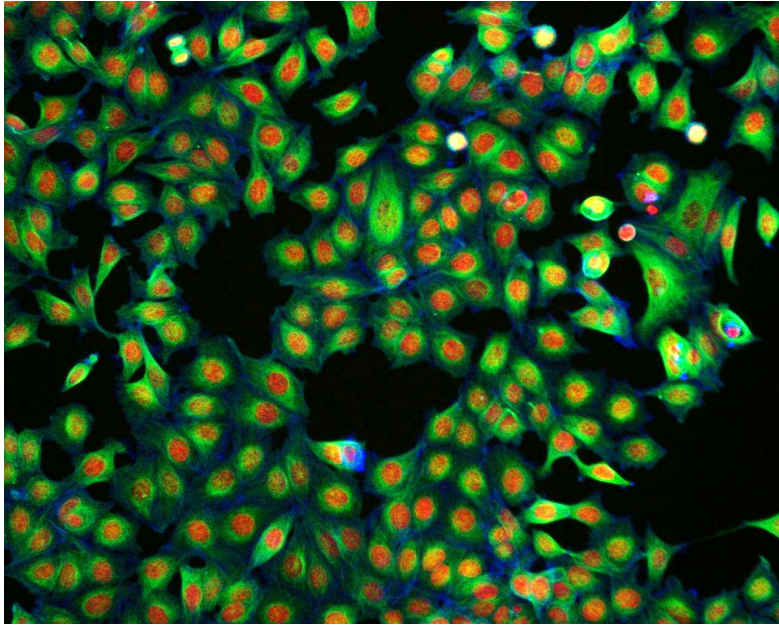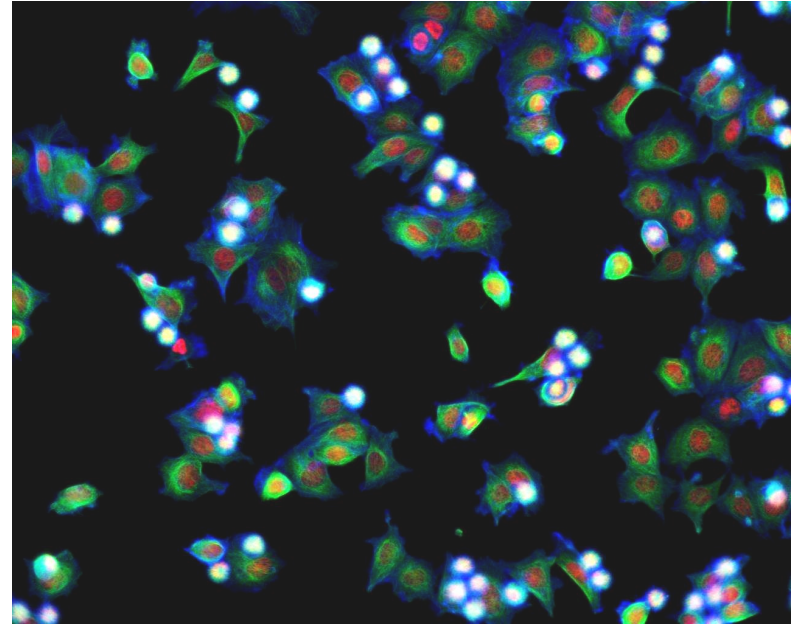# Large Scale Dataset Curation and Model Evaluation

John Peters
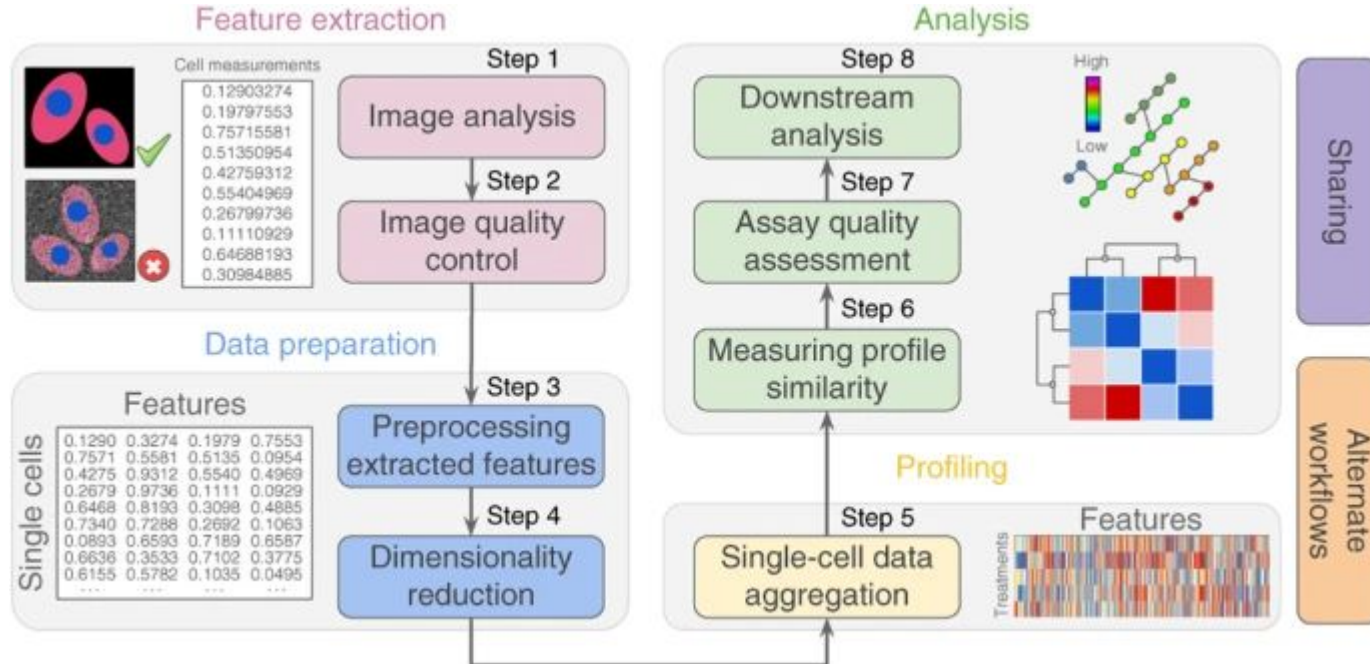Caicedo Lab
6/5/25

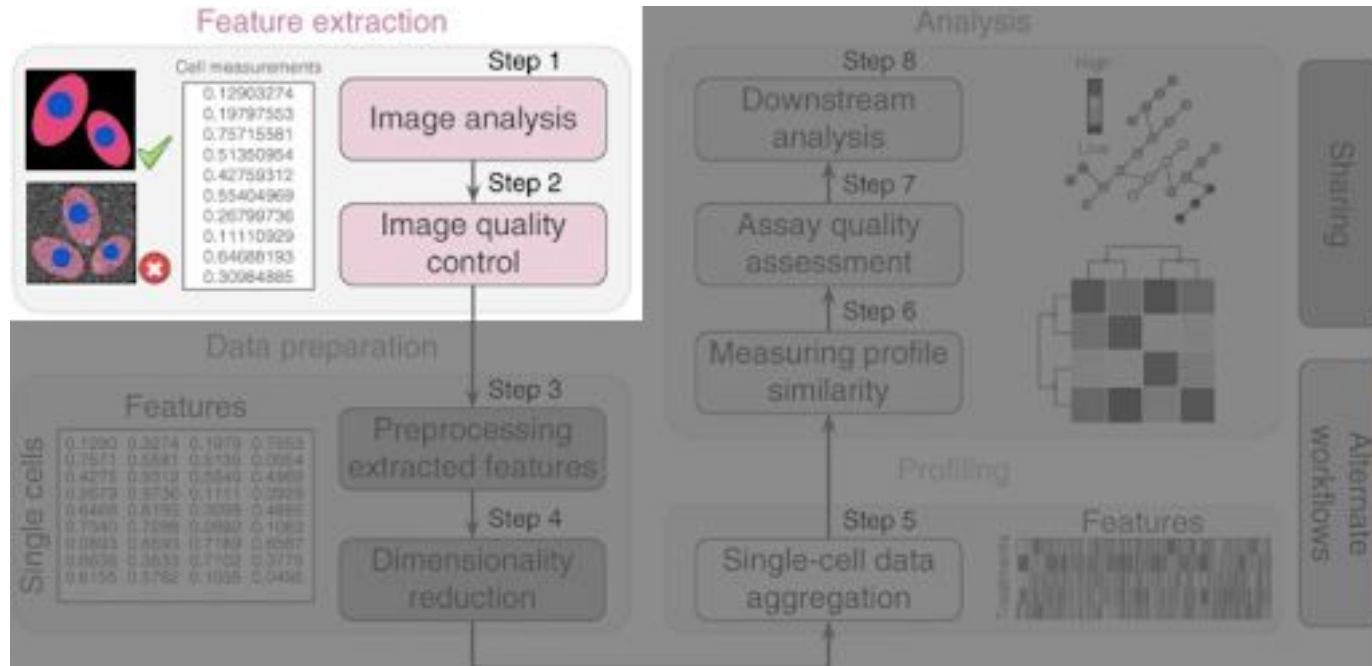# Image-based experiments



**Control condition**

**Treated condition**
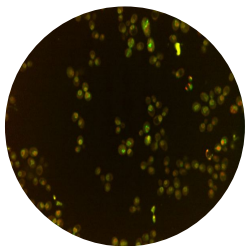
# Image-based profiling workflow



Caicedo et al. "Data-analysis strategies for image-based cell profiling." *Nature methods* 14.9 (2017): 849-863.

# Talk Focus: Feature Extraction



Caicedo et al. "Data-analysis strategies for image-based cell profiling." *Nature methods* 14.9 (2017): 849-863.
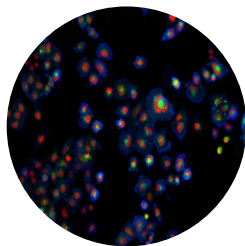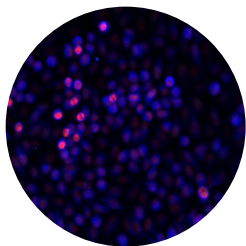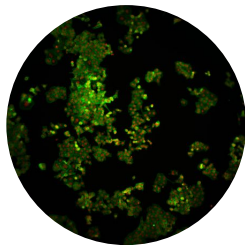
# 5 Studies to Process

idr0007

idr0009

idr0020

idr0017

idr0088

## Directory Organization

**Images**

```
Images
|
|--study
|    |- study-plate_1.zip
|    |
|    |- study-plate_2.zip
|    |        ...
|    |- study-plate_n.zip
|--study-2
     |- study2-plate_1.zip
     |
     |- study2-plate_2.zip
     |        ...
     |- study2-plate_n.zip
```

Data Vault

# Zip Format

Full Image



Cell Body



Protein 1

Protein 2

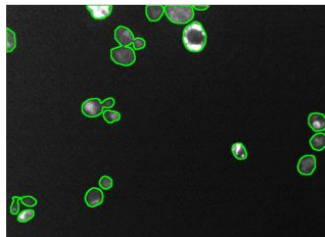| Multi-Channel-Id | Image-Path-In-Zip |
|---|---|
| hpa_image_0001 | .png |
| hpa_image_0001 | .png |
| hpa_image_0001 | .png |

Images saved as individual channel png files within *plate_x.zip files
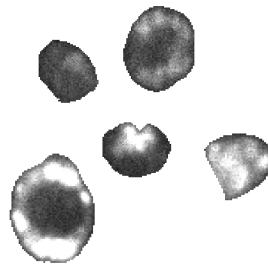
# How do we extract features?



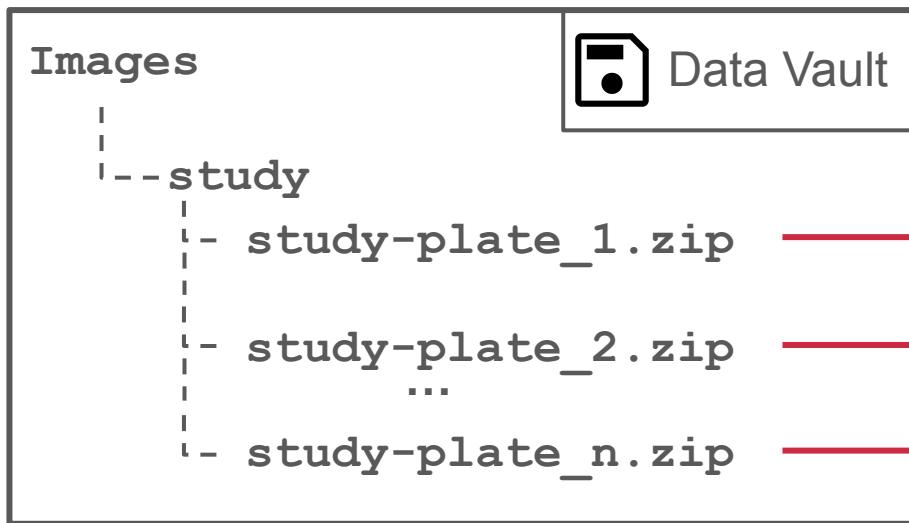**Load Images** → **Segment Cells** → **Isolate Cells** → **Extract Features**

Deep Learning Model

**Single Cell Features**

**4 Steps to Features!**
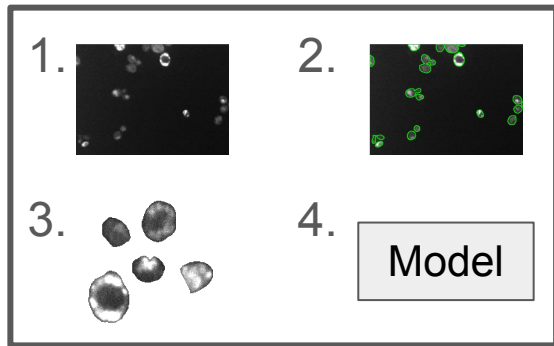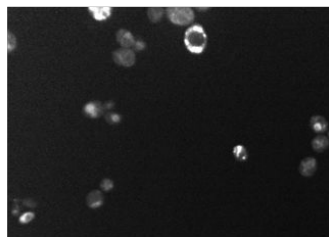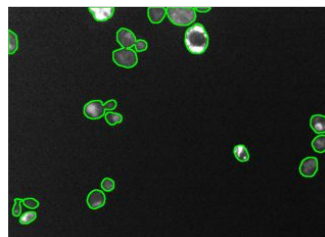
# How can we leverage CHTC?

# Reusable Segmentations – Multiple Models
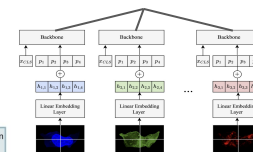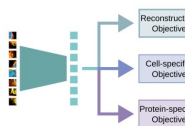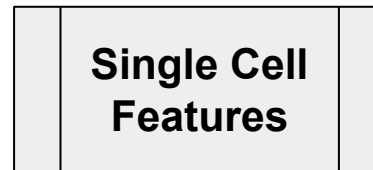


**Load Images**   **Segment Cells**   **Isolate Cells**   **Extract Features**

**Can cache!**

**Single Cell Features**

# Snakemake - Segmentation

Segmentation

CHTC

Parallel!

1.

2.

Masks

# Segmentation in Detail

Reminder Images are Multichannel:

1. Download

2. Unzip

Choose segmentation channel

3. Segment

Segmentation Model

Data Vault

# How large was our data?

Larger studies (GB) with less plates are slower on CHTC

Circle Size = Image Count



Dataset Statistics

dataset
- IDR0007
- IDR0009
- IDR0017
- IDR0020
- IDR0088

# Takeaways

🚀     Small zips work very well on CHTC

🐍     Snakemake allows for local testing before running on CHTC

📦     Caching Saves tons of time with large reruns

# Thanks!



Special thanks to Justin and Ian🙂



Grant #5T15LM007359