# Self-Supervised Representation Learning
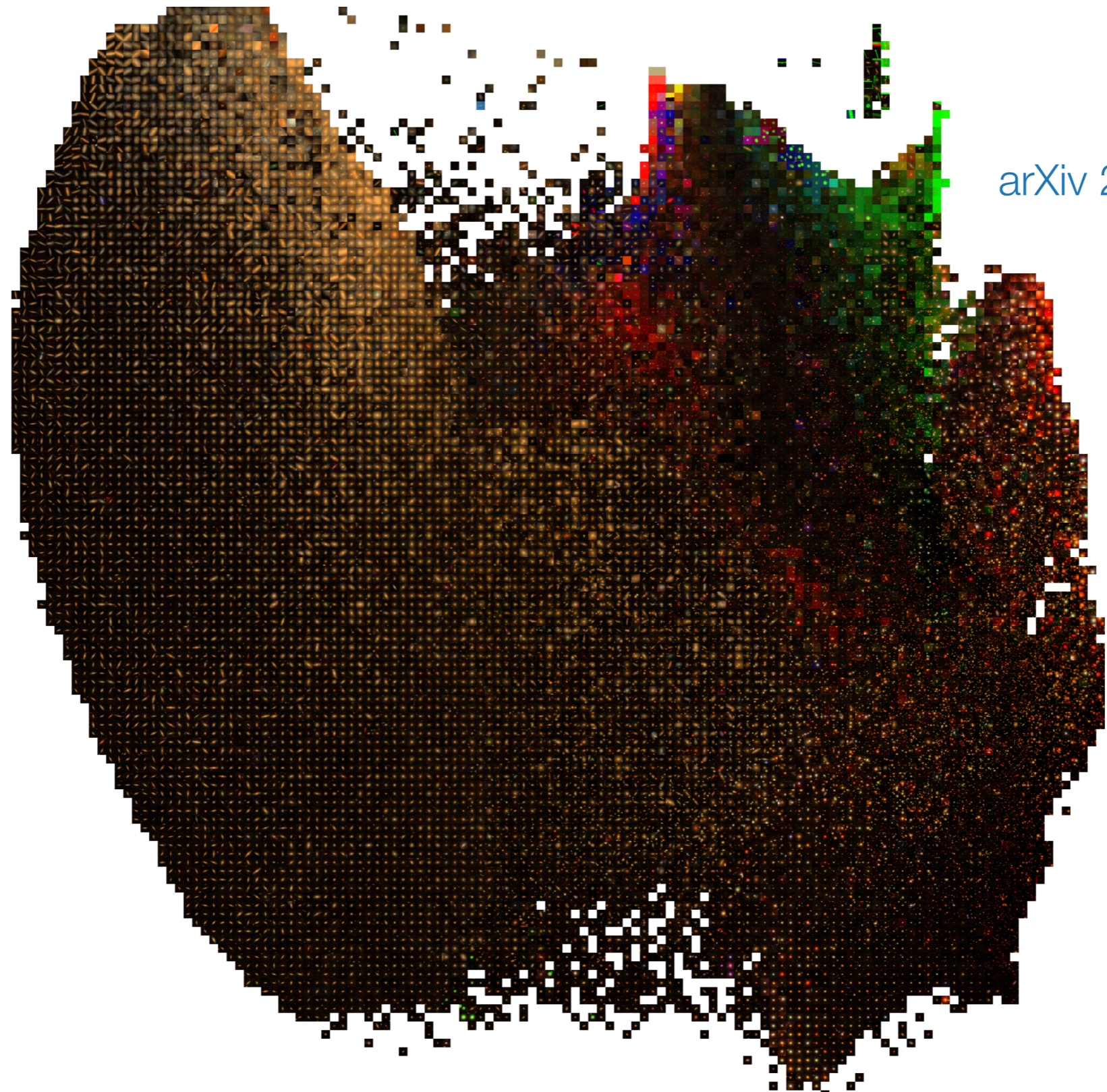## for Astronomical Images

Md Abul Hayat
mahayat@uark.edu

George Stein
gstein@berkeley.edu

Peter Harrington
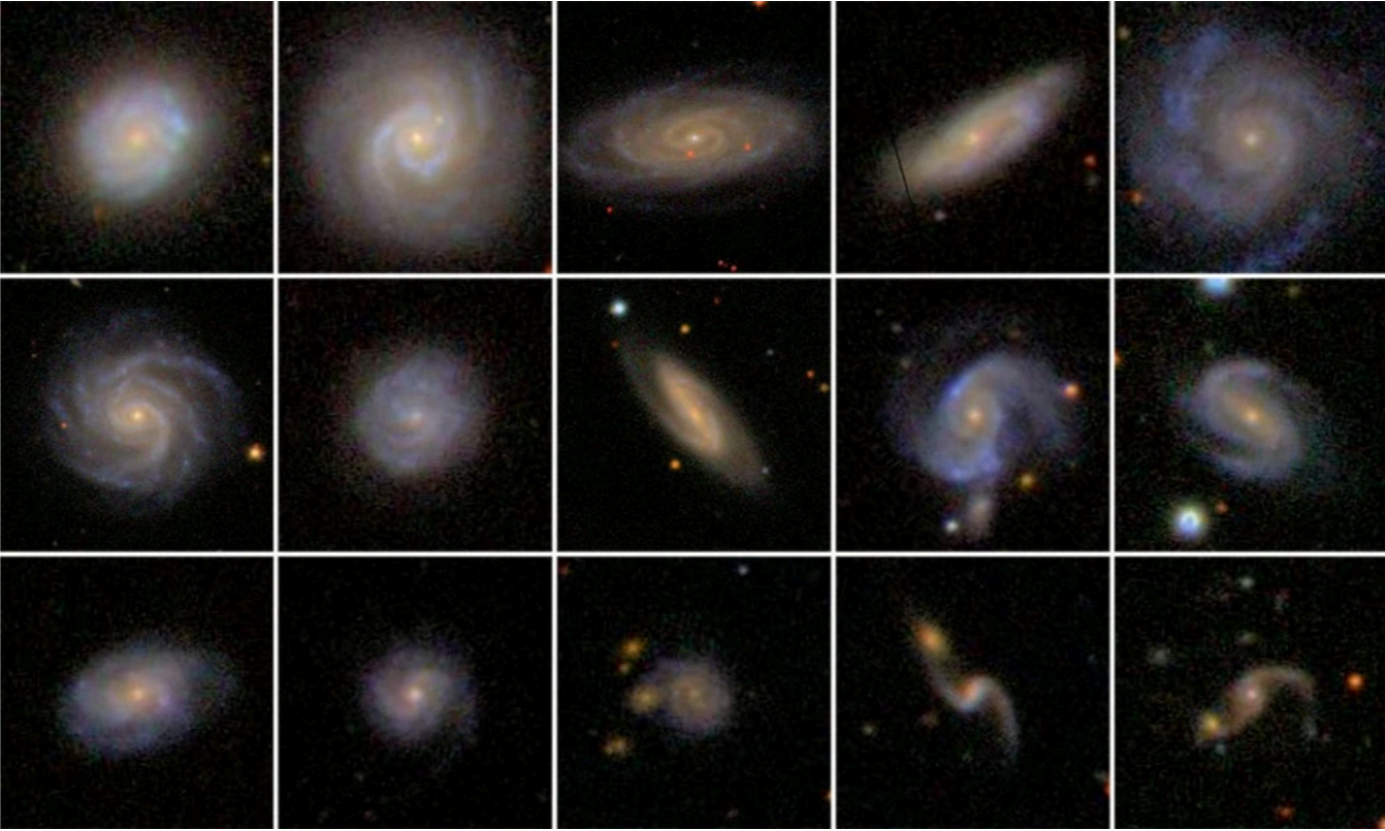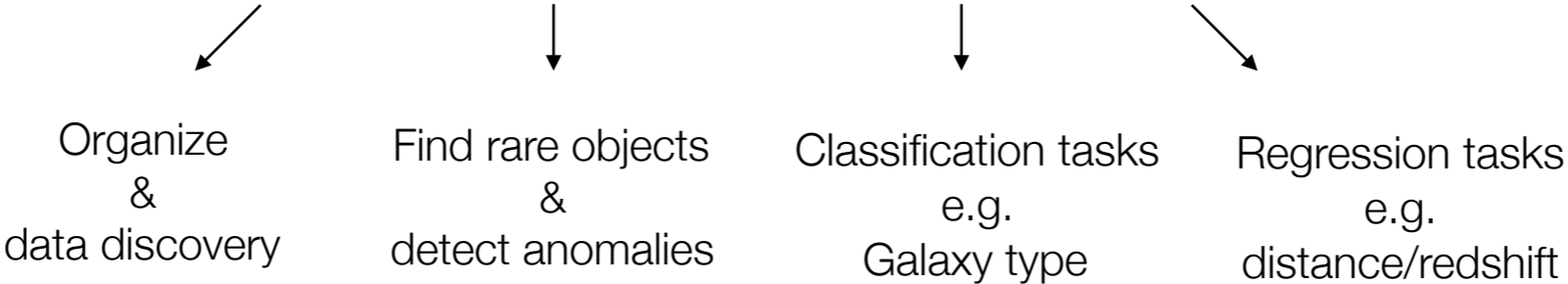pharrington@lbl.gov

Zarija Lukić
zarija@lbl.gov

Mustafa Mustafa
mmustafa@lbl.gov

*blue = clickable links throughout

arXiv 2012.13083

# Sky Surveys

Sky surveys are massive data generators,
imaging ~10 billion galaxies in the near future

| Organize & data discovery | Find rare objects & detect anomalies | Classification tasks e.g. Galaxy type | Regression tasks e.g. distance/redshift |



Galaxies from the Sloan Digital Sky Survey (SDSS)

# Supervised learning

Target a specific task and collect labels

| Organize & data discovery | Find rare objects & detect anomalies | Classification tasks e.g. galaxy type | Regression tasks e.g. distance/redshift |

**Data**

*CNN*

1. Design architecture

2. Split data into train/validate/test

3. Train model weights through back propagation

4. Deploy model, assuming unlabelled data comes from the same distribution

**Labels**

3

# Supervised learning

Target a specific task and collect labels

| Organize & data discovery | Find rare objects & detect anomalies | Classification tasks e.g. galaxy type | Regression tasks e.g. distance/redshift |

**Data**

*CNN*

1. Design architecture
2. Split data into train/validate/test
3. Train model weights through back propagation
4. Deploy model, assuming unlabelled data comes from the same distribution

**Labels**

# Unsupervised learning

Clustering, dimensionality reduction and feature selection, etc..

# Self-supervised representation learning

Without any labels, learn low-dimensional representations of data which preserve semantic information

Then, use representations for "downstream tasks" (regression, classification, etc…)

# Self-supervised representation learning

Without any labels, learn low-dimensional representations of data which preserve semantic information

Then, use representations for "downstream tasks" (regression, classification, etc…)

*Self-supervised Learning: Generative or Contrastive*

**Generative**
- Autoencoder, VAE
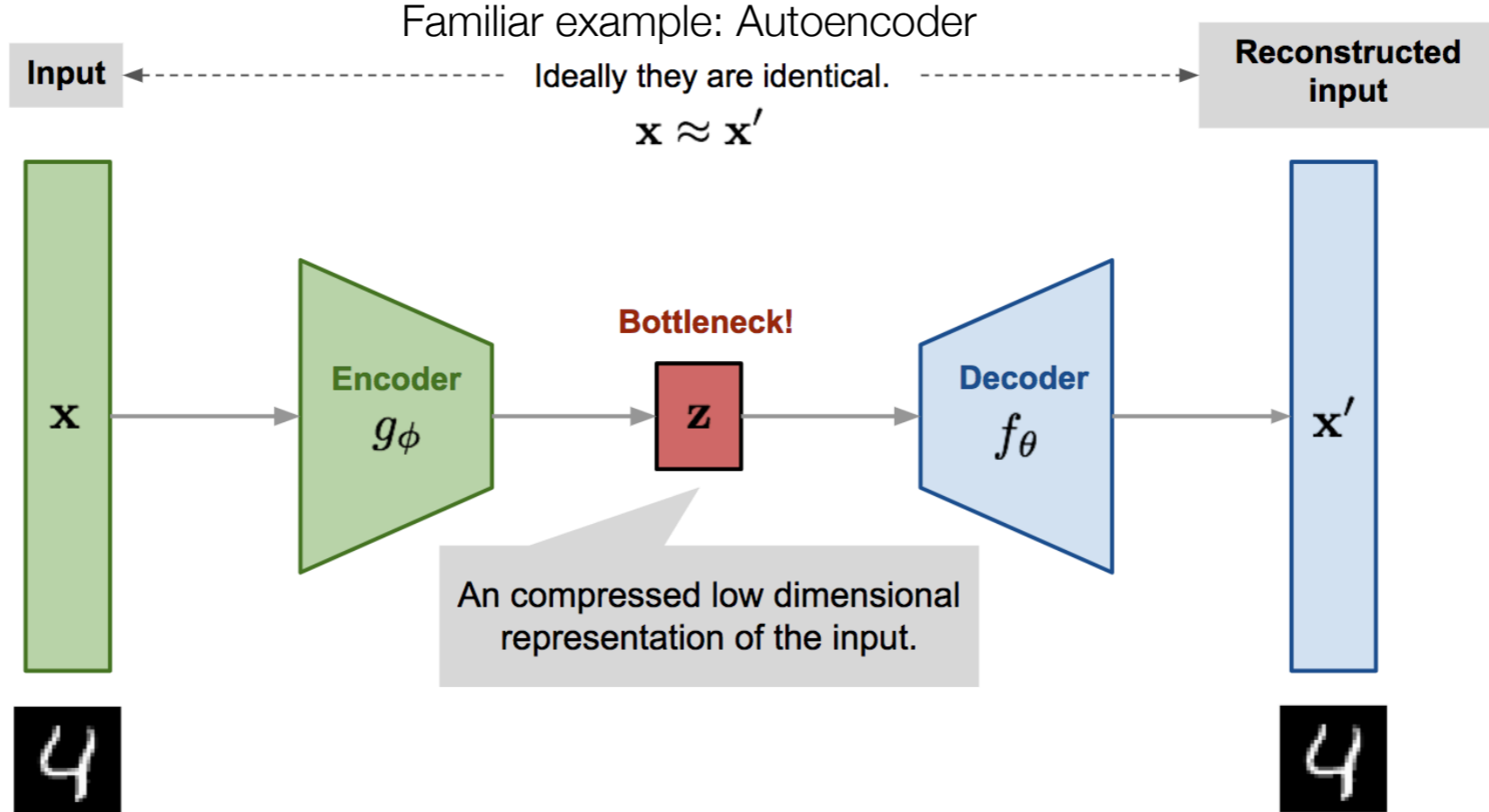- Flow-based
- …

# Self-supervised representation learning

Without any labels, learn low-dimensional representations of data which preserve semantic information

Then, use representations for "downstream tasks" (regression, classification, etc…)

*Self-supervised Learning: Generative or Contrastive*

## Generative

- Autoencoder, VAE
- Flow-based
- …

Familiar example: Autoencoder



Input $\dashleftarrow$ Ideally they are identical. $\dashrightarrow$ Reconstructed input

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

Bottleneck!

$\mathbf{x}$    Encoder $g_\phi$    $\mathbf{z}$    Decoder $f_\theta$    $\mathbf{x}'$

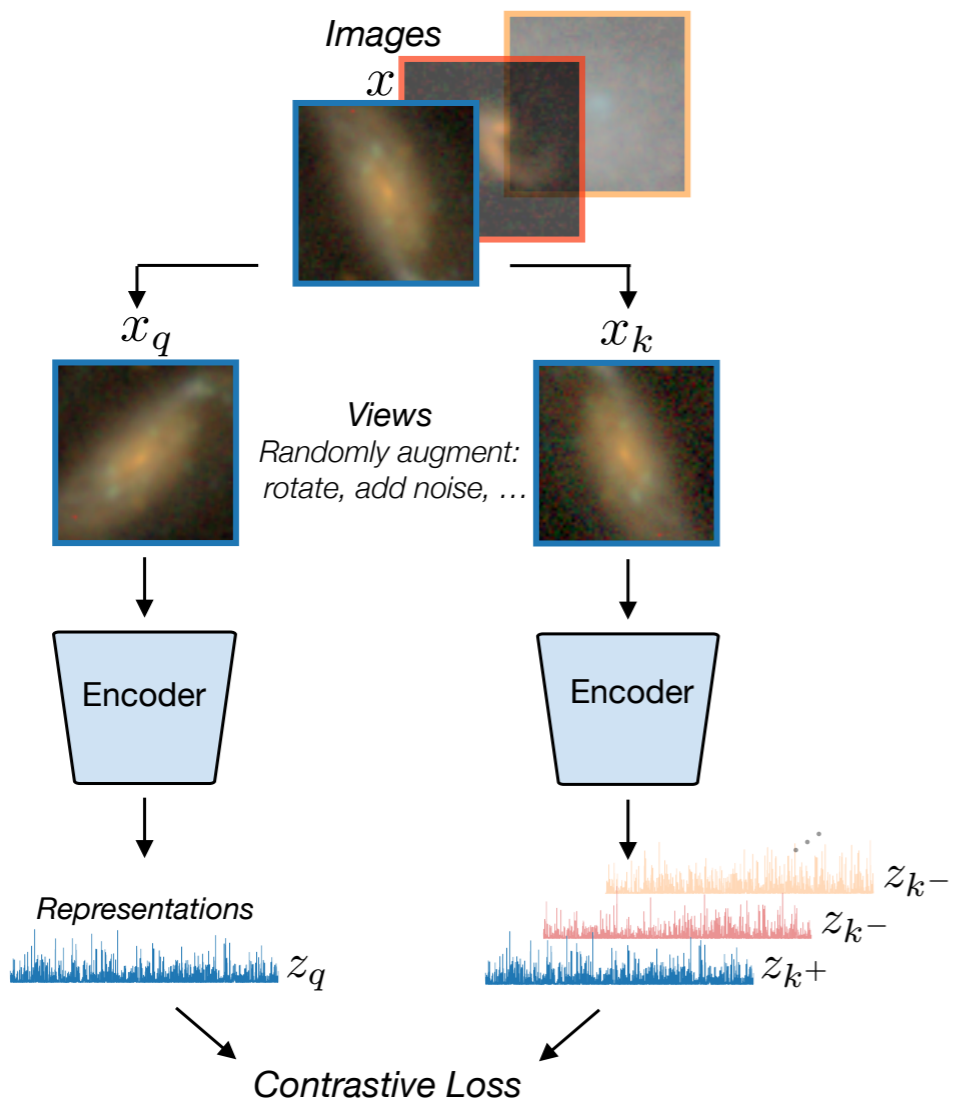An compressed low dimensional representation of the input.

Train using MSE loss or similar

nothing to enforce that small variations in the input image result in similar representations

# Contrastive learning

Learn representations that are invariant to desired augmentations

**1. Self-supervised contrastive representation learning**
Learn representations in an unsupervised manner



*Images*
$x$

$x_q$        $x_k$

*Views*
*Randomly augment:*
*rotate, add noise, …*

Encoder       Encoder

*Representations*

$z_q$

$z_{k^-}$
$z_{k^-}$
$z_{k^+}$

*Contrastive Loss*

Task-agnostic augmentations for galaxy surveys:

-               Rotation
-               Jitter
-          Gaussian noise
-       Galactic extinction
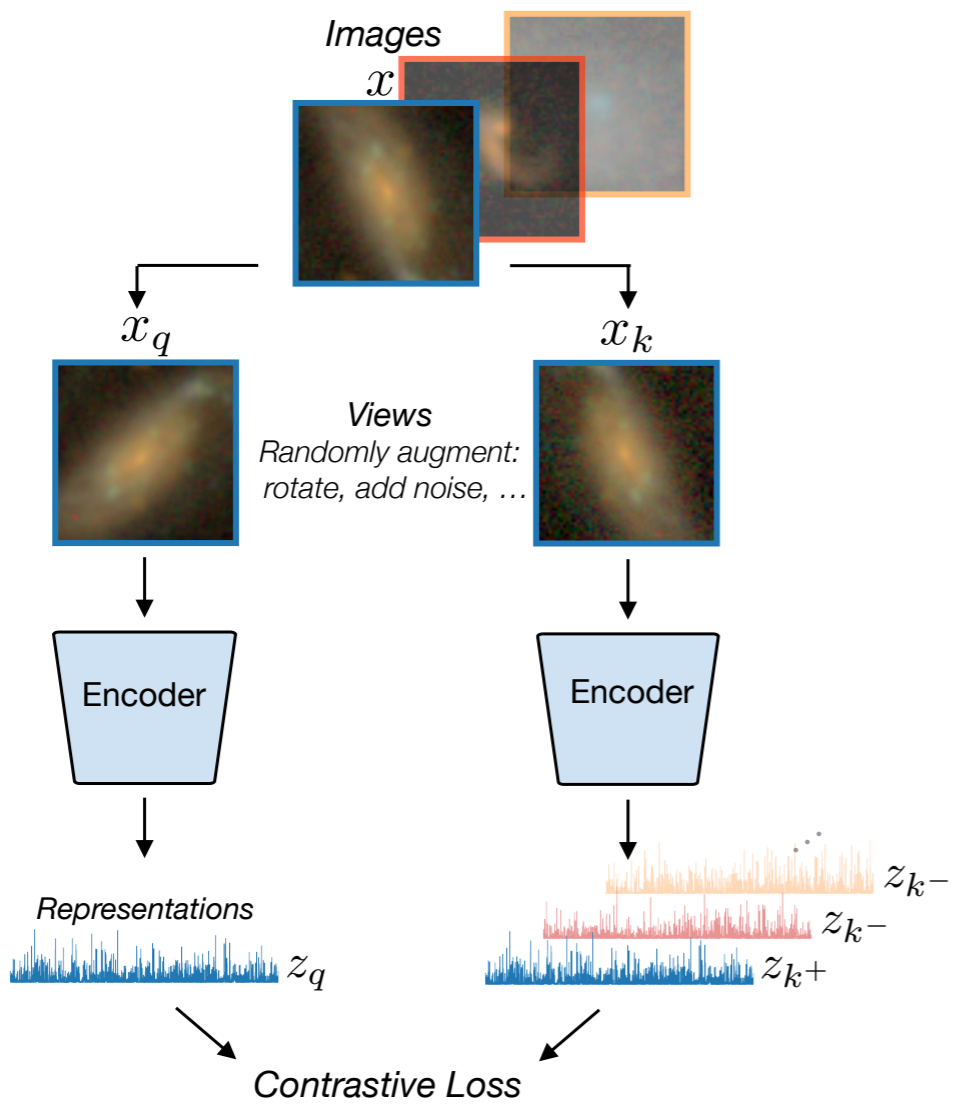-   Point spread function (blur)
-               …

For other data/applications, choose your own
using domain knowledge

Self-supervised representation learning for astronomical images

# Contrastive learning
Learn representations that are invariant to desired augmentations

**1. Self-supervised contrastive representation learning**
Learn representations in an unsupervised manner



Task-agnostic augmentations for galaxy surveys:

-                    Rotation
-                    Jitter
-               Gaussian noise
-             Galactic extinction
-        Point spread function (blur)
-                    …

For other data/applications, choose your own using domain knowledge

$$L_{q,k^+,\{k^-\}} =$$
$$-\log\left(\frac{\exp(\mathrm{sim}(\mathbf{z_q}, \mathbf{z_{k^+}}))}{\exp(\mathrm{sim}(\mathbf{z_q}, \mathbf{z_{k^+}})) + \sum_{\mathbf{k^-}} \exp(\mathrm{sim}(\mathbf{z_q}, \mathbf{z_{k^-}}))}\right), \tag{1}$$
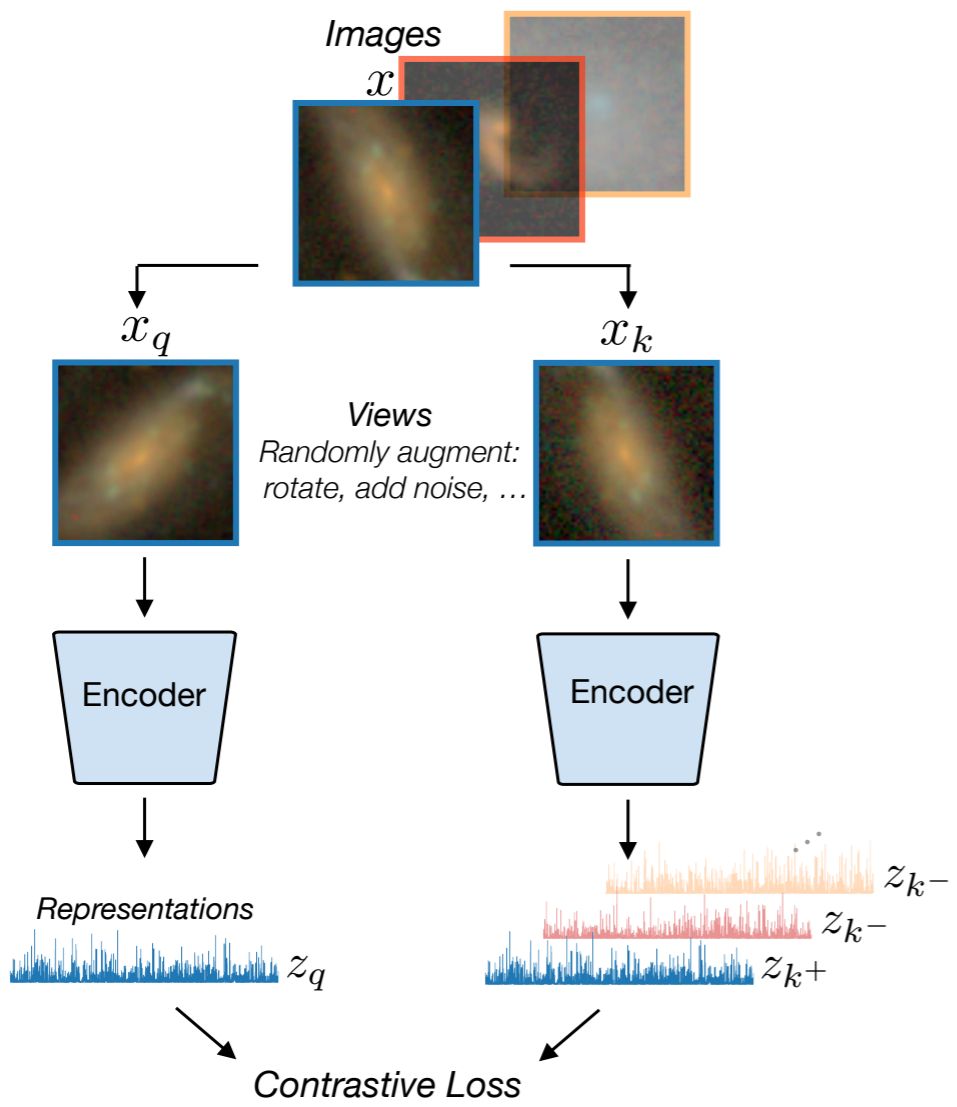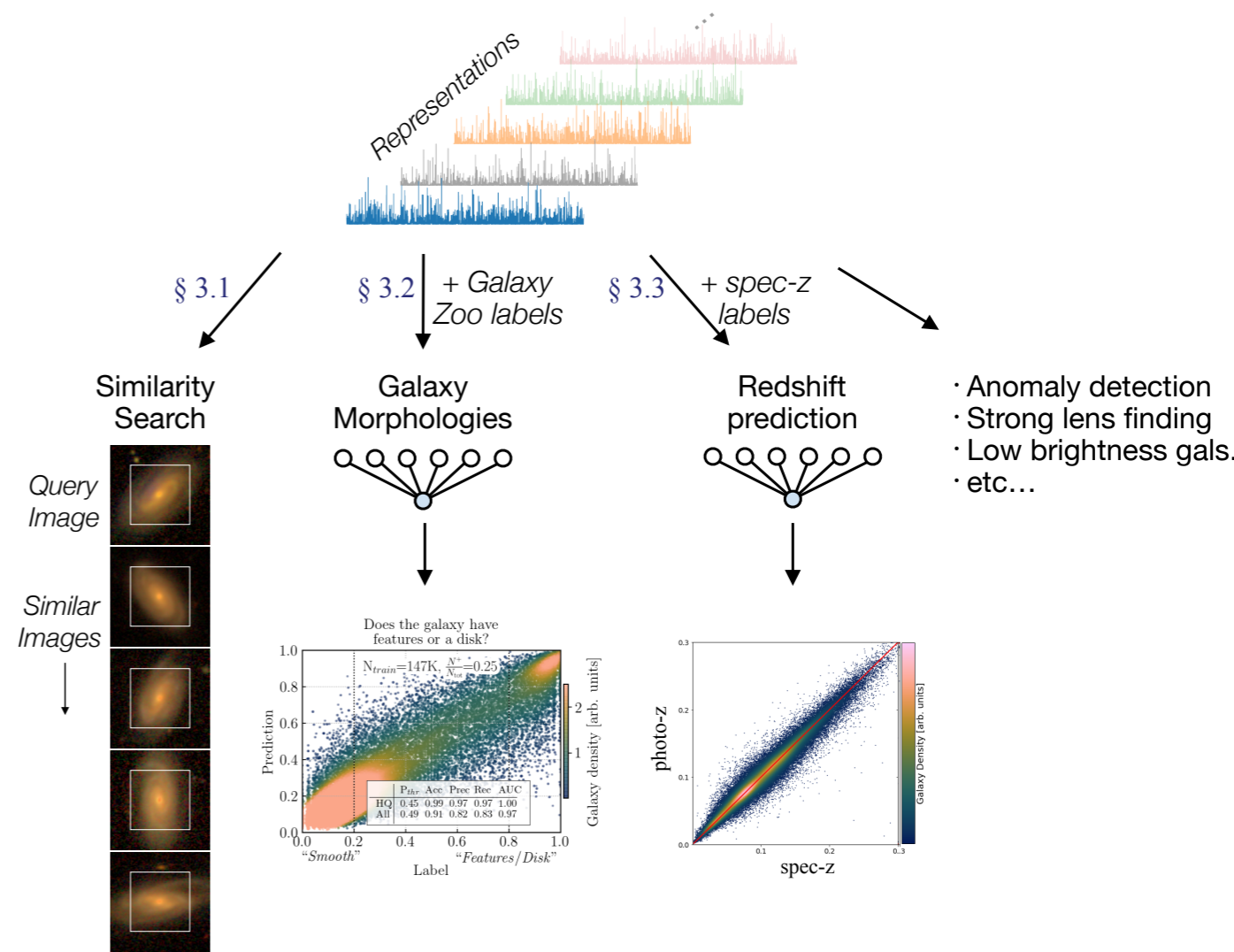
where $\mathrm{sim}(\mathbf{a}, \mathbf{b}) = \mathbf{a} \cdot \mathbf{b}/(\tau\|\mathbf{a}\|\,\|\mathbf{b}\|)$ is the cosine similarity measure between vectors $\mathbf{a}$ and $\mathbf{b}$, normalized by a tunable "temperature" hyper-parameter $\tau$. This loss (InfoNCE, Oord et al. 2018) is minimized when positive pairs have high similarity, while negative pairs have low similarity. We have closely followed Chen, X. et al.

Self-supervised representation learning for astronomical images

# Contrastive learning
Learn representations that are invariant to desired augmentations

**1. Self-supervised contrastive representation learning**
Learn representations in an unsupervised manner

*Images*
$x$

$x_q$     $x_k$

*Views*
*Randomly augment:*
*rotate, add noise, …*

Encoder     Encoder

*Representations*

$z_q$

$z_{k^-}$
$z_{k^-}$
$z_{k^+}$

*Contrastive Loss*

**2. Downstream tasks**
Use representations for a variety of applications

*Representations*

§ 3.1    § 3.2 | *+ Galaxy Zoo labels*    § 3.3 | *+ spec-z labels*

Similarity Search    Galaxy Morphologies    Redshift prediction

· Anomaly detection
· Strong lens finding
· Low brightness gals.
· etc…

*Query Image*

*Similar Images*

Does the galaxy have features or a disk?

$N_{train}$=147K, $\frac{N^+}{N_{tot}}$=0.25

| | $P_{thr}$ | Acc | Prec | Rec | AUC |
|---|---|---|---|---|---|
| HQ | 0.45 | 0.99 | 0.97 | 0.97 | 1.00 |
| All | 0.49 | 0.91 | 0.82 | 0.83 | 0.97 |

Prediction

"Smooth"   Label   "Features/Disk"

Galaxy density [arb. units]

photo-z

spec-z

Galaxy Density [arb. units]

Self-supervised representation learning for astronomical images     arxiv: 2012.13083

# Dataset

SDSS 5 band images (ugriz)

1.3 million images total

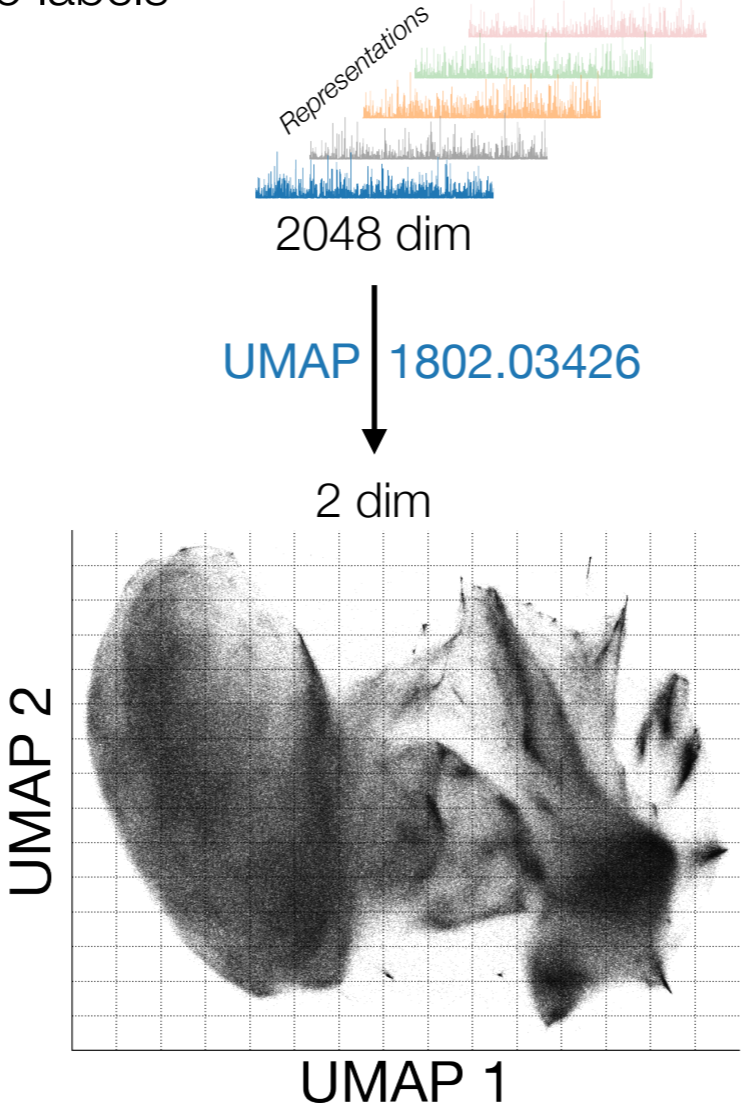500k have redshift labels from spectroscopic followup

180k have crowd sourced morphology classifications

Workflow

1. Train Resnet-50 encoder through self-supervised contrastive representation learning on all 1.3M images using *MoCo v2* framework (see *SimCLR* for alternative)
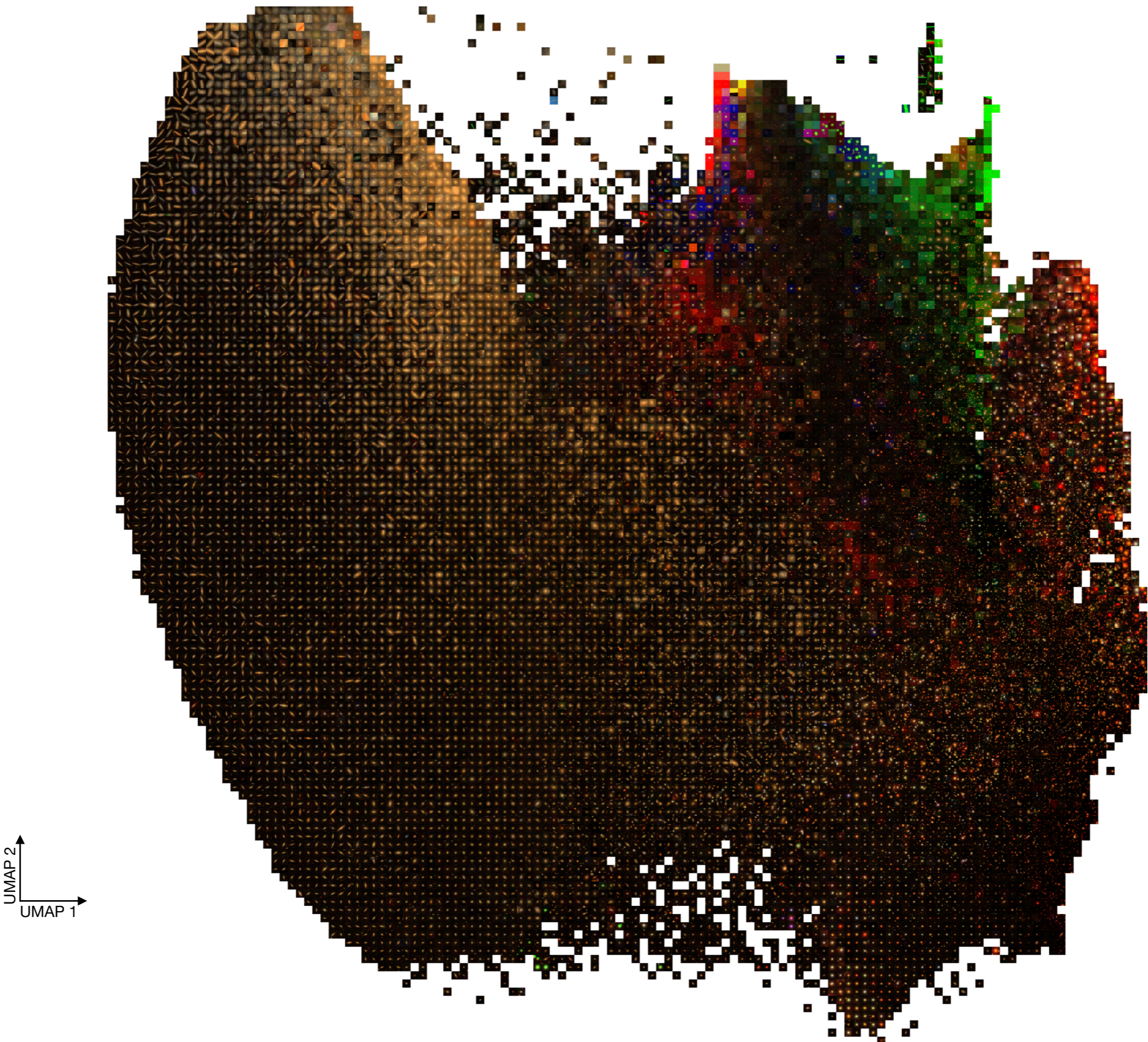
2. Add in labels, and use for downstream tasks

Self-supervised representation learning for astronomical images                                    arxiv: 2012.13083

# Visualize learned representations

Self-supervised = no labels



Representations

2048 dim

UMAP | 1802.03426

2 dim

UMAP 2

UMAP 1

Every data point represents a 5-band galaxy image

Self-supervised representation learning for astronomical images

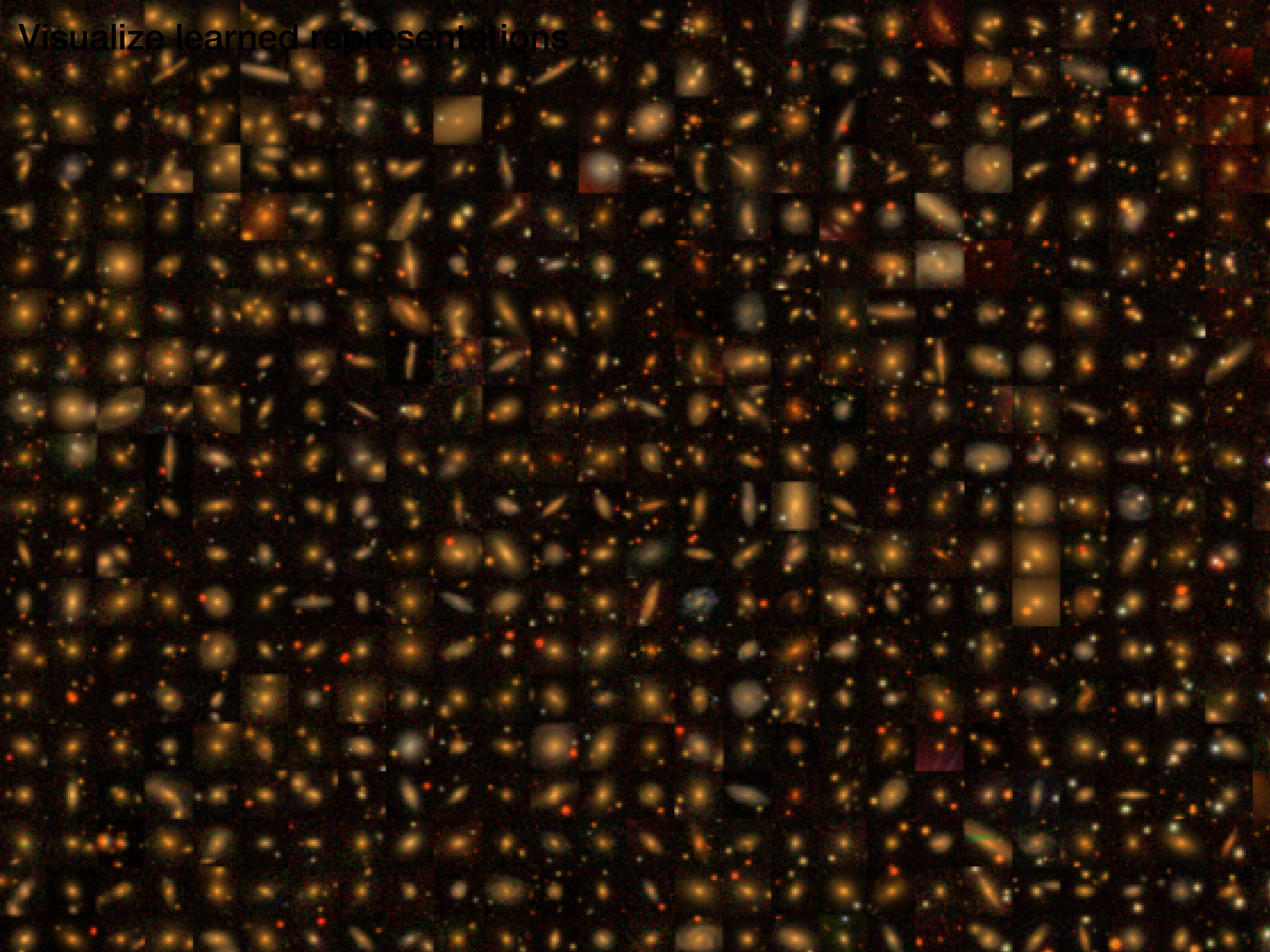# Visualize learned representations

Self-supervised representation learning for astronomical images

Visualize learned representations

Visualize learned representations

# Visualize learned representations

In context of labels



Presence of features or a disk?
● smooth  ● features or disk

Disk viewed edge-on?
● face-on  ● edge-on

Is there anything odd?
● no  ● yes

What is the odd feature?
● ring  ● irregular  ● other  ● merger

Redshift
● low  ● high

Labels
● none  ● only spec-z  ● spec-z & morphology

Labels have not been used during training, yet decision boundaries can almost be drawn by eye

Self-supervised representation learning for astronomical images

arxiv: 2012.13083

# Similarity search for data discovery

Goal:          Given an image, find other "similar" images in the dataset

Approach:
1. Select desired query image
2. Compute similarity metric of query representation with all others
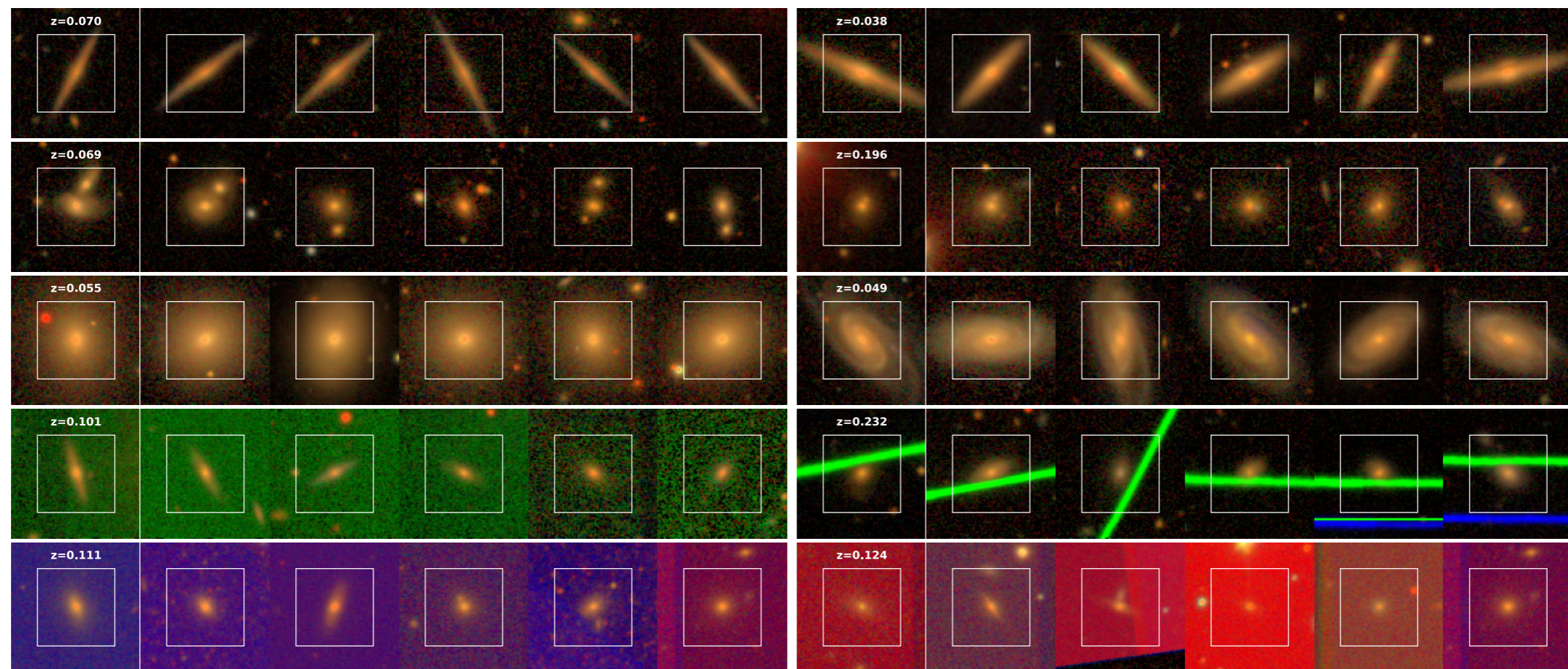3. Sort similarity by decreasing order and return images

2101.04293

# Similarity search for data discovery

Goal: Given an image, find other "similar" images in the dataset

1. Select desired query image

Approach: 2. Compute similarity metric of query representation with all others

3. Sort similarity by decreasing order and return images



Figure 1: Reference SDSS galaxies from the validation set (leftmost panels with redshift labels) and the most similar galaxies from the training set (following 5 panels) identified through a self-supervised similarity search. White squares outline the $64^2$ pixels that are "seen" by the network. 2101.04293

As the contrastive self-supervised loss was similarity-based,
the corresponding representations are by construction organized by their visual similarity

# Morphology classification

Goal: predict Galaxy Zoo 2 crowd sourced morphological labels directly from representations.
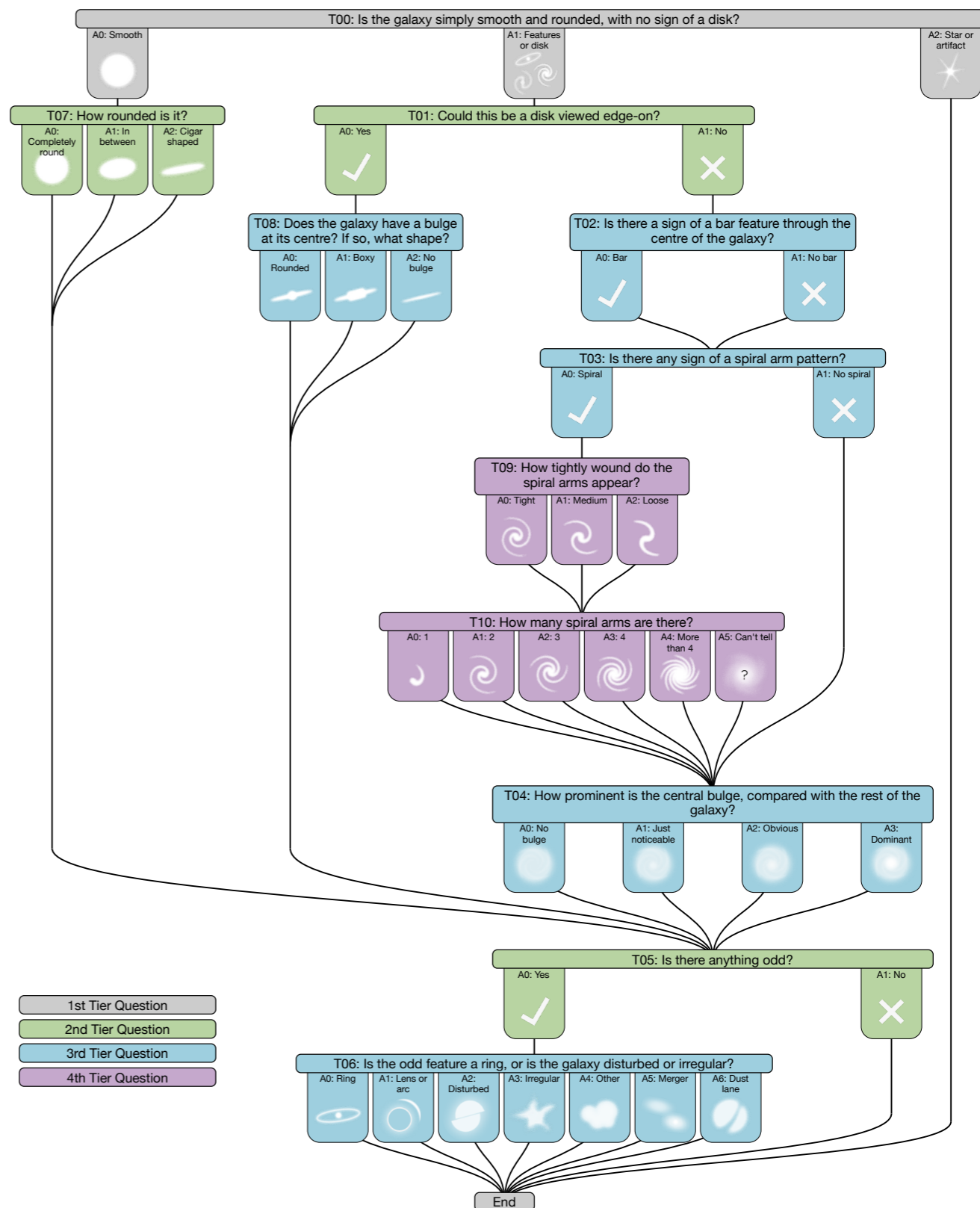
Q: "Does the galaxy have features or a disk?"
A: "smooth" or "features/disk"

Q: "Could this be a disk viewed edge-on?"
A: "Face-on" or "Edge-on"

Approach:

Linear layer from 2048 dimensions to 1,
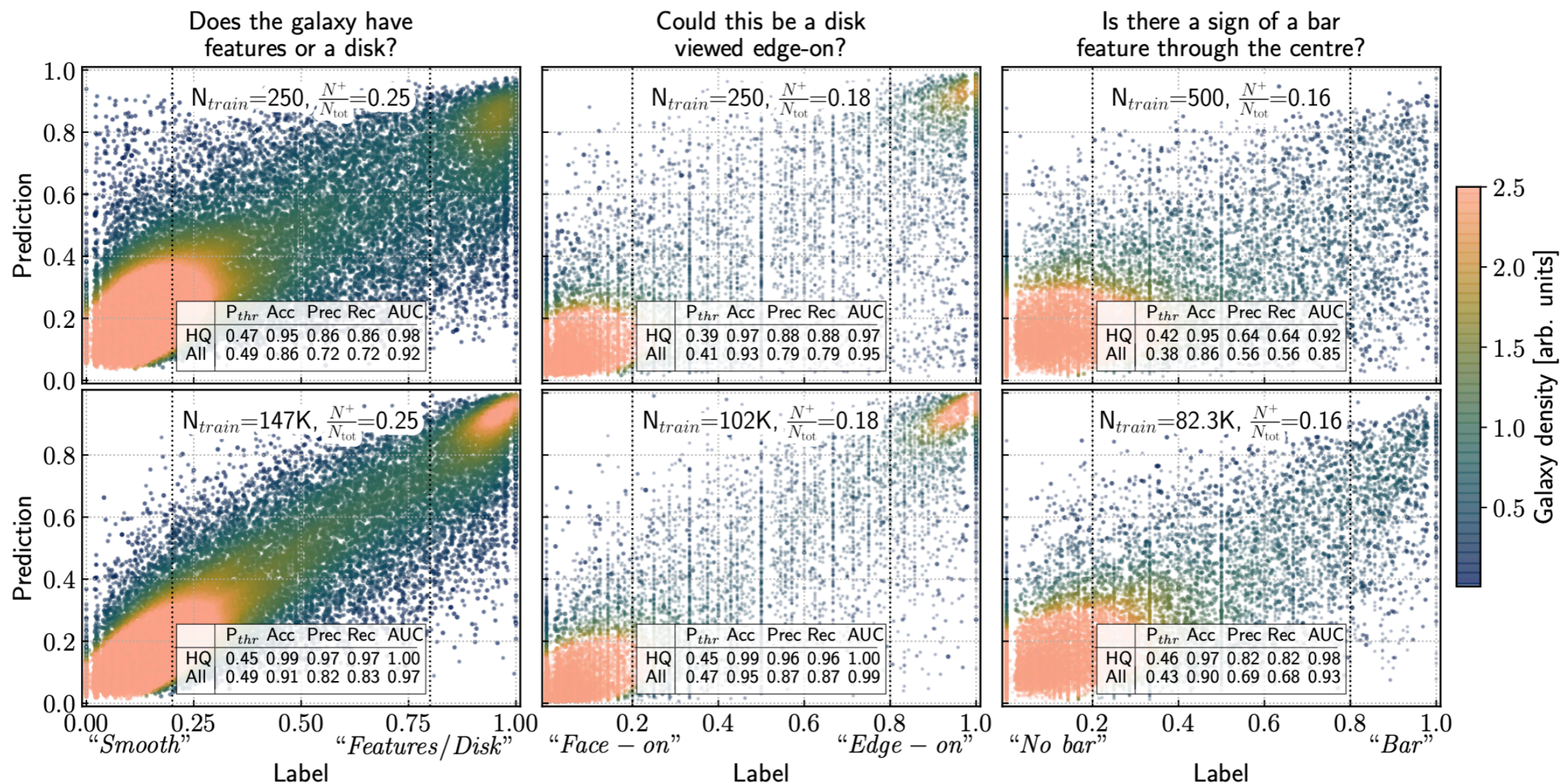followed by sigmoid activation.
Trains in 0.5-10 seconds on a GPU

Complex decision tree



1st Tier Question
2nd Tier Question
3rd Tier Question
4th Tier Question

Self-supervised representation learning for astronomical images

# Morphological classification

Goal: predict Galaxy Zoo 2 crowd sourced morphological labels directly from representations.

"Soft" labels, with high degree of label uncertainty and mislabelling



With highly limited number of training samples achieve high classification performance
Using full set of labels achieve state-of-the-art
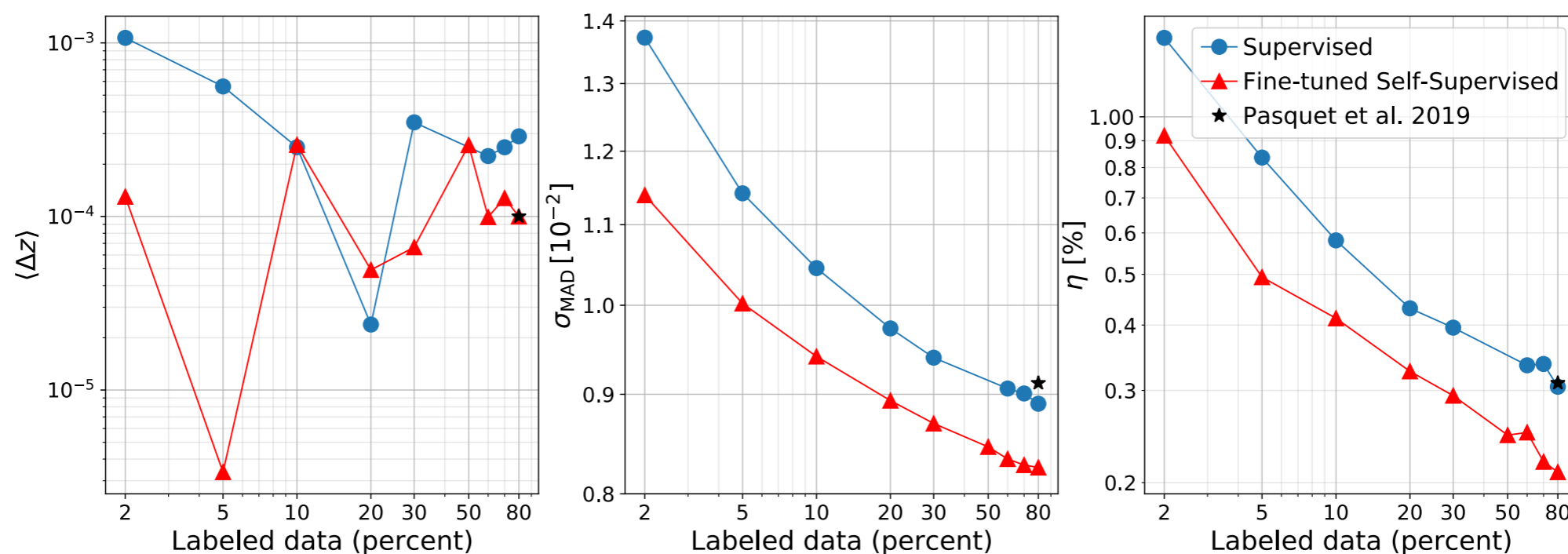No need to train separate networks for each classification task

# Photometric redshift prediction

Goal: Given an image, predict the redshift of the galaxy. ~500,000 labels from spectroscopic followup

Approach: Linear layer, trained as a classifier over a discrete set of 180 redshift bins spanning $0 < z < 0.4$
Fine-tune the encoder during training, using a small learning rate

- The prediction residual $\Delta z = (z_p - z_s)/(1 + z_s)$, where $z_p$ and $z_s$ correspond to the photometric and spectroscopic redshifts, respectively.

- The dispersion or MAD deviation, $\sigma_{\text{MAD}} = 1.4826 \times \text{MAD}(\Delta z)$, where $\text{MAD} = \text{median}(|\Delta z - \text{median}(\Delta z)|)$.

- $\eta$, the percent of "catastrophic" outliers with $|\Delta z| > 0.05$.



Self-supervised pre-training results in accuracy equivalent to 2-4x more labelled samples over supervised learning

Self-supervised representation learning for astronomical images     14     arxiv: 2012.13083

# Summary

**Self-supervised representation learning:**

- yields notable performance gains over supervised learning for multiple tasks, using the same network

- allows for accurate classifications with highly-limited number of training samples

- representations provide a rich avenue for data discovery

- representations allow for similarity search to pull out similar objects (one-shot anomaly detection)

# Future avenues

- Automated anomaly detection, and how this can be done on the representations. Density based models?

- Robustness quantification. Can we use this to determine in-distribution/out-of-distribution for unlabelled samples?
  *Using Self-Supervised Learning Can Improve Model Robustness and Uncertainty*
- Larger model trained on much more data, and serve to the community, much like the operation of existing state-of-the-art language models
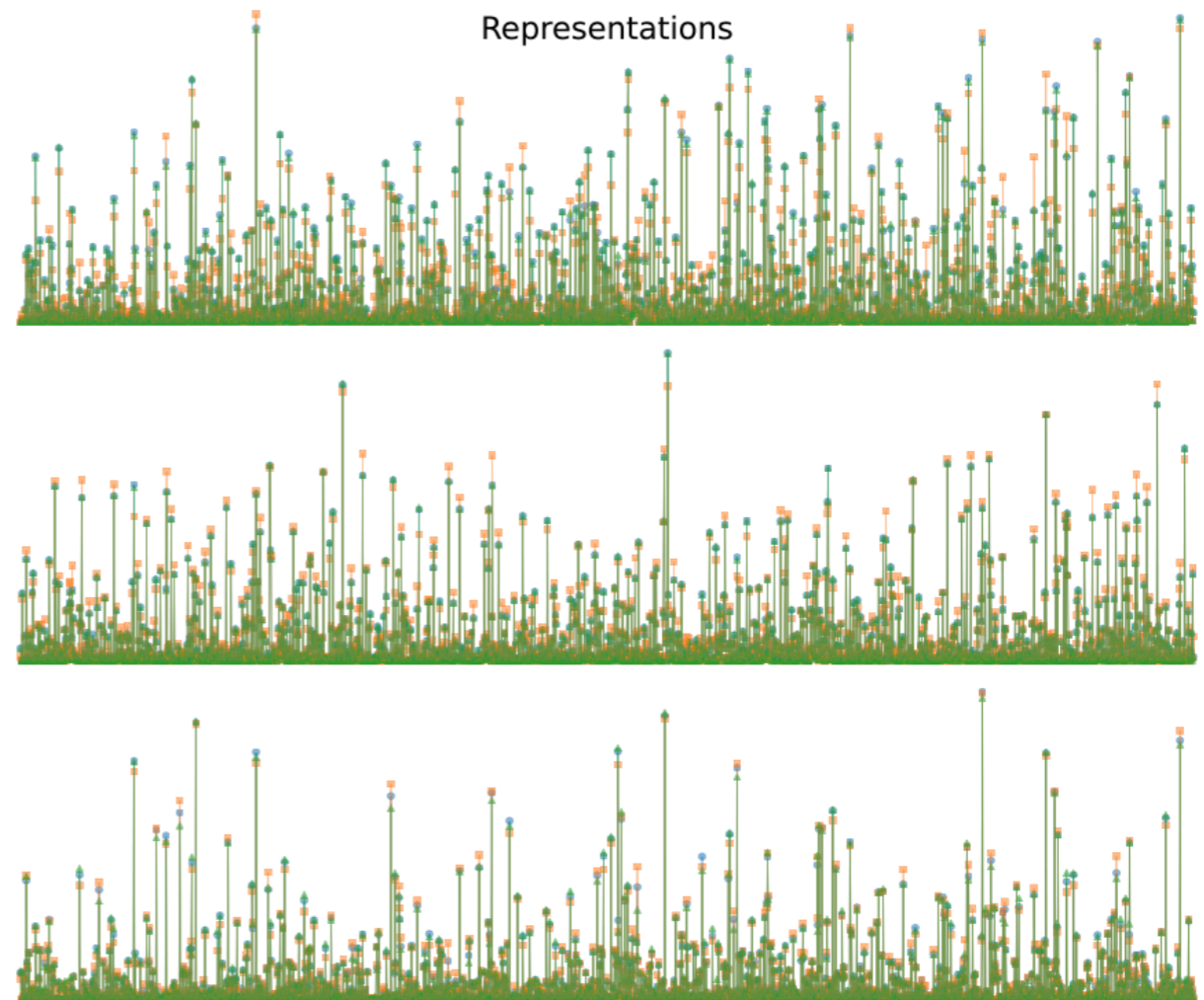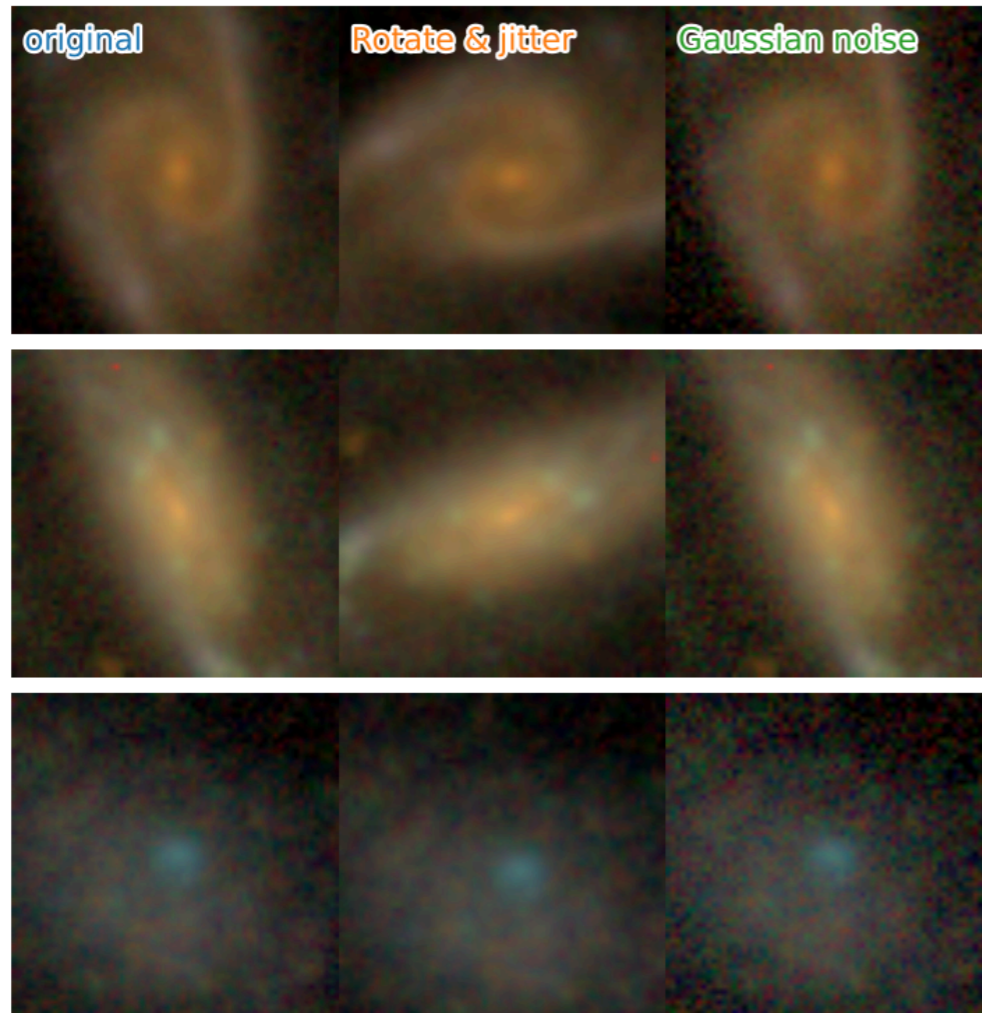
# Sketch for self-supervised learning in other fields

1. For given experiment, take all data with or without labels. Can be 1D, 2D, 3D, …

2. Construct data augmentations that reflect changes in the data you want network to be agnostic to. (don't choose ones that effect semantic information of data for desired downstream tasks)

   - Rotations about certain axes
   - Jitter
   - Masking
   - Smoothing
   - Various types of noise
   - Scaling

3. Learn representations through self-supervised contrastive framework

4. Use representations for downstream tasks

   - Data discovery
   - Anomaly detection/similarity search
   - Classification tasks
   - Regression tasks

5. Improve upon what is capable in a supervised framework

examples:
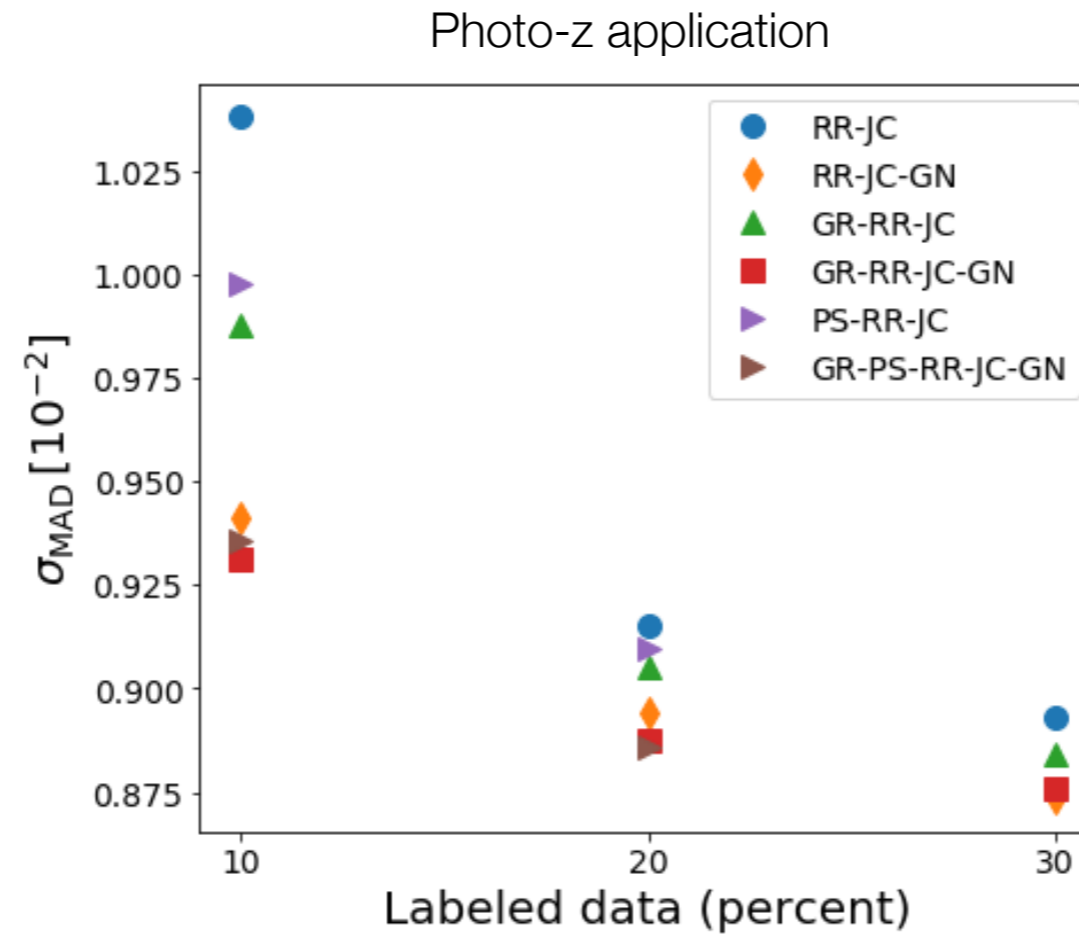*Big Self-Supervised Models Advance Medical Image Classification*

# Extra Slides

## Representations after augmentations



Representations close to invariant under different augmentations

# Extra Slides

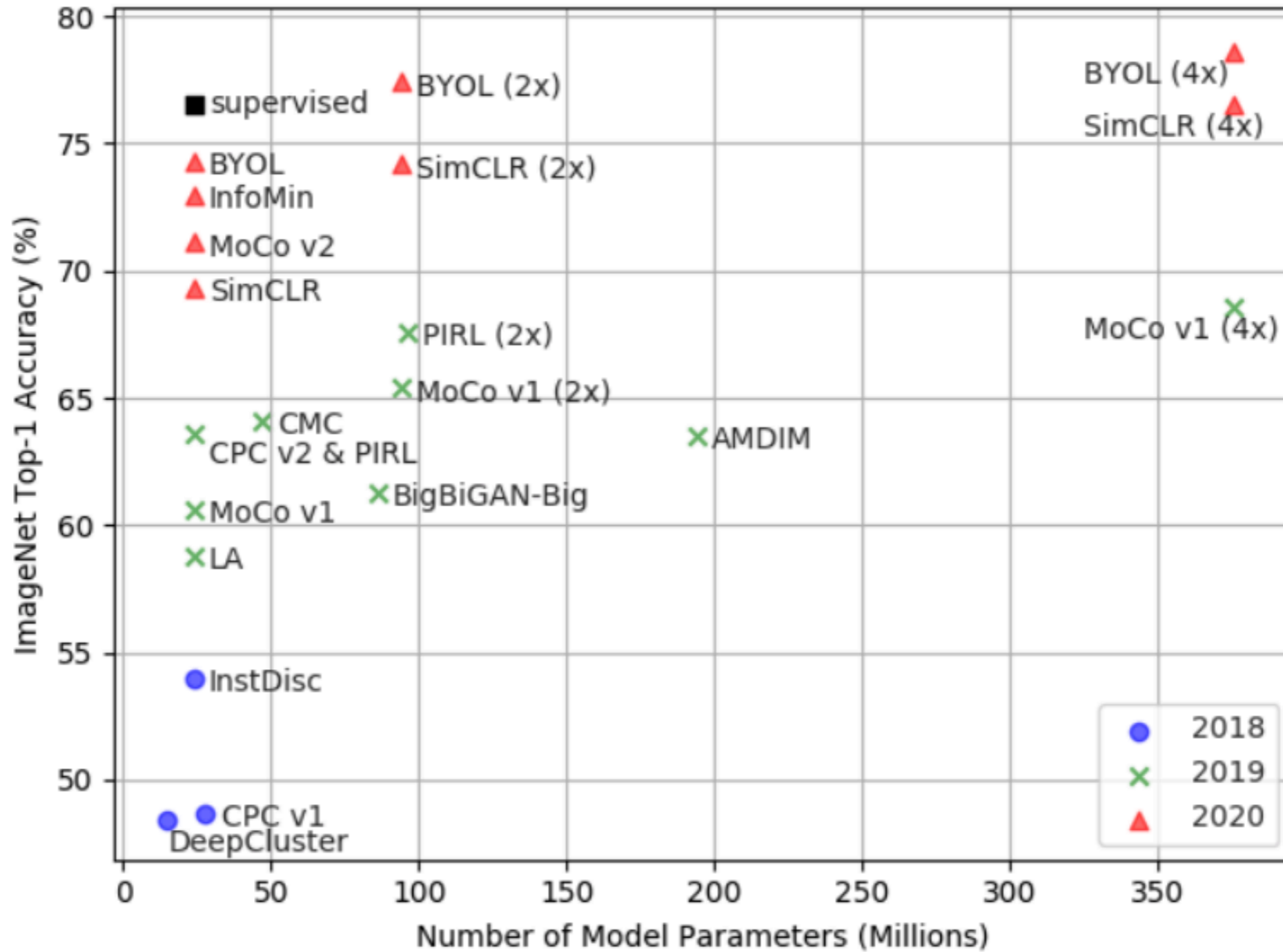## Ablation study

Which augmentations are most powerful?



Photo-z application

RR: Random rotate
JC: Jitter-crop
GN: Gaussian Noise
GR: Galactic reddening
PS: PSF smoothing

Self-supervised representation learning for astronomical images
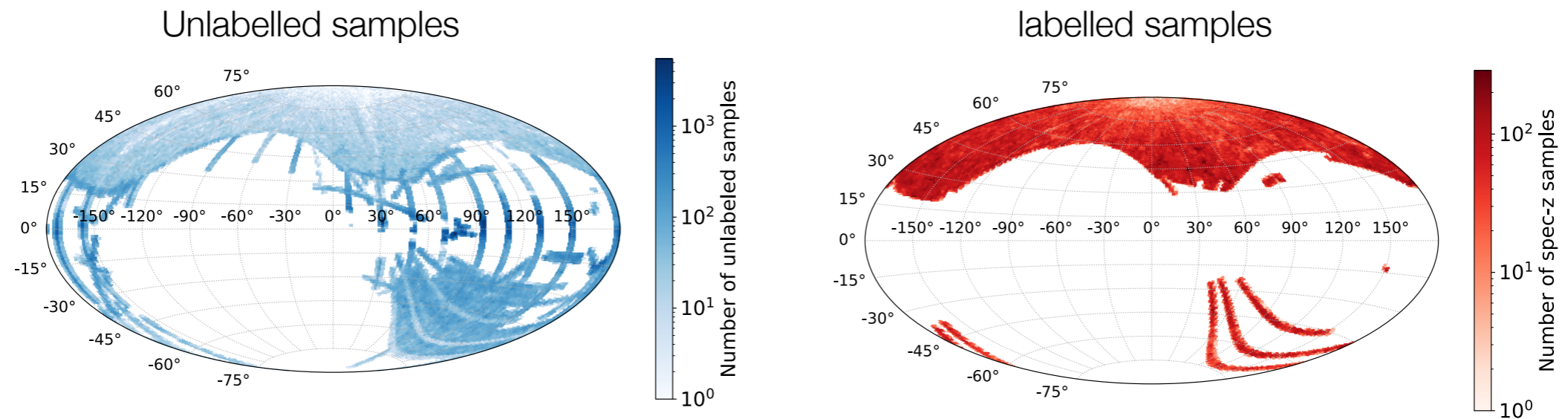
# Extra Slides



## Self-supervised learning frameworks
Self-supervised representation learning performance on ImageNet top-1 accuracy in June, 2020, under linear classification protocol.

from *Self-supervised Learning: Generative or Contrastive*

Self-supervised representation learning for astronomical images

arxiv: 2012.13083

# Extra Slides

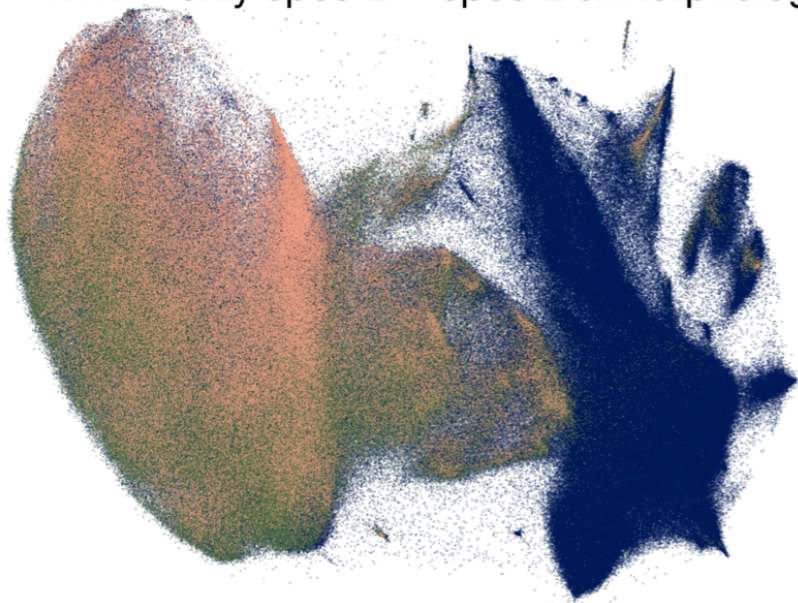## Robustness

Unlabelled samples



labelled samples



A large number of unlabelled samples are from different parts of the sky than those with labels

By design, as (redshift) labels were obtained only for the cleanest images in sky surveys

### Labels
● none  ● only spec-z  ● spec-z & morphology



This is apparent in the learned representation space, i.e. there are unlabelled samples that are not "near" any with labels
(note there are still tonnes of blue dots hidden under green and pink ones)

-> Use for robustness metric of inference on samples that do not have labels:
Unlabelled samples "near" ones with labels = trust inference
Unlabelled samples far from any with labels = don't trust inference

"Near" can be linear classification, density based, iD/OoD, …

Self-supervised representation learning for astronomical images