# How to Collect Segmentations for Biomedical Images?
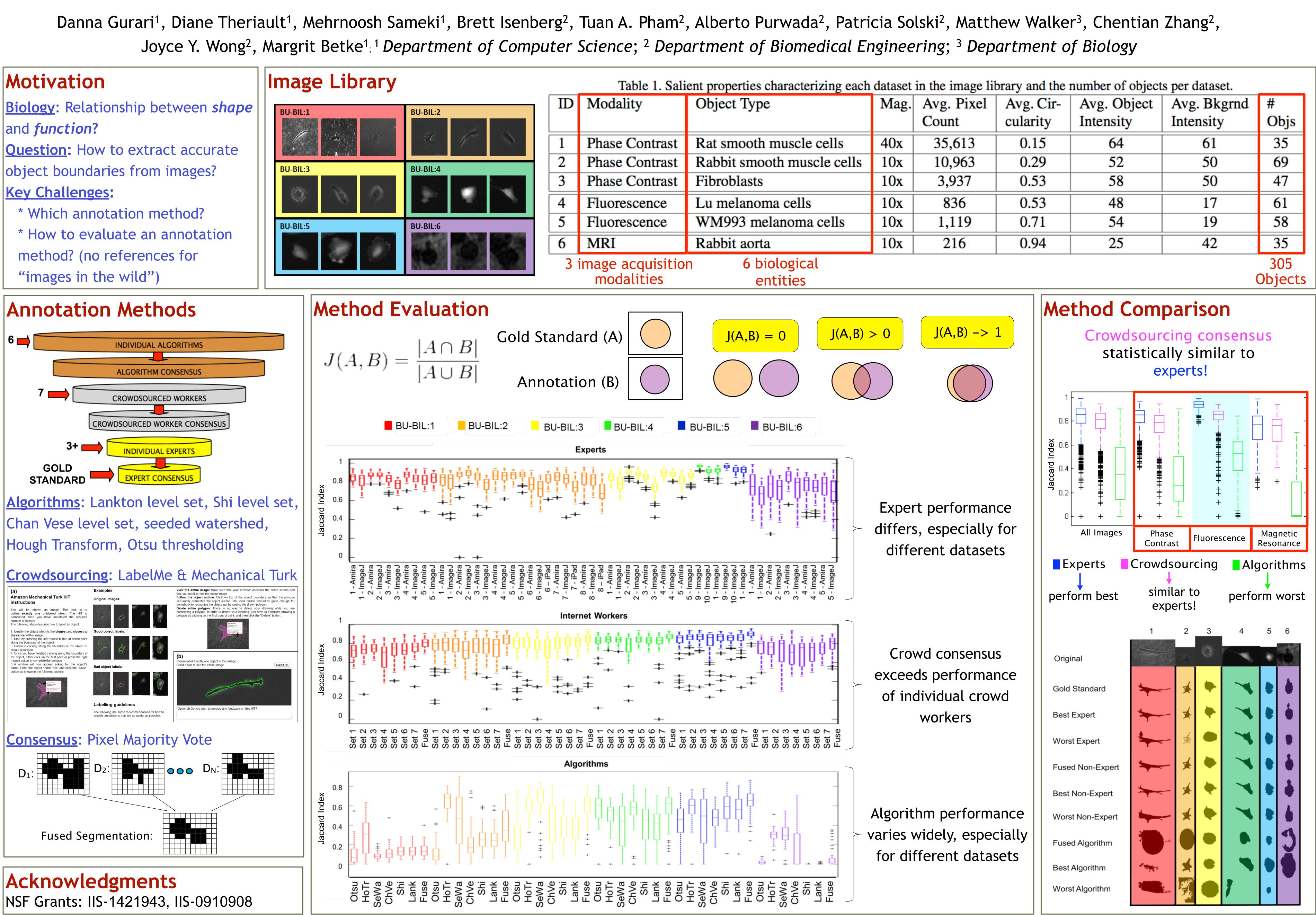# A Benchmark Evaluating the Performance of Experts, Crowdsourced Non-Experts, and Algorithms

Danna Gurari[1], Diane Theriault[1], Mehrnoosh Sameki[1], Brett Isenberg[2], Tuan A. Pham[2], Alberto Purwada[2], Patricia Solski[2], Matthew Walker[3], Chentian Zhang[2],

Joyce Y. Wong[2], Margrit Betke[1]; [1] *Department of Computer Science*; [2] *Department of Biomedical Engineering*; [3] *Department of Biology*

## Motivation

**Biology**: Relationship between *shape* and *function*?

**Question:** How to extract accurate object boundaries from images?

**Key Challenges:**

* Which annotation method?
* How to evaluate an annotation method? (no references for "images in the wild")

## Image Library



BU-BIL:1  BU-BIL:2  BU-BIL:3  BU-BIL:4  BU-BIL:5  BU-BIL:6

Table 1. Salient properties characterizing each dataset in the image library and the number of objects per dataset.

| ID | Modality | Object Type | Mag. | Avg. Pixel Count | Avg. Circularity | Avg. Object Intensity | Avg. Bkgrnd Intensity | # Objs |
|---|---|---|---|---|---|---|---|---|
| 1 | Phase Contrast | Rat smooth muscle cells | 40x | 35,613 | 0.15 | 64 | 61 | 35 |
| 2 | Phase Contrast | Rabbit smooth muscle cells | 10x | 10,963 | 0.29 | 52 | 50 | 69 |
| 3 | Phase Contrast | Fibroblasts | 10x | 3,937 | 0.53 | 58 | 50 | 47 |
| 4 | Fluorescence | Lu melanoma cells | 10x | 836 | 0.53 | 48 | 17 | 61 |
| 5 | Fluorescence | WM993 melanoma cells | 10x | 1,119 | 0.71 | 54 | 19 | 58 |
| 6 | MRI | Rabbit aorta | 10x | 216 | 0.94 | 25 | 42 | 35 |

3 image acquisition modalities    6 biological entities    305 Objects

## Annotation Methods



6 → INDIVIDUAL ALGORITHMS
ALGORITHM CONSENSUS
7 → CROWDSOURCED WORKERS
CROWDSOURCED WORKER CONSENSUS
3+ → INDIVIDUAL EXPERTS
GOLD STANDARD → EXPERT CONSENSUS

**Algorithms:** Lankton level set, Shi level set, Chan Vese level set, seeded watershed, Hough Transform, Otsu thresholding

**Crowdsourcing:** LabelMe & Mechanical Turk



**Consensus:** Pixel Majority Vote



Fused Segmentation:

## Acknowledgments

## Method Evaluation

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Gold Standard (A)    Annotation (B)

J(A,B) = 0    J(A,B) > 0    J(A,B) –> 1

BU-BIL:1  BU-BIL:2  BU-BIL:3  BU-BIL:4  BU-BIL:5  BU-BIL:6



**Experts**

Expert performance differs, especially for different datasets

**Internet Workers**

Crowd consensus exceeds performance of individual crowd workers

**Algorithms**

Algorithm performance varies widely, especially for different datasets

## Method Comparison

**Crowdsourcing consensus** statistically similar to **experts!**



All Images | Phase Contrast | Fluorescence | Magnetic Resonance

■ Experts → perform best
■ Crowdsourcing → similar to experts!
■ Algorithms → perform worst



1  2  3  4  5  6

Original
Gold Standard
Best Expert
Worst Expert
Fused Non-Expert
Best Non-Expert
Worst Non-Expert
Fused Algorithm
Best Algorithm
Worst Algorithm