Simulations to extend manually created reference

- Similarity Requirements

- Estimation of accuracy and sensitivity
- Extraction of cell characteristics

- Statistical significance
- Discovery and decisions

# Computational Science in Biological Metrology

## Biological Metrology

## Computational Science

- Specimen preparation

- Definition of biologically mean

- Requirements

- Statistical significance
- Discovery

- Image Quality
- Image Sampling

- Segmentation & tracking
- Simulations to extend manually created reference

- Estimation of accuracy and cell istics

**BIG DATA**

**LIMITED HUMAN INPUT**

**CONFIDENCE IN CS METROLOGY**

# Central Problem:
# Image Segmentation Accuracy



Labeled Data

Big Data

Data Quality

Assumptions about Data Segmentation

**ACCURACY/ERROR**
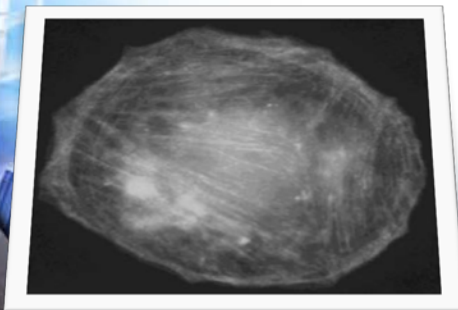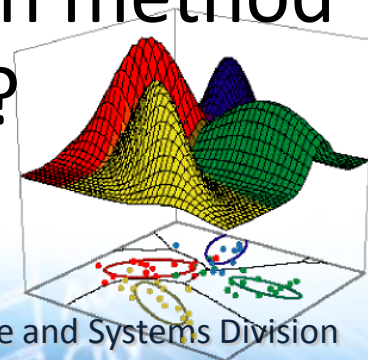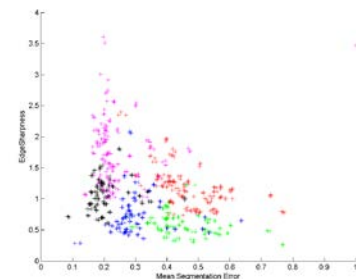
# Image Segmentation Baseline

- What is the baseline for segmentation results in cell microscopy?
  - Disagreement between biologists, statisticians, and computer vision researchers on the baseline criteria
  - Manual versus mathematically grounded approach

# Data Quality

- How does data quality relate to segmentation accuracy?

- What data quality methods are appropriate for driving optimal microscopy settings?

- How does one detect and measure a mismatch between segmentation method assumptions and the input data?

# Data Sampling (Image or Cell)

- What sampling techniques and sample sizes are appropriate for segmentation accuracy evaluations of Big Data?

  – Much work in the signal processing and statistical domains

  – Does the choice of sampling method bias results?

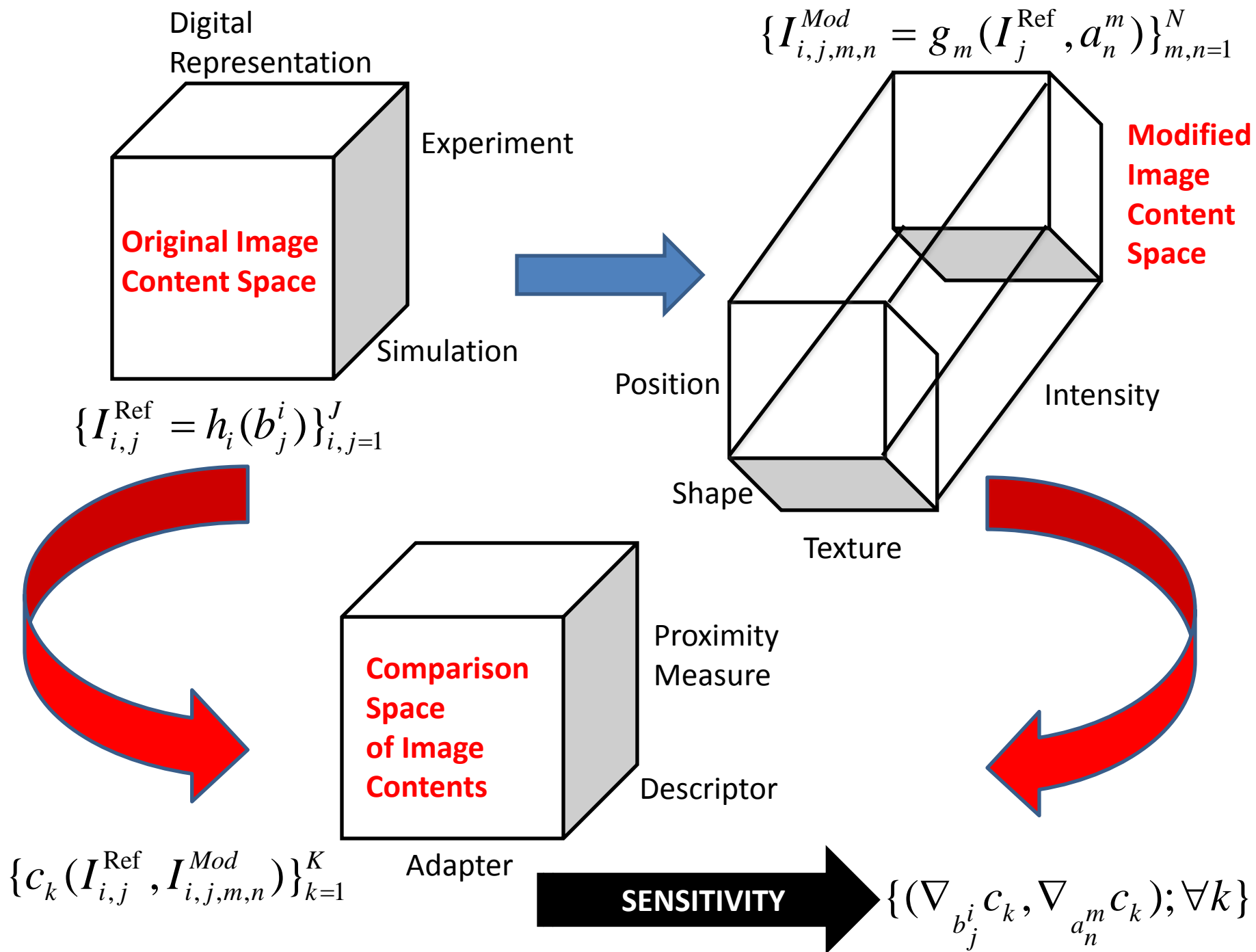  – What uncertainty is acceptable to biologists?

# Data Comparison

- How does one compare two images?
  - Metrics: Euclidean versus Riemannian spaces
  - In general, any two digital data sets?
- How does one choose the most suitable proximity metrics given application requirements?

Digital Representation

$$\{I^{Mod}_{i,j,m,n} = g_m(I^{Ref}_j, a^m_n)\}^N_{m,n=1}$$

**Original Image Content Space**

Experiment

Simulation

**Modified Image Content Space**

$$\{I^{Ref}_{i,j} = h_i(b^i_j)\}^J_{i,j=1}$$

Position

Intensity

Shape

Texture

**Comparison Space of Image Contents**

Proximity Measure

Descriptor

Adapter

$$\{c_k(I^{Ref}_{i,j}, I^{Mod}_{i,j,m,n})\}^K_{k=1}$$

**SENSITIVITY**

$$\{(\nabla_{b^i_j} c_k, \nabla_{a^m_n} c_k); \forall k\}$$

# Sensitivity Signatures of Similarity Metrics

$$\vec{S}_k = \left( \frac{\delta c_k}{\delta I_{i,j}^{\text{Ref}}}, \frac{\delta c_k}{\delta I_{i,j,m,n}^{\text{Mod}}} \right)_{(I_{i,0}^{\text{Ref}}; I_{i,0,m,0}^{Mod})}$$
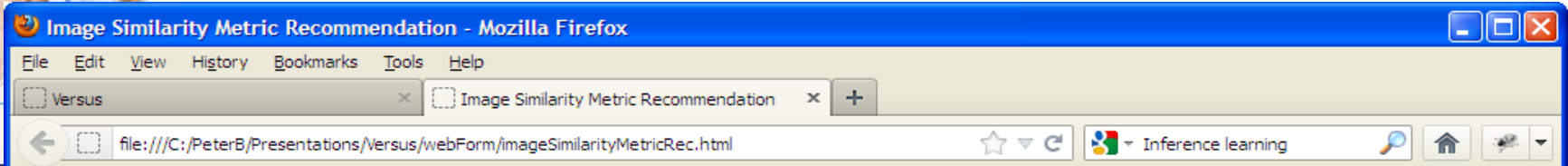
$$\frac{\delta c_k}{\delta I_{i,j}^{\text{Ref}}} = \delta \Phi[c_k(I_{i,j}^{\text{Ref}}, \delta I_{i,j,m,n}^{\text{Mod}}); I_{i,j}^{\text{Ref}}]$$

$$\frac{\delta c_k}{\delta I_{i,j,m,n}^{\text{Mod}}} = \delta \Phi[c_k(I_{i,j}^{\text{Ref}}, \delta I_{i,j,m,n}^{\text{Mod}}); I_{i,j,m,n}^{\text{Mod}}]$$

The differentials are the 1st Frechet derivative of a similarity metric $C_k$ along the direction of reference image content changes (application dependent space of reference image content) and along the direction of any possible modification of reference image content (application task specific) at the referenced point $(I_{i,0}^{\text{Ref}}; I_{i,0,m,0}^{Mod})$

# Similarity Metric Recommendation

# Access, Access, Access

- What is the most efficient Big Data access and retrieval protocol?
  - RESTful web services, Web applications, …

# Visualization

- How to visualize comparisons over Big Data?
  - Human Computer Interfaces for Big Data?

# Data Provenance Gathering

- How do we automate gathering of computational provenance information?
  - Repeatability

- How do we represent provenance information?

- At what granularity should we gather computational provenance?

# Summary

- **Big Data:** Cell microscopy images

- **Basic Challenge:** automated segmentation and its accuracy evaluations over Big Data

- In general, the basic challenge applies to other domains for other automated processing operations

# Summary

- Accuracy evaluations of automated processing over Big Data include at least

  – Data quality, sampling, and comparison measurements

  – Standards for accessing Big Data remotely

  – Standards for computational provenance information representation

# Acknowledgement

**CS**
- Peter Bajcsy
- Adele Peskin
- Mary Brady
- Joe Chalfoun
- Antonio Cardone
- Frederic de Vaulx
- Afzal Godil
- Ben Long
- Ya-Shian Li-Baboud
- Julien Amelot
- Antoine Vandecreme
- Paul Khouri Saba
- Bertrand Stivalet
- Walid Keyrouz
- Doug Foxvog
- Alden Dima

**STATS**
- James Filliben
- Dan Samarov

**BIO**
- Anne Plant
- Kiran Bhadriraju
- Michael Halter
- John Elliott
- TN Bhat

**External Collaborators**
- Marcin Kociolek(1)
- Carol Parent (2)
- Christina Stuelten (2)
- Michael Weiger (2)
- Wolfgang Loster (3)
- Daniel Hoeppner (5)
- Rachel Errington (4)
- Imtiaz Khan (4)
- Justin Martineau (6)
- David Chapman (6)
- Mike Grasso (7)

(1) Technical University of Lodz, Poland ( Marcin Kociolek – image processing
(2) NIH NCI (Carol Parent, Christina Stuelten, Michael Weiger, – cancer cells),
(3) University of Maryland at College Park (Wolfgang Loster – cancer cells),
(4) Broad Institute (Rachel Errington, Imtiaz Khan - bone cells)
(5) The Lieber Institute for Brain Development (Dan Hoeppner - live cell imaging)
(6) University of Maryland, Baltimore County (UMBC) (image processing)
(7) University of Maryland, Baltimore (image processing)

# Questions

- [peter.bajcsy@nist.gov](mailto:peter.bajcsy@nist.gov)