

AI DRIVEN COSMOLOGY

Nesar Ramachandra
nramachandra@anl.gov
Computational Science Division
Argonne National Laboratory



This talk highlights a subset of the AI-for-Cosmology efforts at Argonne.

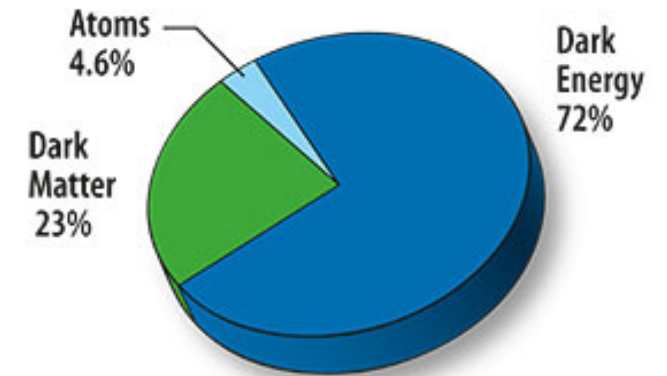
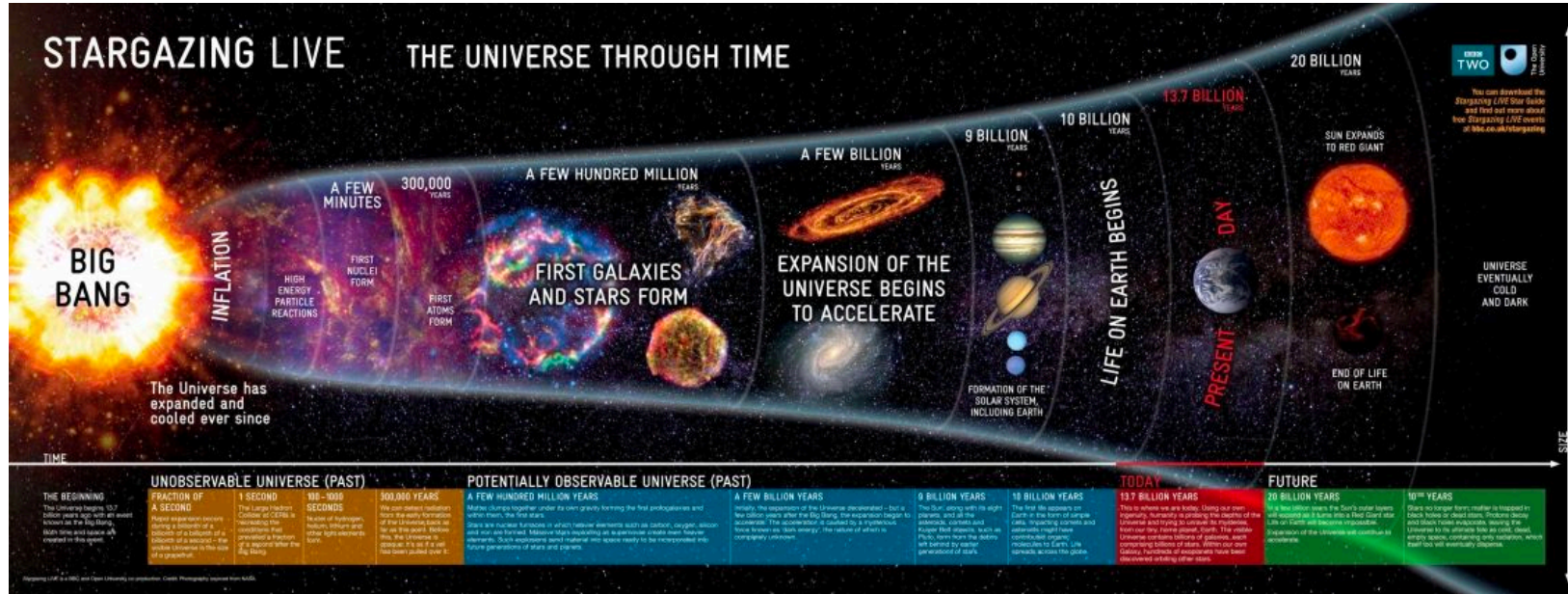
Common themes here are:

- Synthetic/simulation datasets and their connection with real astronomical observations.
- Science requirements:
 - Bayesian/probabilistic schemes rather than point-predictions.
 - Explainability of the AI algorithms.
 - Physics inclusion

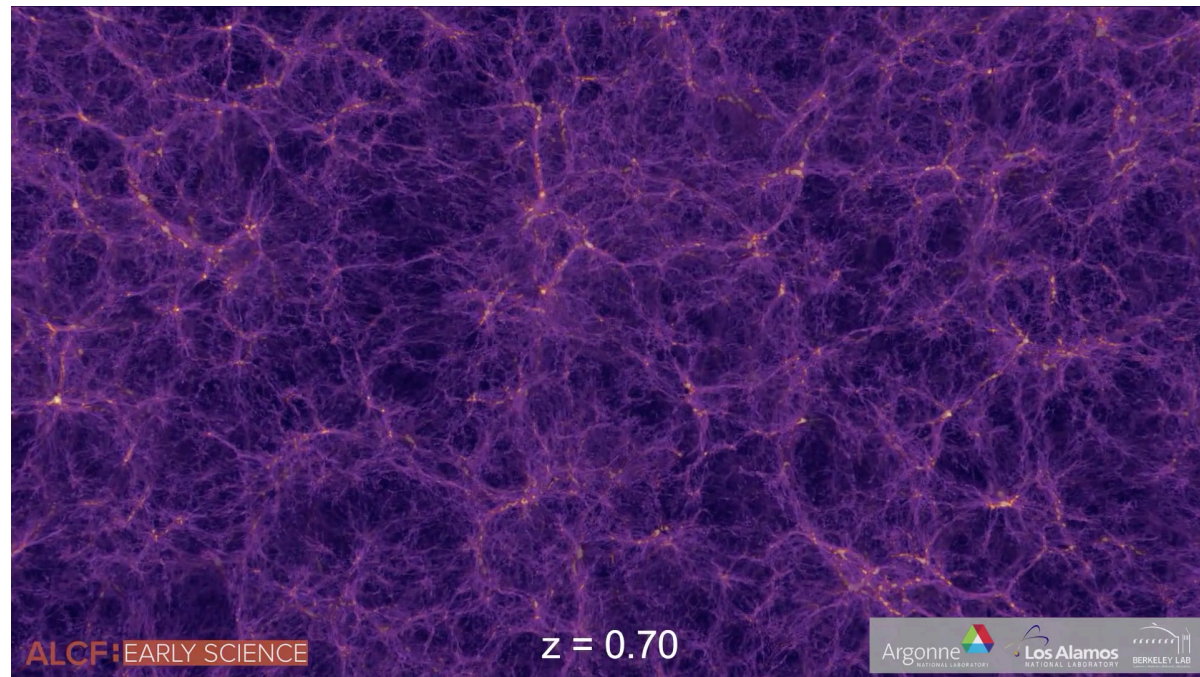
Case study:

- Image processing pipelines for de-noising, de-blending etc.
- Probabilistic classification and regression
- Latent space exploration

SHORT INTRO TO COSMOLOGY



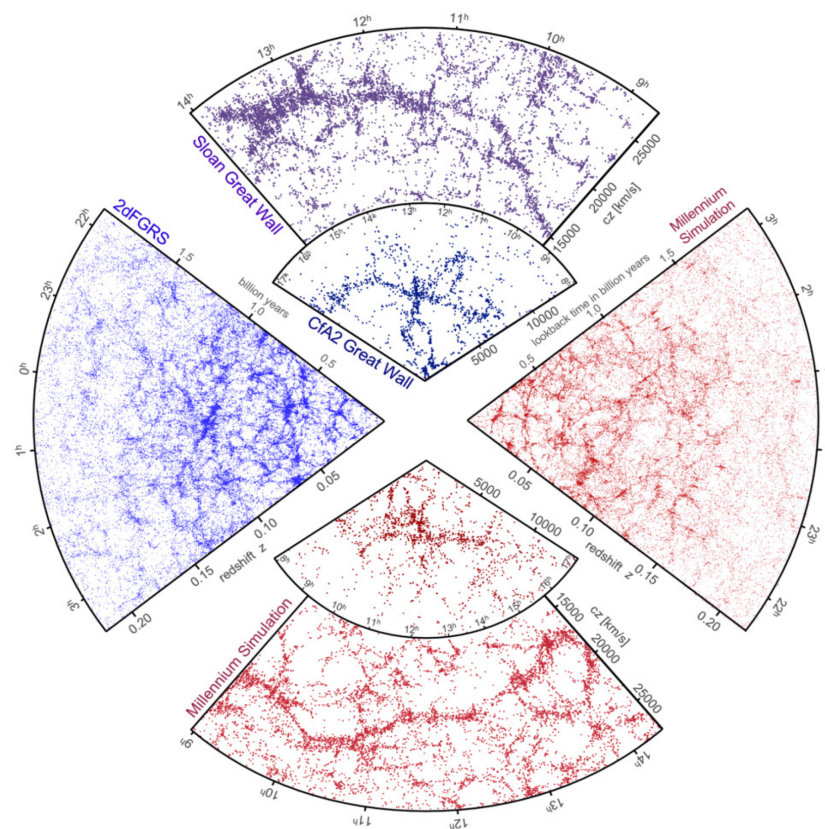
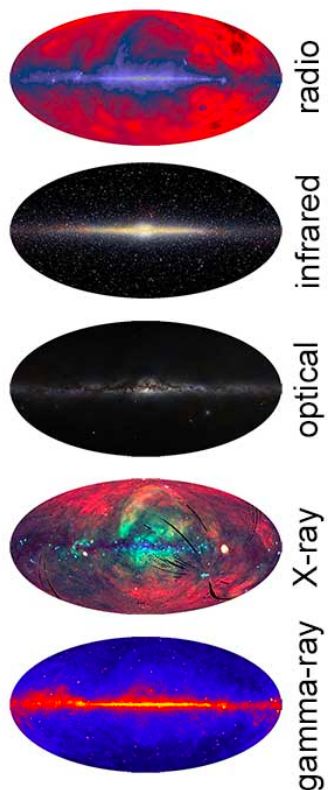
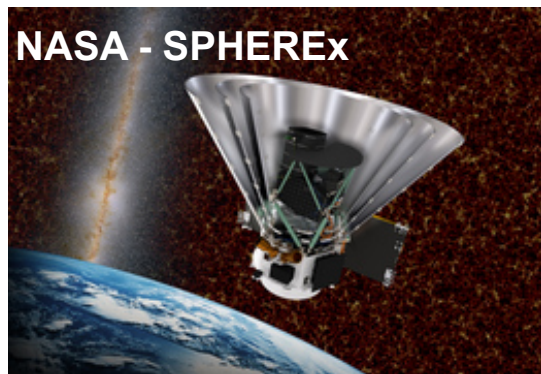
Credit: nasa.gov



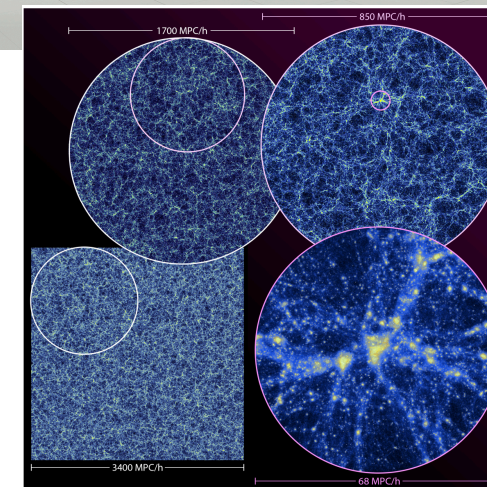
Source: open.edu

Name (Symbol)	Description	Value
Hubble Constant (H_0)	Current rate of expansion of the universe.	67.4 ± 0.5 km/s/Mpc
Cosmological Constant (Λ)	Energy density of space, or vacuum energy.	1.1056×10^{-52} m ⁻²
Dark Energy Density (Ω_Λ)	Fraction of the universe's energy density consisting of dark energy.	0.6847 ± 0.0073
Dark Matter Density (Ω_c)	Fraction of the universe's energy density consisting of dark matter.	0.264 ± 0.013
Baryon Density Parameter ($\Omega_b h^2$)	Density of ordinary matter (baryons) relative to the critical density.	0.0224 ± 0.0001

STUDYING THE UNIVERSE: JOINT EFFORTS



Zavala, J.; Frenk, C.S. Dark Matter Haloes and Subhaloes. *Galaxies* **2019**, 7, 81.



Theoretical and
Computational
Cosmology

Observational
Astronomy

CAVEATS IN COSMOLOGICAL STUDIES

- Multi-modal data:
 - Graphs, images, time-series, summary vectors, scalars, text.
- Expensive:
 - Both simulations and observations are from expensive science campaigns
- Multi-fidelity:
 - Data from different sources have different resolutions, approximations and systematic effects.
 - Transfer of knowledge is not straightforward.
- Data coverage:
 - Gaps, biased datasets are common. No assumption of a 'fair' sampling.

- Analysis requirements
 - Precision cosmology has high error requirements
 - Traditional statistical analyses have been highly successful.
- Prior domain knowledge:
 - Studies assume known physics, conservation laws.

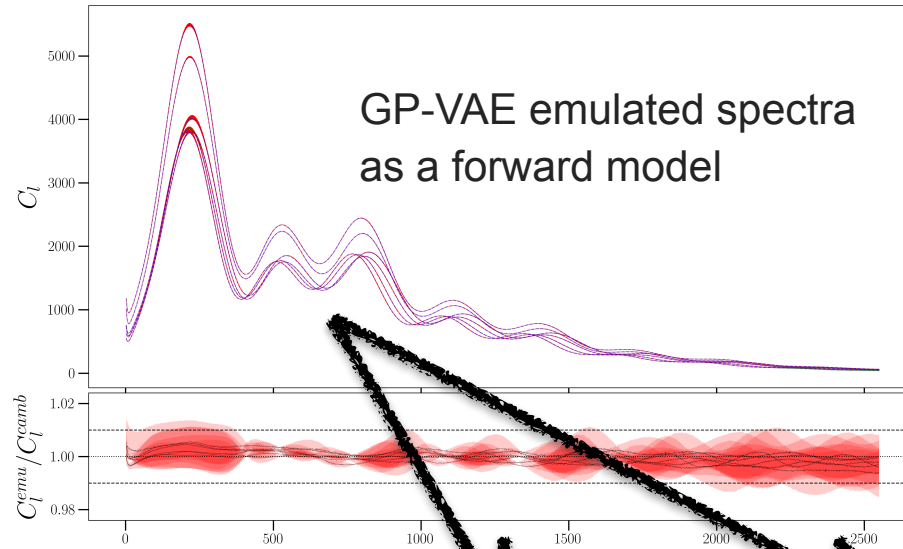


TRAIN



TEST

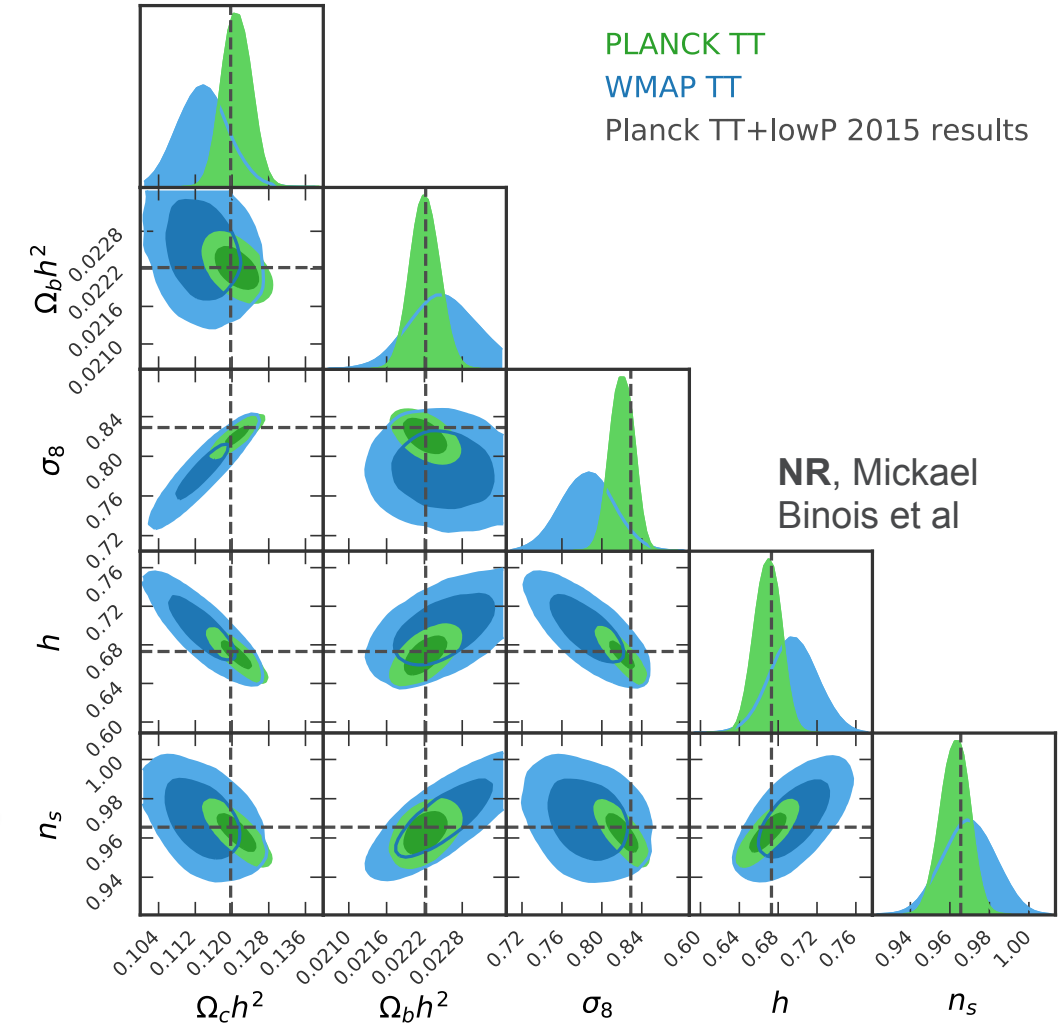
A FEW EXAMPLES: BAYESIAN INFERENCE WITH EMULATORS



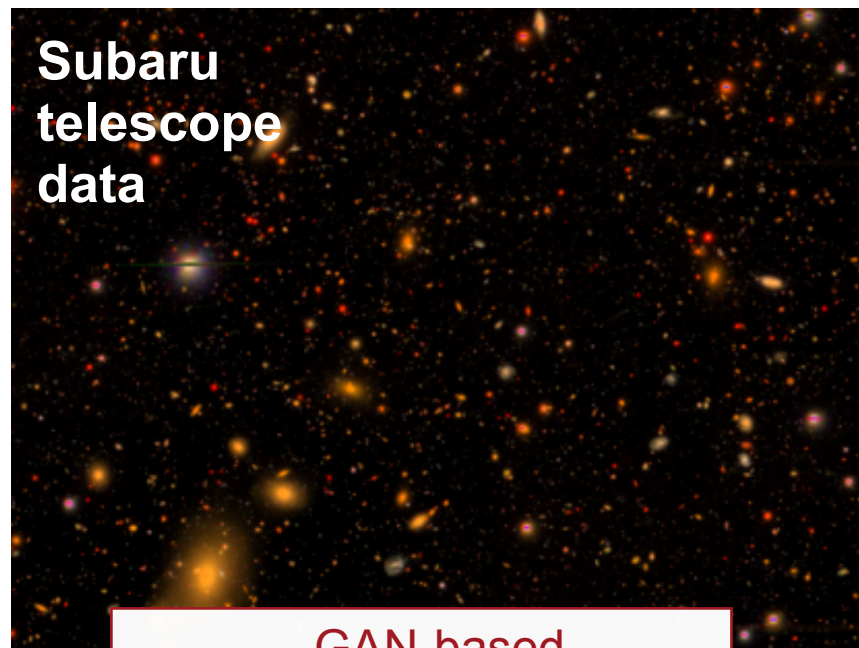
$$\mathcal{L}(D|\theta) \propto \exp \left[-\frac{1}{2} \sum_{i,j} (D - f(\theta))_i C_{ij}^{-1} (D - f(\theta))_j \right]$$

$$P(\theta|D) \propto \mathcal{L}(D|\theta)P(\theta)$$

MCMC sampling
for PLACK/WMAP
data



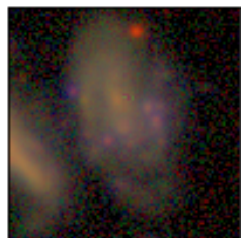
A FEW EXAMPLES: FINDING UNKNOWN UNKNOWN OBJECTS IN THE SKY



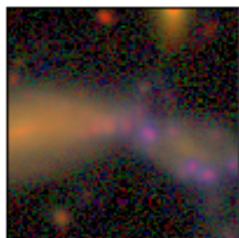
GAN-based
Anomaly finder

<https://arxiv.org/abs/2012.08082>

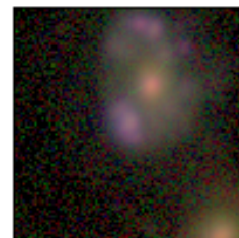
ID: 41631942633848832



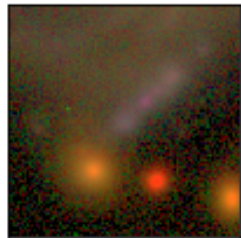
ID: 43774461299654656



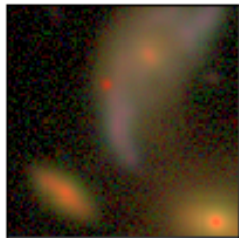
ID: 37480594749259776



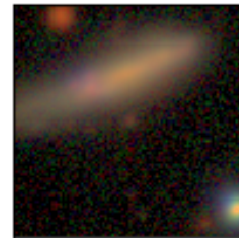
ID: 44781626835599360



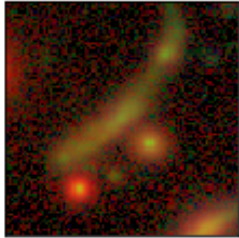
ID: 44223190892806144



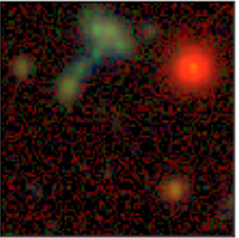
ID: 41196926871273472



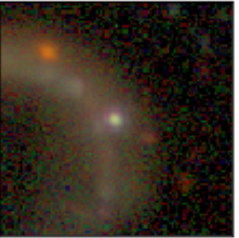
ID: 40669298728894464



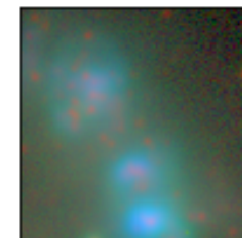
ID: 40665162675388416



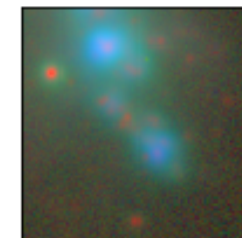
ID: 42182248203550720



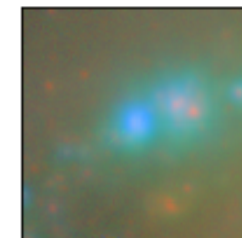
ID: 41012904702509056



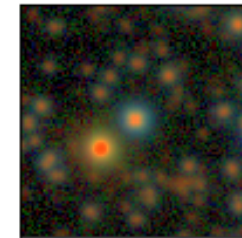
ID: 41012904702509056



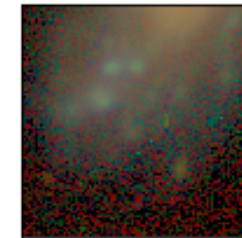
ID: 42086019461283840



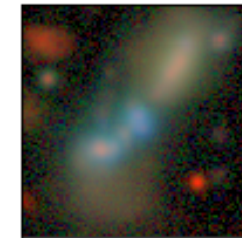
ID: 42093978035683328



ID: 40664896387416064

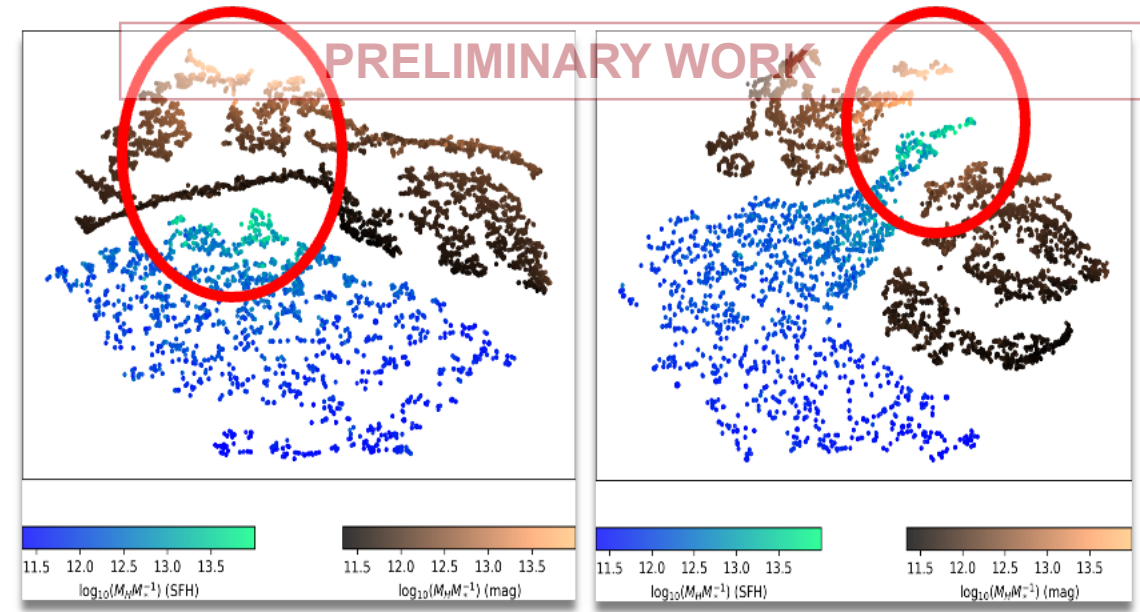


ID: 44763459123937280

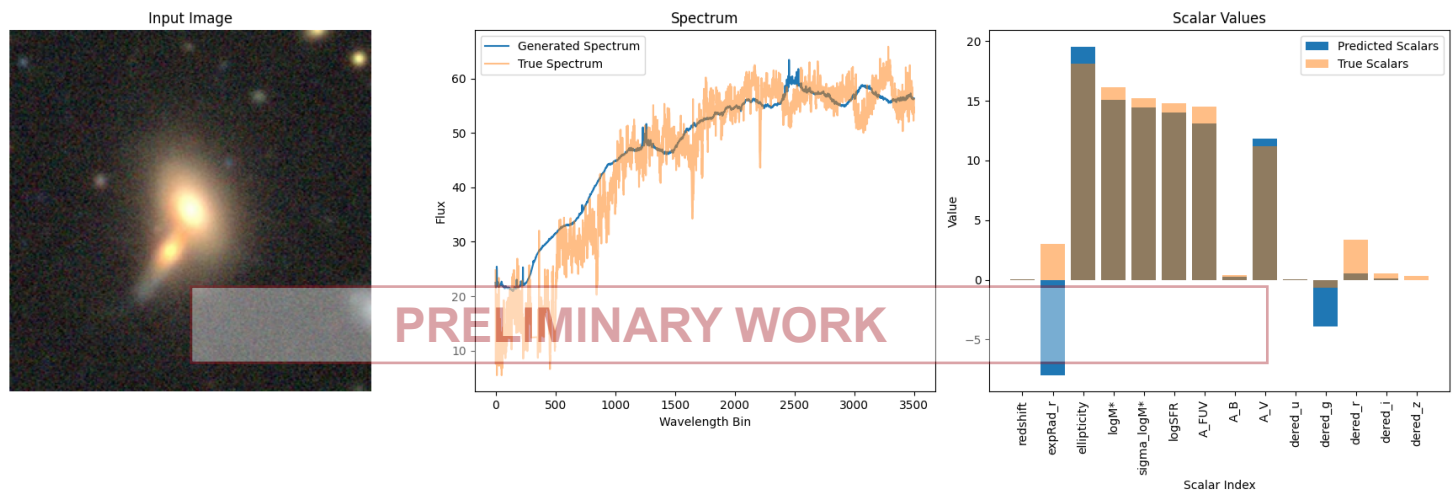


A FEW EXAMPLES: MULTI-MODAL FOUNDATION MODELS

- Foundation models (like Large Language models) encode diverse data into a cohesive representation space
- Going beyond text is the next frontier — unique problems arise in scientific datasets.



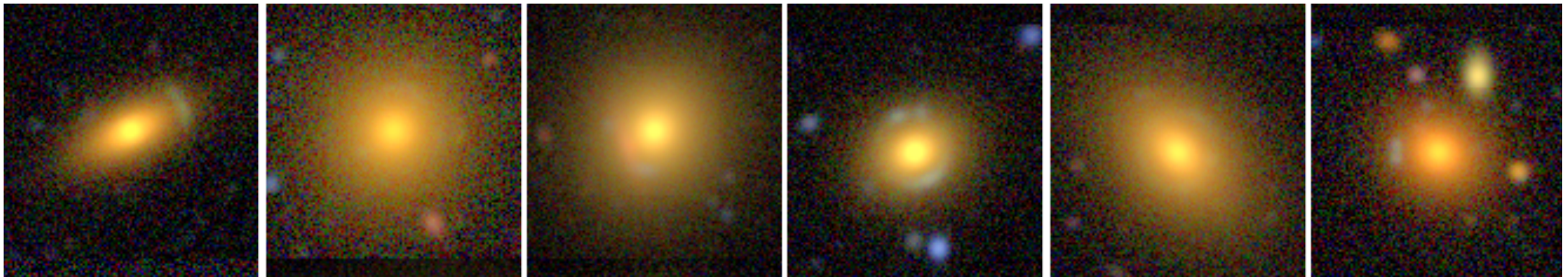
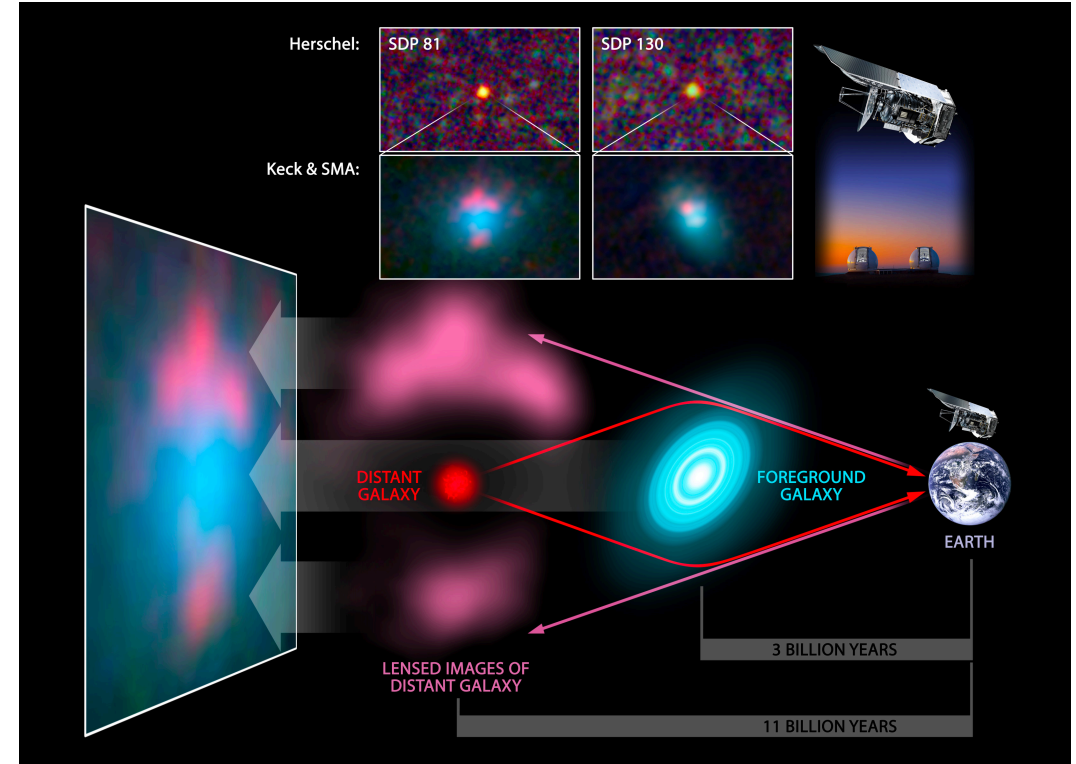
Mapping between different modes



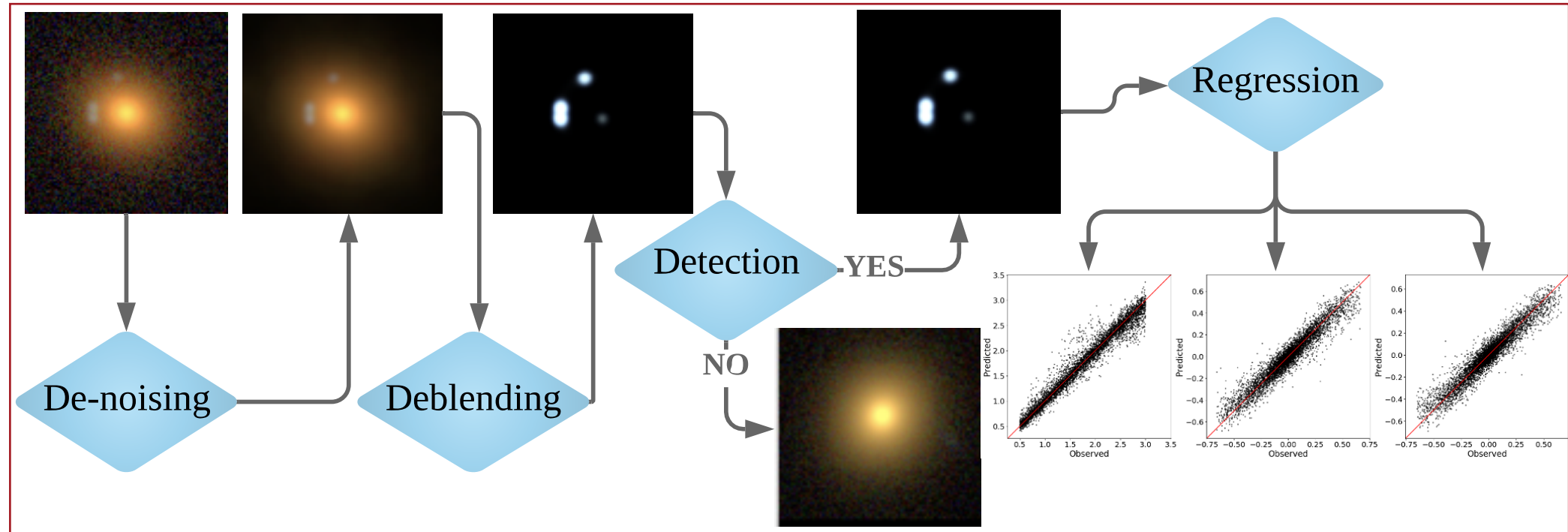
- Foundation models are designed to work on **tasks that are not pre-defined** — a major paradigm shift in AI

CASE STUDY: GALAXY-SCALE STRONG LENSING

- Strong-lenses are rare objects.
 - But understanding them is key to several questions: Distribution of dark matter, expansion of the Universe.
- Discrepancy with current amount of observed data vs future data
 - Observed data is/will be a highly imbalanced dataset
 - Tractable physical models



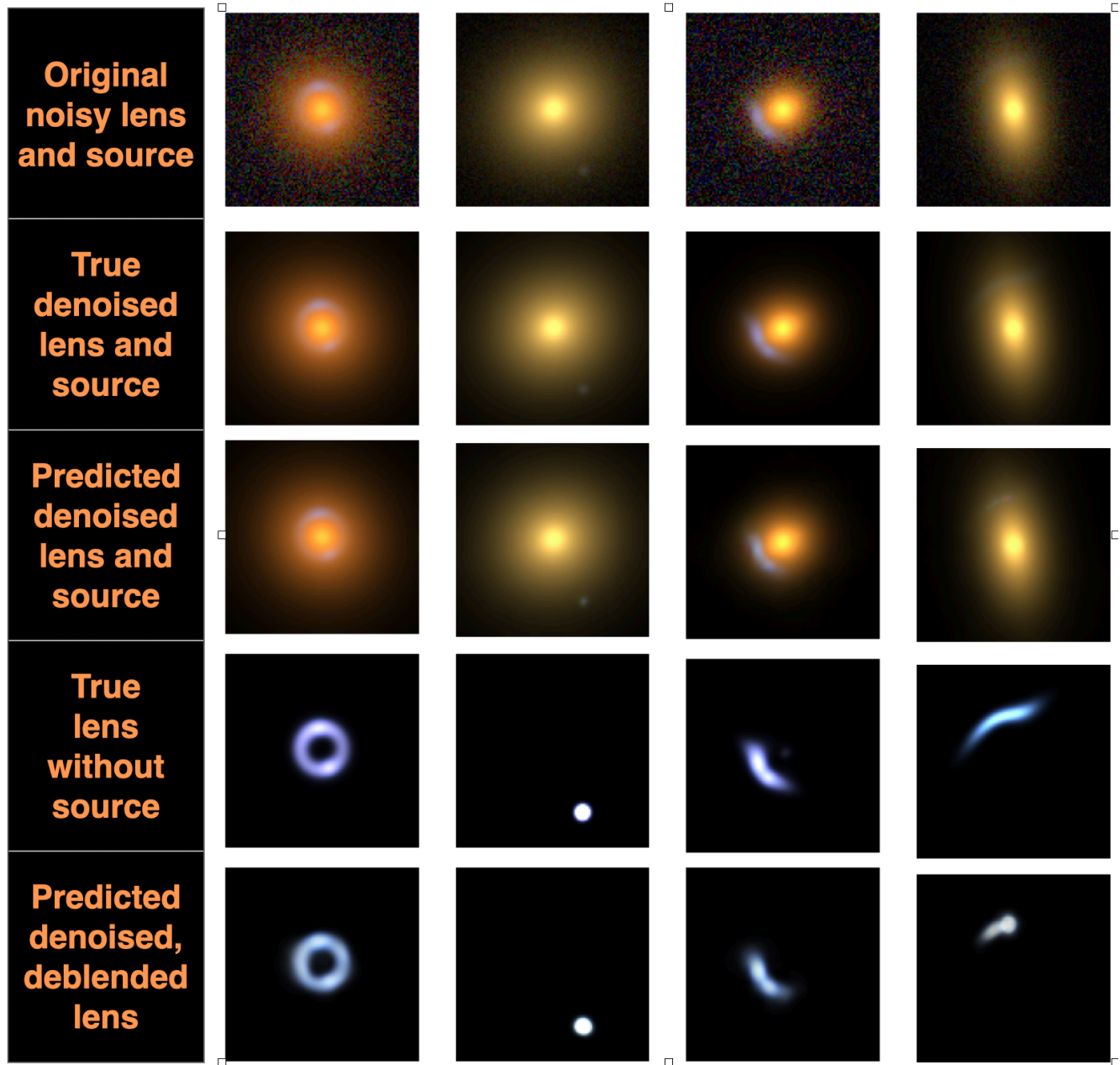
INTERPRETABLE LEARNING PIPELINES

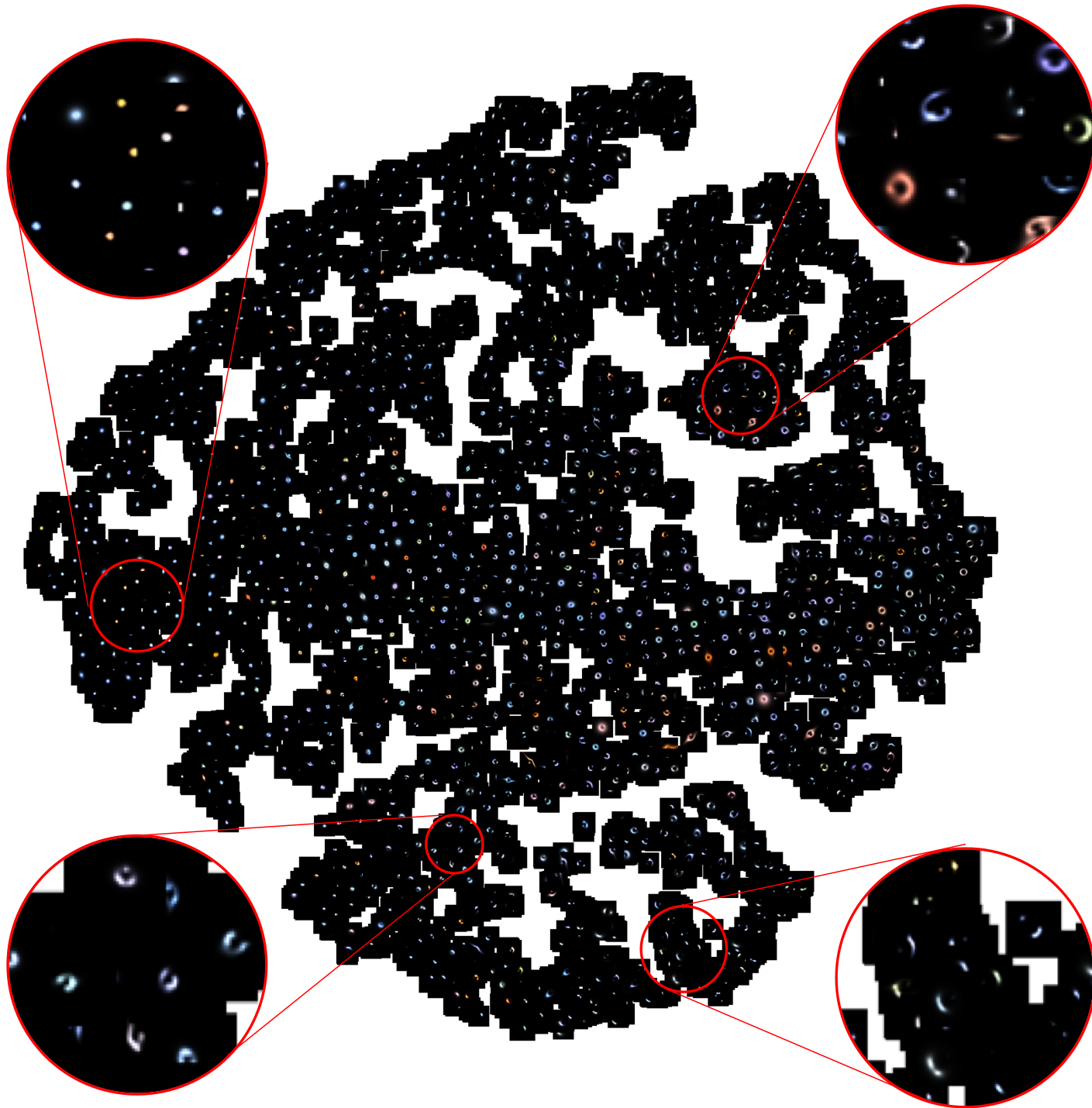


- Each AI-module is independently trained and validated. (Super-resolution modules for Denoising and deblending, Information bottleneck design for detection and regression)
- Synthetic data allows one to train modular pipelines that enable better control over systematics than end-to-end training methods

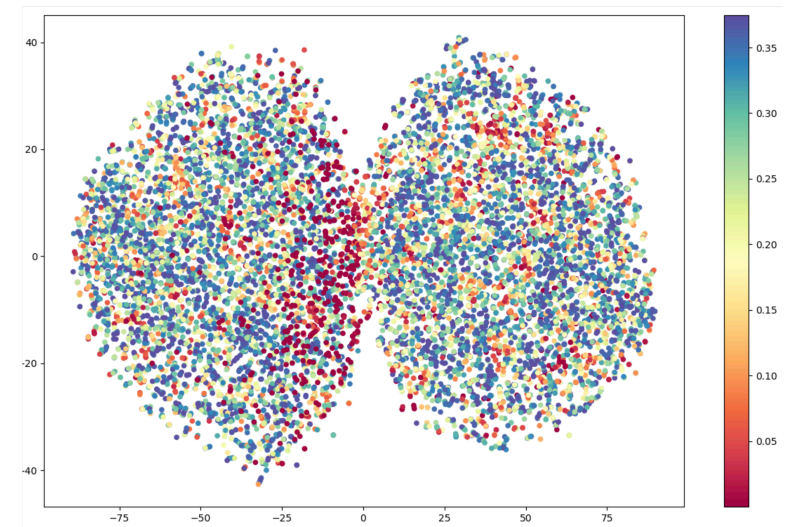
INTERPRETABLE STRONG LENS END-TO-END ANALYSIS PIPELINE

Sandeep Madireddy, NR et al:
[arxiv.org:1911.03867](https://arxiv.org/1911.03867)





Variational Information Bottleneck and representation learning



Uncertainty quantification
for classification

CONCLUSIONS

- Cosmological studies involve variety of data modalities, with vast amount of data. This makes data-driven AI-models extremely valuable.
- Synthetic datasets are often a necessity in Cosmological analysis.
- Careful experimental design, robust data creation, extensive validations are all required while dealing with synthetic data.
- Interpretable, uncertainty quantified models are still very important, probably even more so while using synthetic data in training.