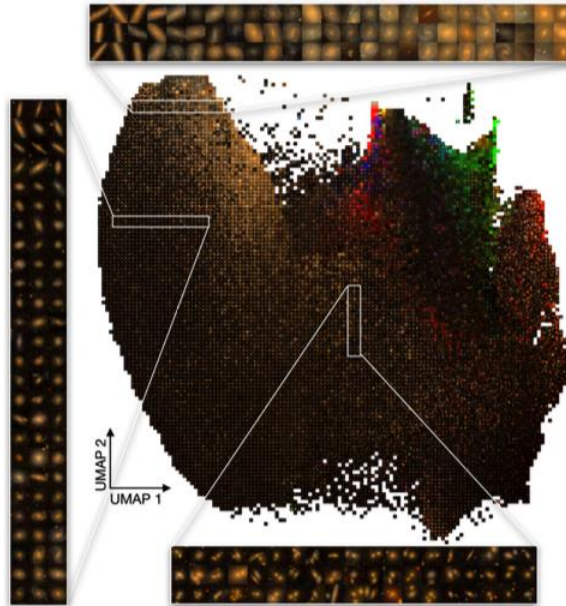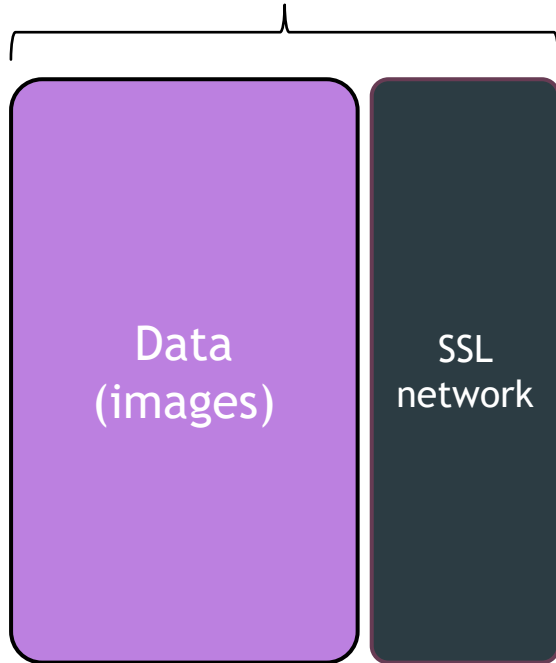# Self-Supervised Learning for MeerKAT Images

E. Lastufka

M. Audard, O. Bait, M. Dessauges-Zavadsky, M. Drozdova, T. Holotyak, V. Kinakh, D. Piras, O. Taran, D. Schaerer, S. Voloshynovskiy

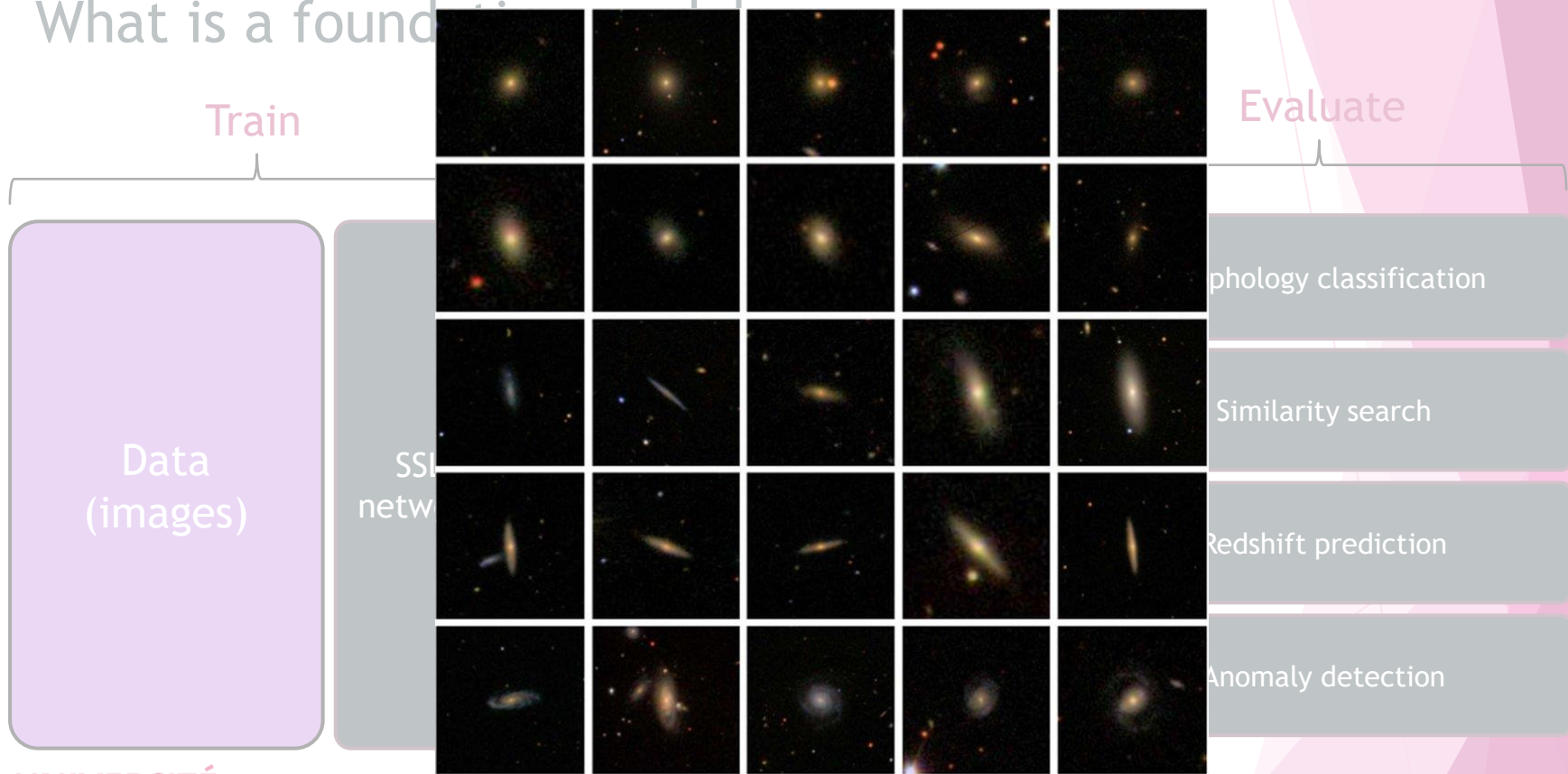# Self-supervised learning is crucial for training foundation models



Hayat et al 2021

Train

Data (images)

SSL network

Evaluate

Morphology classification

Similarity search

Redshift prediction

Anomaly detection

UNIVERSITÉ DE GENÈVE

# What is a foundation model?

Train

Evaluate

**Galaxy Zoo SDSS (scales vary)**



Data
(images)

SSL
network

...phology classification

Similarity search

Redshift prediction

Anomaly detection

UNIVERSITÉ
DE GENÈVE

# Current foundation models in astronomy

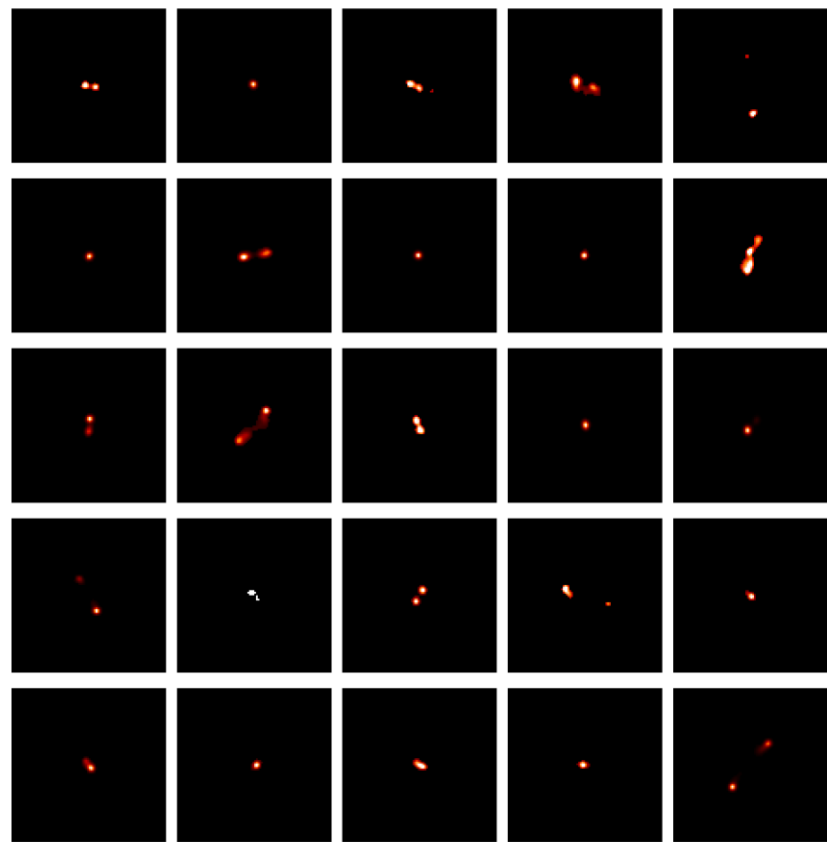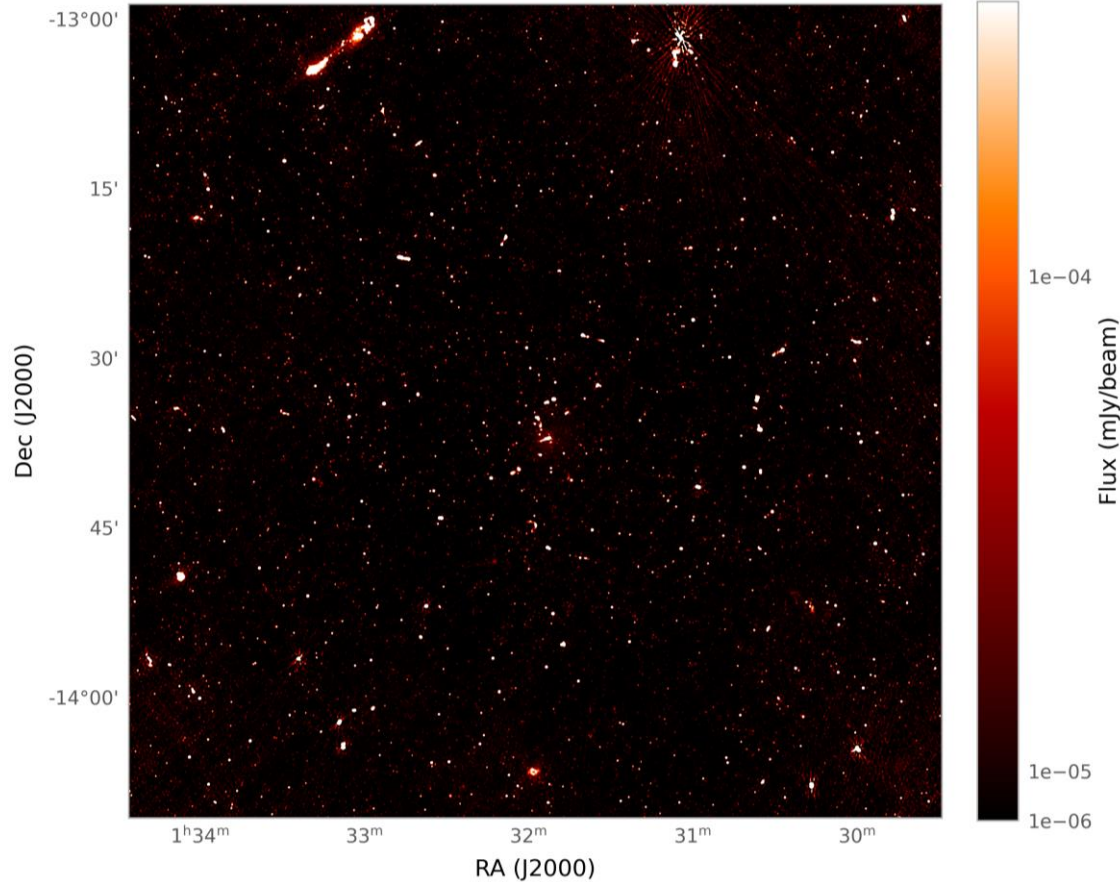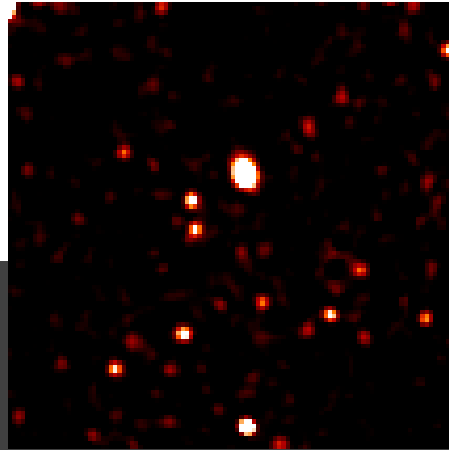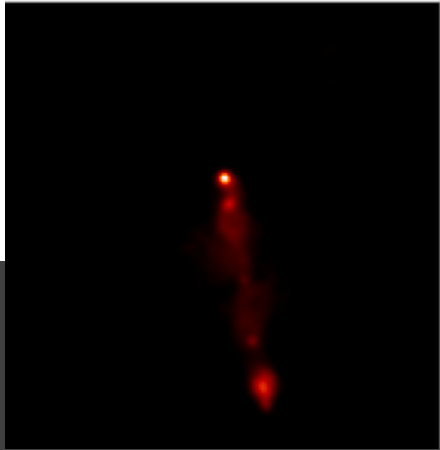| Regime | Work | Dataset | Network |
|--------|------|---------|---------|
| Optical | Hayat et al 2021 | Galaxy Zoo 2 (SDSS) | MoCo v2, Resnet50 |
| Optical | Mohale & Lochner 2023 | Galaxy Zoo DECaLS | BYOL, Resnet50 |
| Radio | Slijepcevic et al 2023 | Radio Galaxy Zoo | BYOL, Resnet18 |
| Radio | Andrianomena & Tang 2023 | Radio Galaxy Zoo | VDVAE, SimCLR BYOL, SimSiam Resnet34 |

Galaxy Zoo SDSS (scales vary)

180"

Radio Galaxy Zoo (VLA)

UNIVERSITÉ DE GENÈVE

MeerKAT L-band   Abell 209

UNIVERSITÉ
DE GENÈVE

# Can we train multi-purpose foundation models with pipeline data?

# MGCLS
## MeerKAT Galaxy Cluster
## Legacy Survey

- ▶ MeerKAT observations of 115 galactic clusters
- ▶ Wide-field 1.2 degree images
  - ▶ ~20,000 crops of pixel size 256 x 256
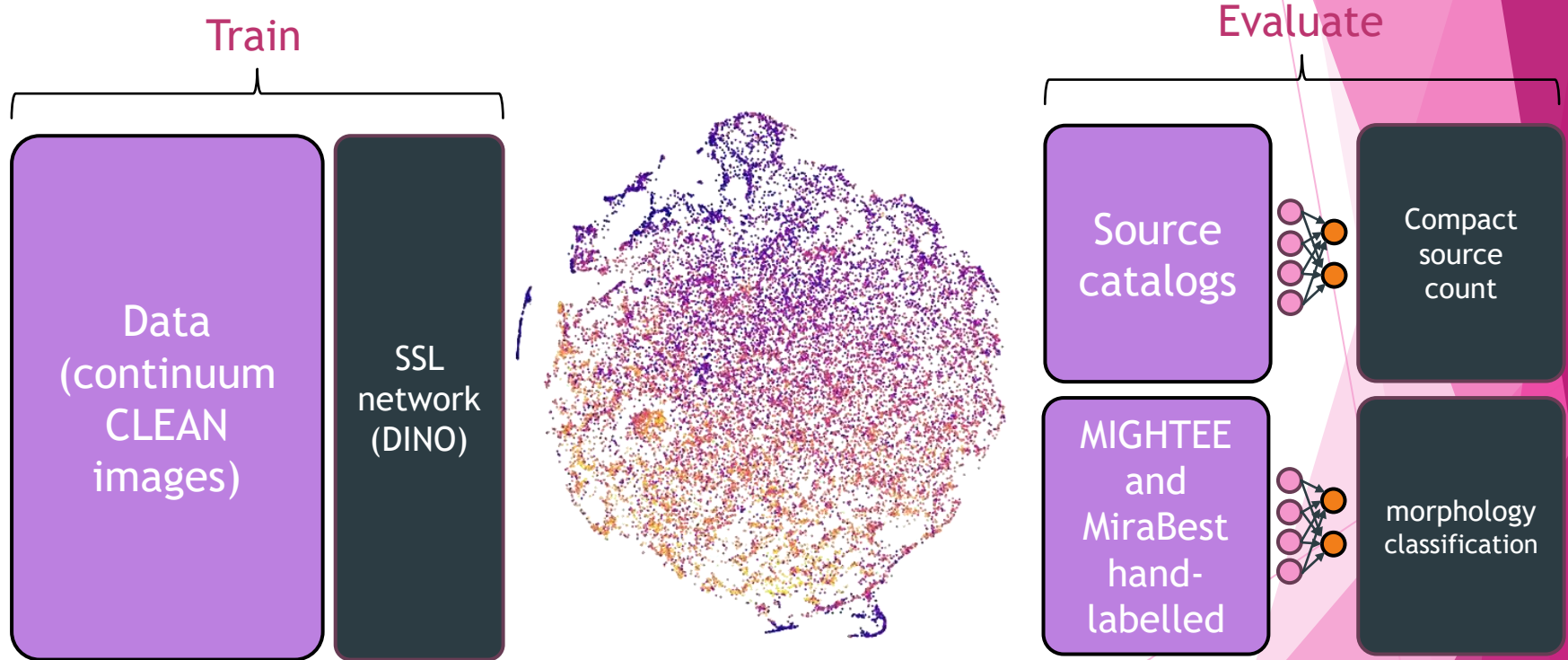- ▶ Continuum CLEAN images

UNIVERSITÉ
DE GENÈVE

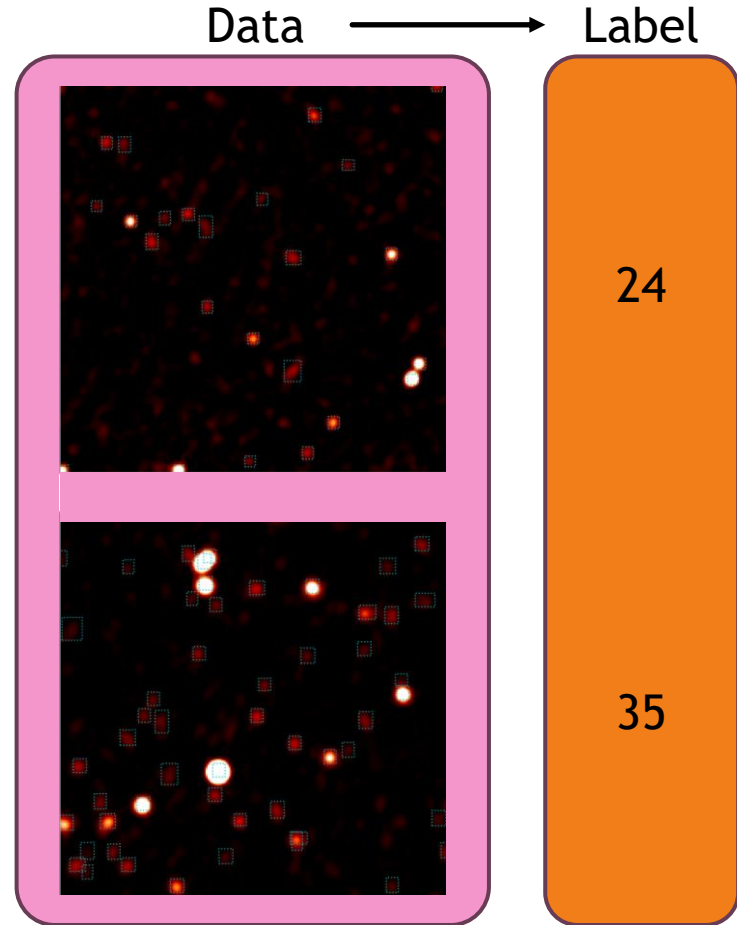# Knowledge distillation with no labels: DINO

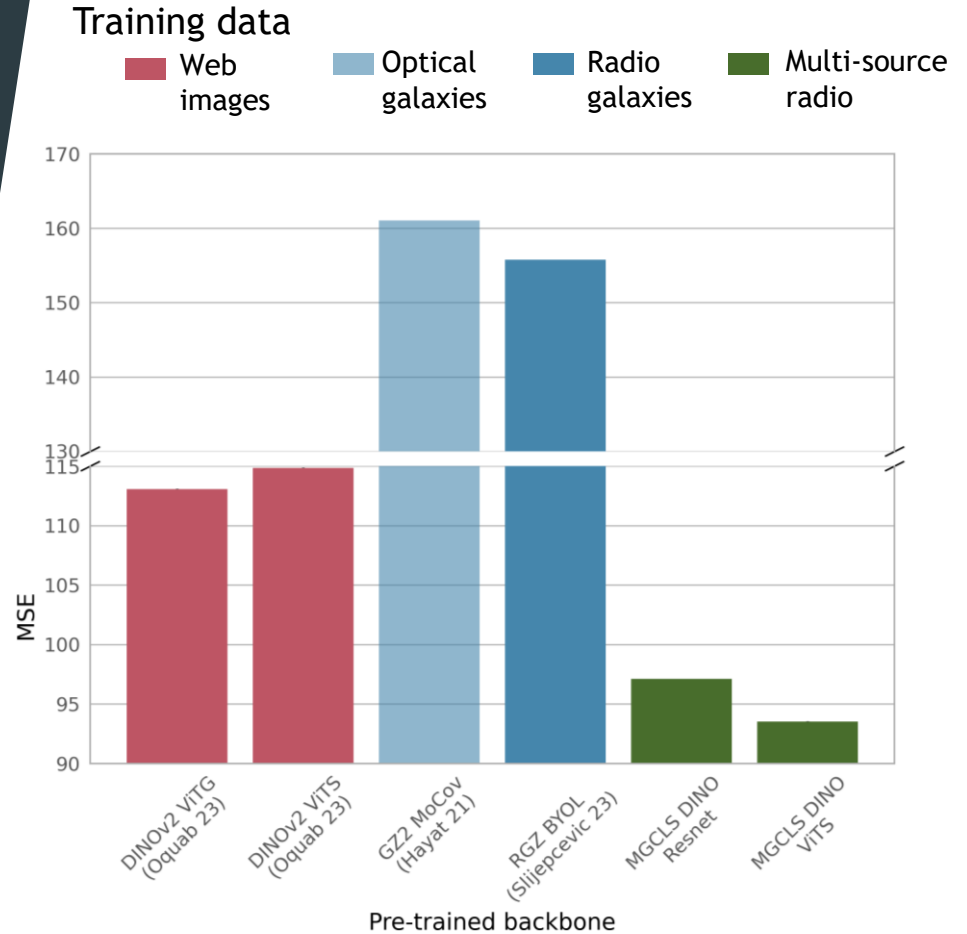# Can SSL create representations that characterize source-rich MeerKAT data?

# Evaluation: compact source count prediction

▶ Labels are from compact source catalog (pyBDSF)

▶ Task: predict source count with image as input
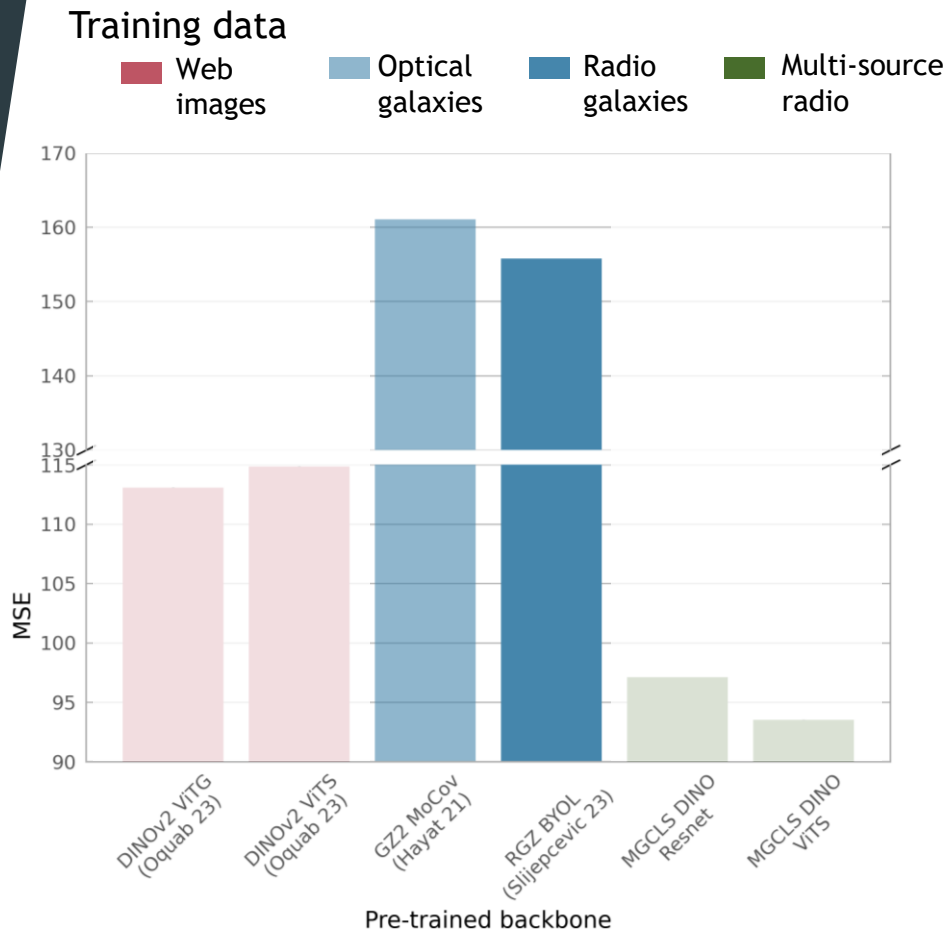
▶ Metric: mean squared error (MSE), the lower the better



Data → Label

24

35

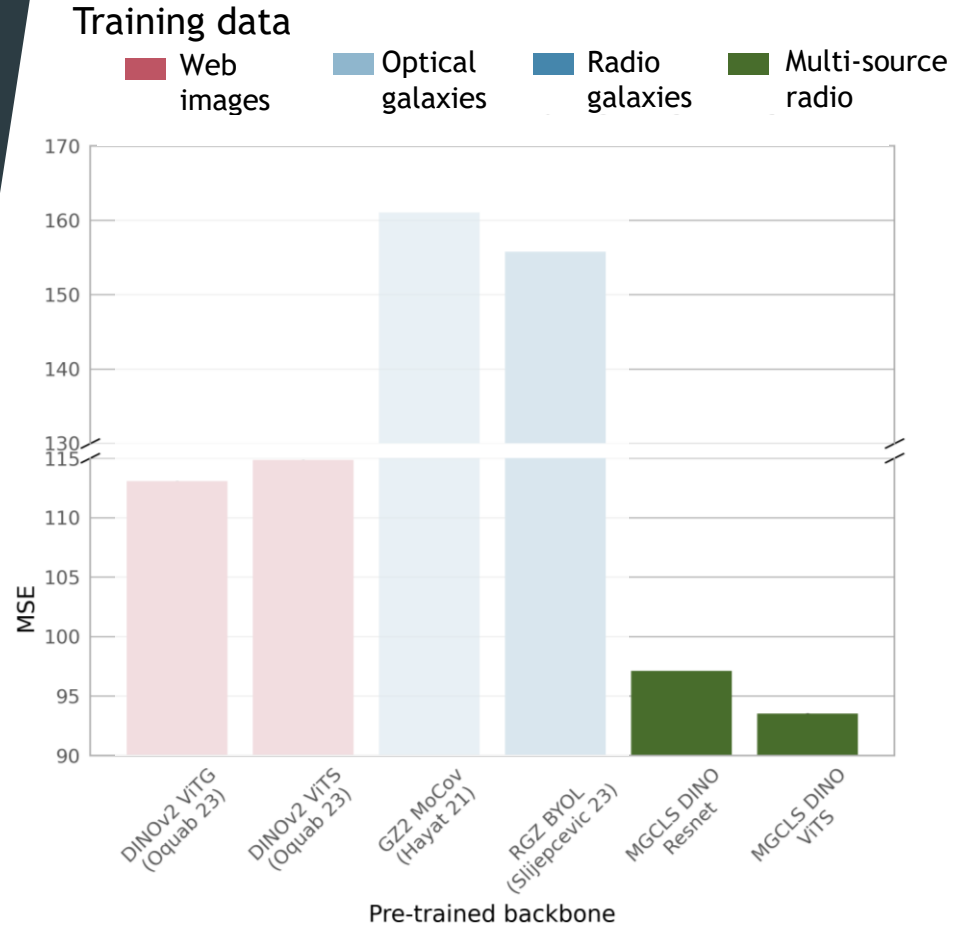# Evaluation: compact source count prediction

SKACH Winter Meeting – E. Lastufka

# Evaluation: compact source count prediction

▶ Backbones trained on curated, single-galaxy images cannot perform this task!



Training data

■ Web images   ■ Optical galaxies   ■ Radio galaxies   ■ Multi-source radio

22-Jan 2024

# Evaluation: compact source count prediction

▶ Domain-specific training data is best for this task



22-Jan 2024

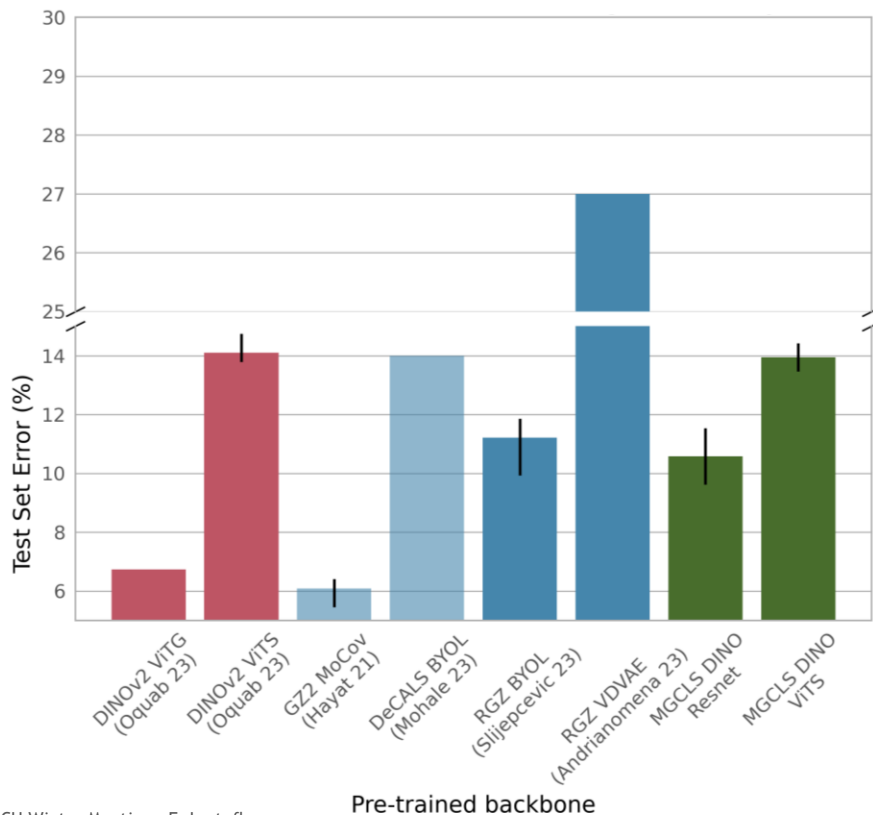# Can the network do more than one thing?

UNIVERSITÉ
DE GENÈVE

# Evaluation: FRI/FRII galaxy morphology classification

▶ Public MiraBest dataset from VLA images
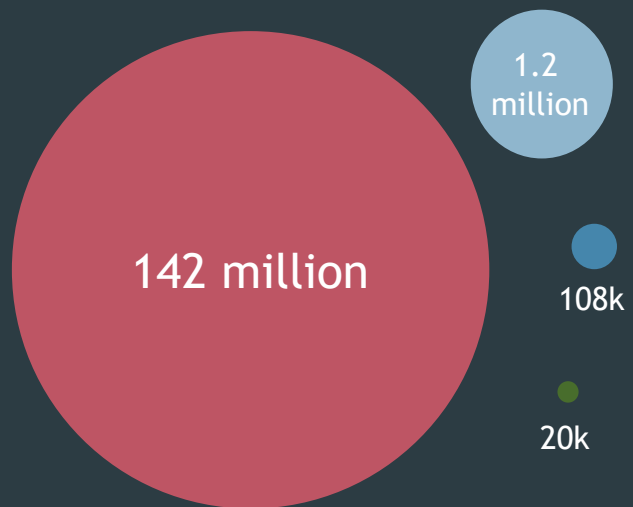
▶ Metric: Test set error (1-accuracy), lower is better

# Evaluation: FRI/FRII galaxy morphology classification

SKACH Winter Meeting – E. Lastufka

# Evaluation: FRI/FRII galaxy morphology classification

▶ **Big training datasets gives good performance!**

1.2 million

142 million

108k

20k

## Training data

- Web images
- Optical galaxies
- Radio galaxies
- Multi-source radio
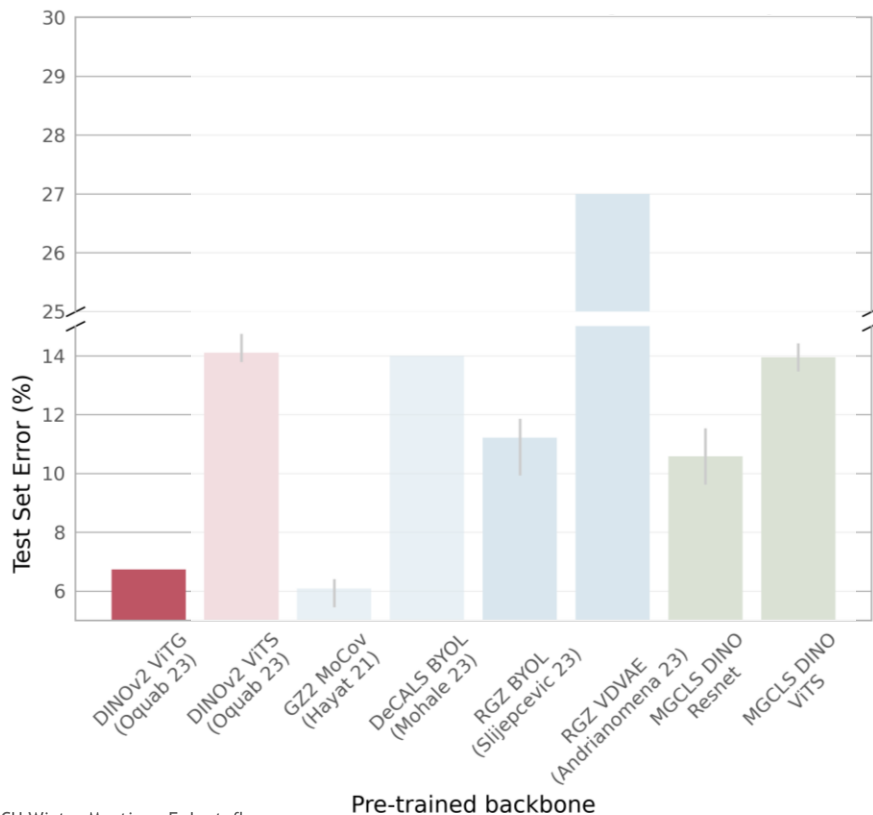
# Evaluation: FRI/FRII galaxy morphology classification

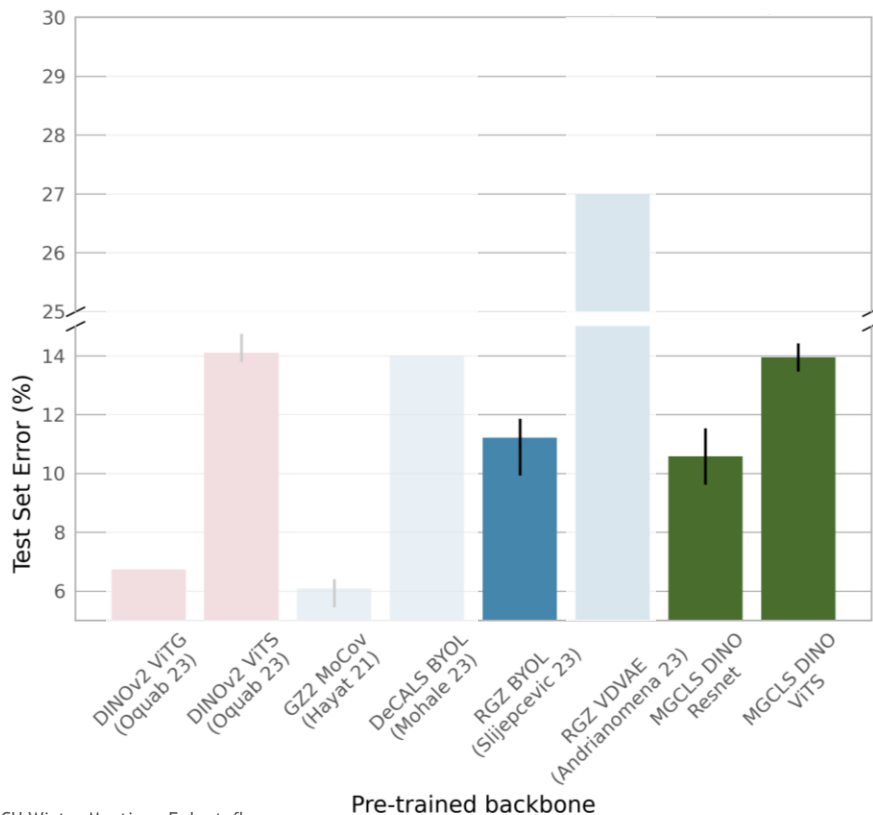▶ Big model makes up for lack of domain-specific training data?



22-Jan 2024

# Evaluation: FRI/FRII galaxy morphology classification

▶ Our model trained on source-rich data can perform as well as academic state-of-the-art trained on highly curated data
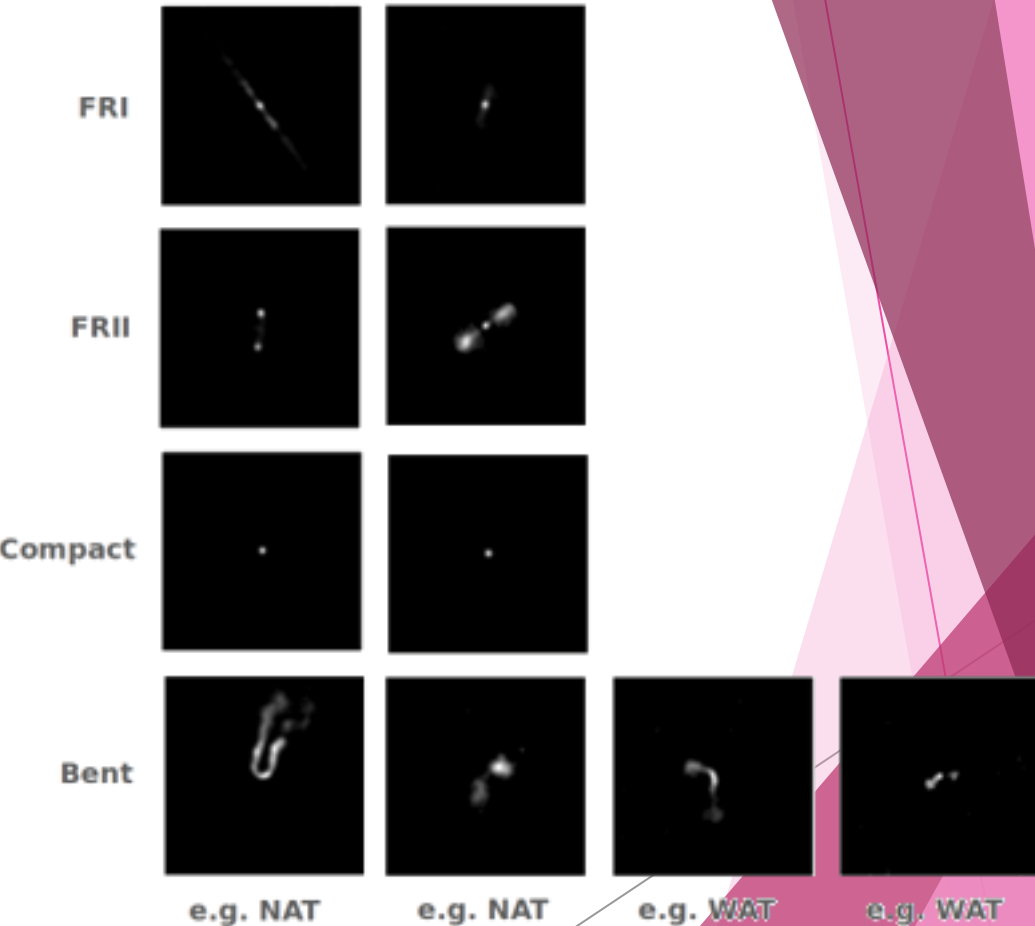


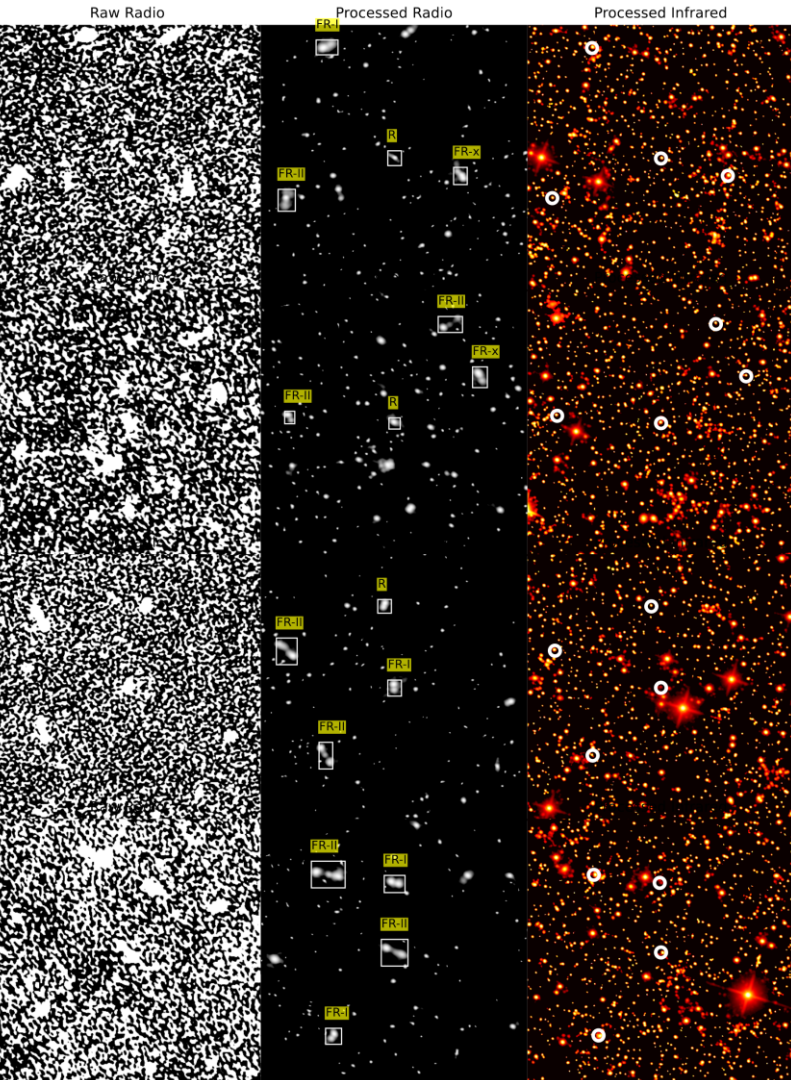22-Jan 2024

# Can the network do complex things?

UNIVERSITÉ
DE GENÈVE

# Multi-class morphology classification

▶ RadioGalaxyDataset

  ▶ VLA images of four morphology classes

E. Lastufka                    22-Jan 2024

# Source detection and segmentation

- [RadioGalaxyNET](#)
  - ASKAP + WISE radio/infrared dataset for object detection

# Future foundation models in astronomy

Training data can come straight from the data processing pipeline!

A mixture of source-rich and single-source?

Evaluated on multiple high-complexity tasks

UNIVERSITÉ DE GENÈVE

# Extra slides

UNIVERSITÉ
DE GENÈVE

# Evaluation: multi-class galaxy morphology classification

SKACH Winter Meeting – E. Lastufka