# Generative approaches to Shapley-based explanations
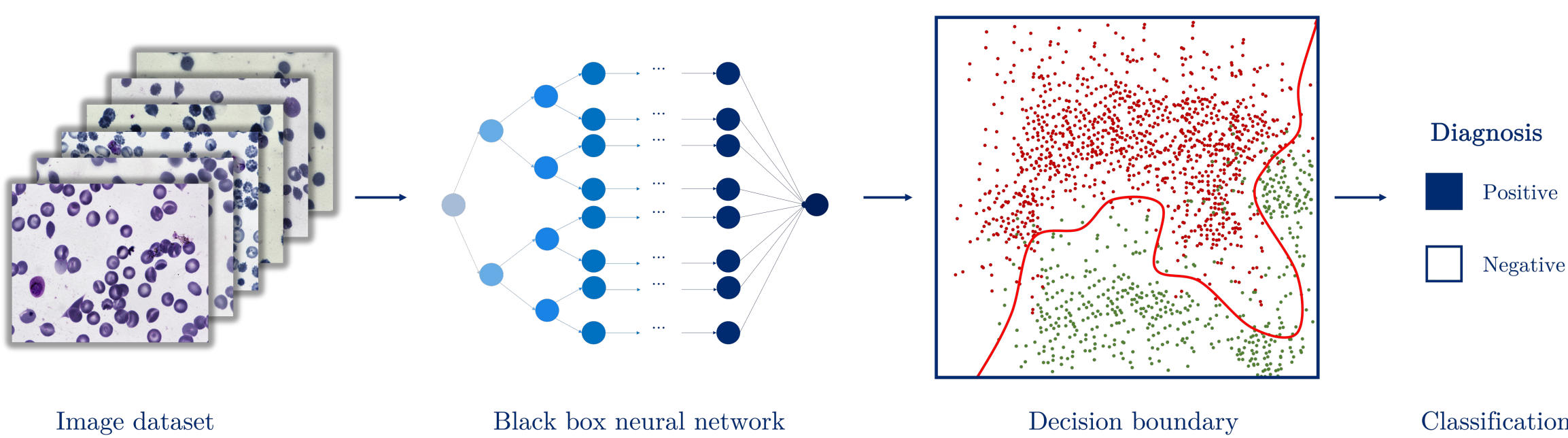
Ishan Kalburge, Siyu Wang, Kuai Yu

Department of Biomedical Engineering, Johns Hopkins University

## Background

### The black box problem

In clinical computer vision settings, deep neural networks (DNNs) can diagnose diseases by creating complex mathematical relationships between image data, sometimes estimating millions of parameters. But how can clinicians trust and verify the conclusions of these DNNs, especially when these mathematical relationships are so complicated?



| | |
|---|---|
| Image dataset | Black box neural network | Decision boundary | Classification |

Diagnosis: Positive / Negative

### *A posteriori* methods: Shapley values

Because DNNs are valuable for their raw predictive power, we can implement methods that can interpret the results of a DNN for us!

**Shapley values**, a tool inherited from cooperative game theory, offer a simple but elegant solution. A Shapley game computes the value added for each member of a team – or, likewise, each feature of an image.



Output **with** Player 3 · Output **without** Player 3 · Player 3 contributions

$10000 - 7000 = 3000$
$6000 - 5000 = 1000$
$4000 - 2000 = 2000$
$2000 - 0 = 2000$
$= 8000$ units

If we treat each feature of an image as players in a Shapley game as shown above, we can summarize the predictive values for a feature using a weighted sum to determine its overall importance:

$$\phi_j(v) = \phi_j = \sum_{S \subseteq M \setminus \{j\}} \frac{|S|!(M-|S|-1)!}{M!}(v(S \cup \{j\}) - v(S)), \quad j = 1, ..., M,$$

## Problems

### But what constitutes a player in a Shapley game?

*Computational feasibility*
- How many players can we have before the problem becomes intractable?

*Semantic Relevance*
- How do we pinpoint features that have conceptual importance? Which **"things,"** rather than pixels, does a DNN latch onto?

### And how should we treat the players we remove?

*Statistical Accuracy*
- If we remove part of the image and replace it with a black space, we are creating an image that does not truly exist in our distribution. How do we fix this?

## Our Approach
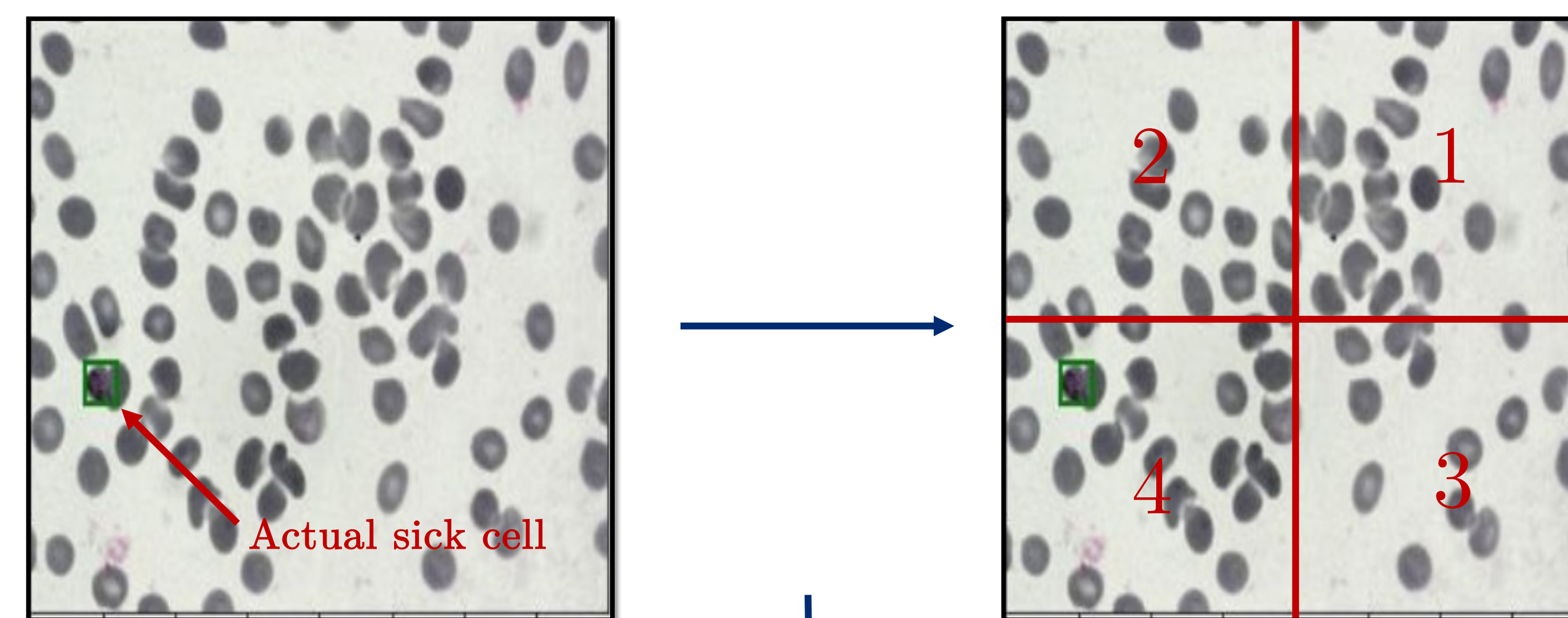
**Reduce** the number of players
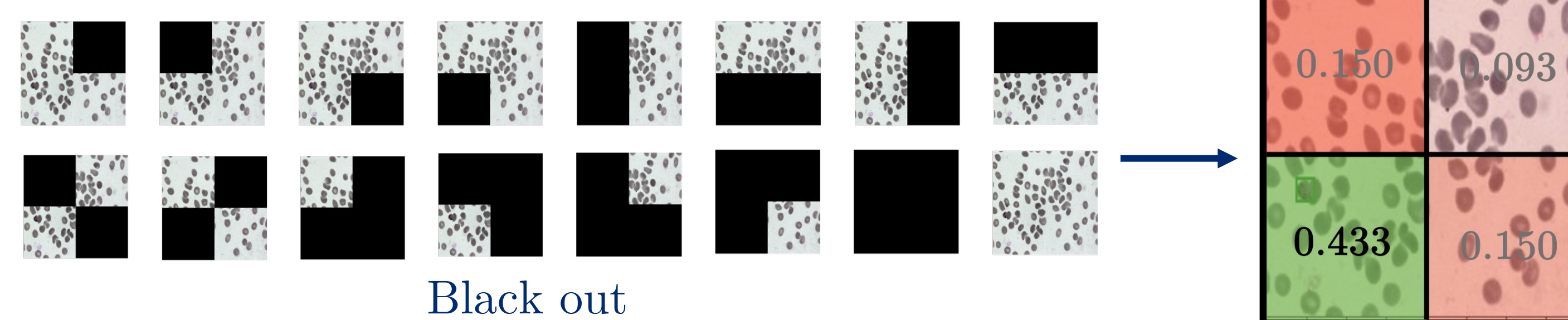
Partition · Segmentation · Unsupervised clustering

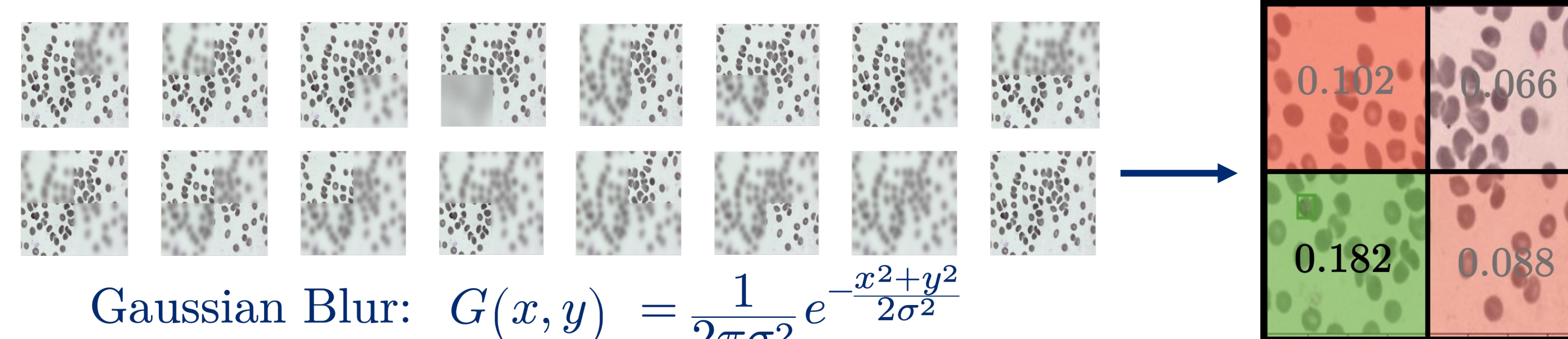**Inpaint** during Shapley games

## Partition Pipeline Comparison



Actual sick cell

Sample Image → Divide into Quadrants

2  1
4  3



**Black out**

| 0.160 | 0.093 |
|---|---|
| **0.433** | 0.150 |



Gaussian Blur: $G(x,y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$

| 0.102 | 0.066 |
|---|---|
| **0.182** | 0.088 |



Generative Adversarial Net

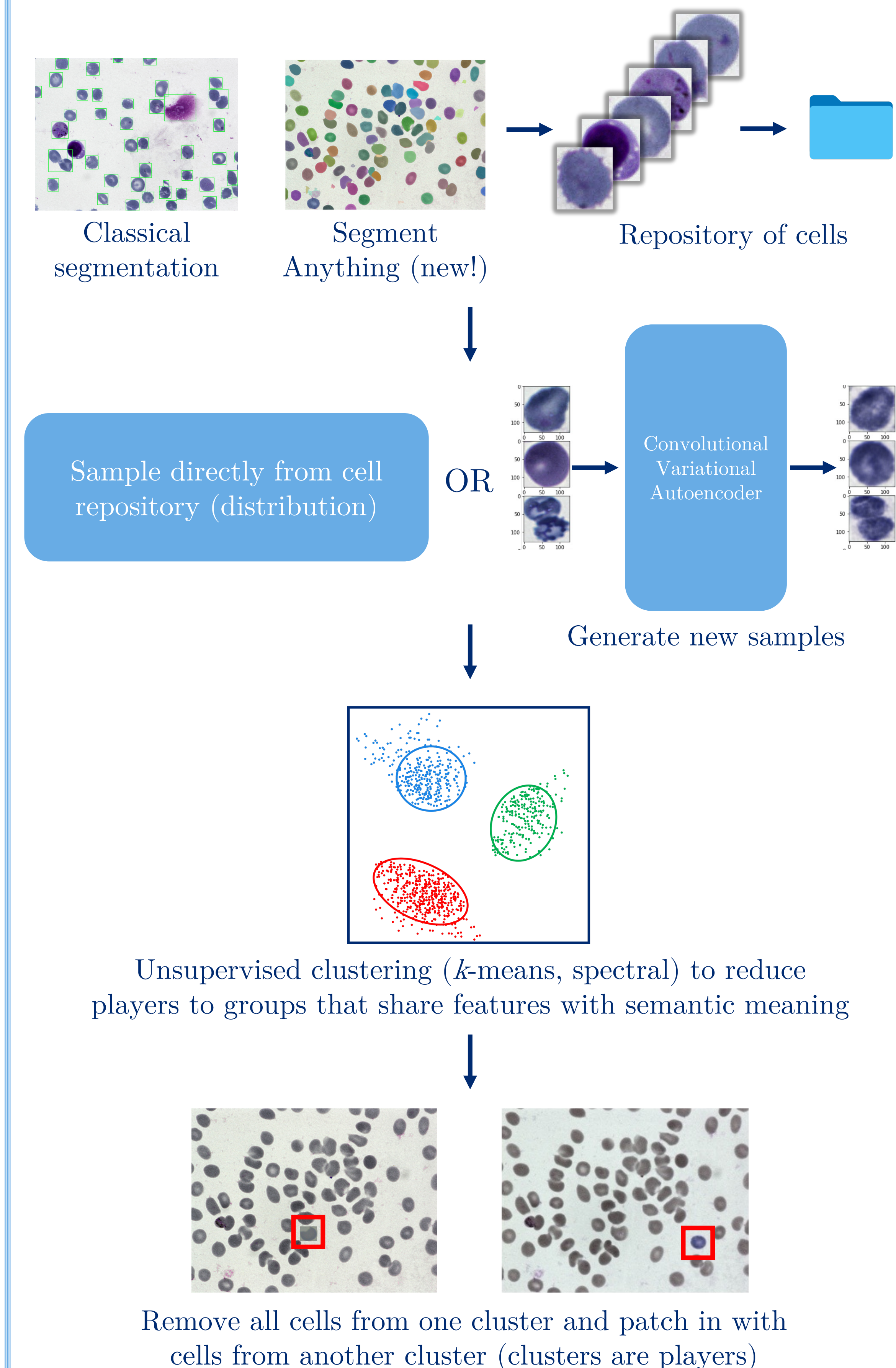| 0.161 | 0.153 |
|---|---|
| **0.270** | 0.048 |

## Segmentation + Clustering Pipeline

### Alternatives to Generative AI

Even when Generative Adversarial Nets (GANs) are trained on our image set, they poorly inpaint removed portions of images – which is not surprising. We offer a pipeline for manual inpainting:



Classical segmentation · Segment Anything (new!) · Repository of cells

Sample directly from cell repository (distribution) · **OR** · Convolutional Variational Autoencoder

Generate new samples



Unsupervised clustering (*k*-means, spectral) to reduce players to groups that share features with semantic meaning



Remove all cells from one cluster and patch in with cells from another cluster (clusters are players)

## Next Steps

*Segmentation and patching process*
- Use Segment Anything? How can we patch for inconsistent dimensions?

*Clustering process*
- Should we use handcrafted features? Should we cluster within or over entire dataset?

*Codebase and documentation*