# Searching for similarities and anomalies in a pool of galaxy images

Matias Carrasco Kind

Senior Research Scientist, National Center for Supercomputing Applications
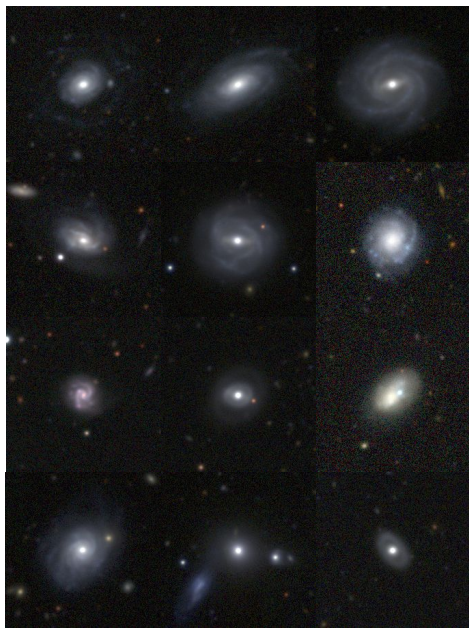
Assistant Research Professor, Astronomy

Data Release Scientist, Dark Energy Survey

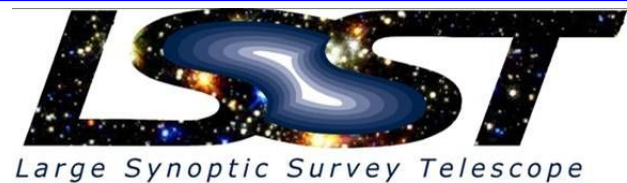University of Illinois at Urbana-Champaign

# Motivation

Astronomy is just one example where image exploration needs to be automated.

Large catalogs, Large number of images, many unexpected objects/problems → <u>Anomaly detection</u>

**LSST**

*Large Synoptic Survey Telescope*

- In operations 2021
- Every night for 10 years
  - 15 TB per night
- 18 billions objects (first year), ~40 billions by the end of survey
- ~1500 images per night
- Stream and static data
- Target to capture new physics (moving and variable objects)

**DARK ENERGY SURVEY**

- More than 500 nights of observation over 5 years, 2TB per night
- 500 millions cataloged galaxies and 100 millions stars
- Many open problems: Systematics, new objects, new physics, etc.
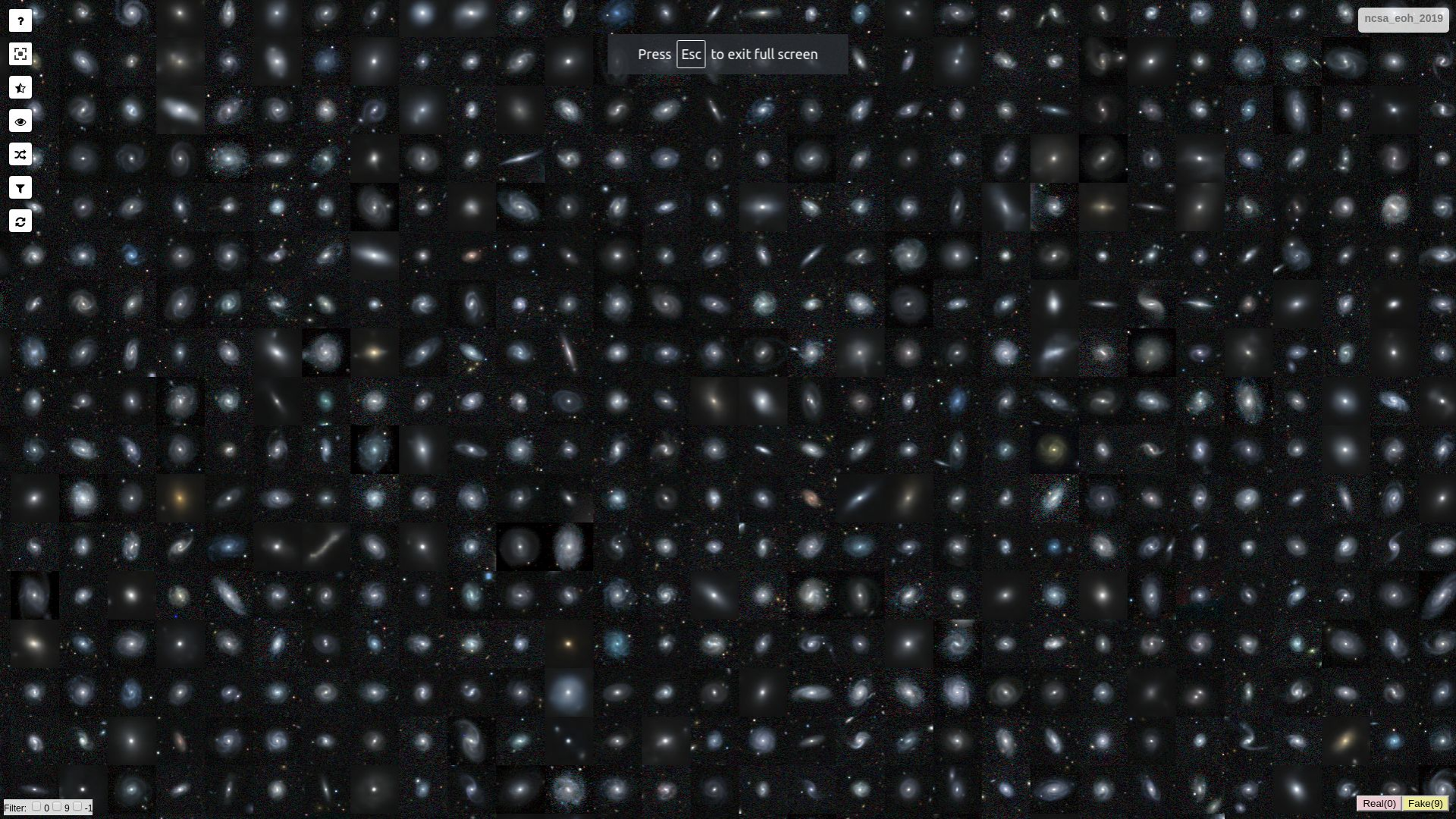- Almost completed

# Current (personal) data discovery challenges

- Visualize large set of galaxy (or other) images
- Quickly classify images for AI using multiple experts
- Compress the important information in a efficient way
- Quickly search images by (dis)similarity (several science cases)
- Find anomalies in a image data set (new phenomena, errors, unrepresentative samples)

Not covered here

- Generate and sample realistic fake images based on a training for modeling and Monte-Carlo Sampling
- Generate and sample realistic fake images based on a training in a controlled manner (with a prior)

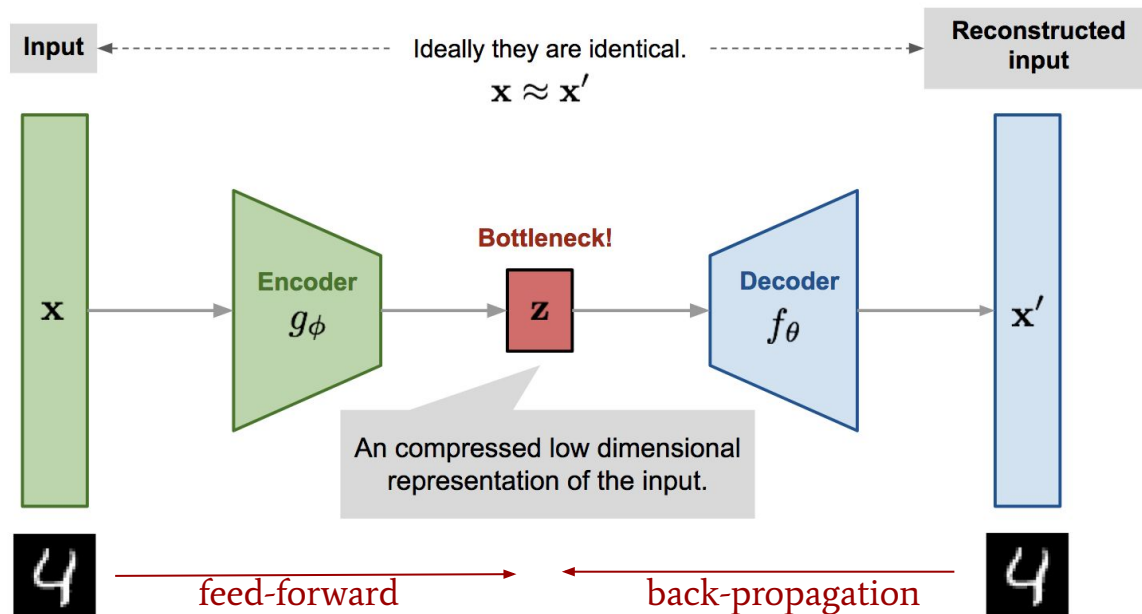Press Esc to exit full screen

# Galaxy Image Exploration and Classification



- Image Exploration
- Resize is done dynamically
- Quick Classification/Label
- Works fine with 10,000 images
- Individual classifications are saved and aggregated
- Keyboard control

https://github.com/mgckind/cutouts-explorer

# Autoencoders review



Input → Ideally they are identical. → Reconstructed input

$$\mathbf{x} \approx \mathbf{x'}$$

$\mathbf{x}$ — Encoder $g_\phi$ — **Bottleneck!** $\mathbf{z}$ — Decoder $f_\theta$ — $\mathbf{x'}$

An compressed low dimensional representation of the input.

feed-forward ← back-propagation

- Around since the 80's
- Data compression
- Anomaly detection
- Denoising
- Regular Machine Learning
- PCA

Many variants: Sparse AE, Contractive AE, Stacked AE, etc...

$$z = g_\phi(x)$$
$$x' = f_\theta(g_\phi(x))$$
$$\mathcal{L}(x, x') + reg.$$
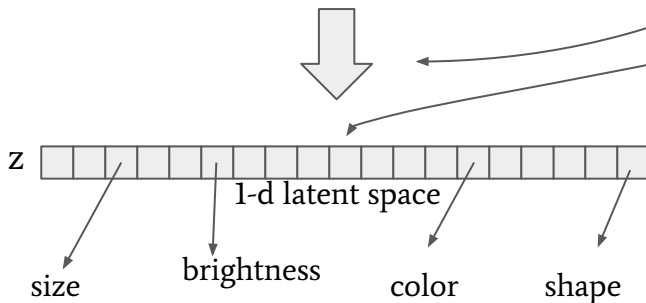
$$\ell_1 = \lambda \sum_i |a_i^{(h)}|$$

$$\mathcal{L}(x, x') = ||x - x'||^2$$

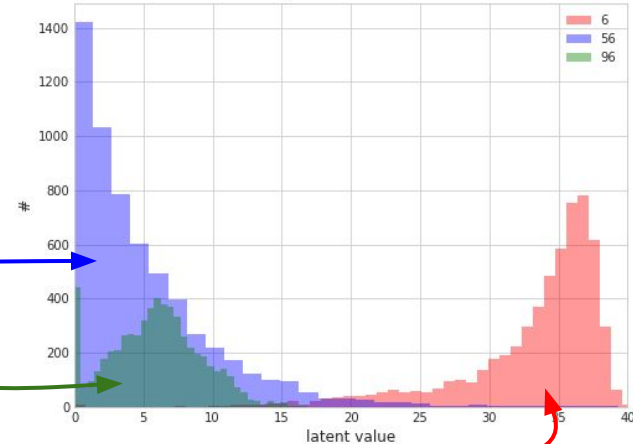# The latent space



Credit: esa

# The latent space



Credit: esa



Input ⟵ · · · · · · · · · · Ideally they are identical. · · · · · · · · · · ⟶ Reconstructed input

$$\mathbf{x} \approx \mathbf{x}'$$

**Encoder** $g_\phi$

**Bottleneck!**

$\mathbf{z}$

**Decoder** $f_\theta$

$\mathbf{x}$   $\mathbf{x}'$

An compressed low dimensional representation of the input.

z

1-d latent space

# The latent space

Credit: esa

flatten

encoding

z

1-d latent space

Input — Ideally they are identical. — Reconstructed input
$$\mathbf{x} \approx \mathbf{x}'$$

**x**

Encoder $g_\phi$

**Bottleneck!**

**z**

An compressed low dimensional representation of the input.

Decoder $f_\theta$

**x'**

# The latent space

flatten

Input ← ⋯⋯⋯⋯⋯ Ideally they are identical. ⋯⋯⋯⋯⋯ → Reconstructed input

$$\mathbf{x} \approx \mathbf{x}'$$

Bottleneck!

$\mathbf{x}$

Encoder $g_\phi$

$\mathbf{z}$

Decoder $f_\theta$

$\mathbf{x}'$

An compressed low dimensional representation of the input.

Credit: esa

encoding

z

1-d latent space

size       brightness       color       shape

# The latent space



Credit: esa

flatten

$x \approx x'$

Bottleneck!

Encoder $g_\phi$

$z$

Decoder $f_\theta$

$x$

$x'$

An compressed low dimensional representation of the input.

Each galaxy image is encoded to a unique 1-d vector

encoding

z

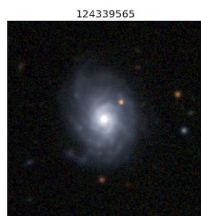1-d latent space

size    brightness    color    shape

$z_i$ distribution

# Similarity ranking

Using standard ML in latent space to look for neighbors, outliers, etc.

Compress images from 220x220x3 pixels to 100-vector (2000x), for fast similarity search, anomaly detection, etc...
No need decoder (only for Loss)

# Similarity ranking
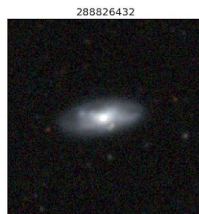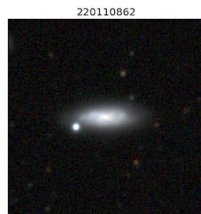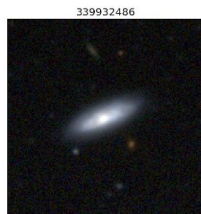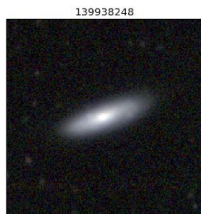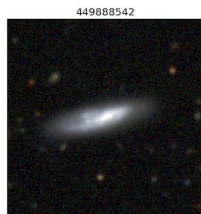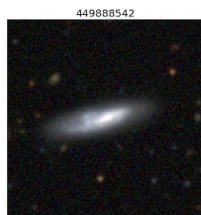
Using standard ML in latent space to look for neighbors, outliers, etc.

Compress images from 220x220x3 pixels to 100-vector (2000x), for fast similarity search, anomaly detection, etc...
No need decoder (only for Loss)

# Similarity ranking

Using standard ML in latent space to look for neighbors, outliers, etc.

Compress images from 220x220x3 pixels to 100-vector (2000x), for fast similarity search, anomaly detection, etc...
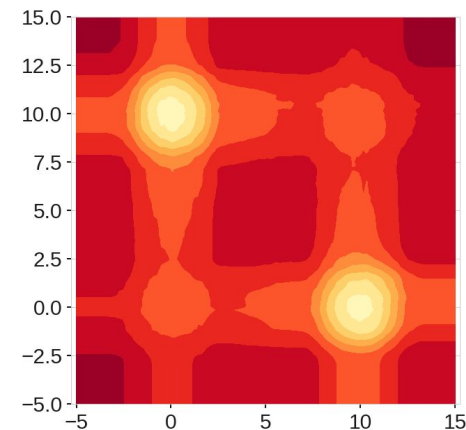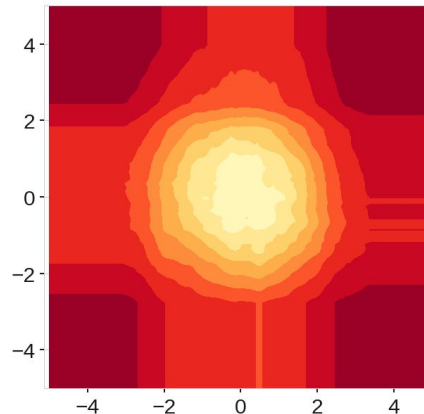No need decoder (only for Loss)

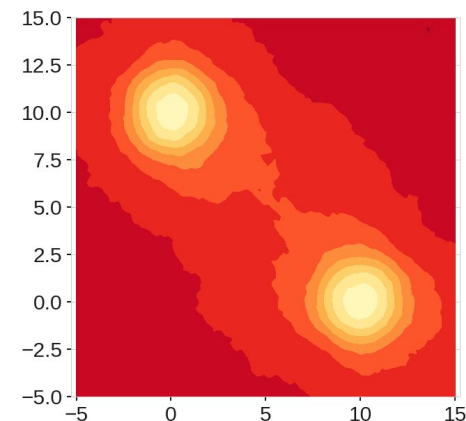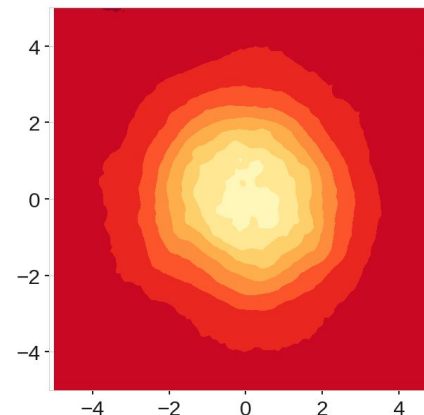# Anomaly detection with Extended Isolation Forest

**Isolation Forest:**

- ✔ Model free
- ✔ Computationally efficient
- ✔ Readily application to high dimensional data
- ✘ Inconsistent scoring seen in score maps

**Extended Isolation Forest:**

- ✔ Model free
- ✔ Computationally efficient
- ✔ Readily application to high dimensional data
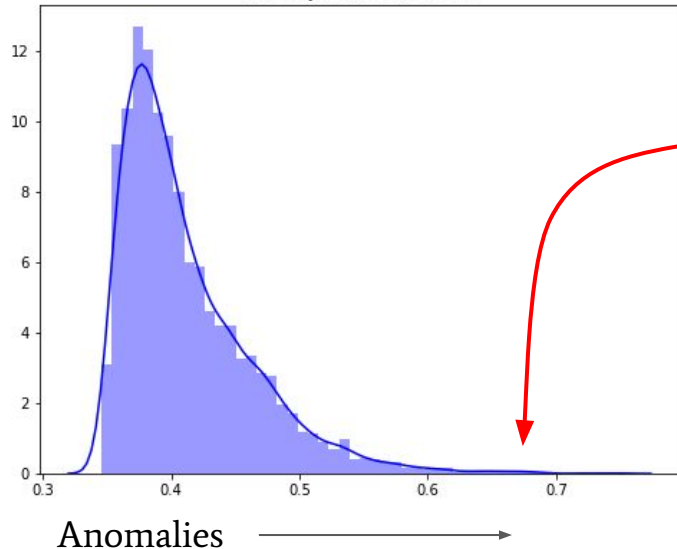- ✔ Consistent scoring

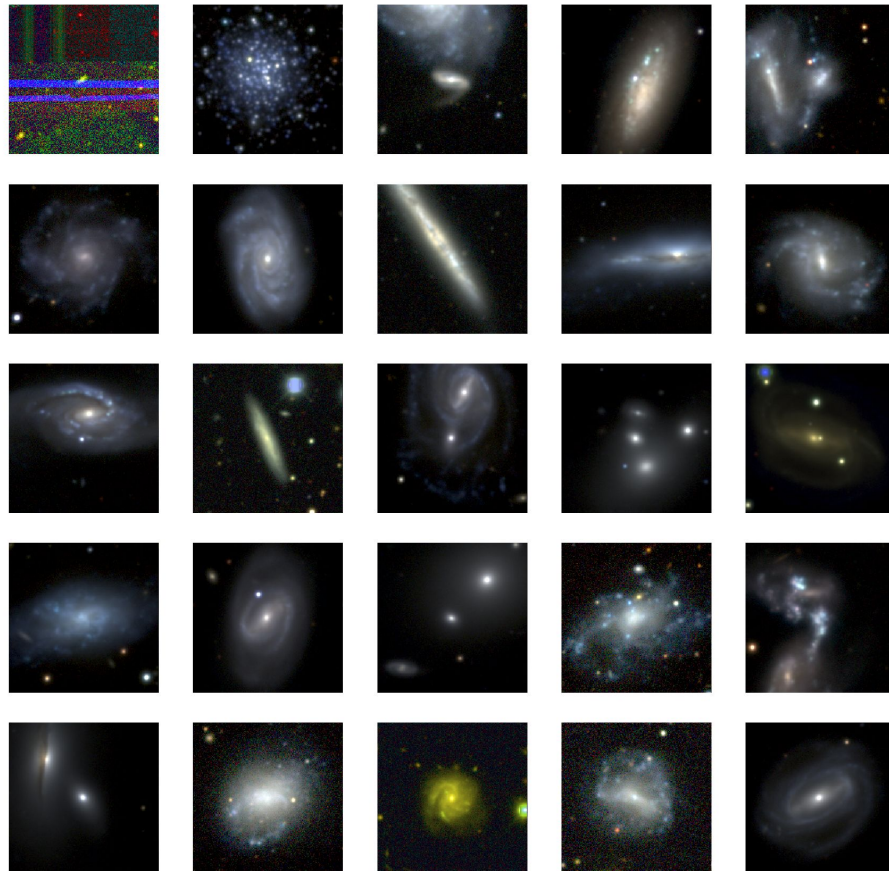Hariri, Carrasco-Kind, Brunner, 2019 , arXiv: 1811.02141

https://github.com/sahandha/eif

# Anomaly detection with EIF in high dimensional latent space

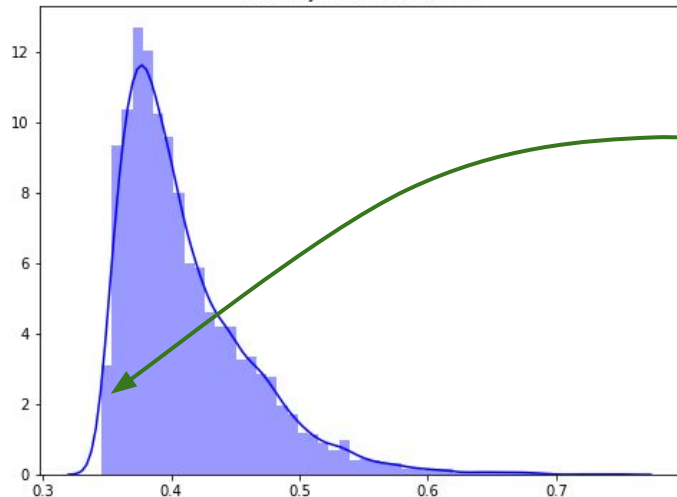Anomaly Score Distribution

Anomalies →

Top 25 anomalies

The EIF algorithm produces a anomaly score that can be used to select outliers galaxies.

Errors, special cases, unrepresented galaxies

# Anomaly detection with EIF in high dimensional latent space
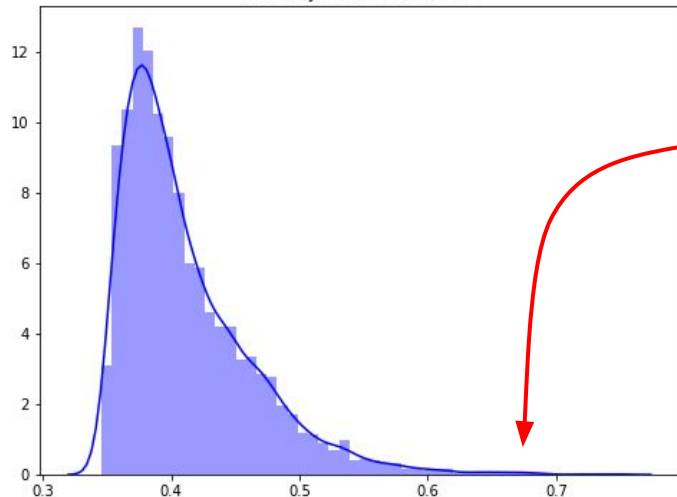


Anomaly Score Distribution

Anomalies

Top 25 nominals

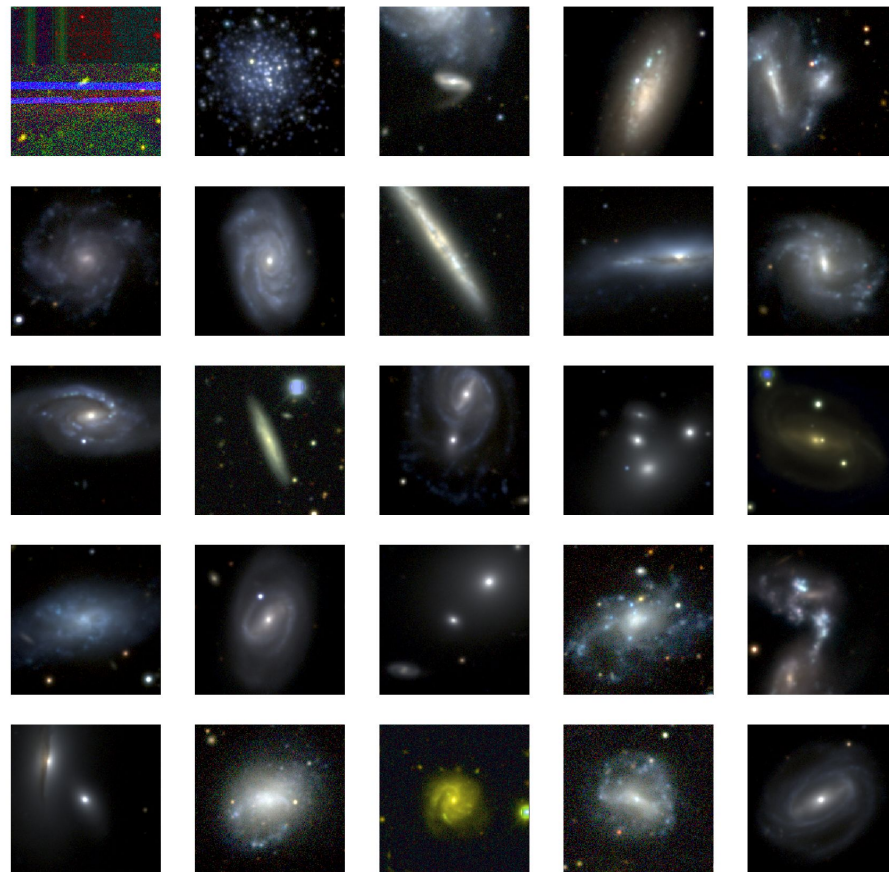But also can tell us about the repeated and common cases as shown here

These are all different sources!

# Anomaly detection with EIF in high dimensional latent space

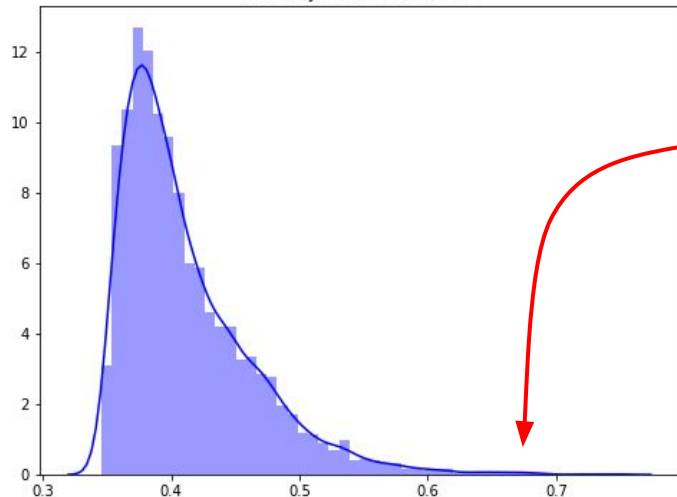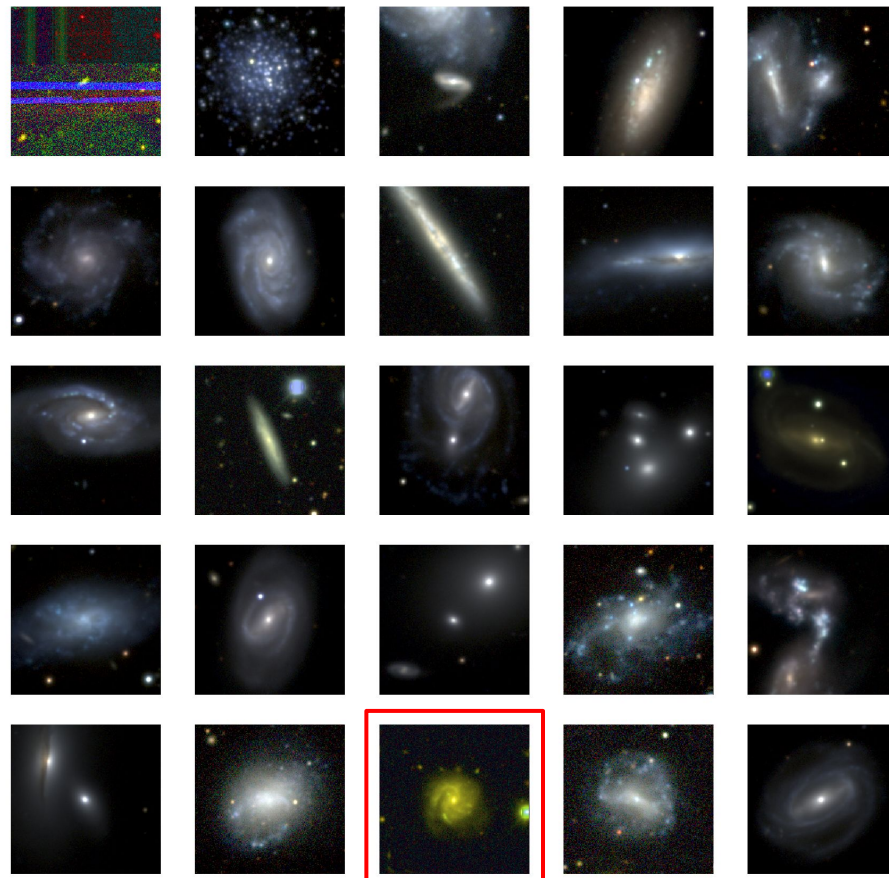Anomaly Score Distribution

Anomalies →

Top 25 anomalies

We can find other errors by looking at similarities

# Anomaly detection with EIF in high dimensional latent space



Anomaly Score Distribution

Anomalies

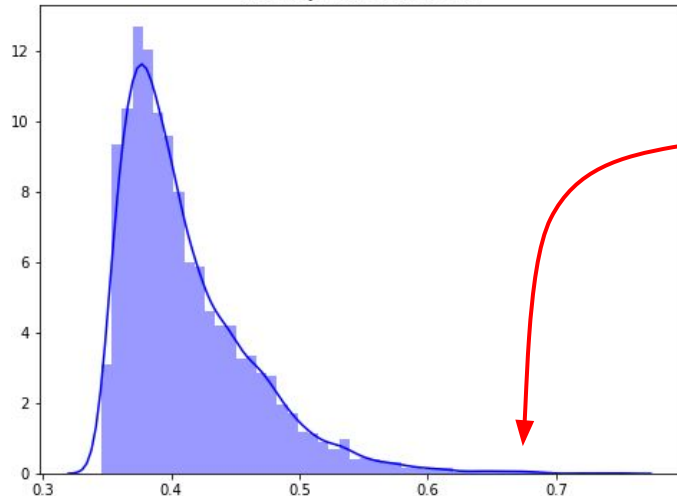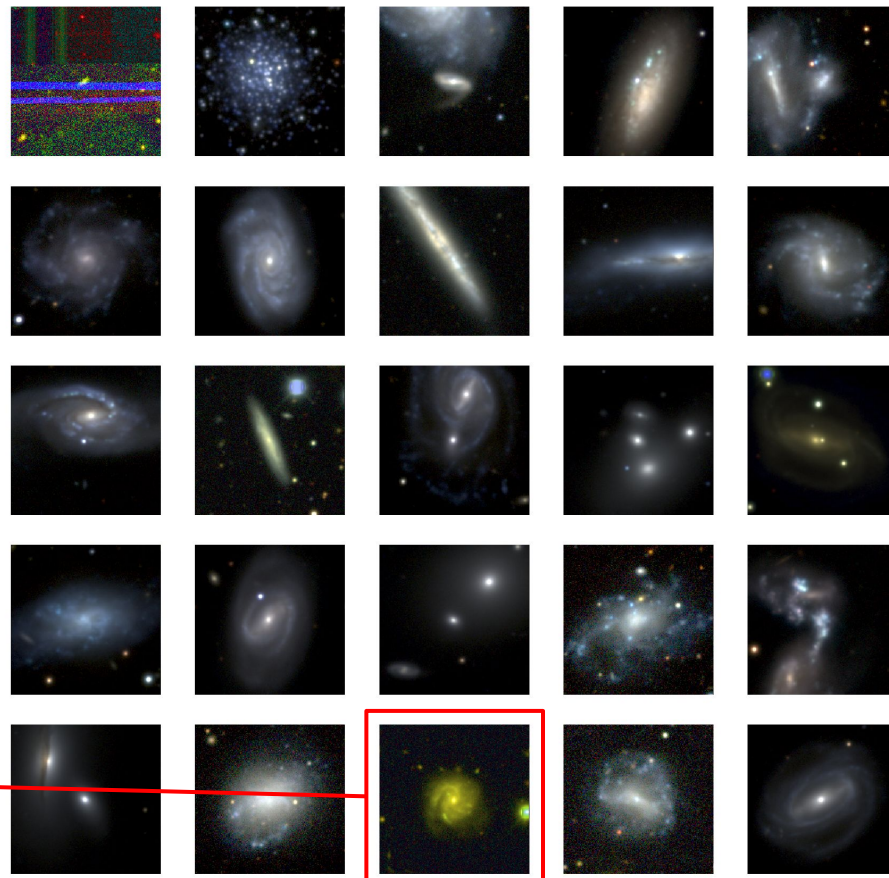Top 25 anomalies

We can find other errors by looking at similarities

# Anomaly detection with EIF in high dimensional latent space



Anomaly Score Distribution
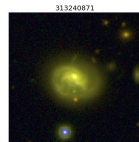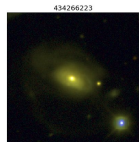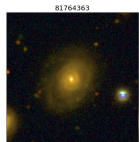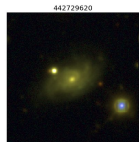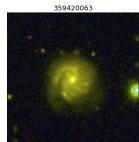
Anomalies →
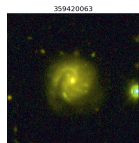
Top 25 anomalies

We can find other errors by looking at similarities

# Can we keep reducing dimensions?

- Apply clustering and unsupervised techniques to latent space to find patterns
- Self Organizing Maps and T-SNE are perfect candidates
- We have used Uniform Manifold Approximation and Projection (UMAP)

Credit: esa

2d visualization

z

1-d latent space

# UMAP Representation (6000 galaxies)

# UMAP Representation (6000 galaxies)

# UMAP Representation (6000 galaxies)

# UMAP Representation (6000 galaxies)

# UMAP Representation (6000 galaxies)

# Conclusions

- We developed a visualization and classification tool for multiple images
- Using Autoencoders we can compress images to small (but high-n) latent space
- Look for similarities and anomalies in that space
- Represent even more in a 2d graph using t-SNE, SOM or UMAP
- Scientific driven cases
- State-of-the-art models allows a bayesian manipulation of the latent space

# Thank you!

## Questions?

Matias Carrasco Kind -- NCSA
mcarras2@illinois.edu
github.com/mgkind
matias-ck.com

Original

Model 1

Model 2

Blurry images and structure is lost, but angular sizes, radial profiles and brightness are a match.

What if we can make the model learn properties at the same time as images.? What if can sample from the latent space?

**Input** ⟵- - - - - - - - - - - - - - - - Ideally they are identical. - - - - - - - - - - - - - -⟶ **Reconstructed input**

$$\mathbf{x} \approx \mathbf{x}'$$

**Probabilistic Encoder**

$$q_\phi(\mathbf{z}|\mathbf{x})$$

Mean $\boldsymbol{\mu}$

$\mathbf{x}$

Std. dev $\boldsymbol{\sigma}$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}$$
$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{I})$$

**Sampled latent vector**

$\mathbf{z}$

An compressed low dimensional representation of the input.

**Probabilistic Decoder**
$$p_\theta(\mathbf{x}|\mathbf{z})$$

$\mathbf{x}'$

We can generate samples from z, next step is can we constrain what's being sampled?

**Multimodal Generative Models for Scalable Weakly-Supervised Learning**

Mike Wu
Department of Computer Science
Stanford University
Stanford, CA 94025
wumike@stanford.edu

Noah Goodman
Departments of Computer Science and Psychology
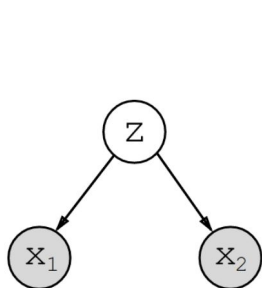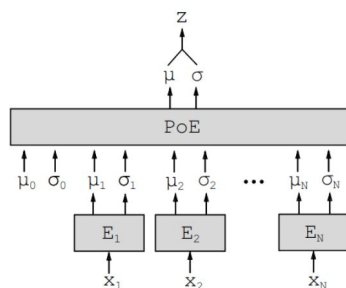Stanford University
Stanford, CA 94025
ngoodman@stanford.edu

multimodal-generative-models-for-scalable-weakly-supervised-learning



(a)          (b)          (c)

Learning joint representation of conditionally independent modalities using product of experts.

We can:
- Conditional sample with certain attributes
- Sample without any limitations
- Change the attribute of an existing input data
- Similarity search and anomaly detection
- Predict one modality from the others
- Sample and train with missing modalities

Samples with changing brightness (increasing downwards)



Bowen, Carrasco-Kind, et al. in prep.

Samples with changing area
(increasing downwards)



encoder

X    Y

$\mu_X, \sigma_X$    $\mu_0, \sigma_0$    $\mu_Y, \sigma_Y$

$\mu_z, \sigma_z$

PoE

Z

X'    Y'

decoder

Bowen, Carrasco-Kind, et al. in prep.

36

Bowen, Carrasco-Kind, et al. in prep.
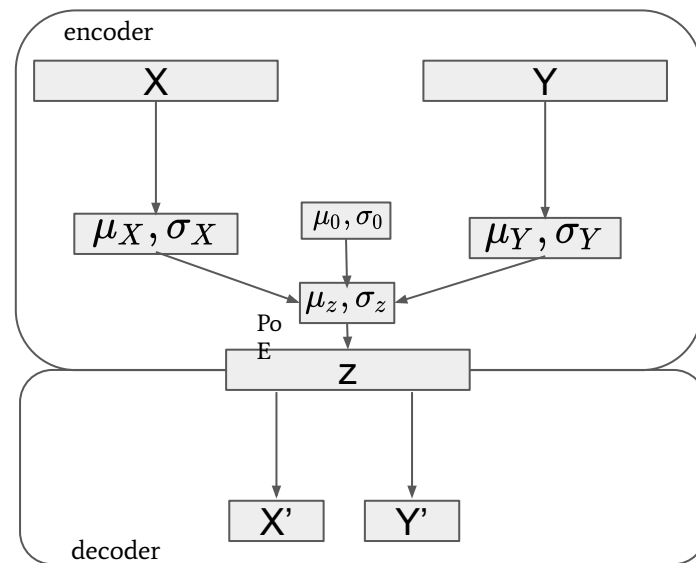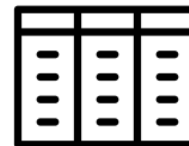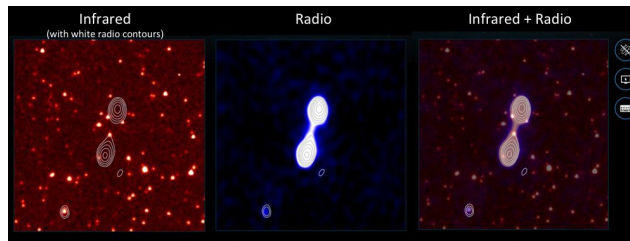
# Can we sample z to generate fake data?



Exist $\theta$ for max the likelihood

$$p(x) = \int p(x|z, \theta)p(z)dz$$

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

<span style="color:red">Too expensive</span>

We need an approximate posterior (prob. encoder)

$$q_\lambda(z|x) \approx p(z|x)$$

And we can use q to be Gaussian (there are other alternatives)

$$q_\lambda(z|x) = \mathcal{N}(z; \mu_\lambda(x), \sigma_\lambda(x))$$

$$p(z) = \mathcal{N}(0, I)$$

- Map x to a distribution $\quad p(z|x)$
- Sample from distribution $\quad z_i \sim p(z)$
- Generate fake data $\quad x'_i \sim p_\theta(x'|z)$
- Probabilistic approach

$$p(x, z) = p(x|z)p(z)$$

**Variational "Autoencoder"**

# Generative Adversarial Networks (GAN)

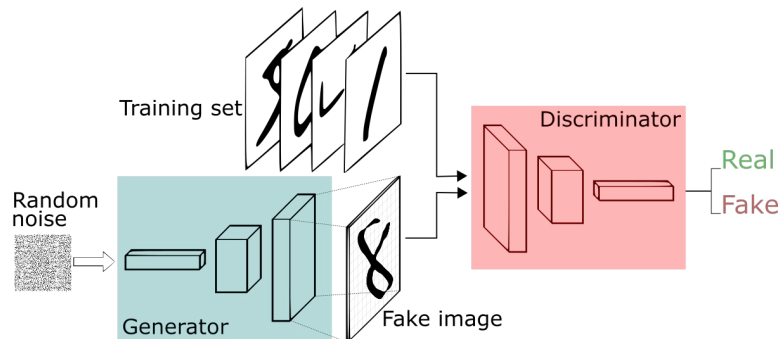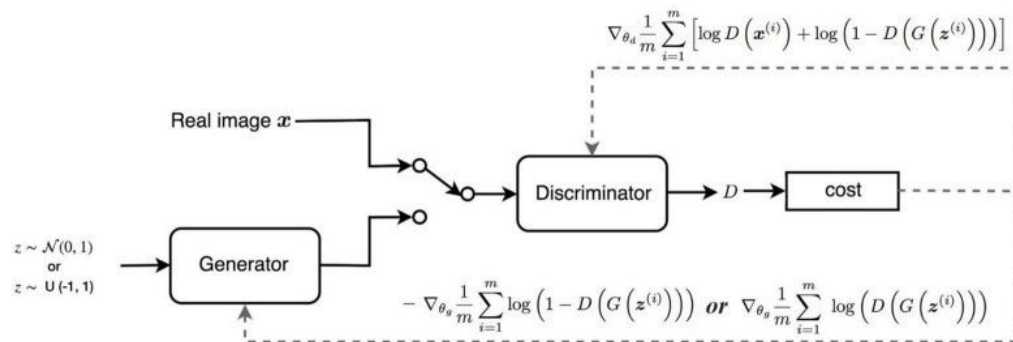Gan can be VERY good to specific image generation and create realistic images. Very powerful discriminator



But:
- Very hard to train (unstable)
- Not really sampling methods
- Hard to evaluate likelihood of data p(x)
- Tend to underfit data distribution
- Main goal is to fool the discriminator

Very powerful if combined with VAE



$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(\boldsymbol{x}^{(i)}\right) + \log \left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \right]$$

$z \sim \mathcal{N}(0, 1)$
or
$z \sim U(-1, 1)$

$$-\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log \left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \ \ or \ \ \nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log \left(D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right)$$