Aleksandra Ćiprijanović
(she/her/hers)

Fermilab, DSSL
aleksand@fnal.gov

**Towards flexible domain adaptation methods for cross-datasets studies of galaxies**

KITP
March, 2023

# Vision of the Future



~2023.

**Rubin LSST**

~ 20 TB / day
~ 100 PB total by DR11



- **Real-time:**
  - data handling,
  - decision making
  - detection of interesting events
  - inference
- **Automated experiments**
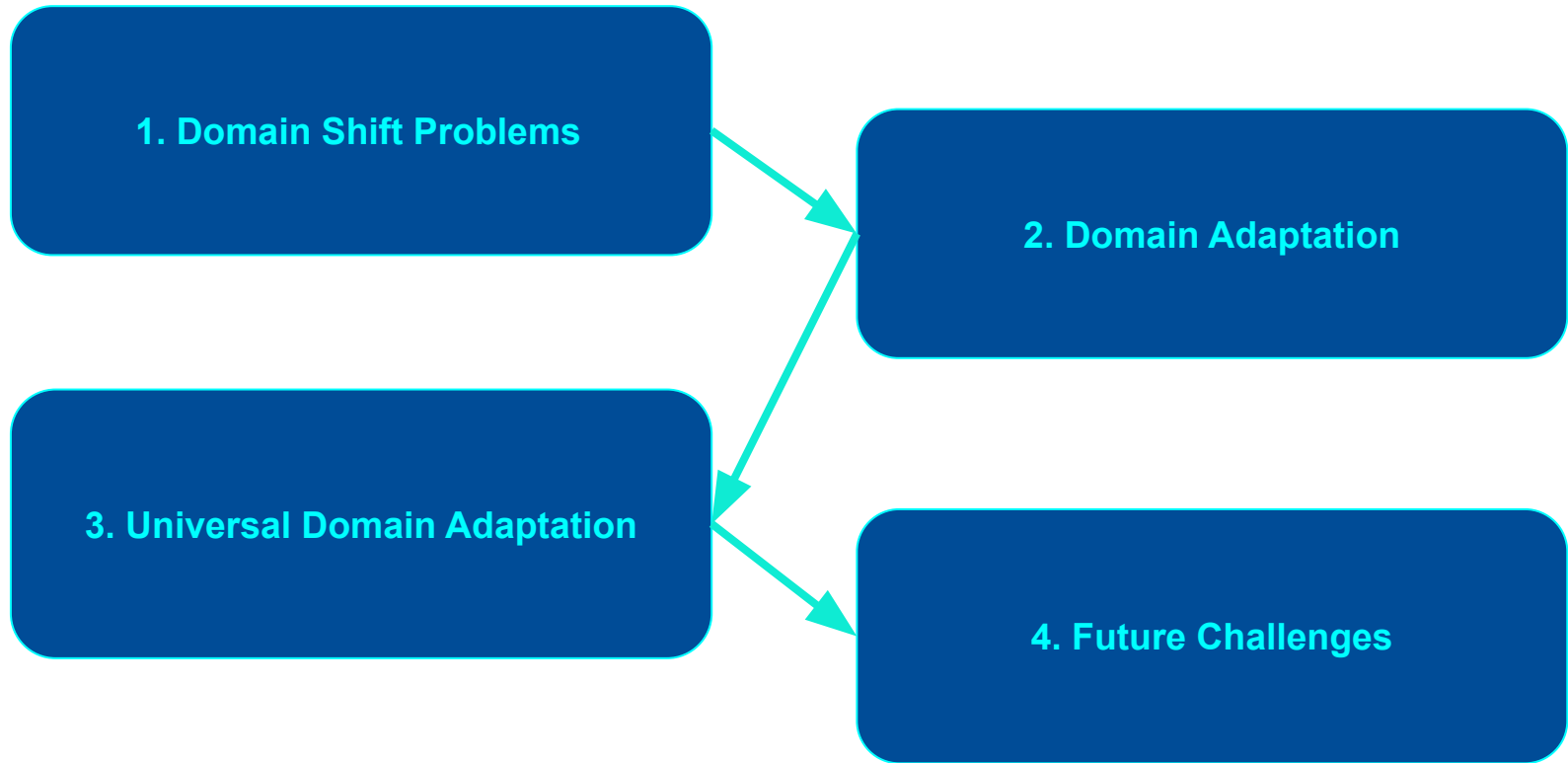- **Working with big data** later in the process



**Often research starts from simulations…**
**…but when AI gets involved…**
**…we face some challenges.**

🟦 **Fermilab**

# Talk Outline

1. Domain Shift Problems

2. Domain Adaptation

3. Universal Domain Adaptation

4. Future Challenges

**Fermilab**

# Talk Outline

1. Domain Shift Problems

2. Domain Adaptation

3. Universal Domain Adaptation

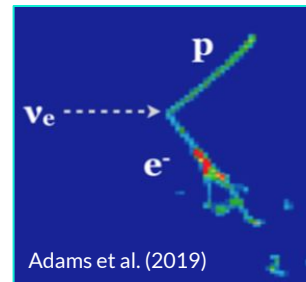4. Future Challenges

Fermilab

# Combining Datasets

All areas of Fermilab science often need to create **model trained on simulated data, that also work on real detector data**!
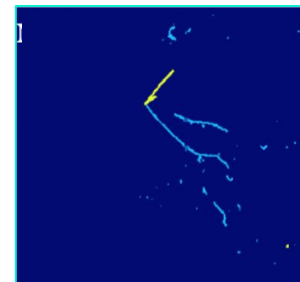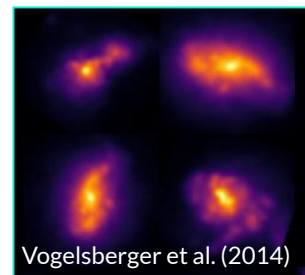
**DATASET SHIFT**



SIMULATED — REAL

MicroBooNE (neutrinos)

Adams et al. (2019)

Illustris / Hubble (merging galaxies)

Vogelsberger et al. (2014)
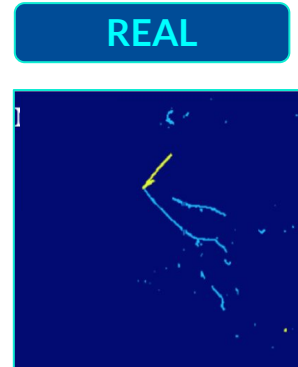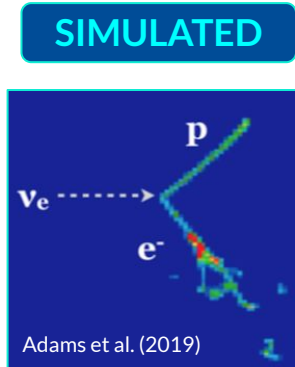
Hubble

# Combining Datasets

All areas of Fermilab science often need to create **model trained on simulated data, that also work on real detector data**!
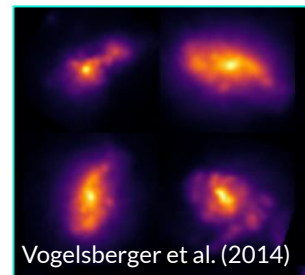
**DATASET SHIFT**

Missing and unknown physics, wrong geometry, background levels

Computational constraints for simulations



SIMULATED

REAL

MicroBooNE (neutrinos)

Adams et al. (2019)

Illustris / Hubble (merging galaxies)

Vogelsberger et al. (2014)

Hubble

🌊 **Fermilab**

# Combining Datasets

All areas of Fermilab science often need to create **model trained on simulated data, that also work on real detector data**!
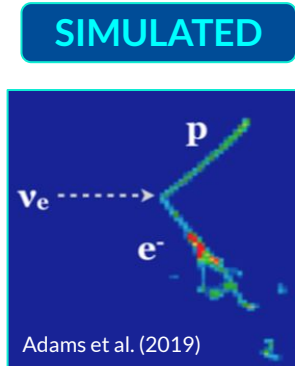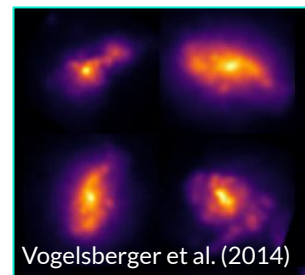
## DATASET SHIFT

| | |
|---|---|
| Missing and unknown physics, wrong geometry, background levels | Computational constraints for simulations |
| Detector problems, transients, errors, data compression | Imperfect addition of observational effects |



SIMULATED · REAL

MicroBooNE (neutrinos)

Adams et al. (2019)

Illustris / Hubble (merging galaxies)

Vogelsberger et al. (2014)   Hubble

‡ Fermilab

# Combining Datasets

All areas of Fermilab science often need to create **model trained on simulated data, that also work on real detector data**!

## DATASET SHIFT

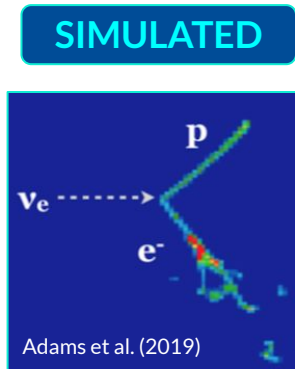Missing and unknown physics, wrong geometry, background levels

Computational constraints for simulations

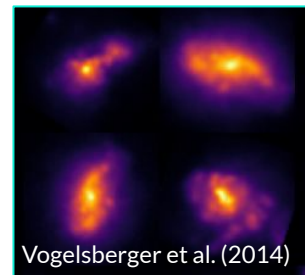Detector problems, transients, errors, data compression

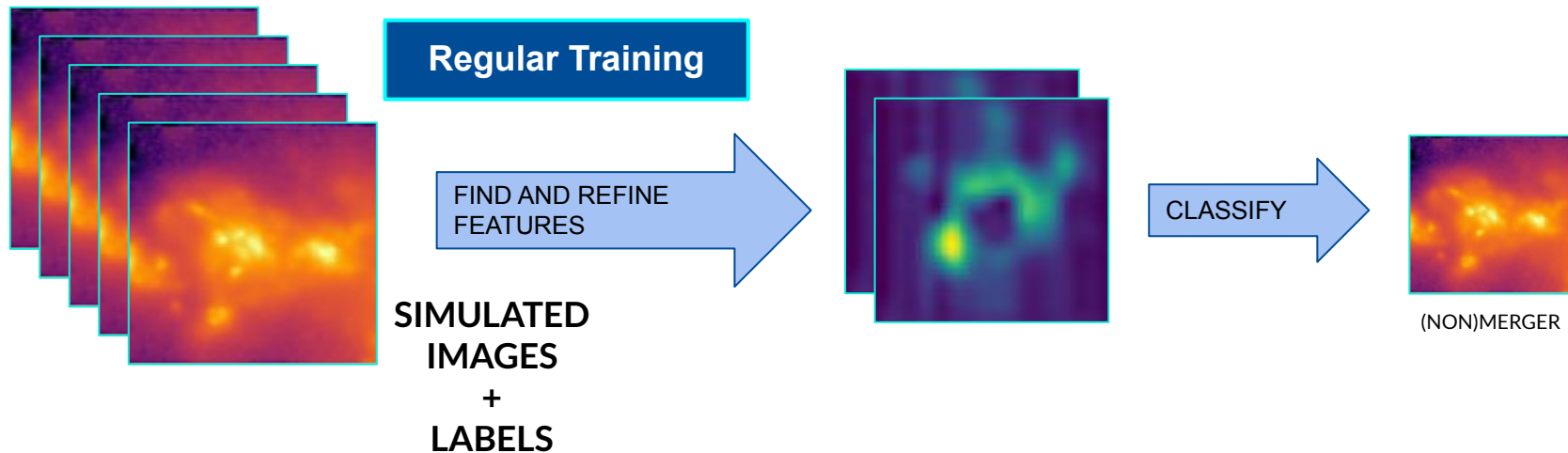Imperfect addition of observational effects
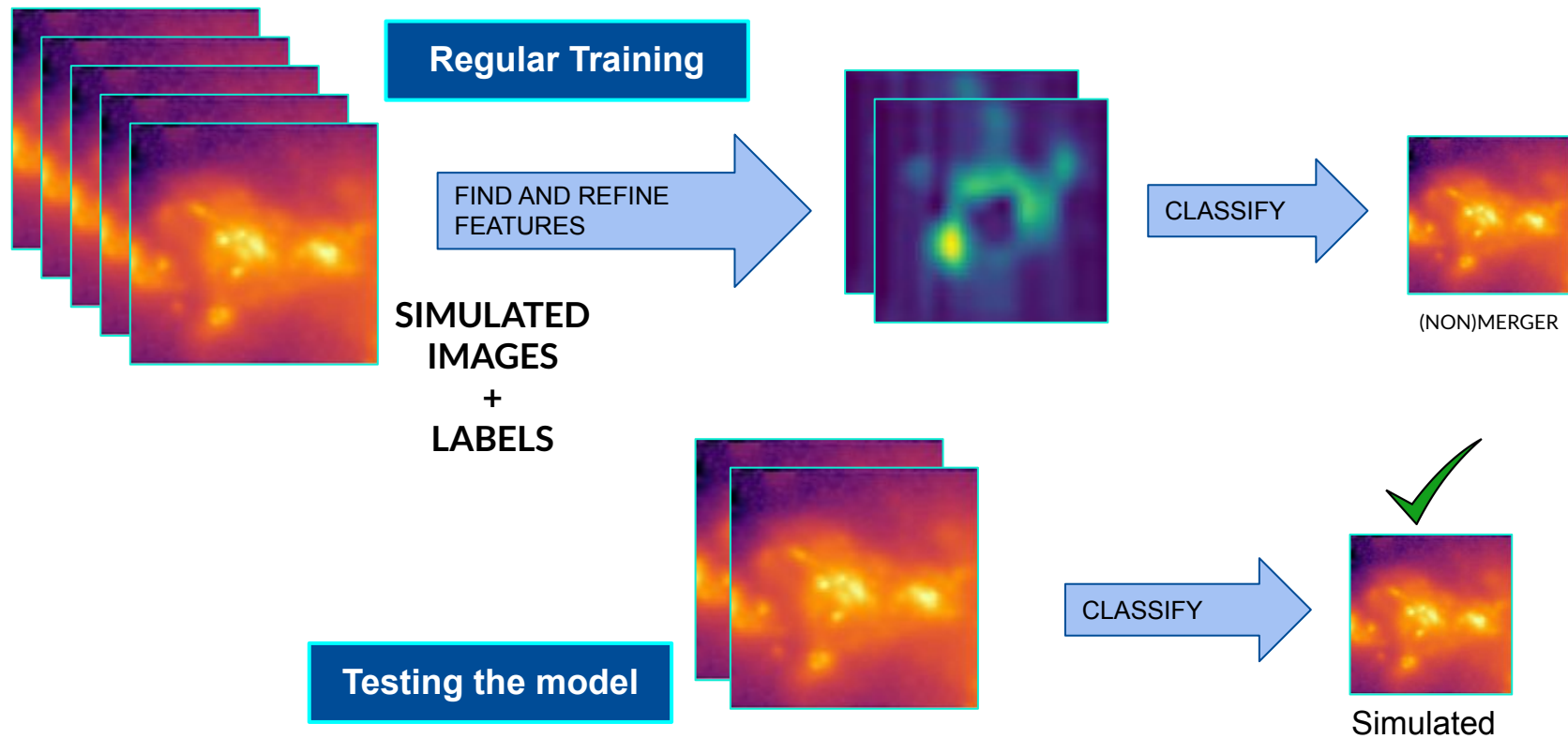
Different detectors or telescopes

**SIMULATED**

**REAL**

MicroBooNE (neutrinos)



Adams et al. (2019)

Illustris / Hubble (merging galaxies)



Vogelsberger et al. (2014)

Hubble

🔷 **Fermilab**

# Combining Datasets



**Regular Training**

FIND AND REFINE FEATURES

**SIMULATED
IMAGES
+
LABELS**

CLASSIFY

(NON)MERGER

🎗️ **Fermilab**

# Combining Datasets

**Regular Training**

FIND AND REFINE FEATURES

**SIMULATED IMAGES + LABELS**

CLASSIFY

(NON)MERGER

**Testing the model**

CLASSIFY

Simulated

Fermilab

# Combining Datasets



**Regular Training**

FIND AND REFINE FEATURES

CLASSIFY

(NON)MERGER

SIMULATED
IMAGES
+
LABELS

**Testing the model**

CLASSIFY

Simulated   Observed

🎗 Fermilab

# Combining Datasets

**Why does this happen?**

**Fermilab**

# Combining Datasets

**Why does this happen?**

Source Domain



Train the model on source dataset and find the decision boundary.

Fermilab

# Combining Datasets

## Why does this happen?



New domain is shifted, learned decision boundary doesn't work.
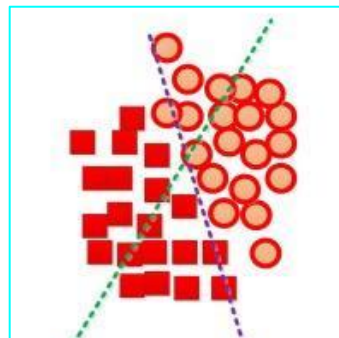
Source Domain

Target Domain

🎔 Fermilab

# Combining Datasets
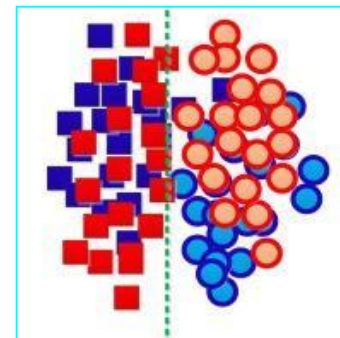
## Why does this happen?

We need to align the data during training!

| Source Domain | Target Domain | Domain Alignment |

Fermilab

# Talk Outline

1. Domain Shift Problems

2. Domain Adaptation

3. Universal Domain Adaptation

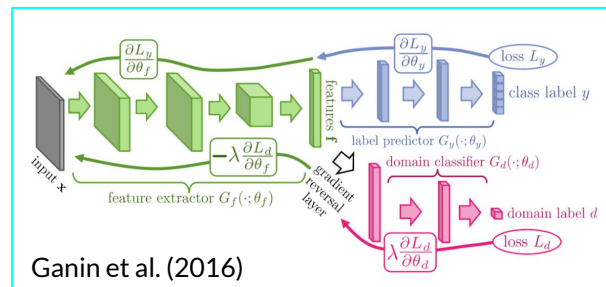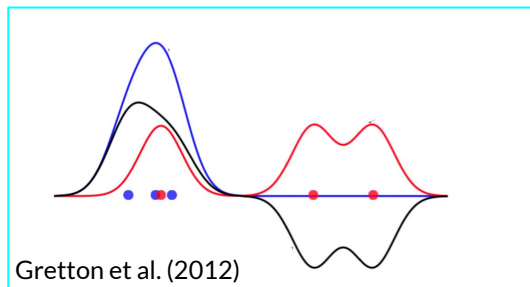4. Future Challenges

🎇 Fermilab

# Combining Datasets

## DOMAIN ADAPTATION

Align data distributions in the latent space of the network by forcing the network to **find more robust domain-invariant features.**

**Fermilab**
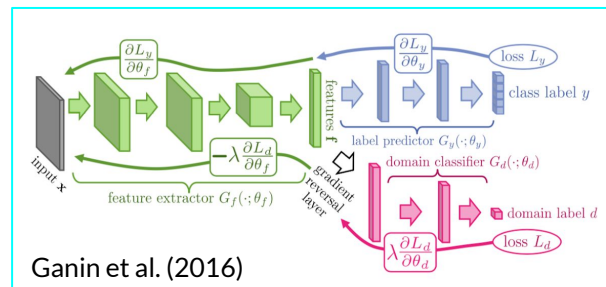
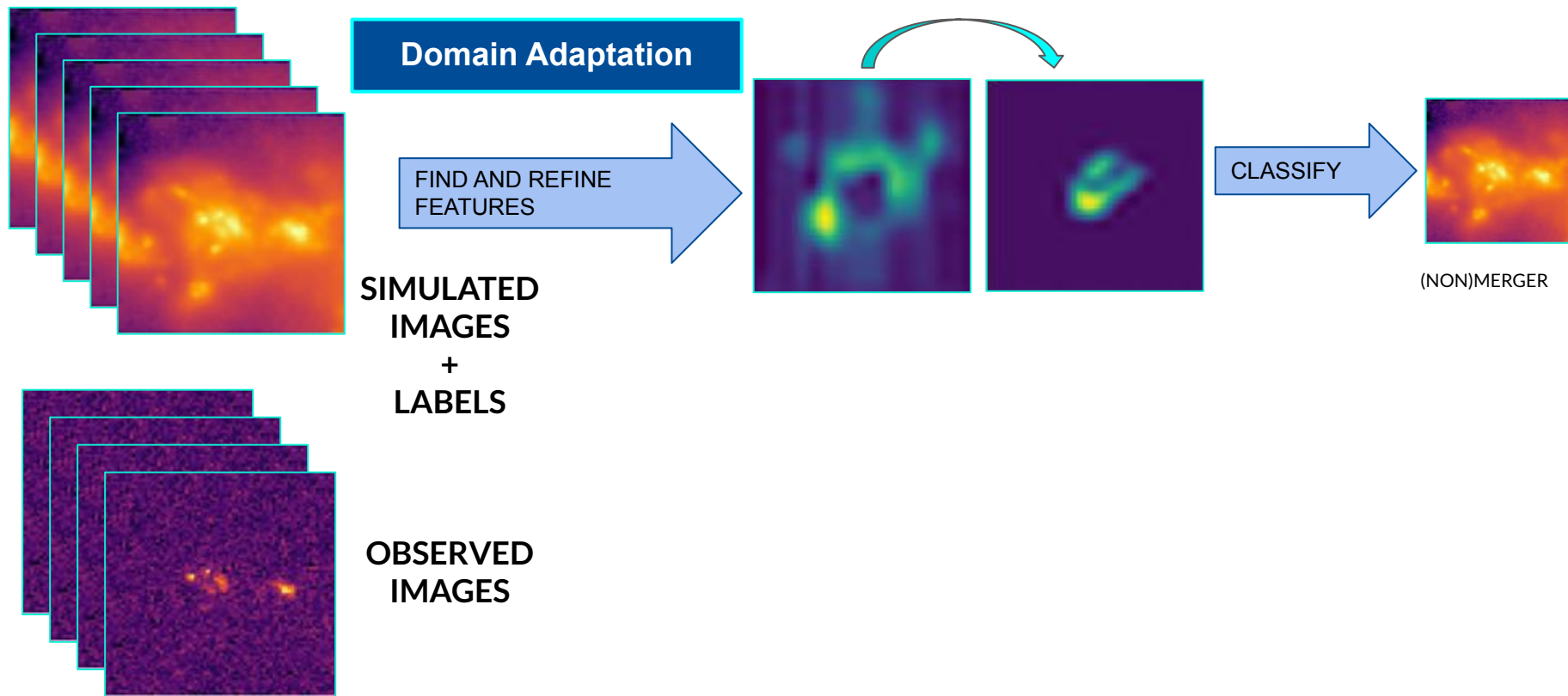# Combining Datasets

**DOMAIN ADAPTATION**

Align data distributions in the latent space of the network by forcing the network to **find more robust domain-invariant features.**
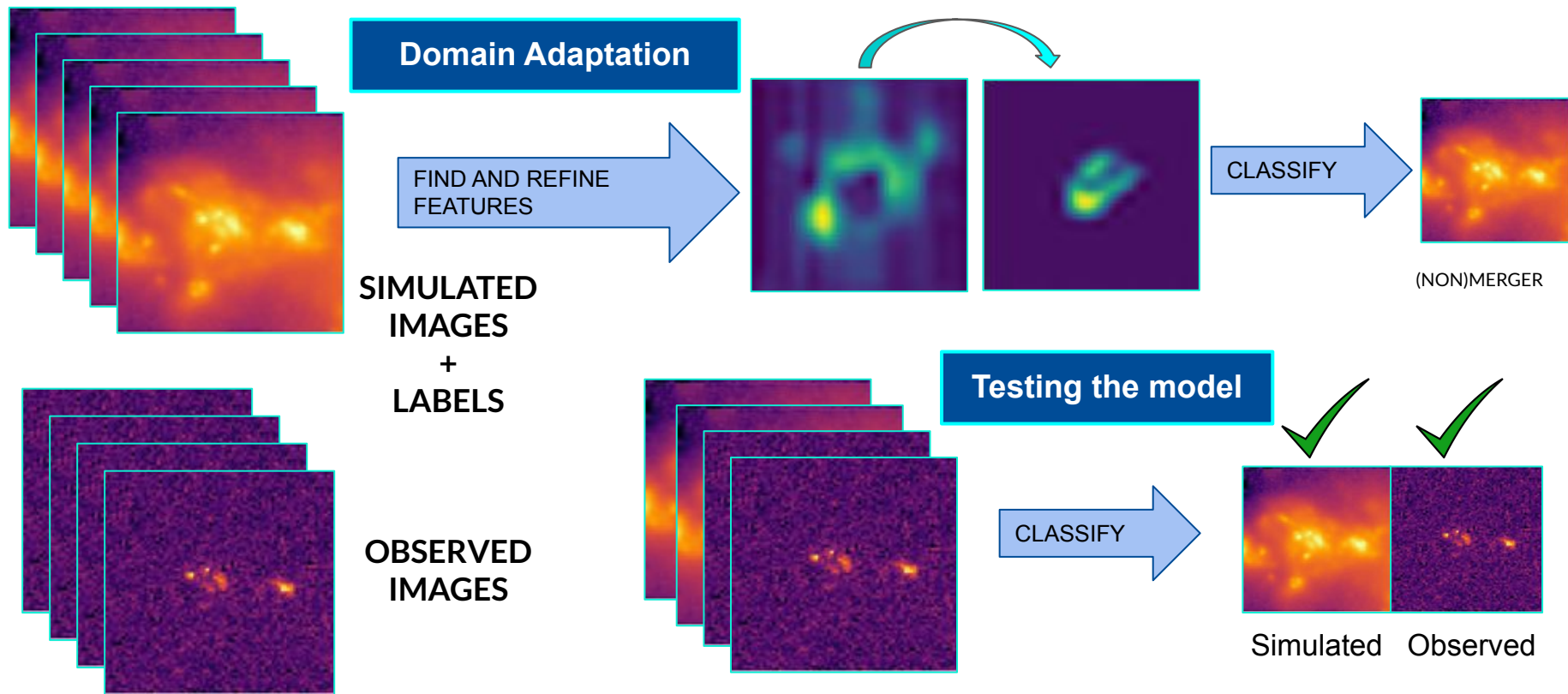
**Distance-based methods**

**Adversarial methods**



Gretton et al. (2012)



Ganin et al. (2016)

🦬 **Fermilab**

# Combining Datasets

**DOMAIN ADAPTATION**

Align data distributions in the latent space of the network by forcing the network to **find more robust domain-invariant features.**

Distance-based methods

Adversarial methods

Training
=
Task Loss
+
DA Loss


Gretton et al. (2012)


Ganin et al. (2016)

🎄 **Fermilab**

# Combining Datasets

**DOMAIN ADAPTATION**

Align data distributions in the latent space of the network by forcing the network to **find more robust domain-invariant features.**

Distance-based methods

Adversarial methods

Works on **unlabeled target domain**!
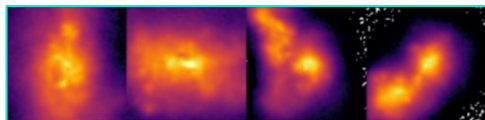Can be applied to **new data**, no need for scientists to label anythin
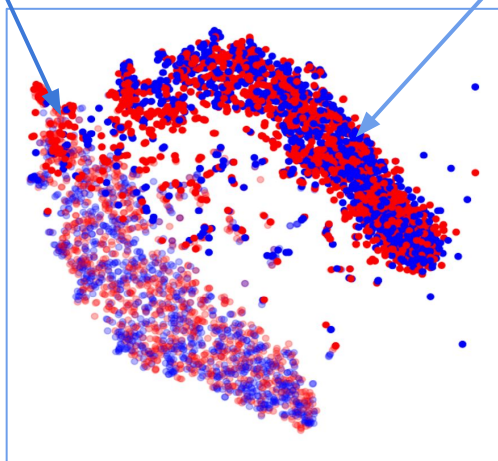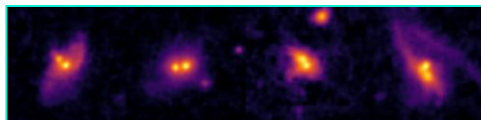
🐝 **Fermilab**

# Combining Datasets

# Combining Datasets



Domain Adaptation

FIND AND REFINE FEATURES

CLASSIFY

(NON)MERGER

SIMULATED IMAGES + LABELS

OBSERVED IMAGES

Testing the model

CLASSIFY

Simulated    Observed

Fermilab

# Combining Datasets
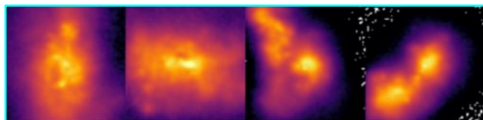
**Source - Illustris**　　**Target - SDSS observations**



This is how the network sees the data.
2D representation of network's latent space.

Ćiprijanović et al. 2020b.
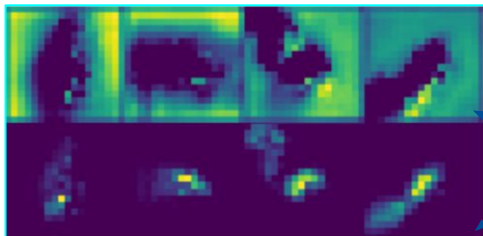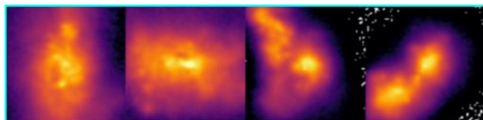Ćiprijanović et al. 2021a.

Fermilab

# Combining Datasets

**Source - Illustris**

M

NM

Important regions are highlighted!

**Regular Training**

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 202

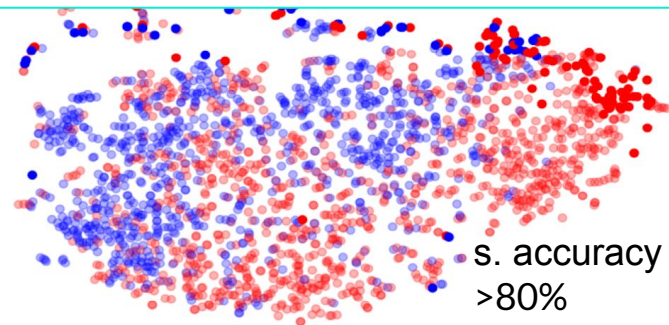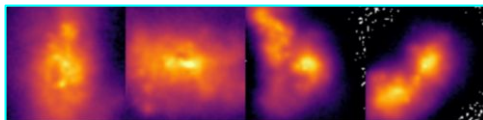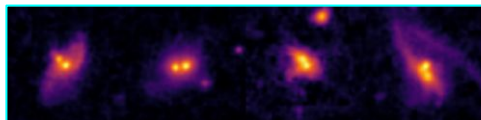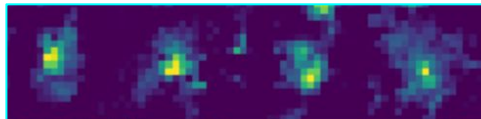🎗️ **Fermilab**

# Combining Datasets

**Source - Illustris**



M

NM

Important regions are highlighted!

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 202

**Regular Training**



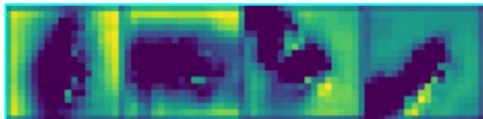s. accuracy >80%

🎶 **Fermilab**

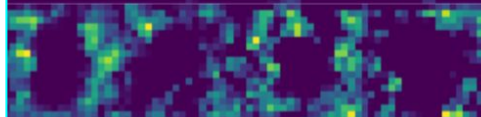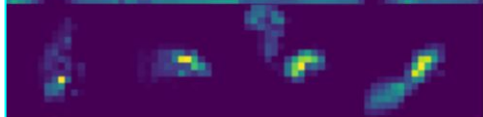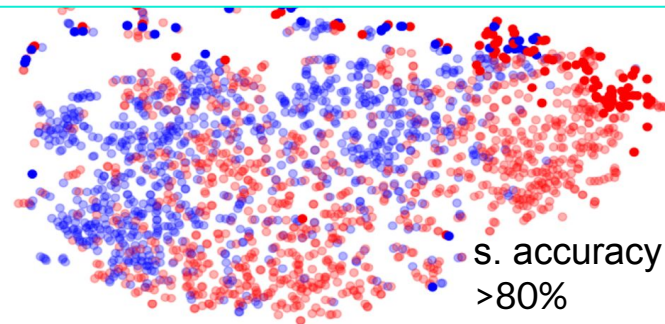# Combining Datasets

**Source - Illustris**

**Target - SDSS observations**



M

NM

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 202

**Regular Training**



s. accuracy
>80%

Fermilab

# Combining Datasets

**Source - Illustris**

**Target - SDSS observations**

M

NM



t. accuracy ~50%

s. accuracy >80%

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 202

🔶 **Fermilab**

# Combining Datasets

**Source - Illustris**

**Target - SDSS observations**

M

NM

M

NM

**Domain Adaptation**

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 202

🟰 Fermilab
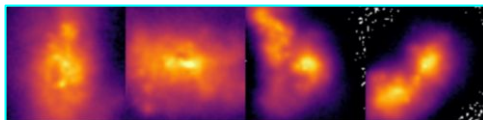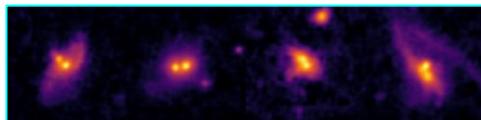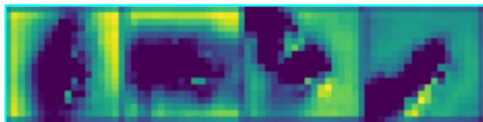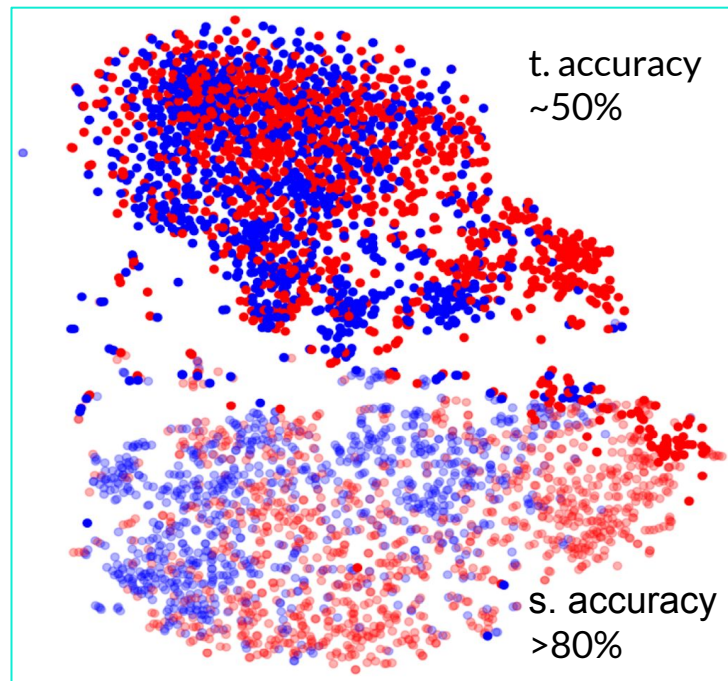
# Combining Datasets



Source - Illustris    Target - SDSS observations

M

NM

M

NM

Domain Adaptation

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 202

🪒 Fermilab

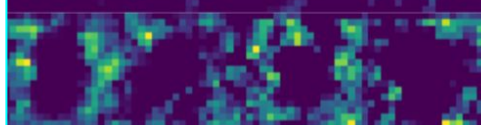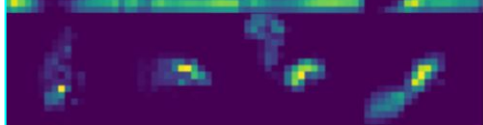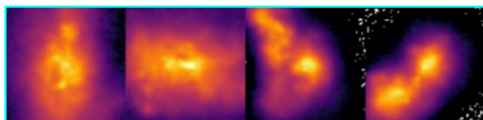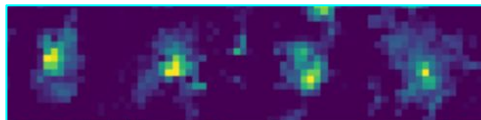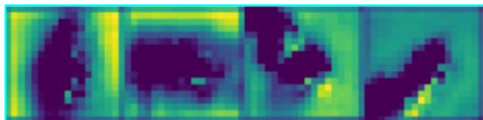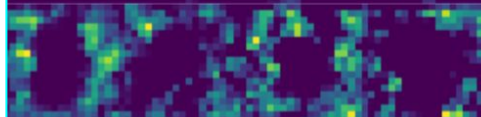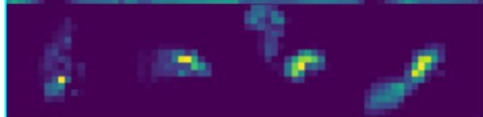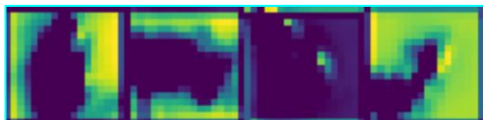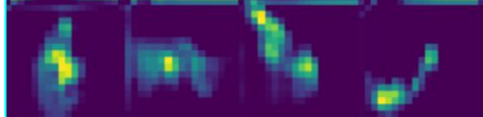# Combining Datasets

**Source - Illustris**　　　**Target - SDSS observations**



M

NM

M

NM

Ćiprijanović et al. 2020b.
Ćiprijanović et al. 2021a

Up to 30% increase!



t. accuracy
~80%

s. accuracy
~90%

Fermilab

# Talk Outline

1. Domain Shift Problems

2. Domain Adaptation

3. Universal Domain Adaptation

4. Future Challenges

Fermilab

# Bridging between observations - Much Harder!

The gap between observational datasets is much larger:
- Noise, PSF
- Pixel scale
- Depth of the survey
- Magnitude limit
- Perhaps different filters
- Different data distributions….

How do we build something flexible enough to handle any kind of data distributions and distribution overlaps?



SDSS to DECaLS?

**⚛ Fermilab**

# Types of Dataset Shift Problems

- Overall distribution per class can be different between datasets.
  - Overlapping classes should be aligned independently instead of aligning the entire data distribution.
- We can even have classes present in only one of the datasets - old labeled data or even new unlabeled data (so we won't even know it's there!)
  - Non-overlapping classes should not be aligned with anything.



source   target

Closed

Partial

Open

Open-Partial

Fermilab

# Universal Domain Adaptation (DeepAstroUDA)

Classification of known classes

**+**

Clustering of similar known and unknown samples

**+**

Separation of different (anomalous) unknown samples

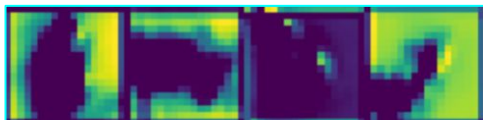Ćiprijanović et al. 2022.
Ćiprijanović et al. 2023.



source    target

Entropy Separation
Adaptive Clustering
Cross-Entropy Clustering

Fermilab

# Universal Domain Adaptation (DeepAstroUDA)

Classification of known classes

**+**

Clustering of similar known and unknown samples

**+**

Separation of different (anomalous) unknown samples
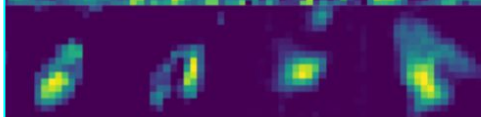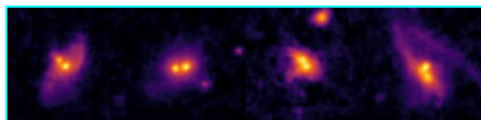
Ćiprijanović et al. 2022.
Ćiprijanović et al. 2023.

source   target

Entropy Separation
Adaptive Clustering
Cross-Entropy Clustering

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

Output vector **p**

**Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)

Classification of known classes

$$\mathcal{L}_{CE} = \frac{-\sum_{k=1}^{K} w_k y_k \log \hat{y}_k}{\sum_{k=1}^{K} w_k},$$

Using true and predicted labels



source   target        ⇨ Entropy Separation
▲  (△)                   ⇨ Adaptive Clustering
                         ⇨ Cross-Entropy Clustering

0  1  2  3  4  5  6  7  8  9

Output vector **p** ⟹ compare predicted y' with true label y

🎅 **Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)

Clustering of similar known and unknown samples

Via self-supervision:
comparing pairs of output features
between all samples from both domains

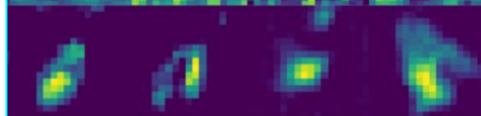| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

Output vector **p**

**Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)

Ćiprijanović et al. 2022.
Ćiprijanović et al. 2023.

Clustering of similar known and unknown samples

Via self-supervision:
comparing pairs of output features
between all samples from both domains

$$\mathcal{L}_{AC} = -\sum_{i \in B}\sum_{j \in b_t} s_{ij}\log(\mathbf{p}_i^\top \mathbf{p}_j) + (1 - s_{ij})\log(1 - \mathbf{p}_i^\top \mathbf{p}_j),$$

(1)



source    target

Entropy Separation
Adaptive Clustering
Cross-Entropy Clustering

0  1  2  3  4  5  6  7  8  9

Output vector p ➡ rank order to create similarity labels

Fermilab

# Universal Domain Adaptation (DeepAstroUDA)

Ćiprijanović et al. 2022.
Ćiprijanović et al. 2023.

Separation of different (anomalous) unknown samples

Pushing away samples with high entropy of outputs features



source    target
● ▲    (◯) (△)

⟹ Entropy Separation
⟸ Adaptive Clustering
◀ Cross-Entropy Clustering

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |

Output vector p

**Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)

Ćiprijanović et al. 2022.
Ćiprijanović et al. 2023.

Separation of different (anomalous) unknown samples

Pushing away samples with high entropy of outputs features

$$\mathcal{L}_{\text{ES}}(\mathbf{p}_i) = \begin{cases} -|H(\mathbf{p}_i) - \rho| & |H(\mathbf{p}_i) - \rho| > m, \\ 0 & \text{otherwise.} \end{cases} \qquad \mathcal{L}_{\text{ES}} = \frac{1}{|b_t|} \sum_{i \in b_t} \mathcal{L}_{\text{ES}}(\mathbf{p}_i).$$

$$H(X) = -\sum_{x \in X} p(x) \log p(x)$$

source   target

Entropy Separation
Adaptive Clustering
Cross-Entropy Clustering

0  1  2  3  4  5  6  7  8  9

Output vector p ⟹ calculate entropy of each output

🌟 Fermilab

# Universal Domain Adaptation (DeepAstroUDA)

Ćiprijanović et al. 2022.
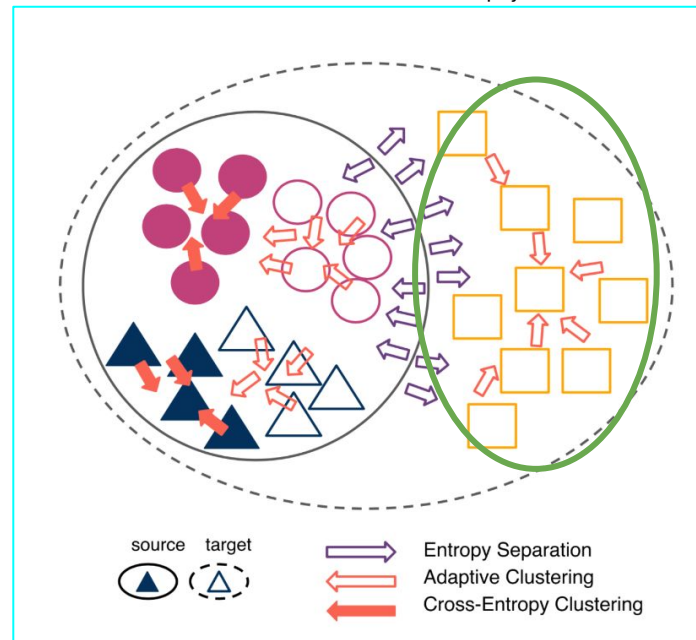Ćiprijanović et al. 2023.

Separation of different (anomalous) unknown samples

Pushing away samples with high entropy of outputs features

$$\mathcal{L}_{ES}(\mathbf{p}_i) = \begin{cases} -|H(\mathbf{p}_i) - \rho| & |H(\mathbf{p}_i) - \rho| > m, \\ 0 & \text{otherwise.} \end{cases} \qquad \mathcal{L}_{ES} = \frac{1}{|b_t|} \sum_{i \in b_t} \mathcal{L}_{ES}(\mathbf{p}_i).$$

$$H(X) = -\sum_{x \in X} p(x) \log p(x)$$



$$\mathcal{L} = \mathcal{L}_{CE} + \lambda(\mathcal{L}_{AC} + \mathcal{L}_{ES}).$$

0  1  2  3  4  5  6  7  8  9

Output vector p ⟹ calculate entropy of each output

Fermilab

# Universal Domain Adaptation (DeepAstroUDA)

DA tests we ran:
- **Two data releases from the same telescope**
  - LSST mocks Y1 and Y10
- **Different surveys**
  - SDSS and DECaLS
- **Wide and deep fields in the same survey**
  - SDSS wide and Stripe 82 deep field

**Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)

DA tests we ran:
- **Two data releases from the same telescope**
  - LSST mocks Y1 and Y10
- **Different surveys**
  - SDSS and DECaLS
- **Wide and deep fields in the same survey**
  - SDSS wide and Stripe 82 deep field

Class labels are from Galaxy Zoo 2 & 3 (crowdsourcing labels ~10^5 volunteers).

Known classes:
Disturbed (0)
Merging (1)
Round smooth (2)
Cigar shaped smooth (3)
Barred spiral (4)
Unbarred tight spiral (5),
Unbarred loose spiral (6)
Edge-on without bulge (7),
Edge-on with bulge (8),

Unknown anomaly class (only in DECaLS):
Strong gravitational lens (9)

‡ Fermilab

# Universal Domain Adaptation (DeepAstroUDA)

**SDSS**



**DECaLS**



Class labels are from Galaxy Zoo 2 & 3 (crowdsourcing labels ~10^5 volunteers).

Known classes:
Disturbed (0)
Merging (1)
Round smooth (2)
Cigar shaped smooth (3)
Barred spiral (4)
Unbarred tight spiral (5),
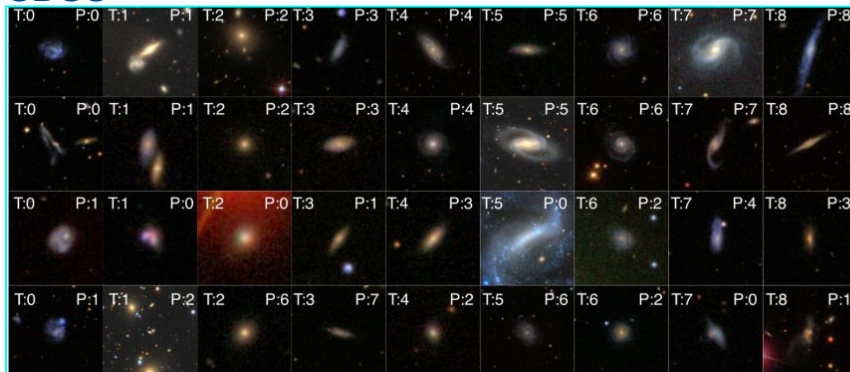Unbarred loose spiral (6)
Edge-on without bulge (7),
Edge-on with bulge (8),
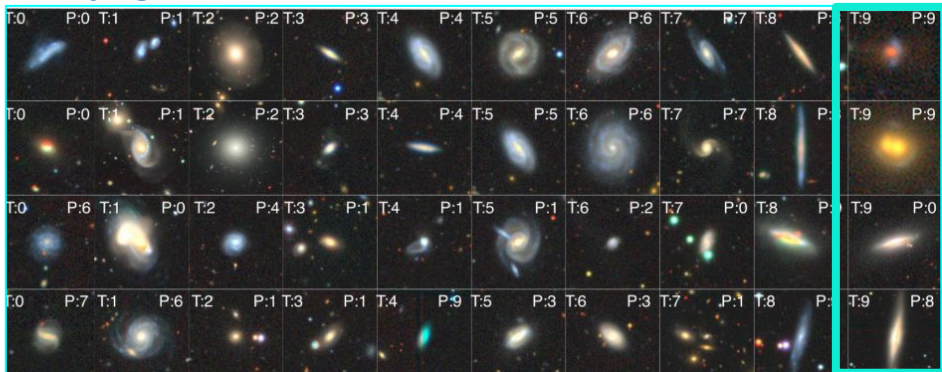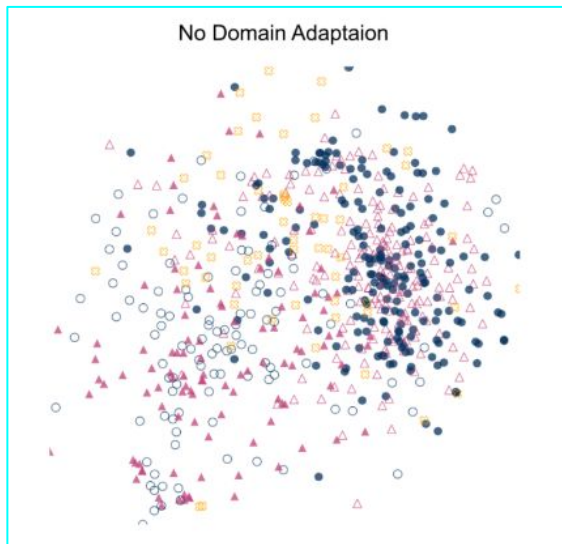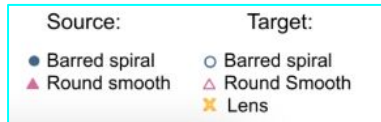
Unknown anomaly class (only in DECaLS):
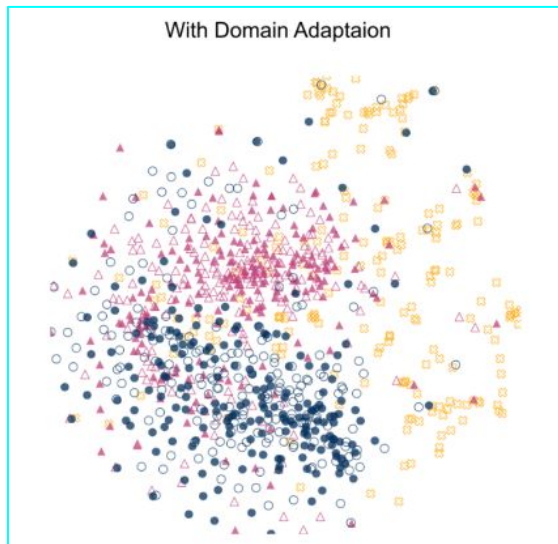Strong gravitational lens (9)

**Fermilab**

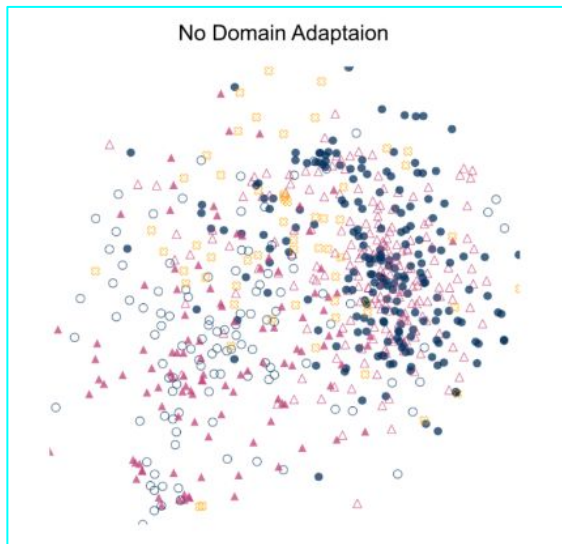# Universal Domain Adaptation (DeepAstroUDA)



No Domain Adaptaion

Classes are mixed!

Source:
- Barred spiral
- Round smooth

Target:
- Barred spiral
- Round Smooth
- Lens

**Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)



No Domain Adaptaion

With Domain Adaptaion

Classes are mixed!

Source:
- ● Barred spiral
- ▲ Round smooth

Target:
- ○ Barred spiral
- △ Round Smooth
- ✕ Lens

Known classes overlap,
unknown is pushed to the side.

‡ Fermilab
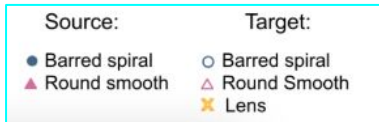
# Universal Domain Adaptation (DeepAstroUDA)



Classes are mixed!

Source: Target:
- Barred spiral (filled) / Barred spiral (open)
- Round smooth (filled triangle) / Round Smooth (open triangle)
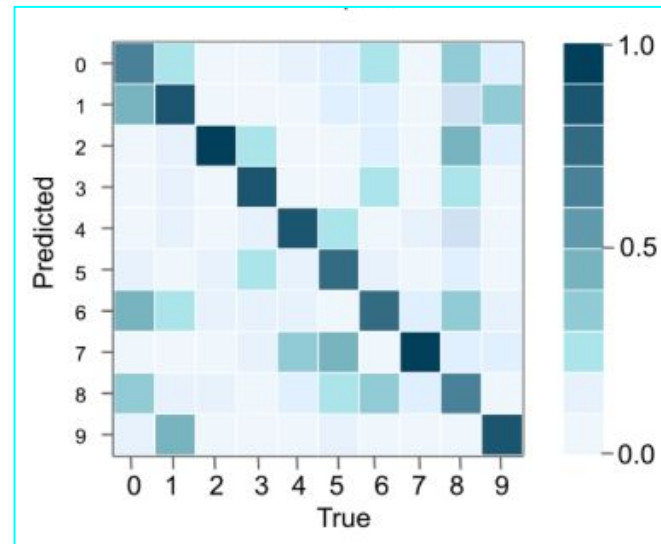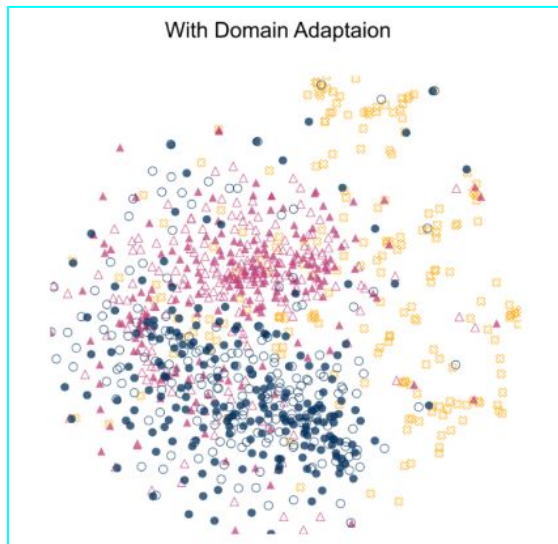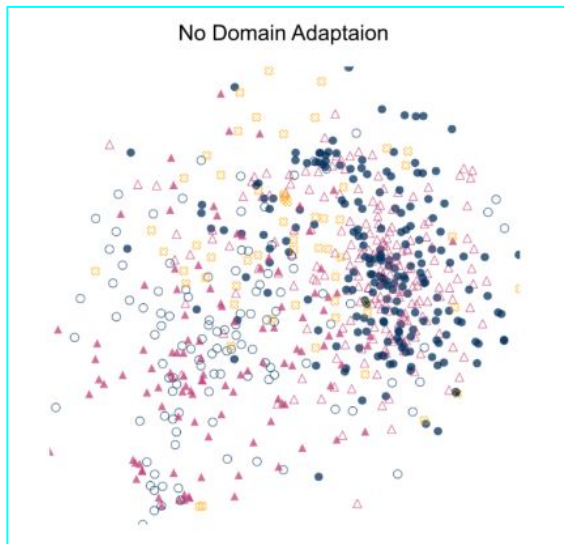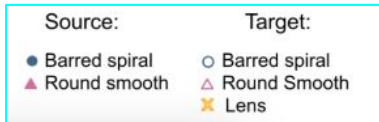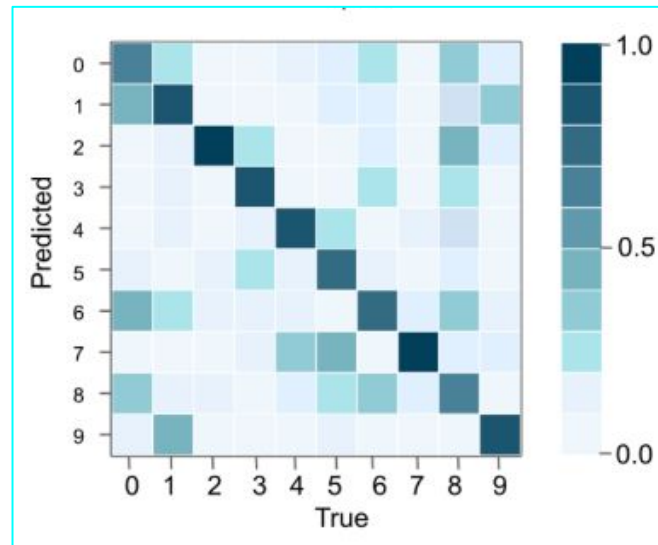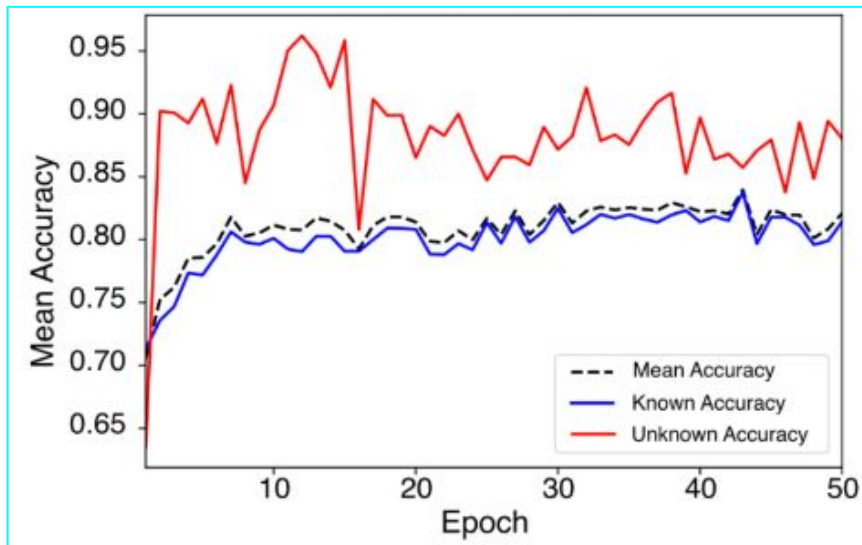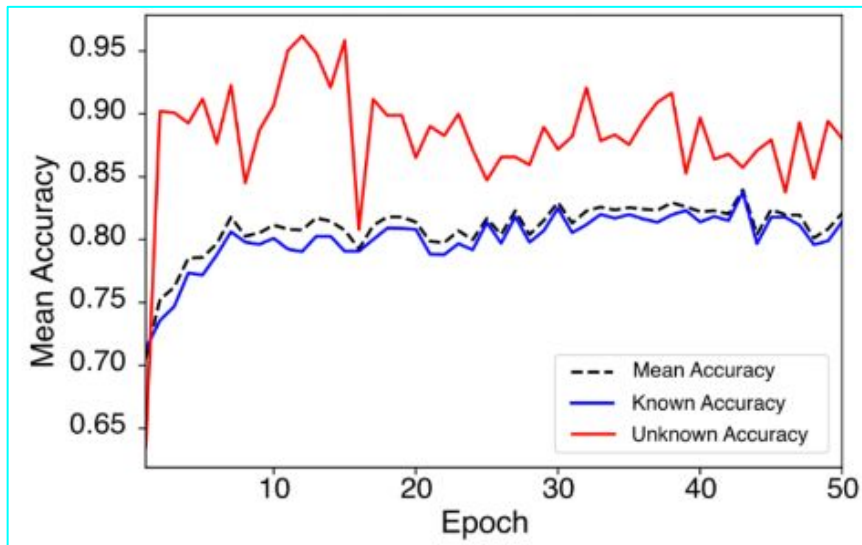- Lens (X)

Known classes overlap,
unknown is pushed to the side.

# Universal Domain Adaptation (DeepAstroUDA)



- Most confusion between classes is for truly morphologically similar classes, like disturbed and merging.
- Model is very sure about the unknown lens class - it can recognize these object look different than all other known classes.

🎇 **Fermilab**

# Universal Domain Adaptation (DeepAstroUDA)
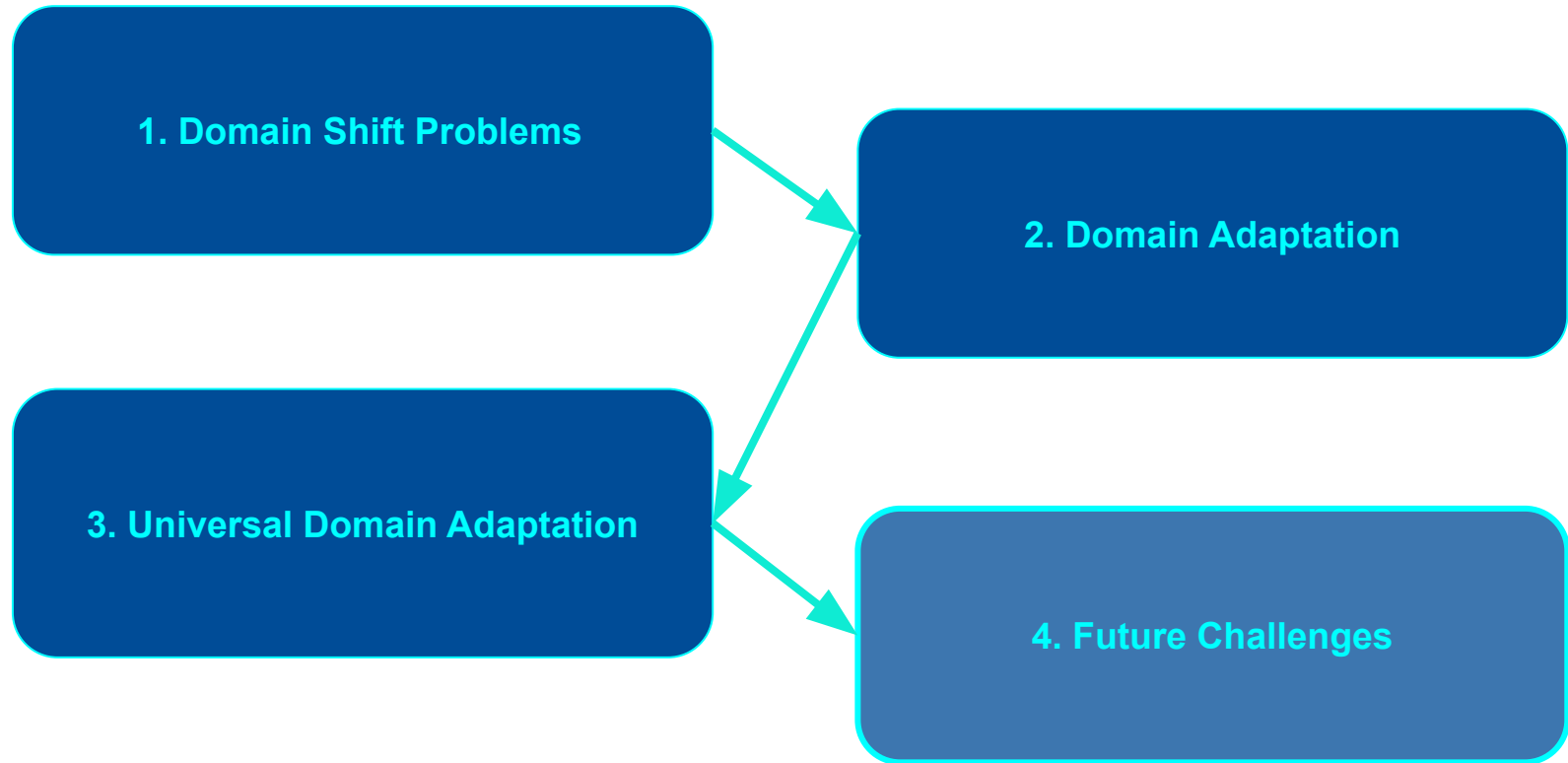


- Most confusion between classes is for truly morph[...] and merging.
- Model is very sure about the unknown lens class than all other known classes.

Fermilab

# Talk Outline

**1. Domain Shift Problems**

**2. Domain Adaptation**

**3. Universal Domain Adaptation**

**4. Future Challenges**

- Simulation and observations

  ```
  Ćiprijanović et al. 2021.
  ```

- Increase robustness to data perturbations

  ```
  Ćiprijanović et al. 2022.
  ```

- Different data releases from the same survey
- Different surveys
- Wide and deep fields of the same survey

  ```
  Ćiprijanović et al. 2022.
  Ćiprijanović et al. 2023.
  ```

**🟁 Fermilab**

- Simulation and observations

  Ćiprijanović et al. 2021.

- Increase robustness to data perturbations

  Ćiprijanović et al. 2022.

- Different data releases from the same survey
- Different surveys
- Wide and deep fields of the same survey
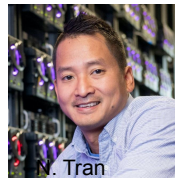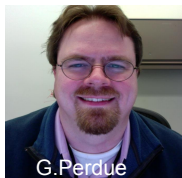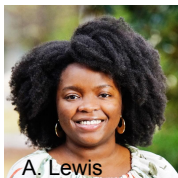
  Ćiprijanović et al. 2022.
  Ćiprijanović et al. 2023.
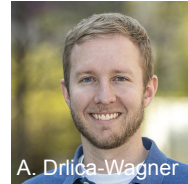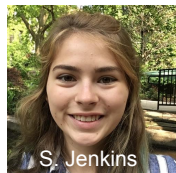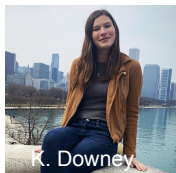
# FUTURE CHALLENGES

- Connecting extracted features to physical properties.

- Guiding the model to use some preferred physical features and discover the rest.

- Understanding and exploring the latent space.

- Can we get any new insights from AI?

- What if the domain shift is physical not instrumental/computational?

**🐝 Fermilab**

# Big thanks to all my amazing collaborators



**Fermilab**

A. Lewis  G.Perdue  D.Kafkes  B. Nord  N. Tran  Pedro

**University of Chicago**

K. Downey  S. Jenkins  J. Poh  A. Drlica-Wagner  D. Tanoglidis

**Argonne, Oakridge**

S. Madireddy  T. Johnston  ●●● and many more!

**Space Telescope Science Institute**

G. Snyder  J. Peek

THANK YOU!

KITP
March, 2023

Aleksandra Ćiprijanović
(she/her/hers)

Fermilab, DSSL
aleksand@fnal.gov