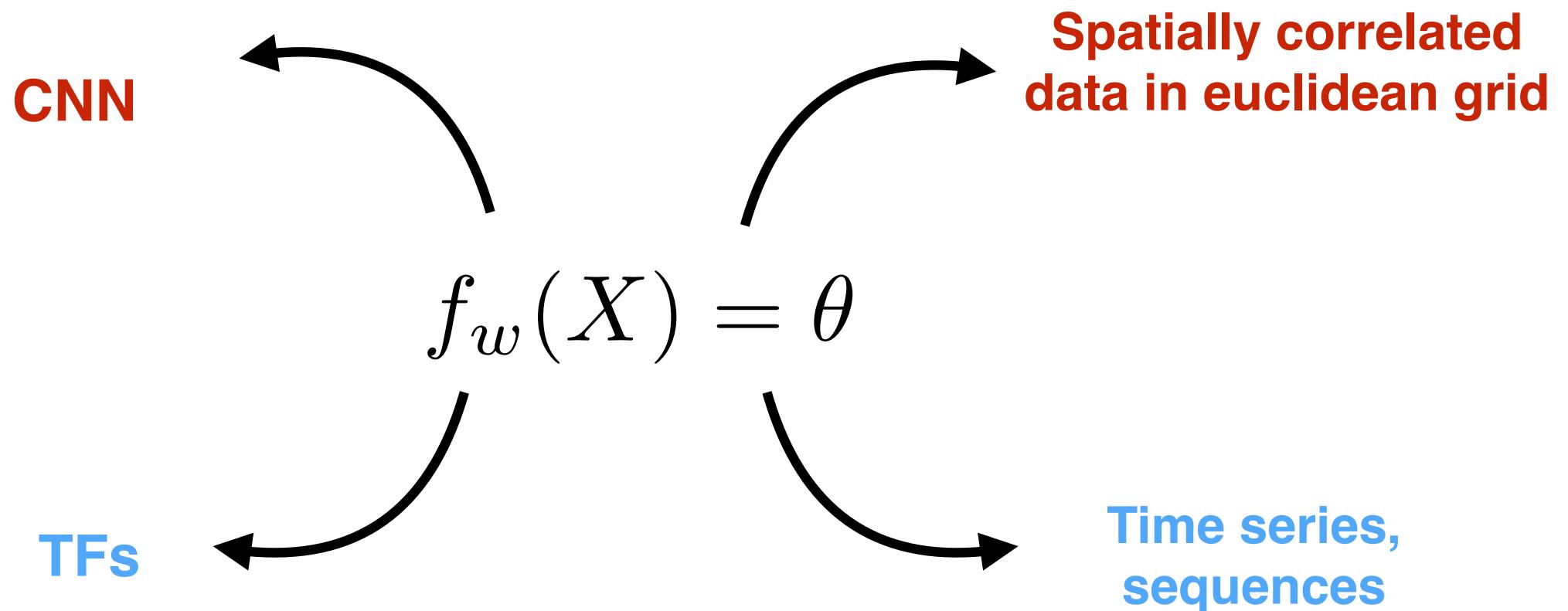


Foundation models and representation learning

Marc Huertas-Company & Hubert Bretonnière

Neural networks as universal feature extractors



Supervised learning pipeline

1. **Gather labeled dataset** (X, θ)
2. **Train** a network to map between labels and data $f_w(x_i) = \theta_i$ for $(x_0, \dots, x_N) \in X_{tr}$
3. **Validate on a test set** X_t
4. **Apply to new unlabelled dataset** X_u

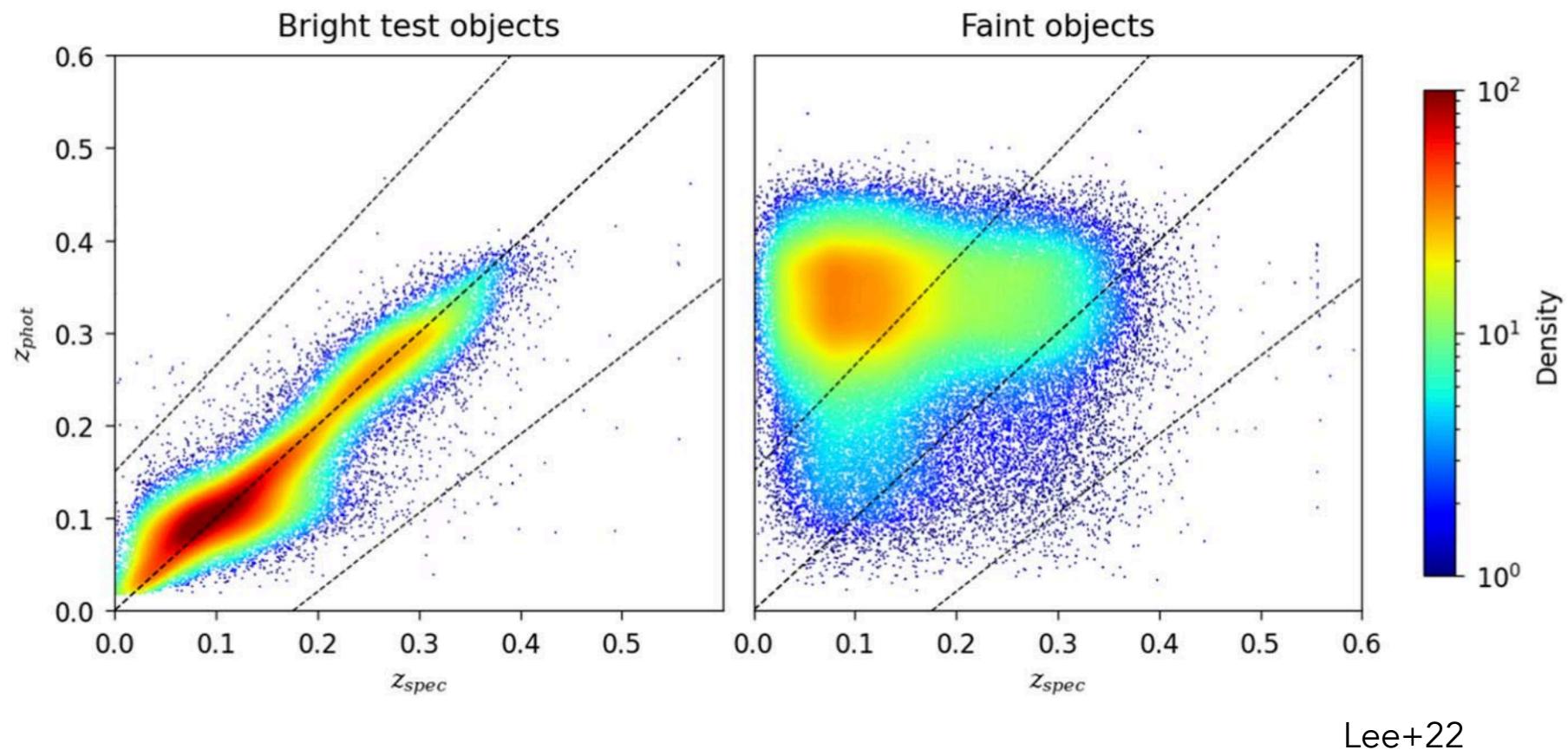
The limits of supervised learning

requires large amount of labeled data
(e.g. galaxy zoo) - the majority of the data is not labelled

limited representation of the information content;
network needs to be retrained for new dataset /
new labels even for similar tasks

limited discovery potential - rare or novel objects
have by definition few labeled examples

photometric redshift estimation is a paradigmatic example

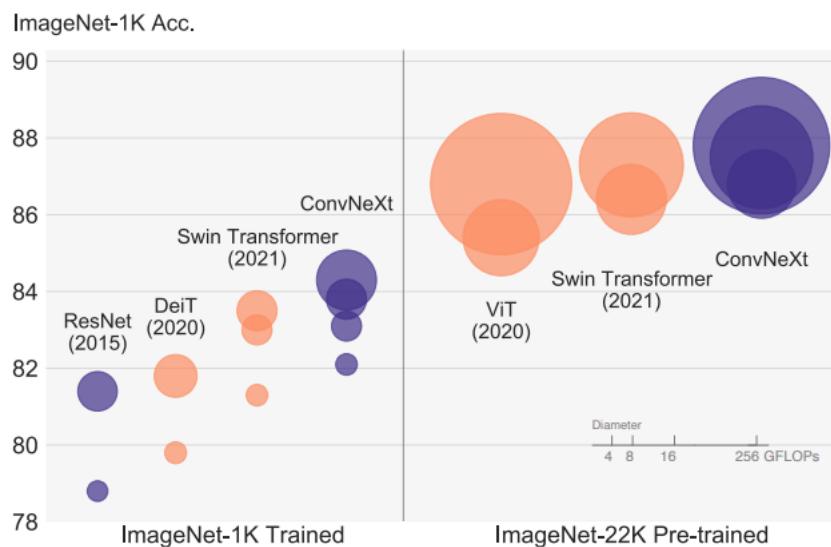


*almost all the extensive literature is focused on mitigating the effects of biased labelling

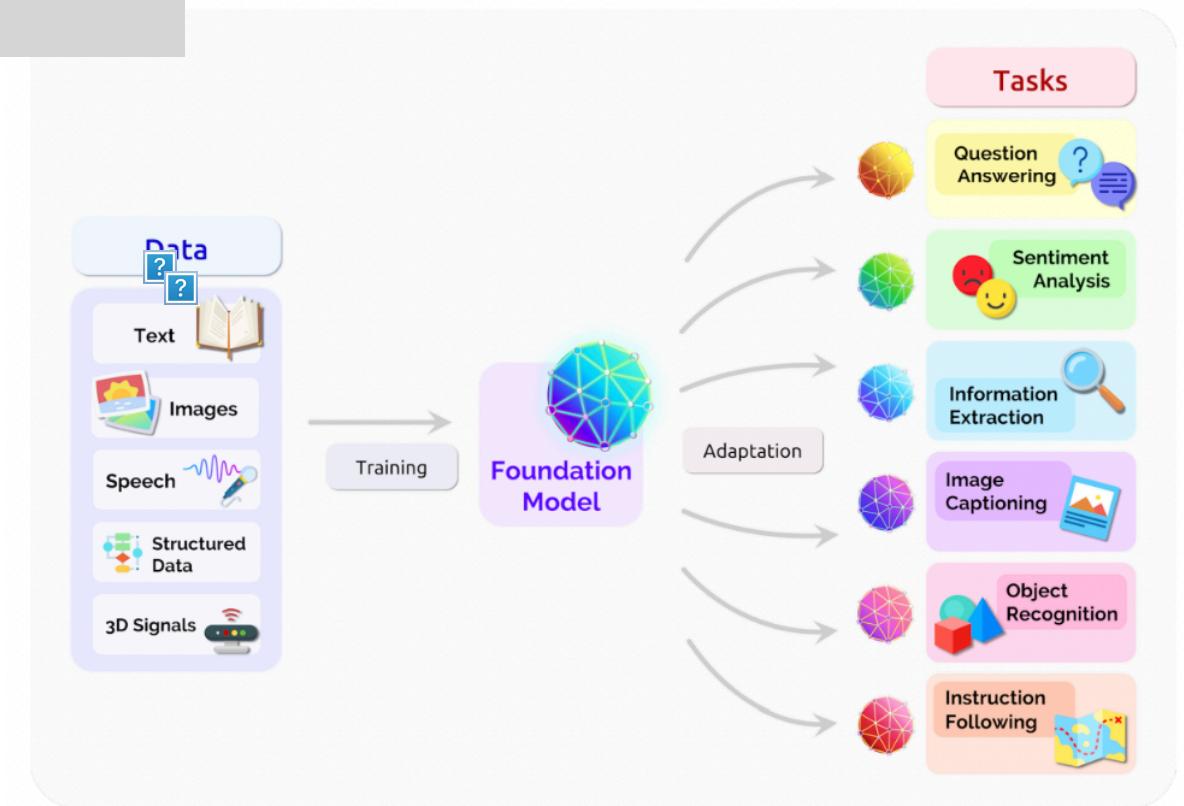
Foundation Models: One model to rule them all

Foundation Model approach:

1. **Pretrain** models on **pretext tasks**, without labels, on very large scale datasets.
2. **Adapt** pretrained models to **downstream tasks**.

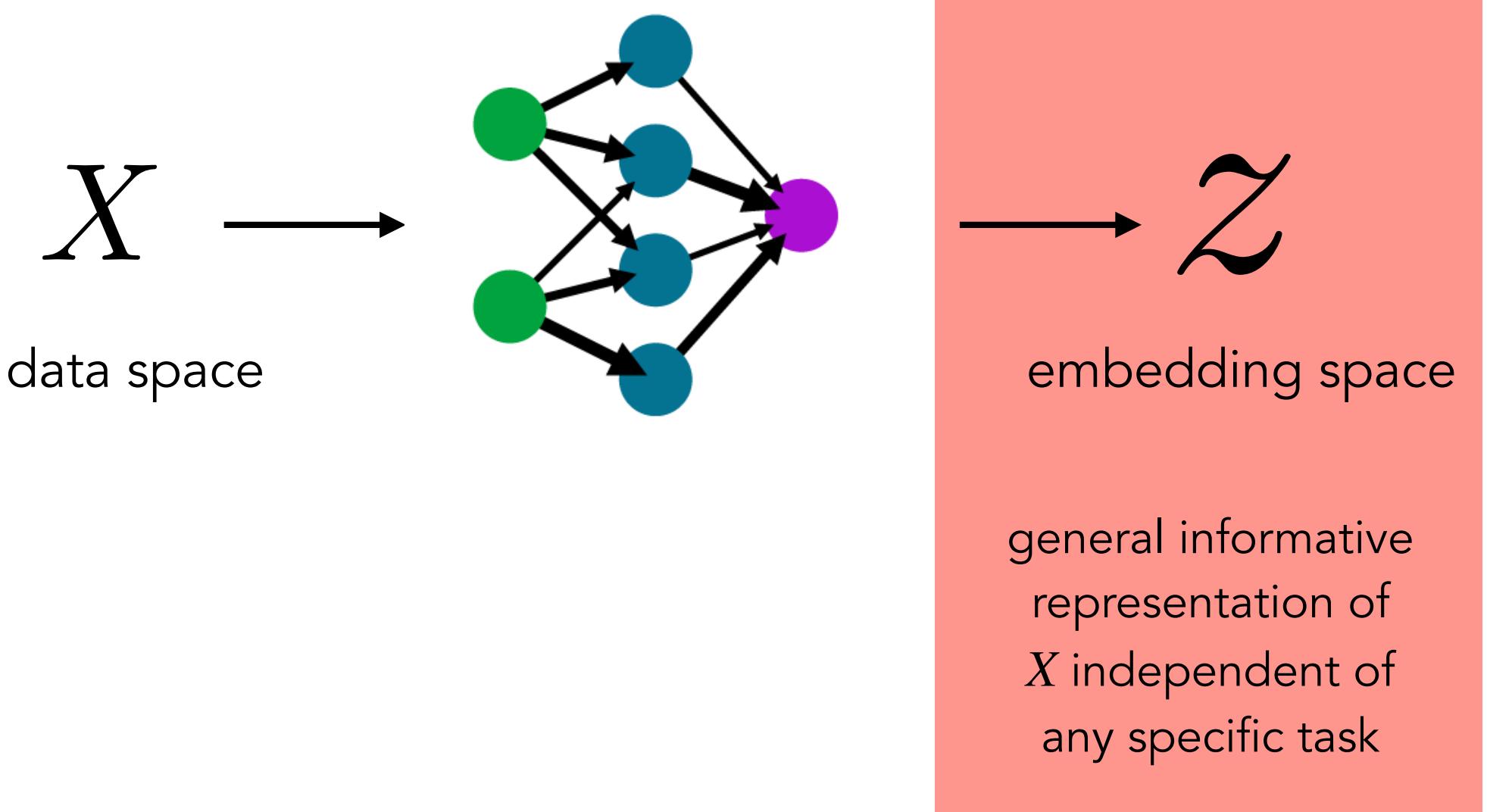


Liu+20

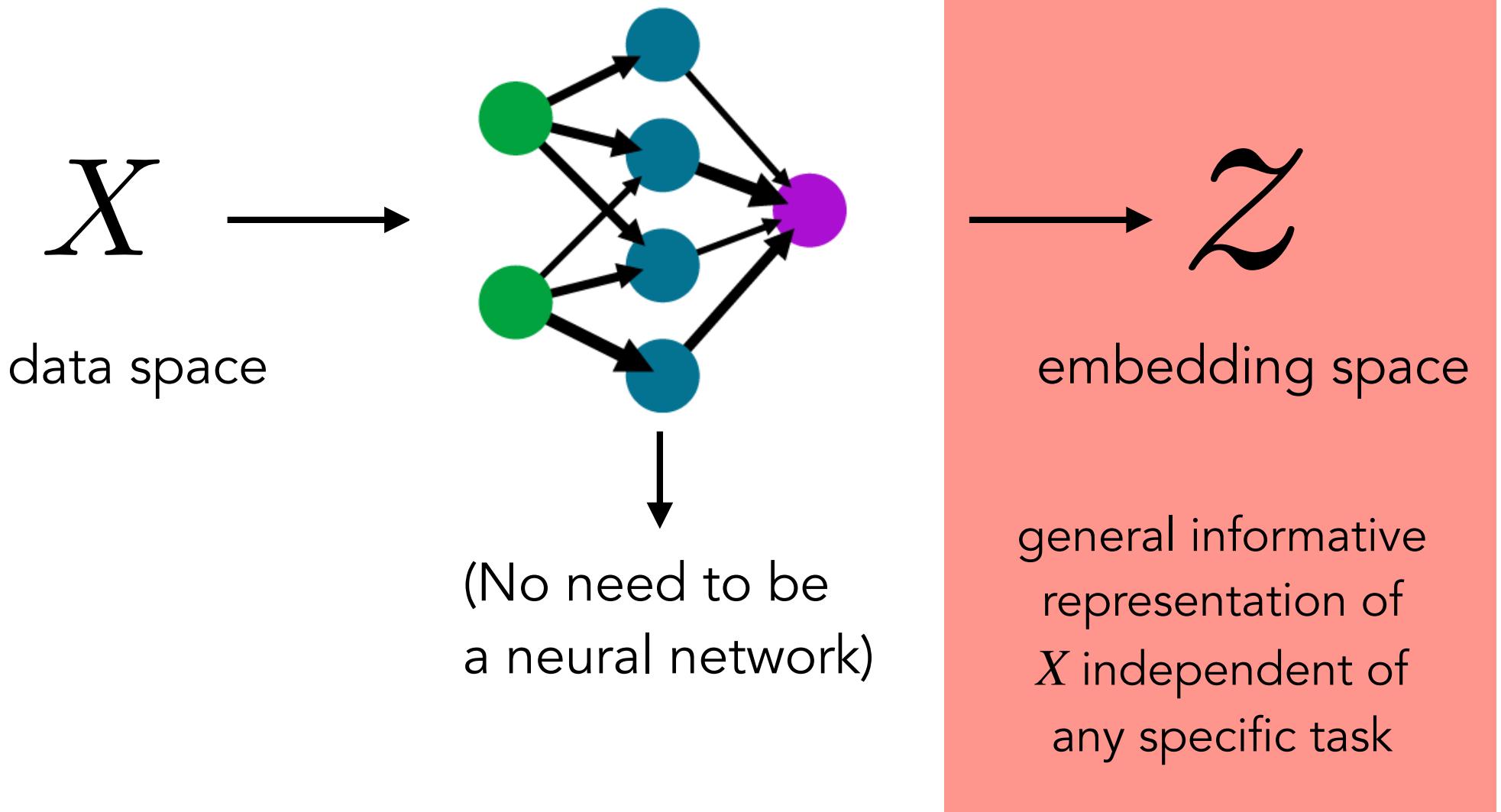


Bommasani+21

representation learning - use neural networks as a universal embedder

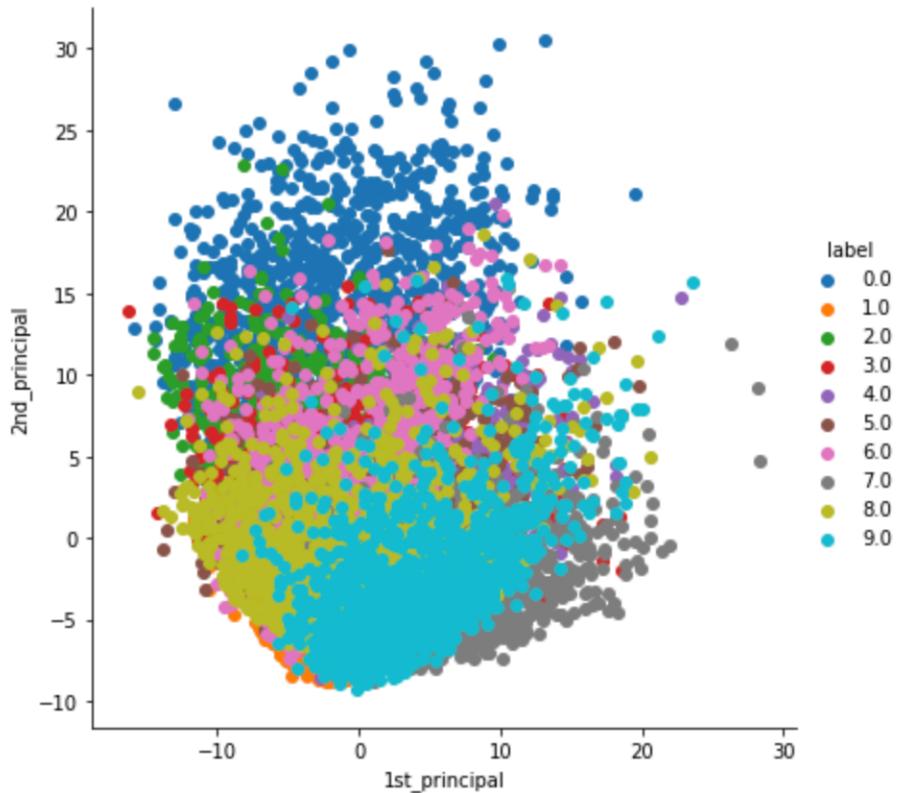
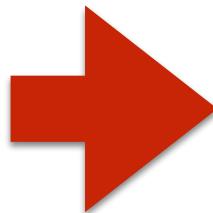
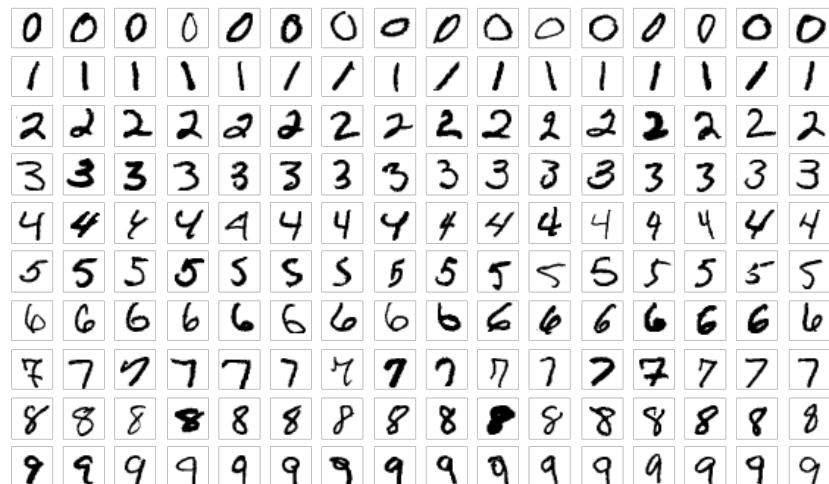


representation learning - use neural networks as a universal embedder

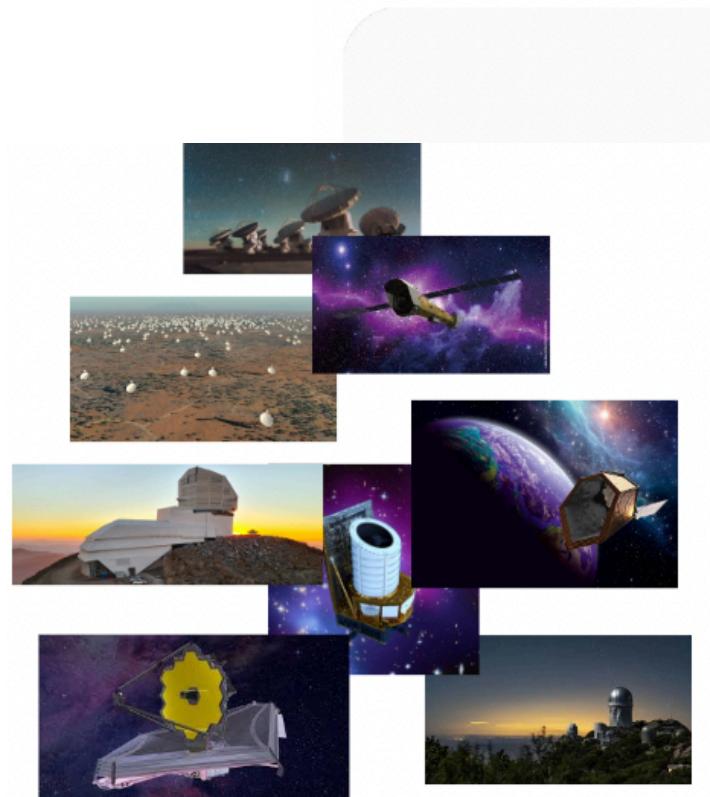


limitations of PCA

assumes data follows a multivariate gaussian distribution



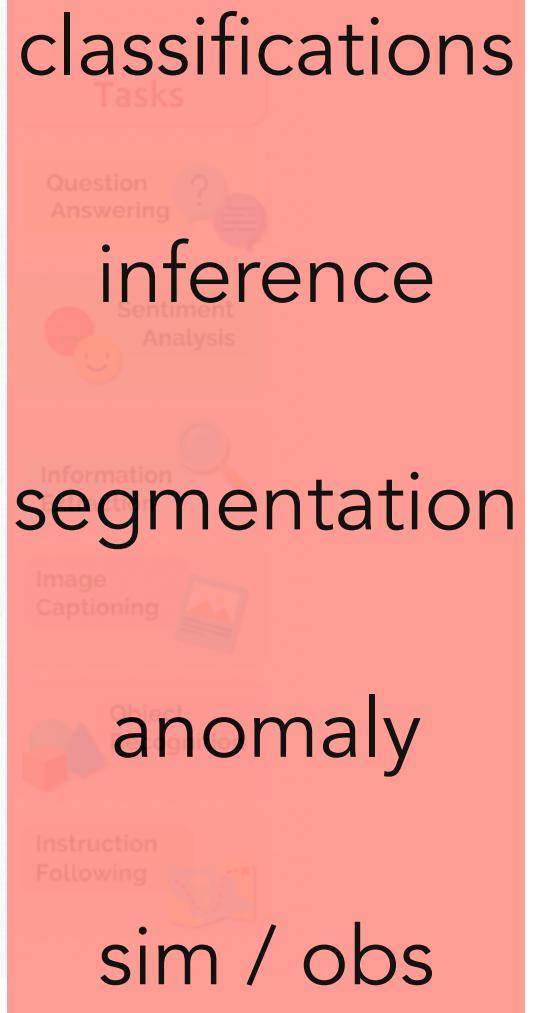
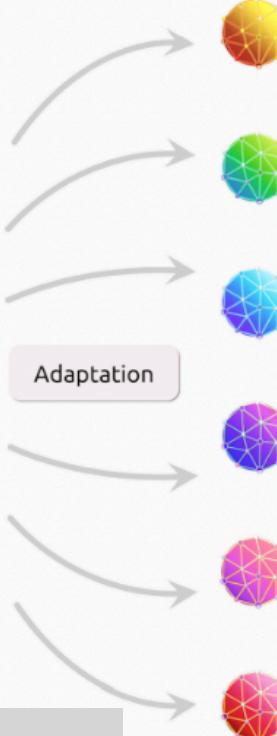
what does this mean for astro?



Training

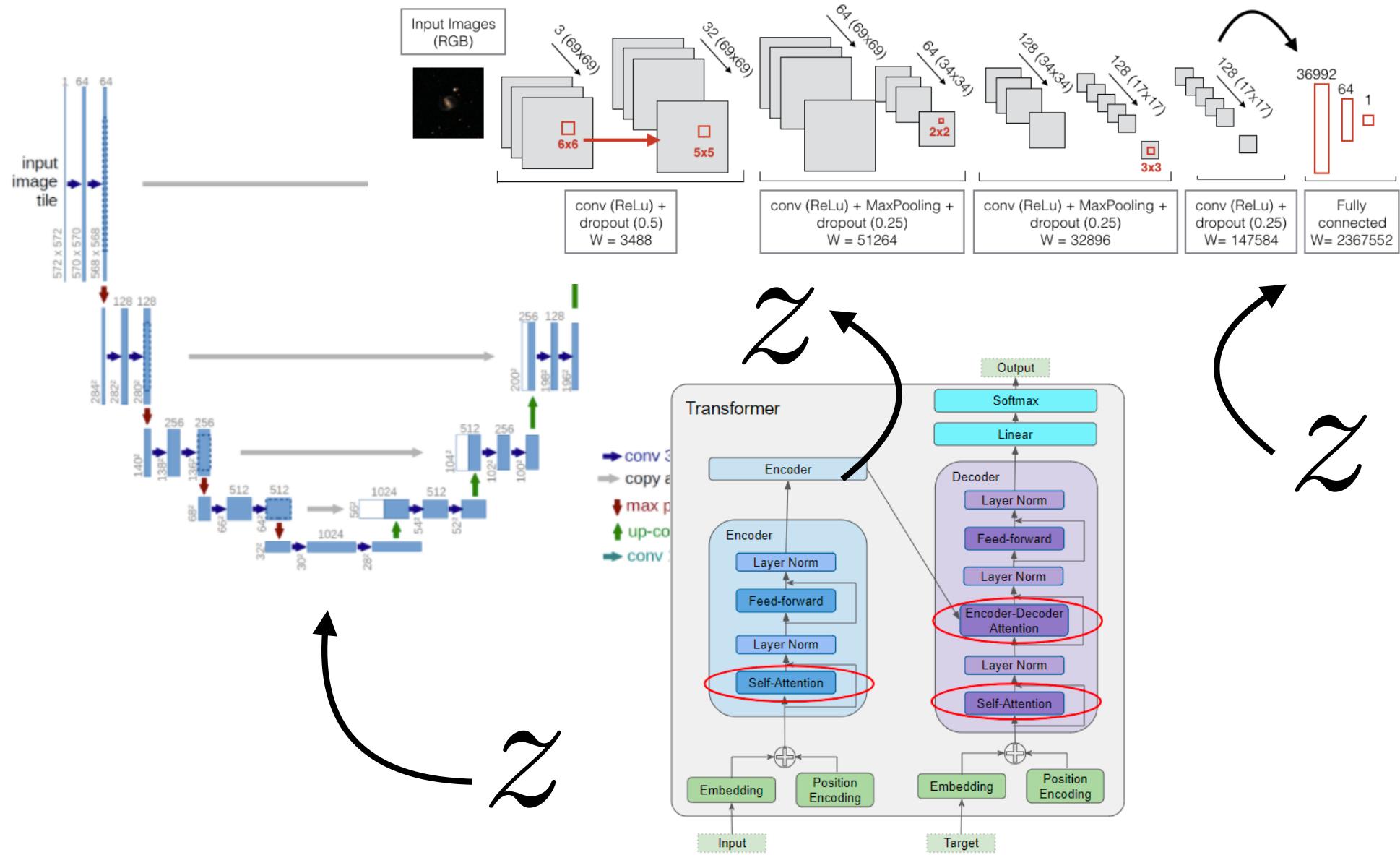


Foundation
Model

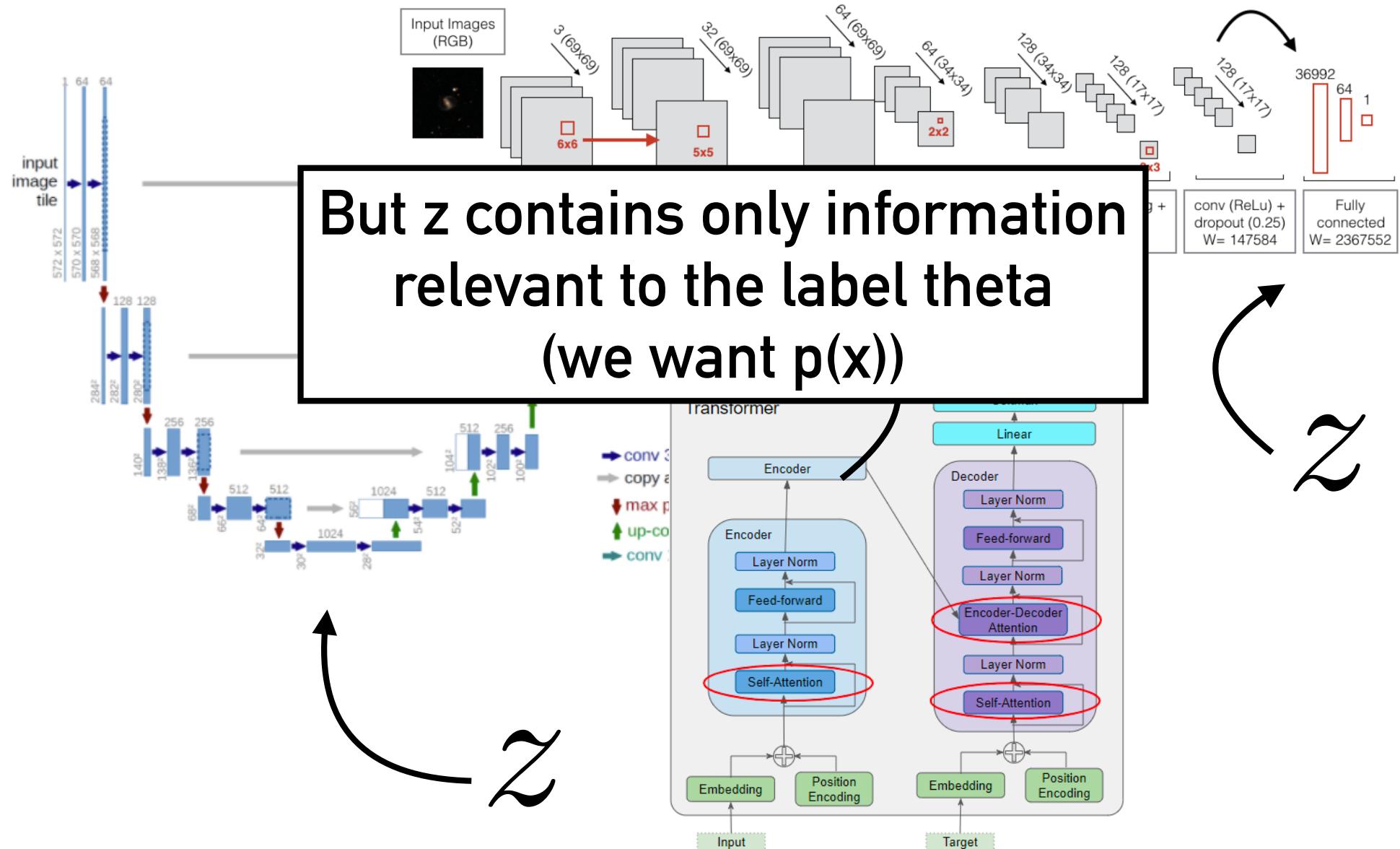


1. never again train a neural network from scratch
2. no labels required (or very few)
3. integrates all data beyond images
4. very simple analysis tools for downstream tasks

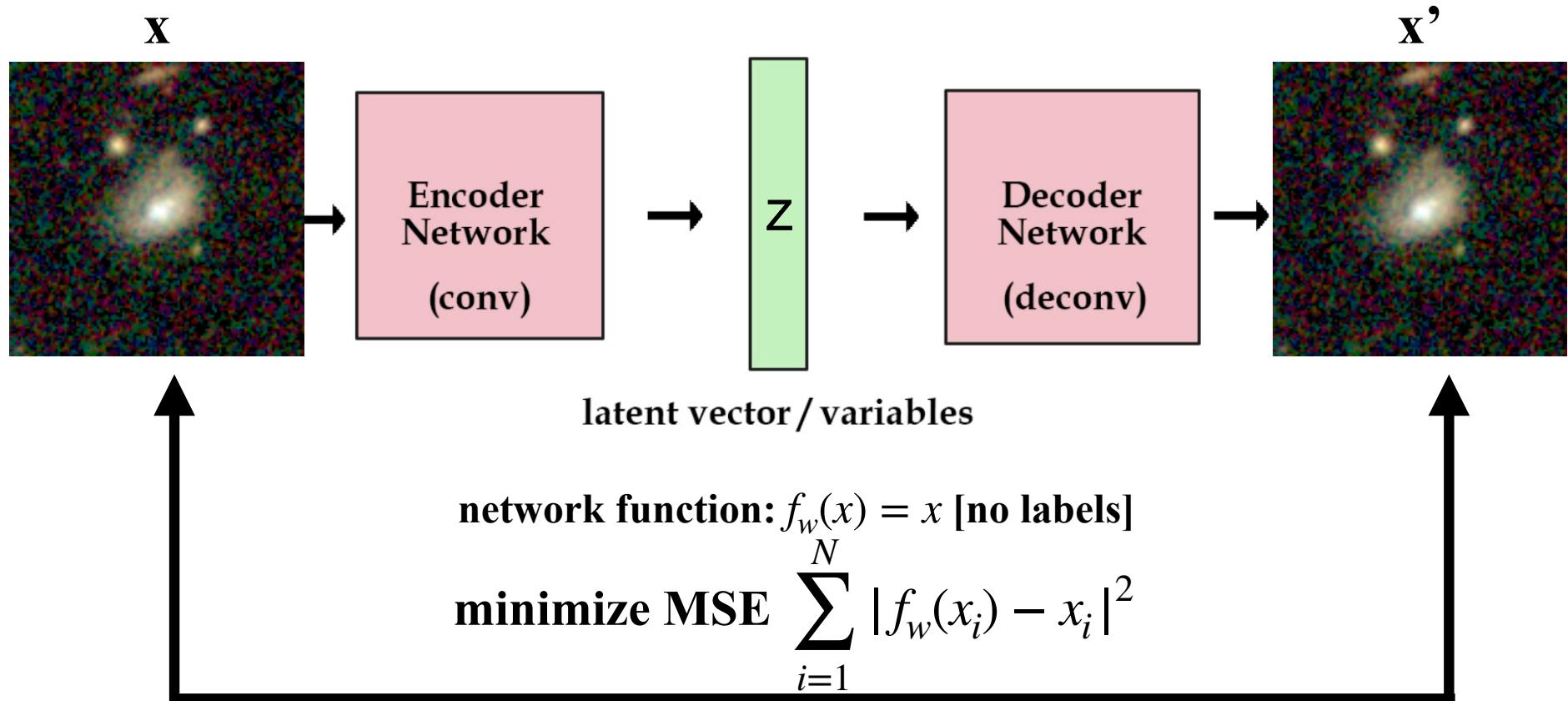
any supervised network obtains a non-linear representation of the data



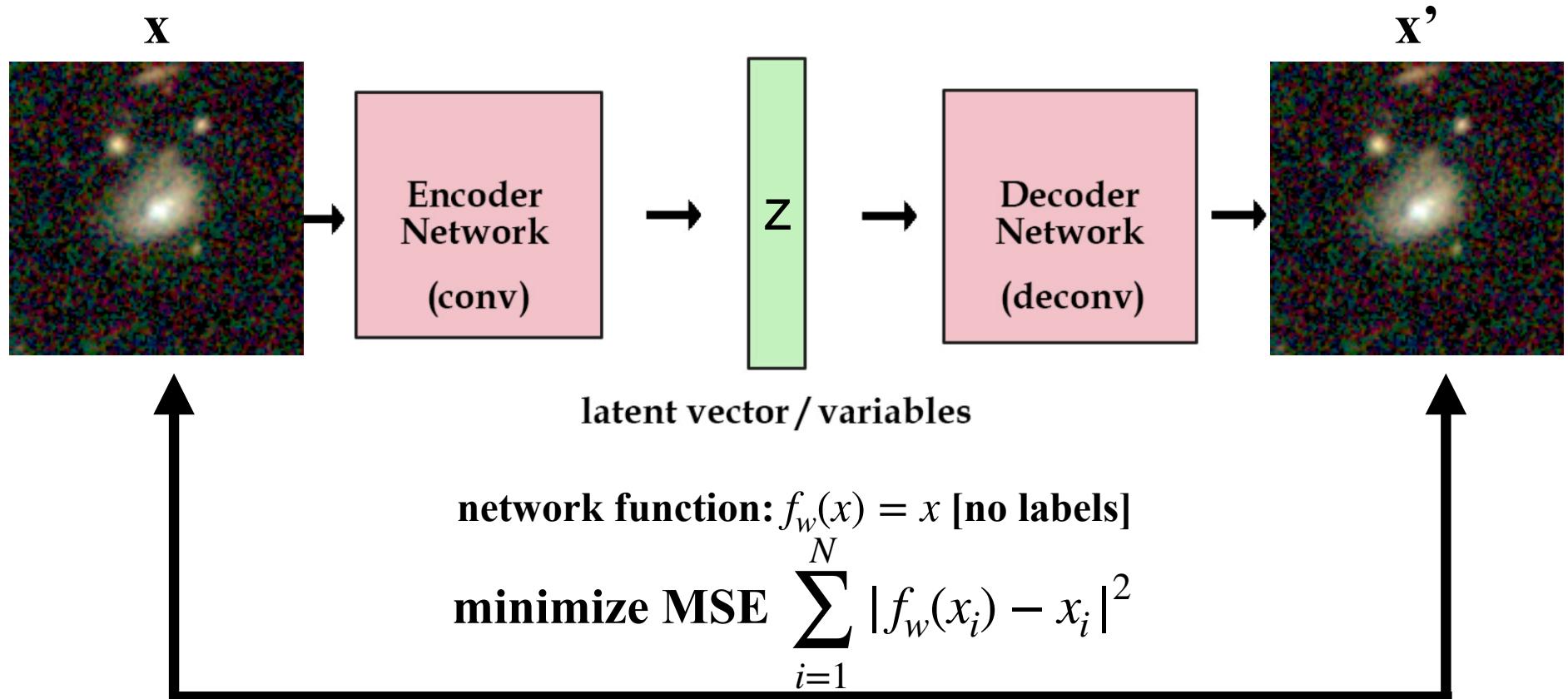
any supervised network obtains a non-linear representation of the data



auto-encoder



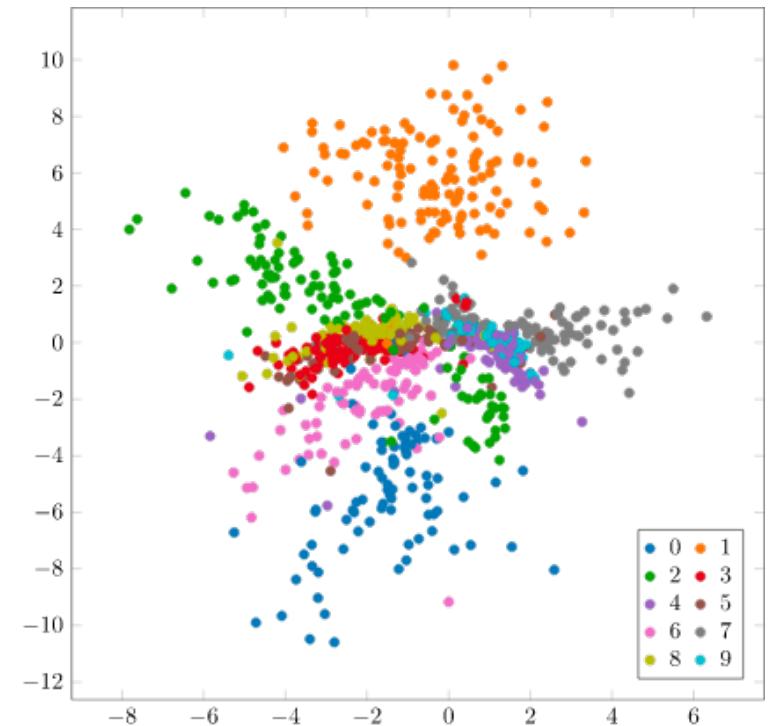
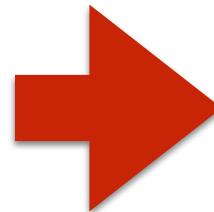
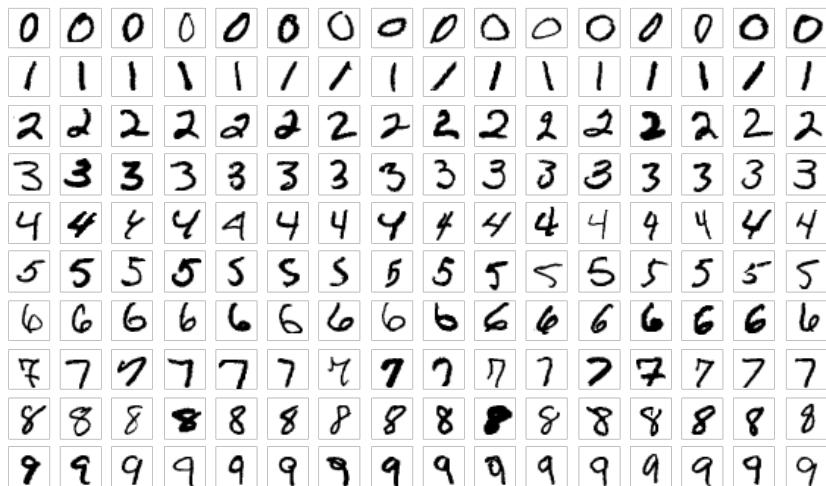
auto-encoder



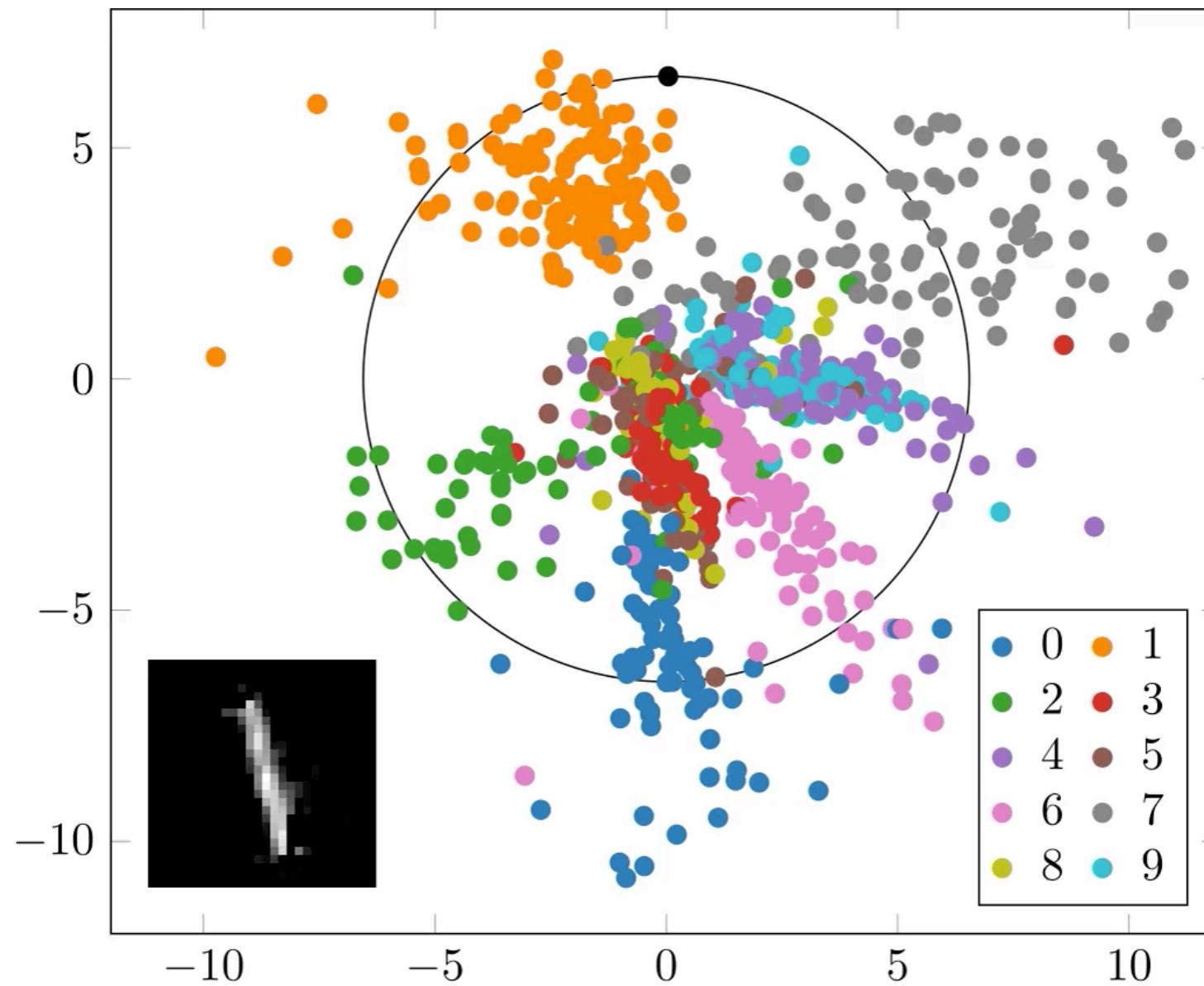
By creating a bottleneck, we encourage z to encode a lower dimensional informative representation of the data

* **question:** what would happen if we train an auto encoder with no activation functions - only linear operations?

AUTOENCODER REPRESENTATION OF MNIST



AUTOENCODER REPRESENTATION OF MNIST



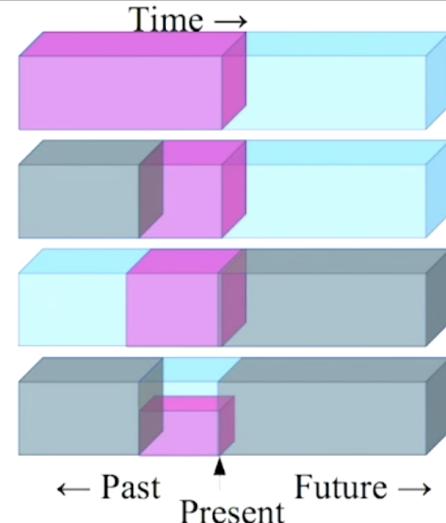
self-supervised learning

Another alternative is to invent a target to obtain z

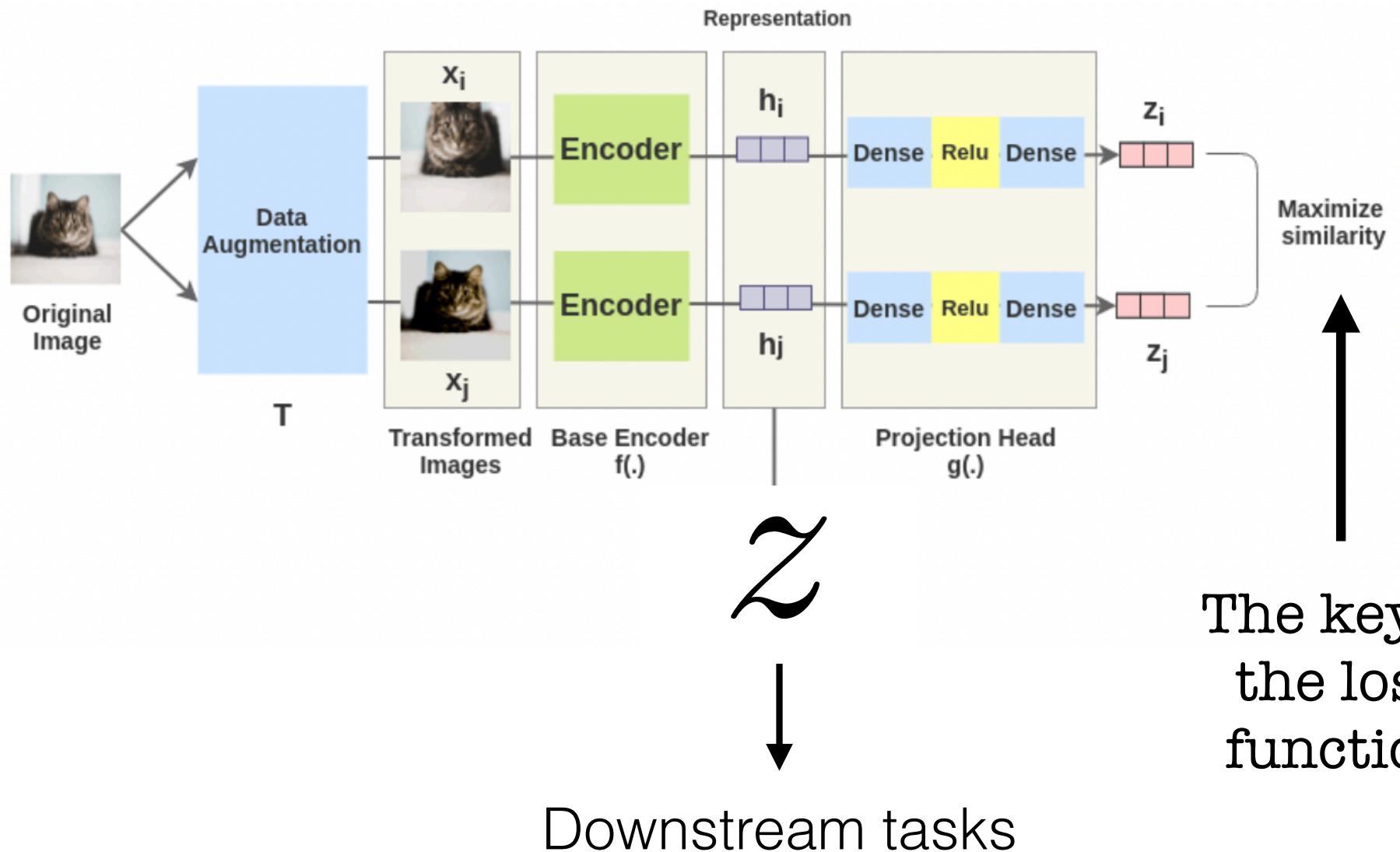
$$f_w(X) = \hat{\theta}$$

(Still using a discriminative model, with a general target, and use the latent z as a representation of the data)

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ Pretend there is a part of the input you don't know and predict that.



contrastive self-supervised learning



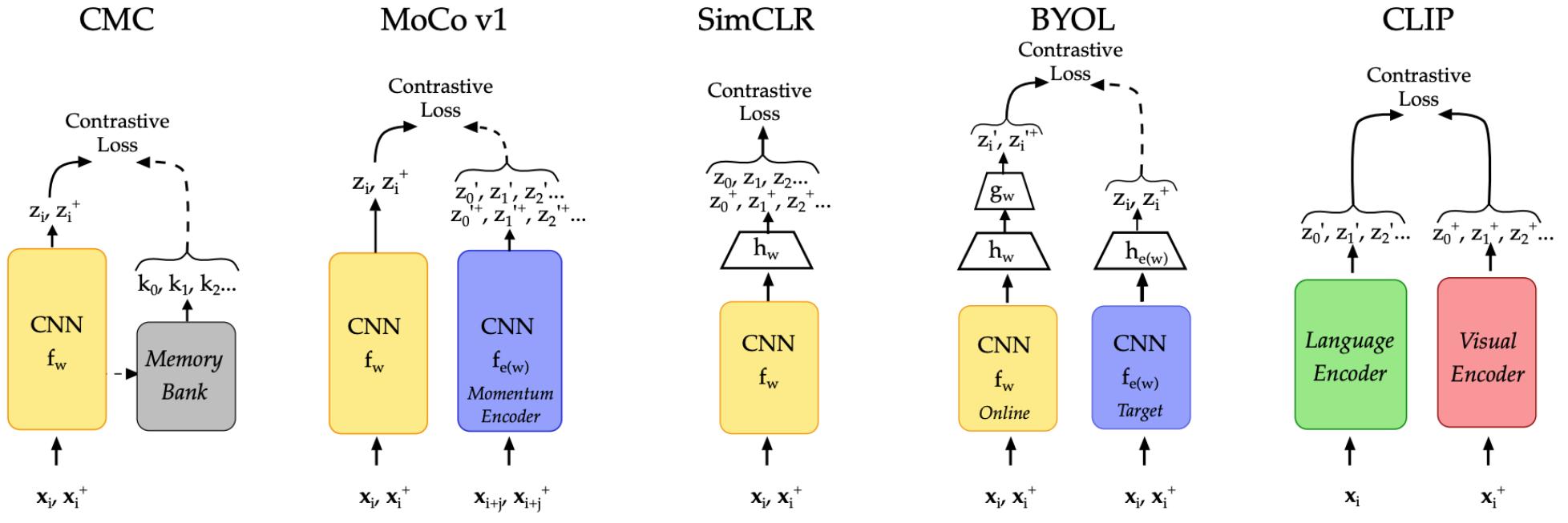
Similarly between
two representations of positive pairs

The contrastive loss:

$$l_{i,j} = - \log \frac{\exp(\langle z_i, z_j \rangle / h)}{\sum_{k=1, k \neq i}^{2N} \exp(\langle z_i, z_k \rangle / h)},$$

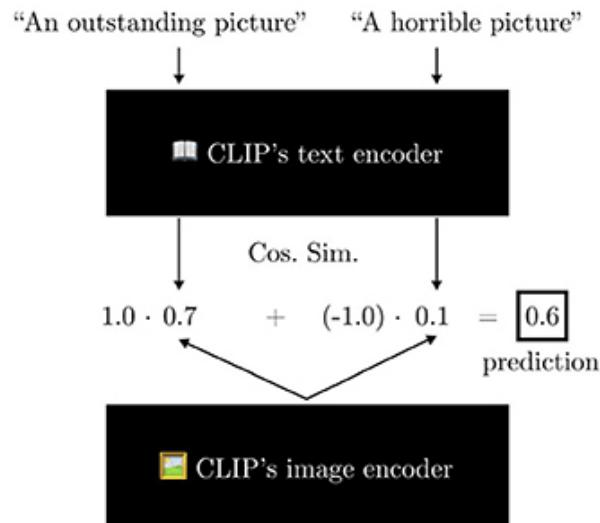
Sum of all similarities between
negative pairs

contrastive learning

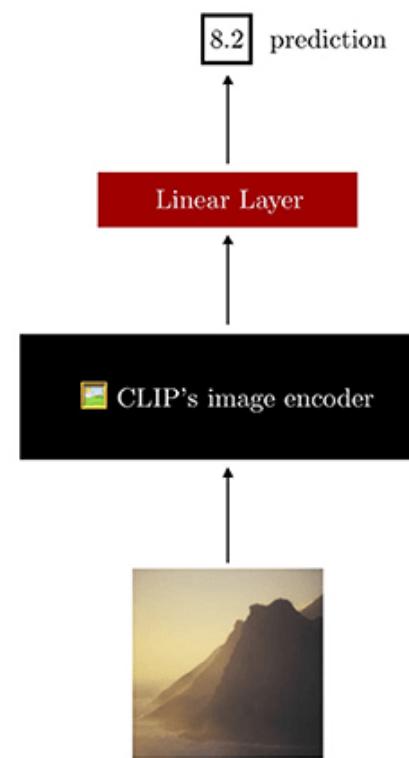


MHC+23

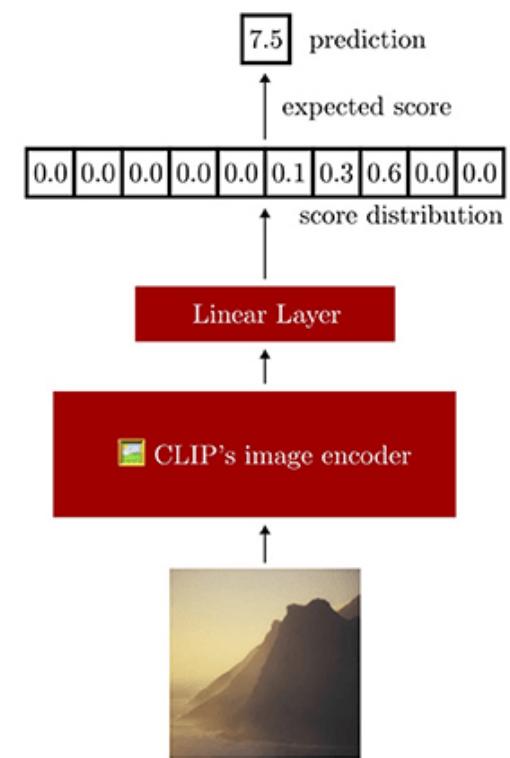
Prompting



Linear Probing



Fine-Tuning



█ component with **frozen weights**
█ component with **trainable weights**

computational complexity

Zhai+22

autoregressive

$$P_{\theta}(x) = \prod_{t=2}^T P_{\theta}(x_t | \{x_{<t}\})$$

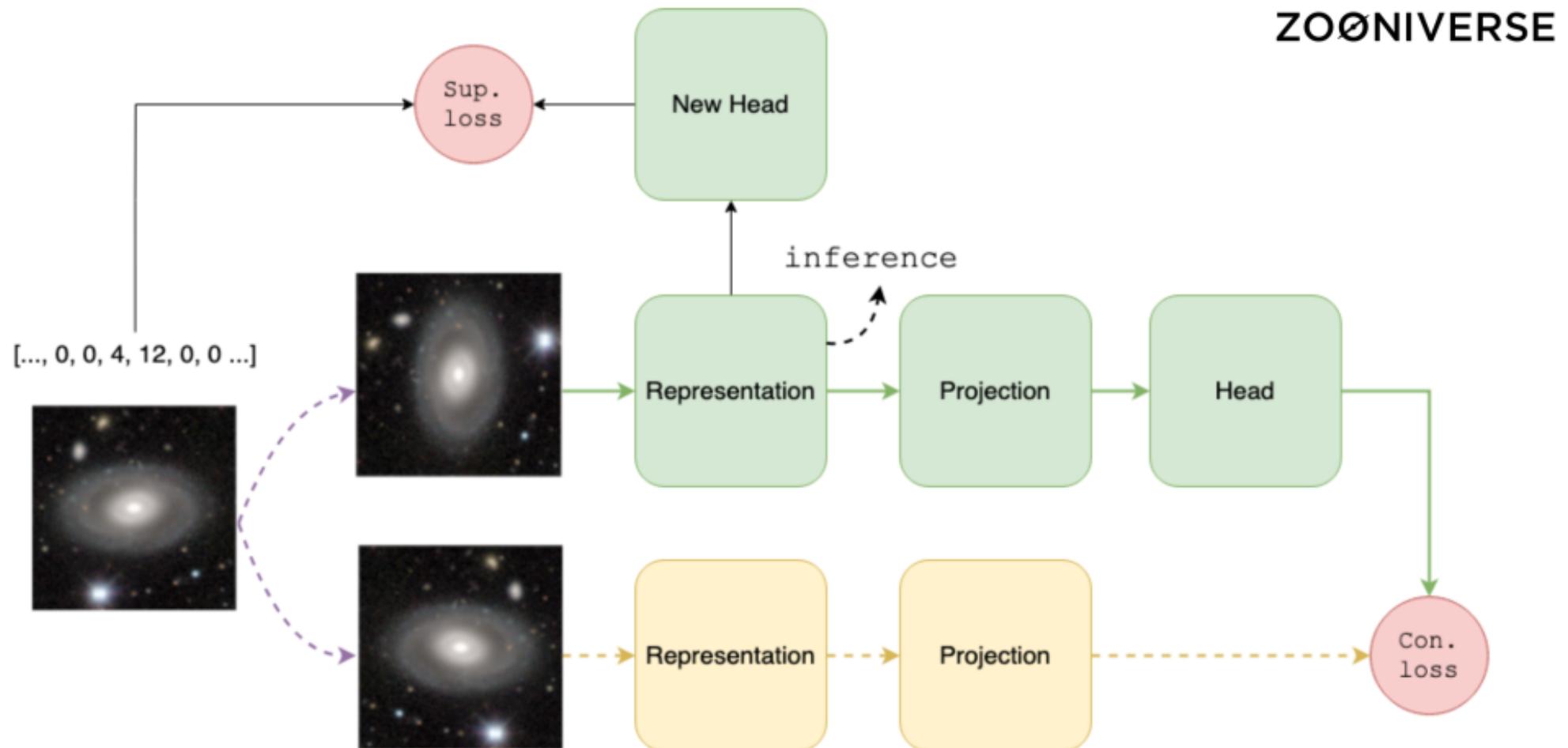
$$= P_{\theta_2}(x_2 | \{x_1\}) P_{\theta_3}(x_3 | \{x_1, x_2\}) \dots P_{\theta_T}(x_T | \{x_1, \dots, x_{T-1}\})$$

$$\begin{matrix} o \leftarrow \{\bullet\} \\ x_2 \quad x_1 \end{matrix}$$

$$\begin{matrix} o \leftarrow \{\bullet \quad \bullet\} \\ x_3 \quad x_1 \quad x_2 \end{matrix}$$

$$\begin{matrix} o \leftarrow \{\bullet \quad \dots \quad \bullet\} \\ x_T \quad x_1 \quad x_{T-1} \end{matrix}$$

ZOOBOT: an hybrid foundation model

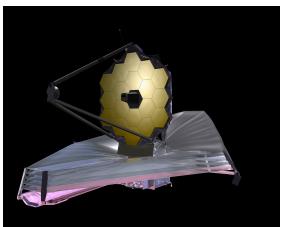


Adding a new supervised head to guide the contrastive representation

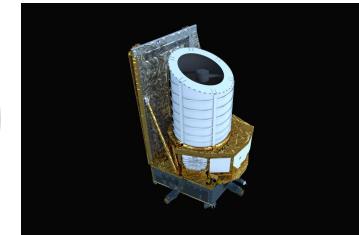
walmsley+



4 years



3 months



0 months



Credit: ESA/Euclid/Euclid Consortium/NASA, image processing by M. Walmsley, M. Huertas-Company, J.-C. Cuillandre



EC: M. Walmsley. MHC+25

Everything else

Smooth or featured?

Answer

Smooth

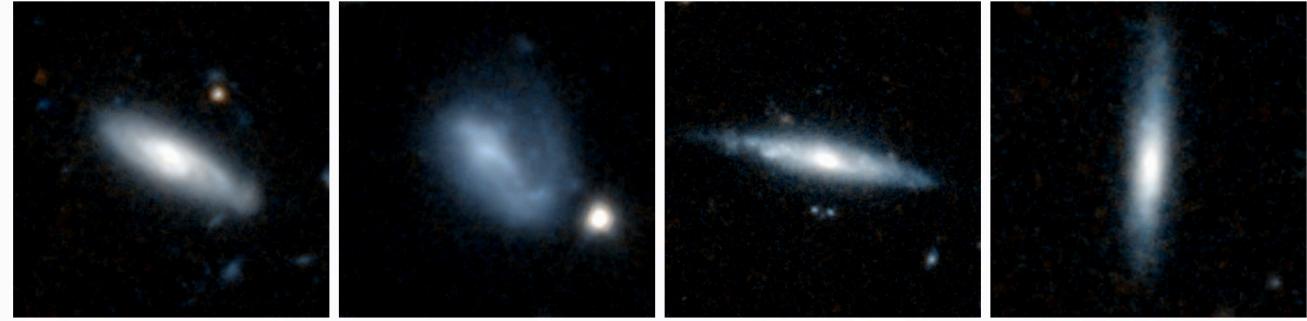


Vote Fraction



0.00

1.00



Disk edge on?

Answer

Yes

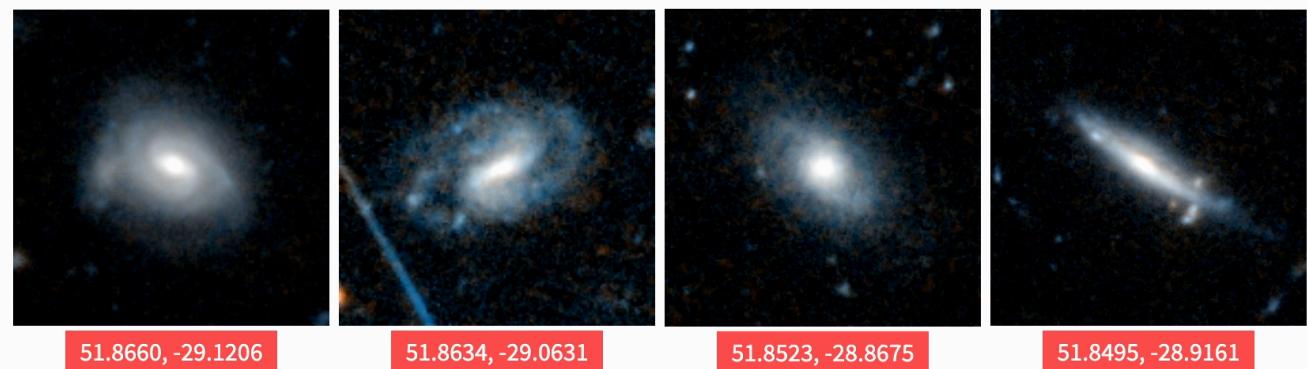


Vote Fraction



0.00

1.00



Bar?

Answer

Strong

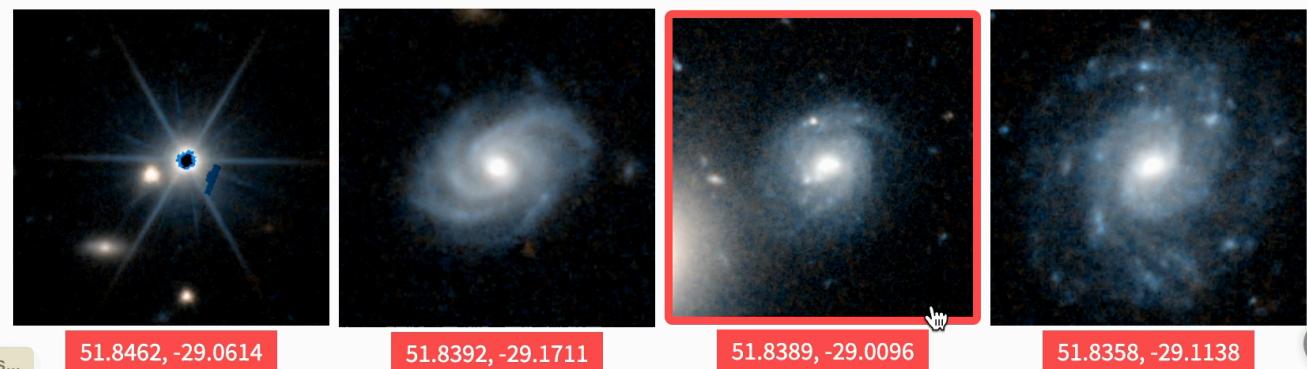


Vote Fraction



0.00

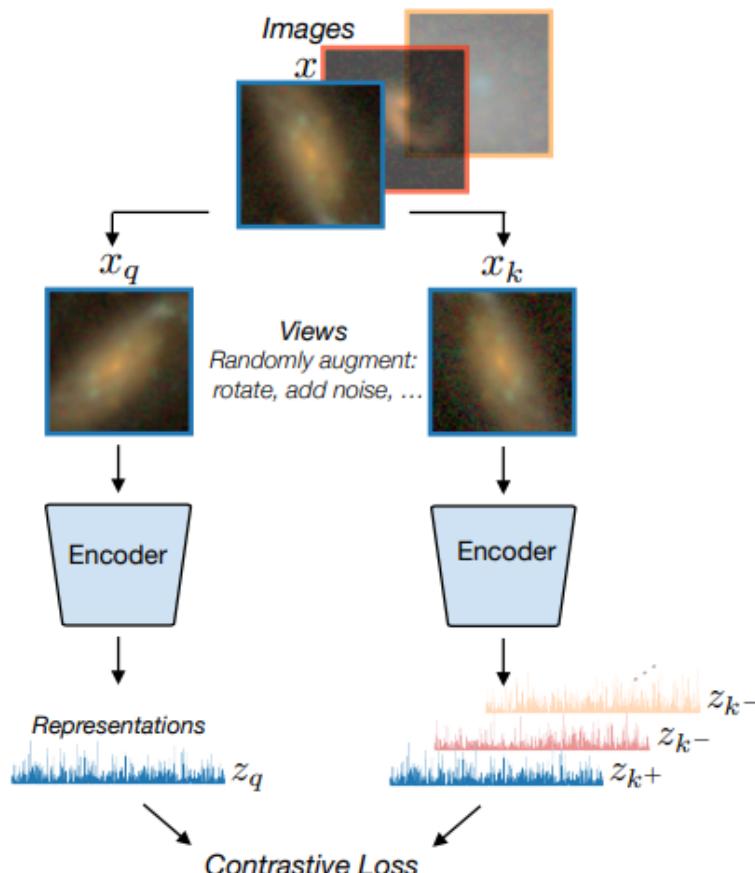
1.00



first applications show good generalisation to downstream tasks

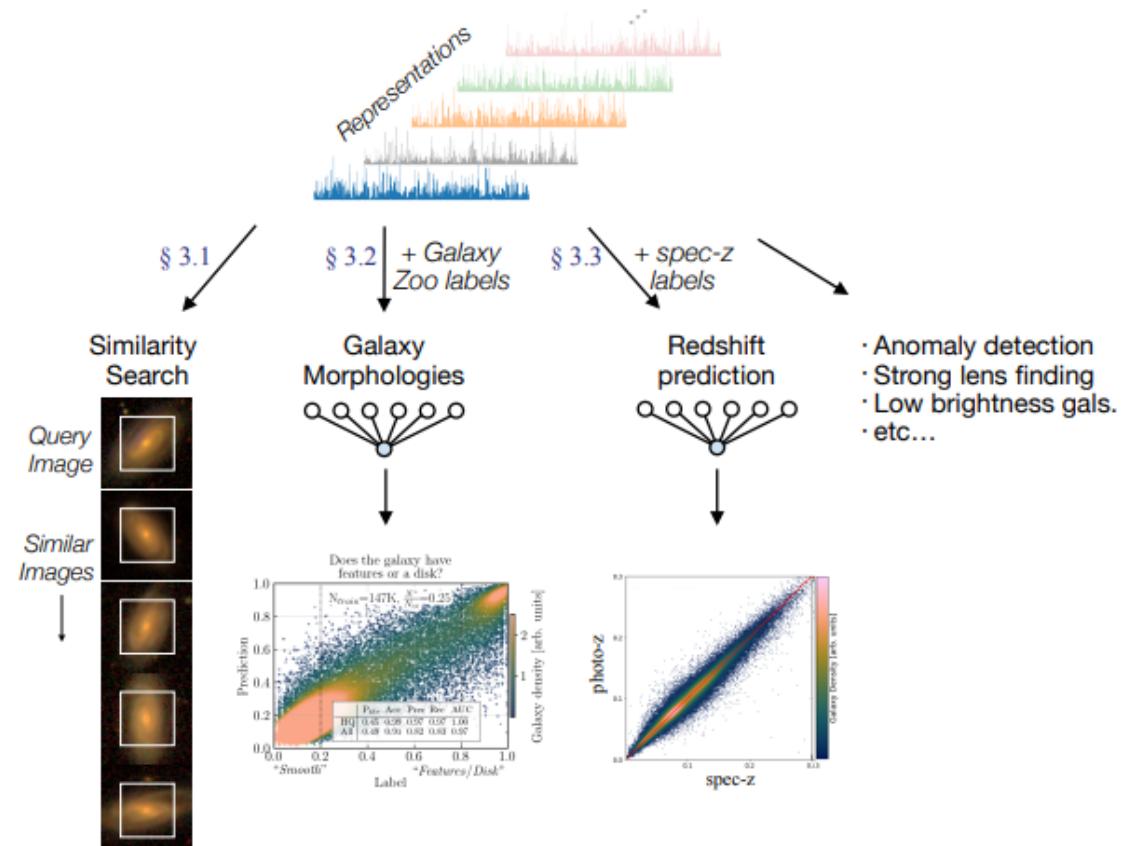
1. Self-supervised contrastive representation learning

Learn representations in an unsupervised manner

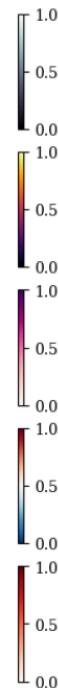
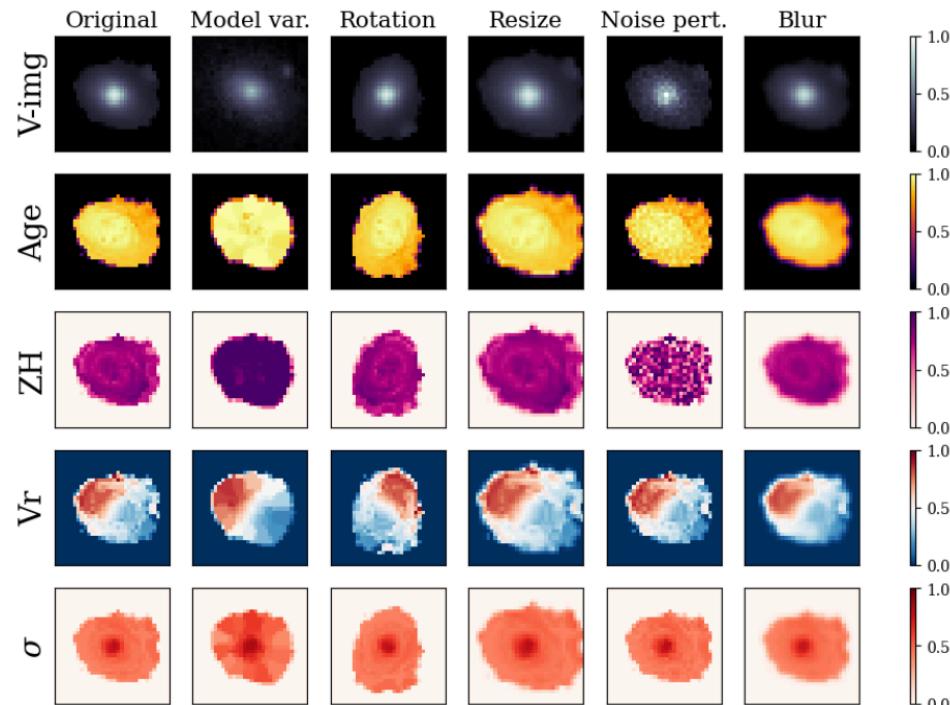


2. Downstream tasks

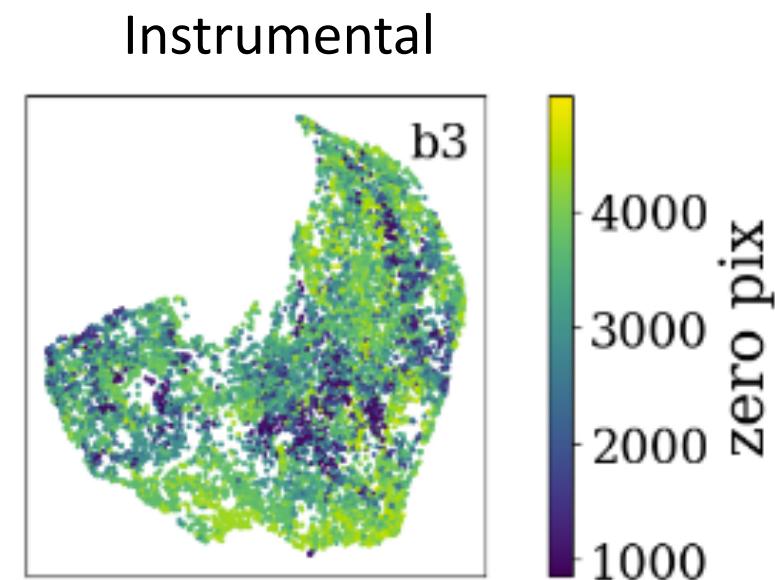
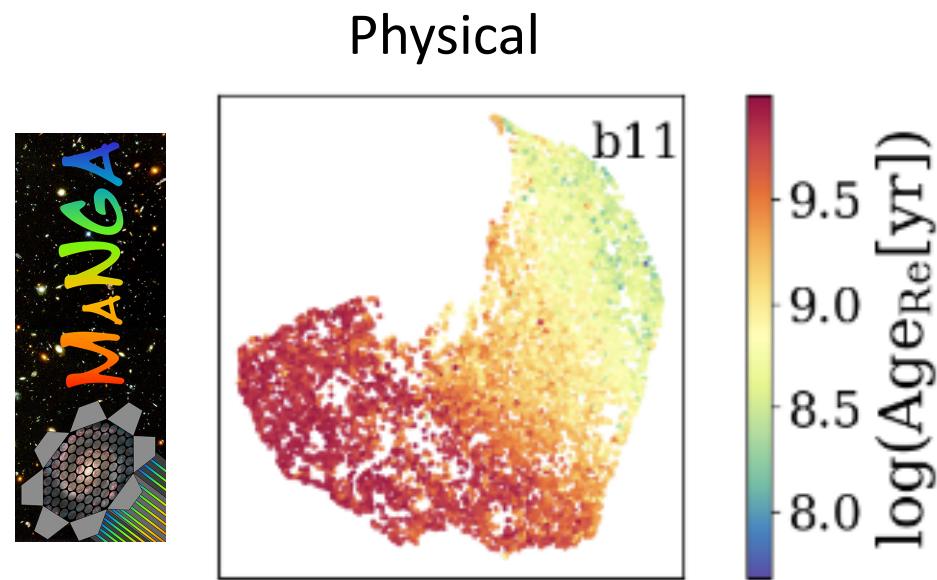
Use representations for a variety of applications



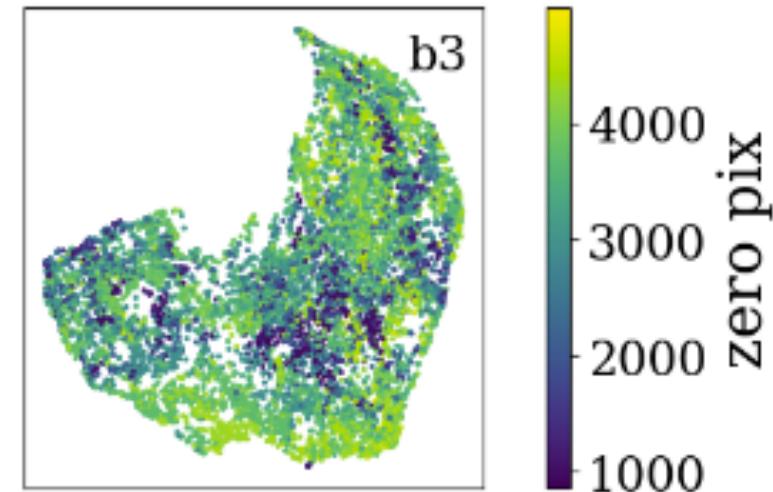
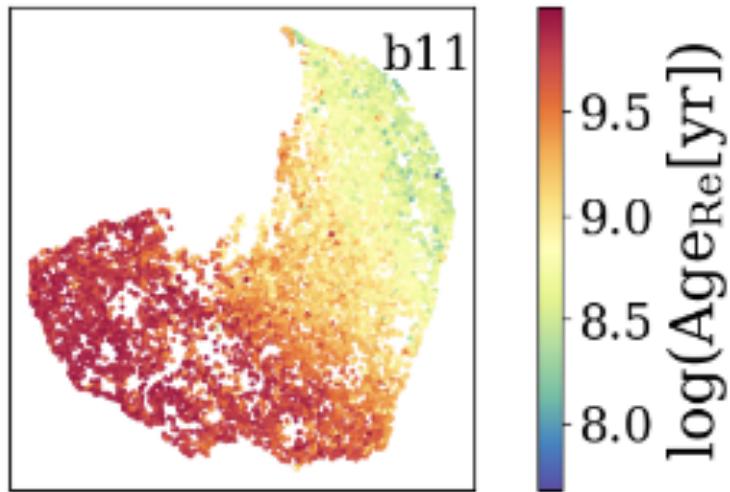
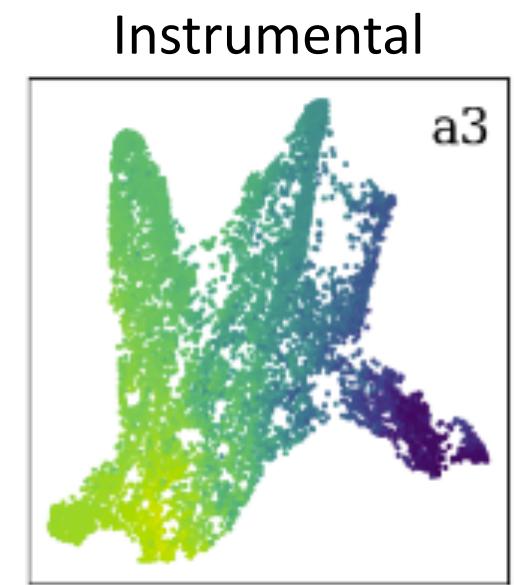
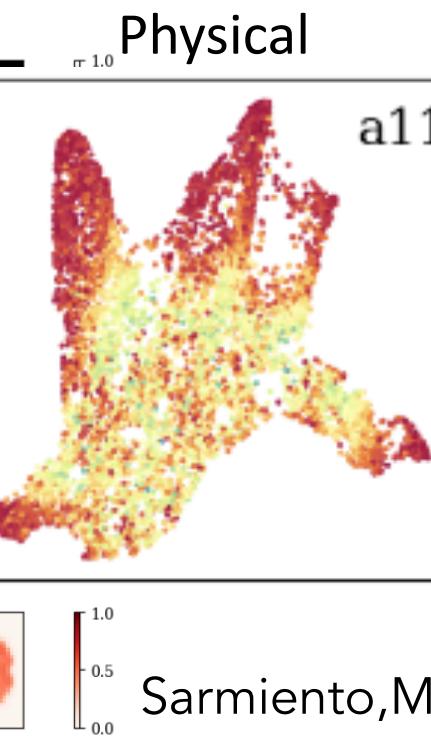
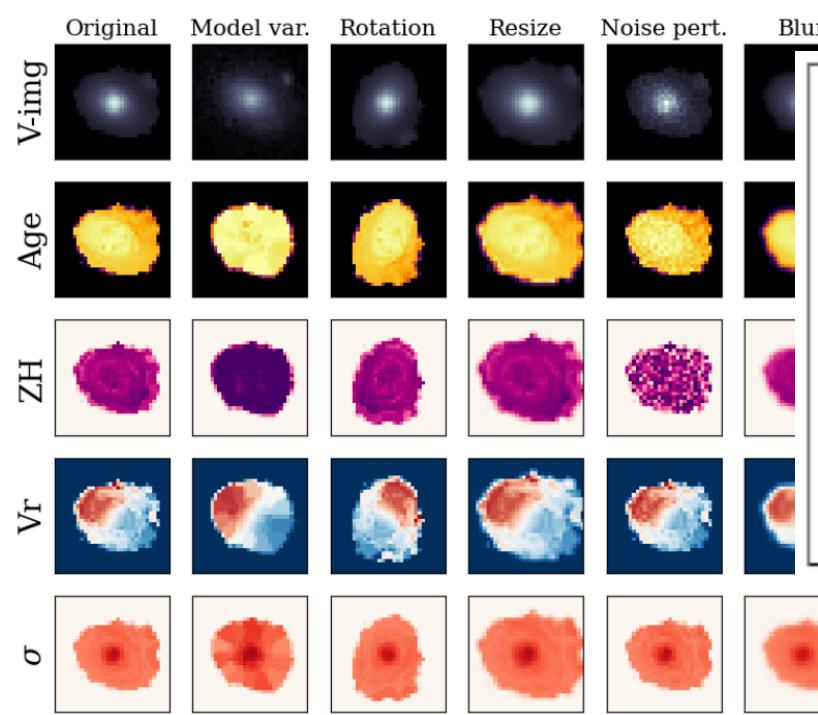
first applications show informative embeddings



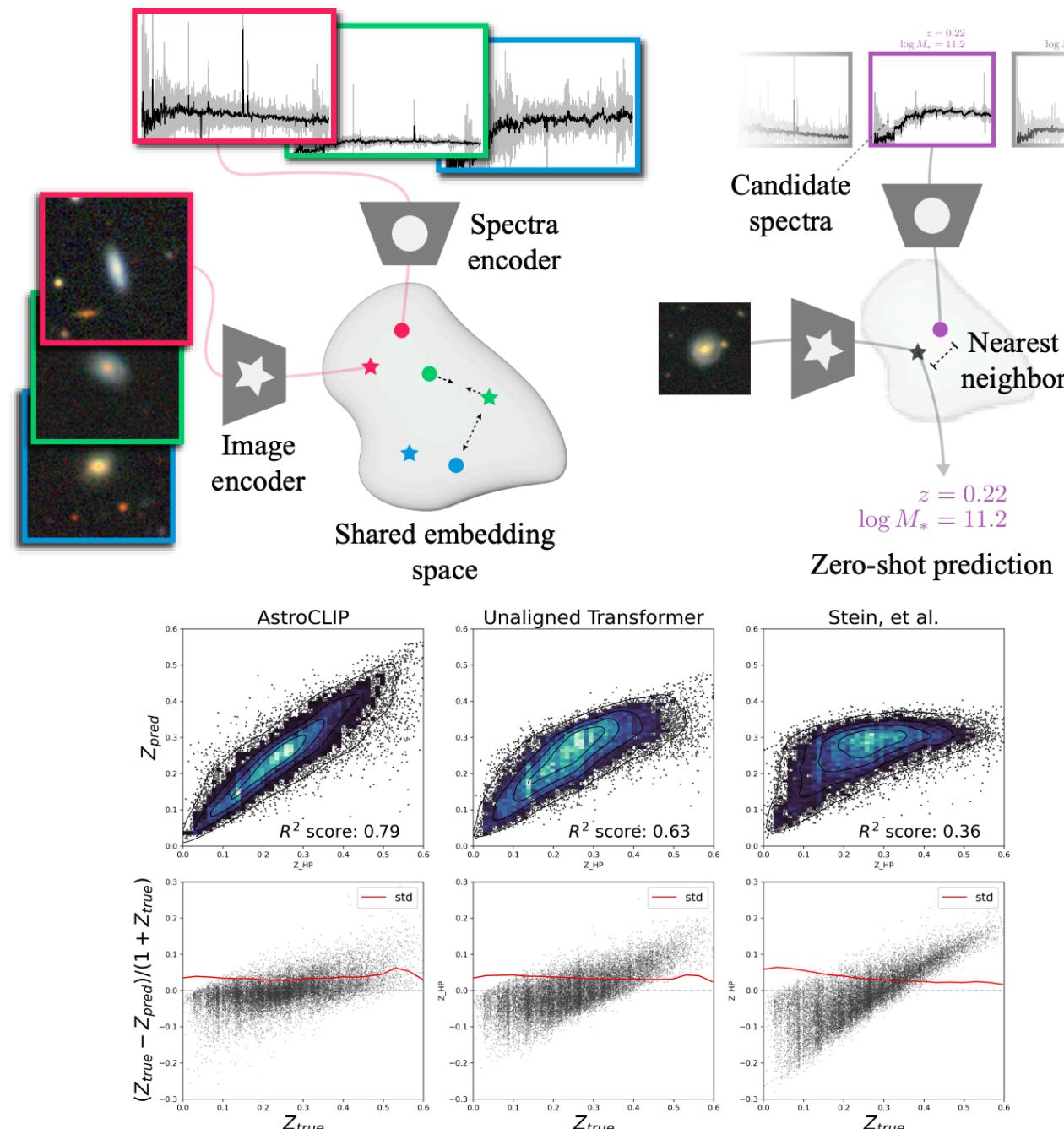
Sarmiento,MHC+21



first applications show informative embeddings

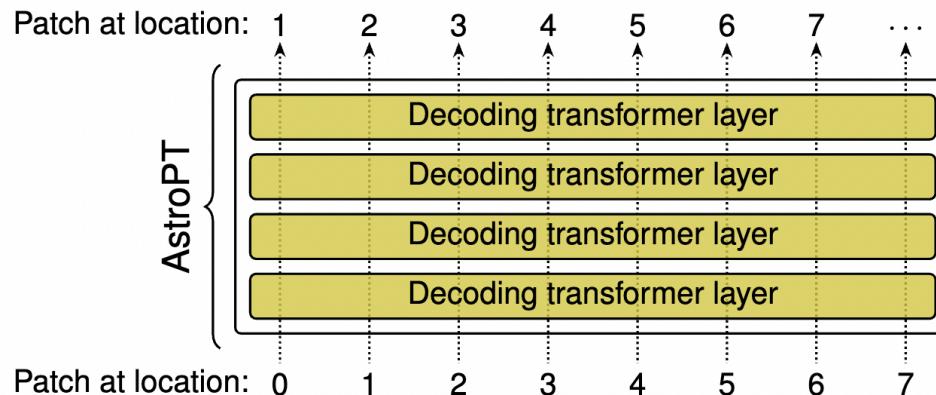
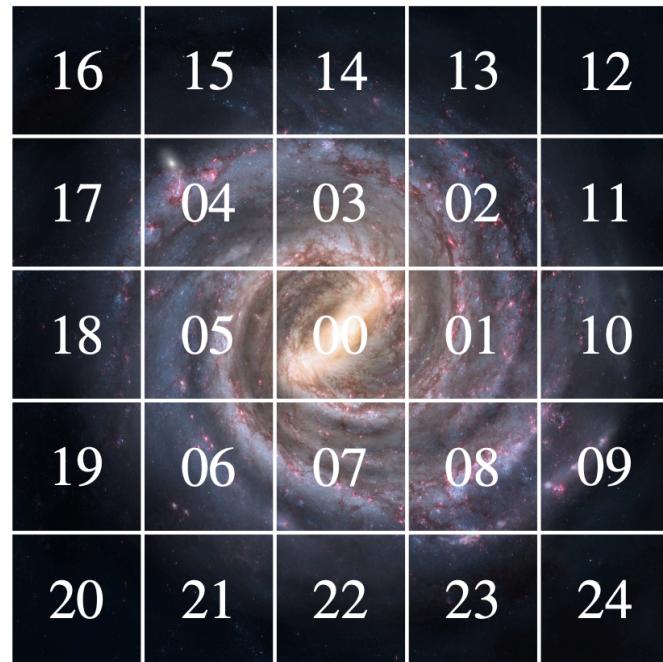


straightforward generalization to two modalities

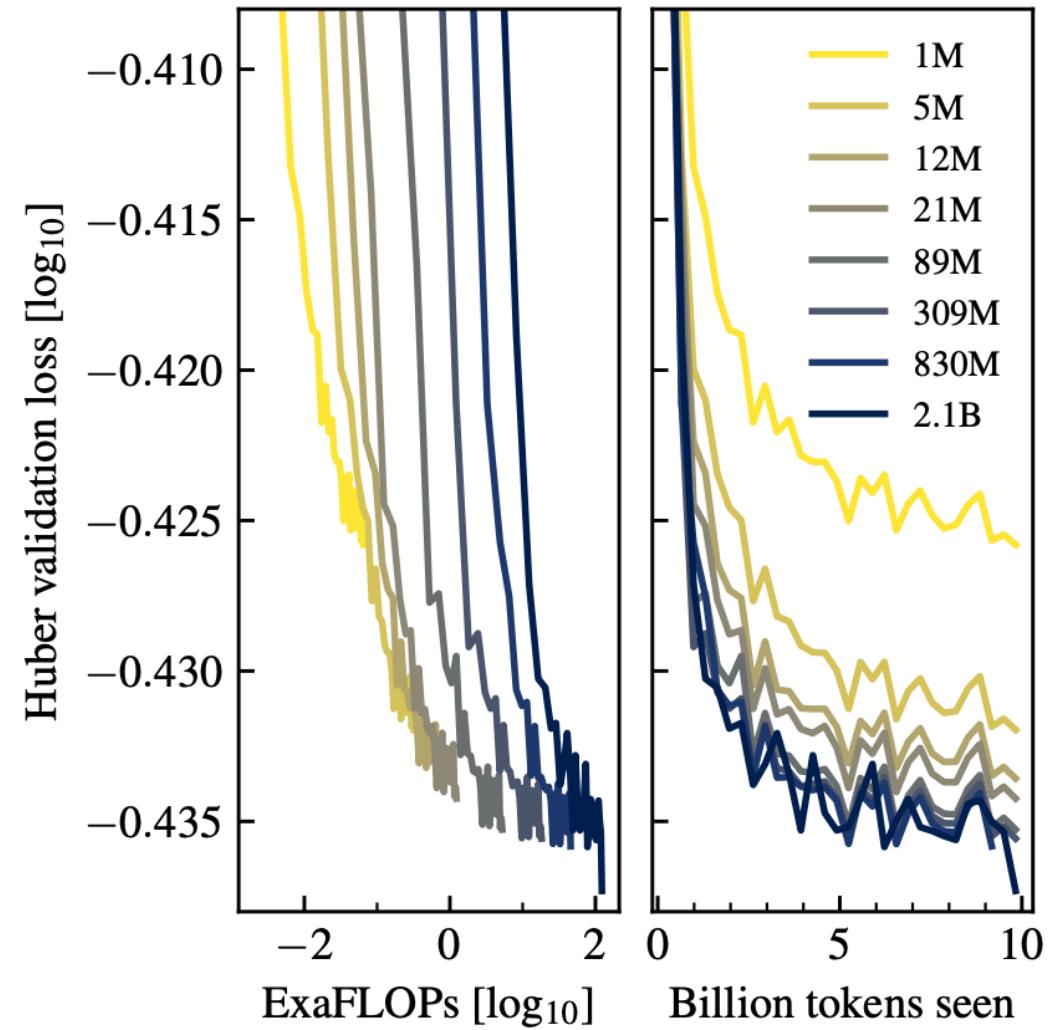


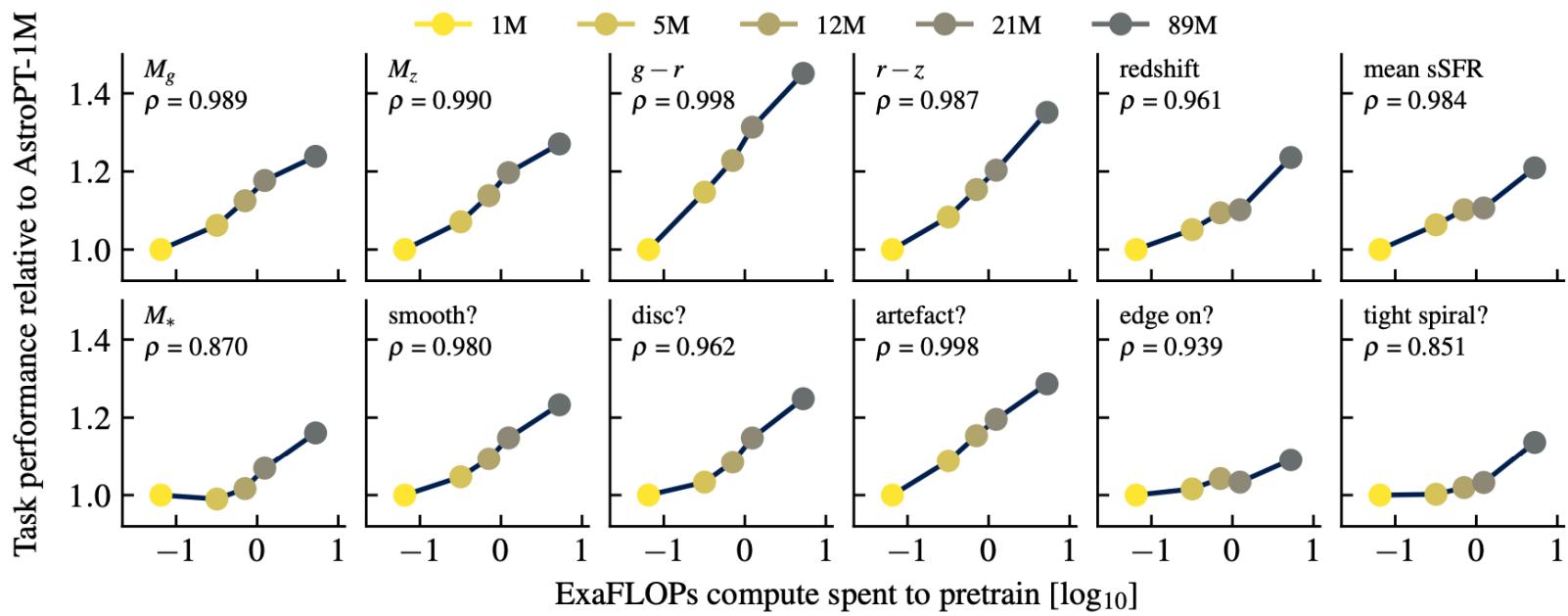
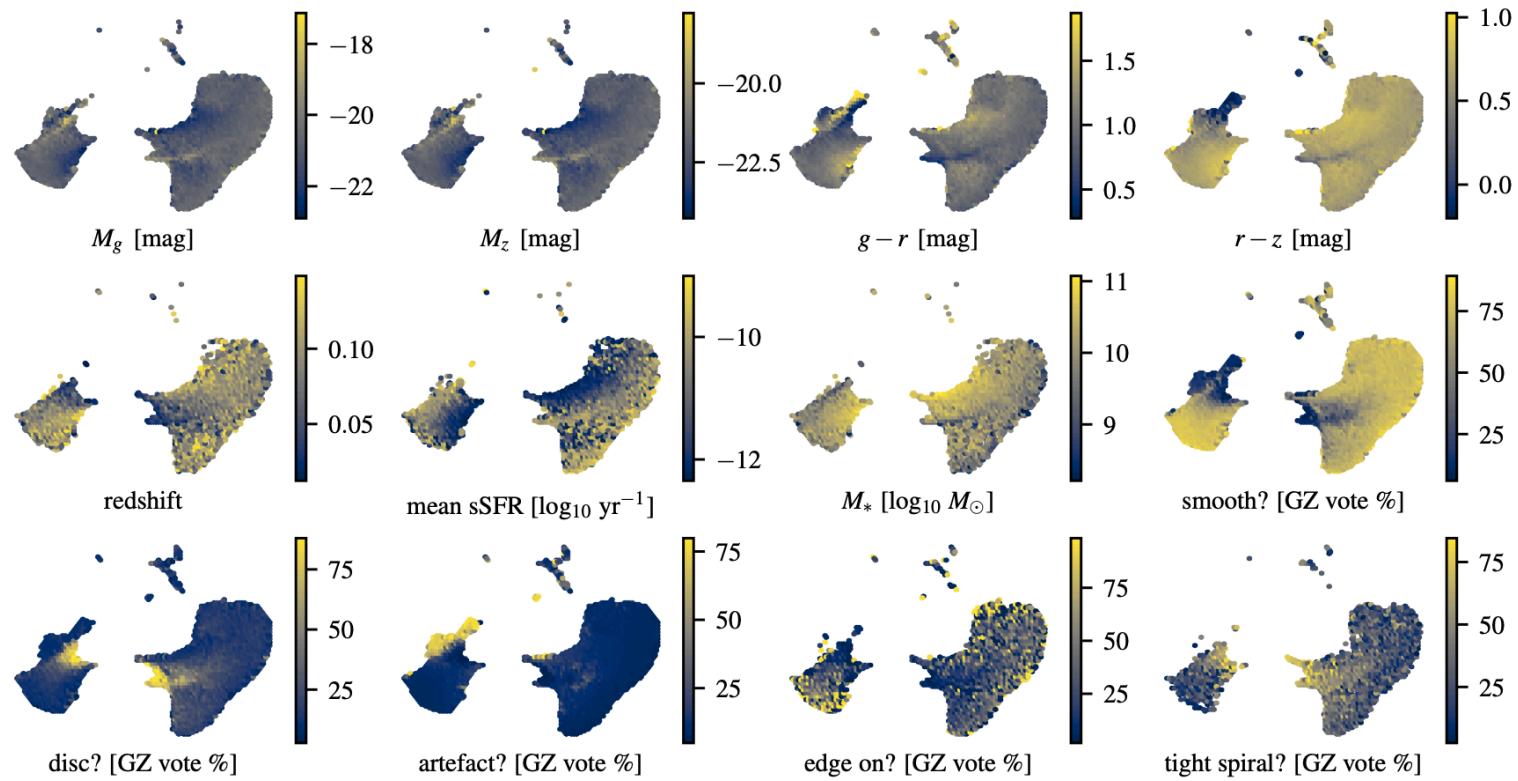
Parker+23

astro-PT: auto-regressive encoding

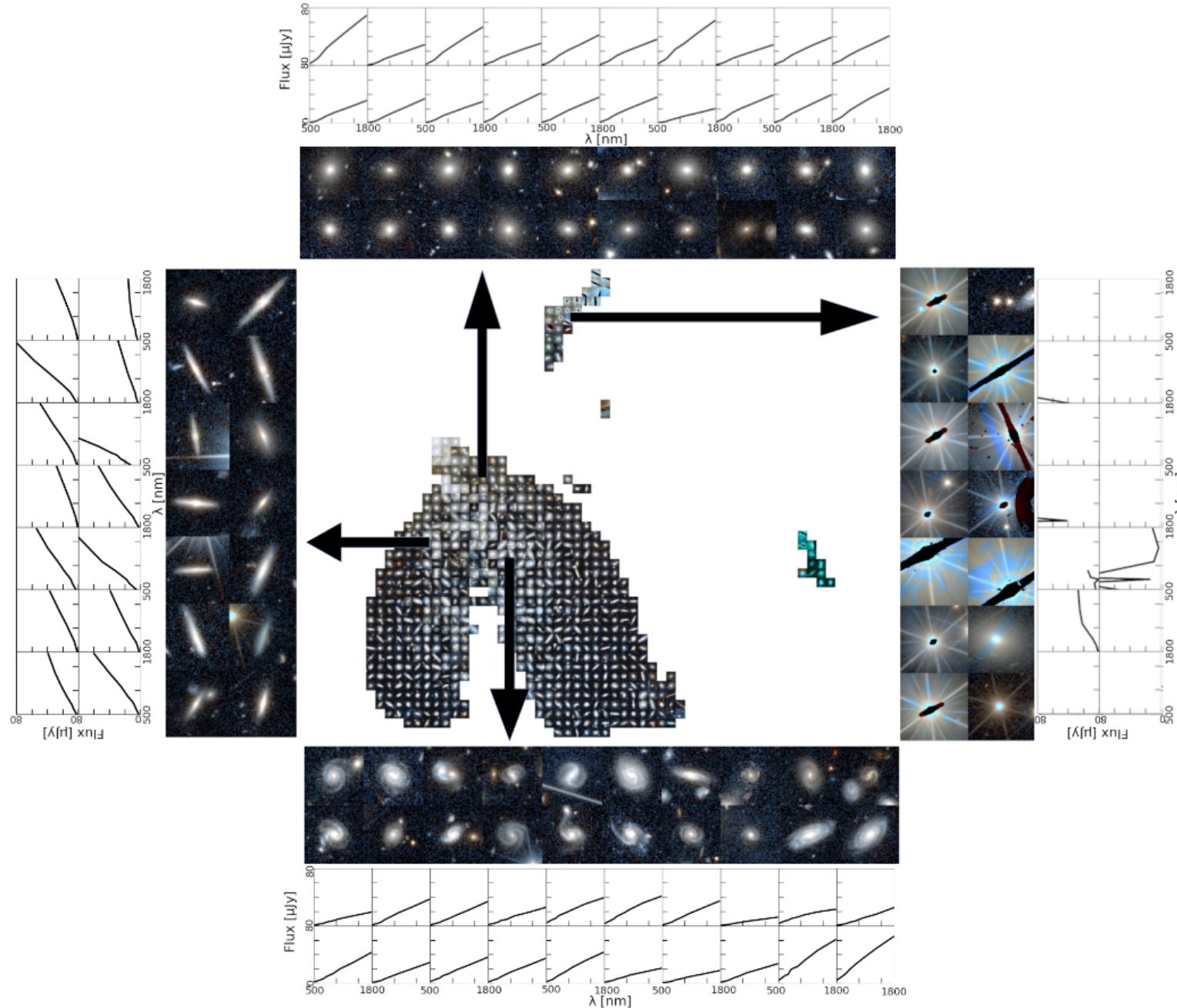


Information content of galaxy images.

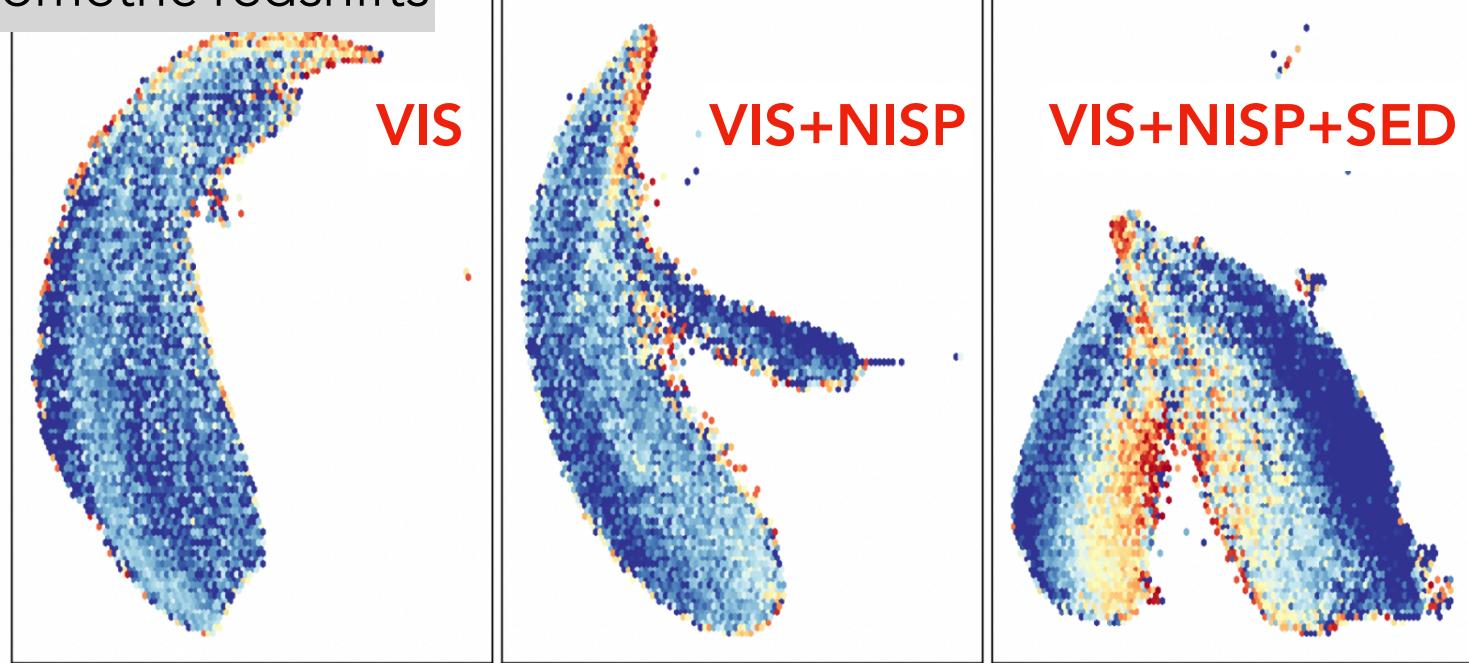




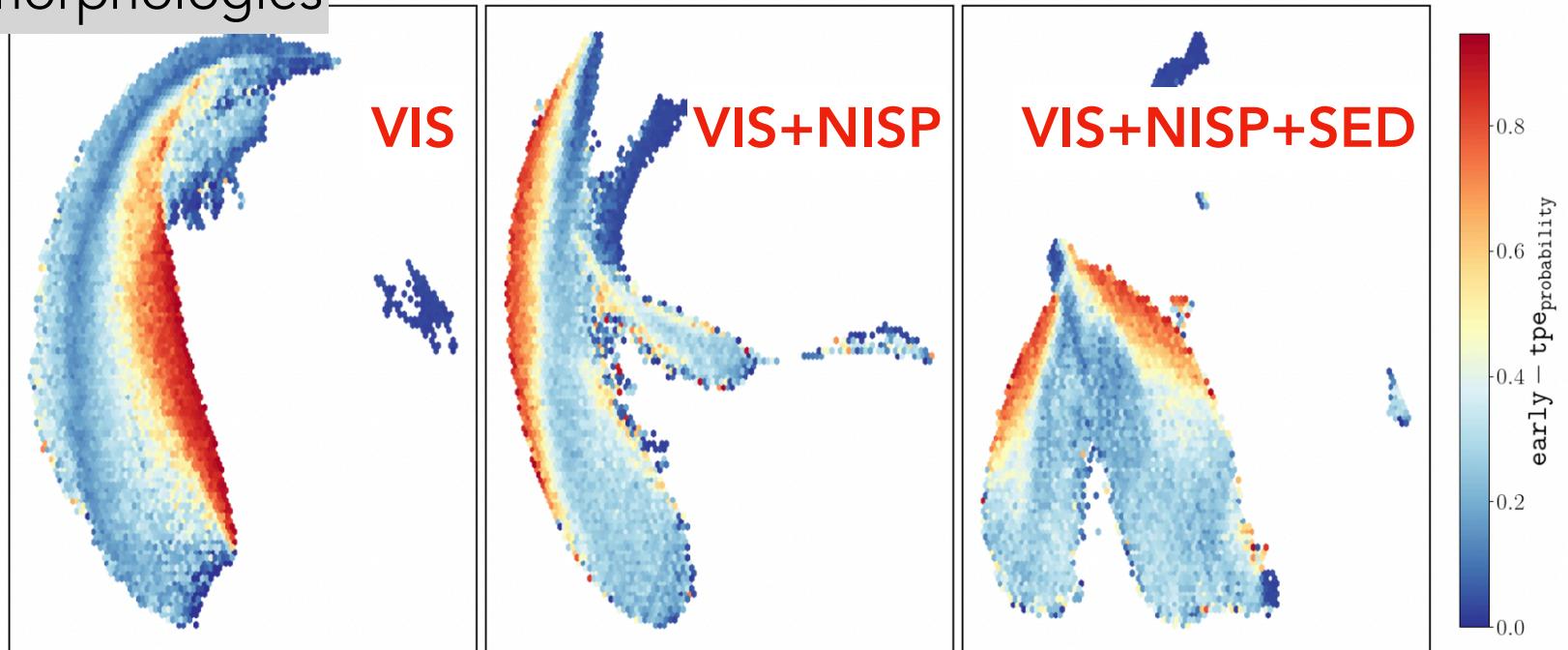
astro-PT joint embedding of SEDs and images



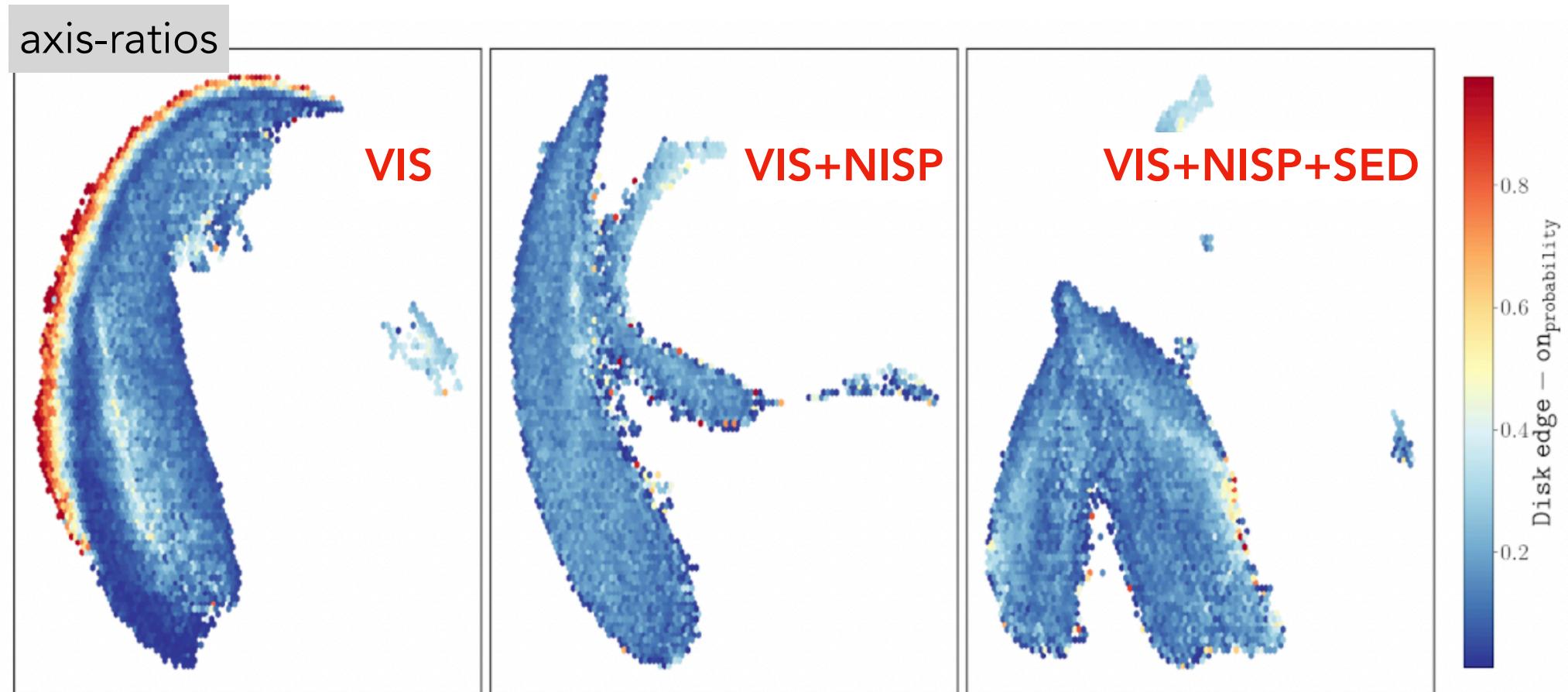
photometric redshifts



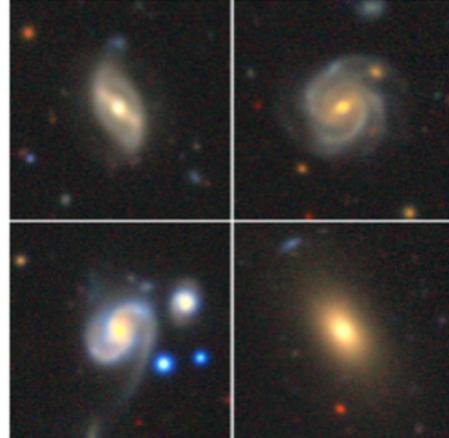
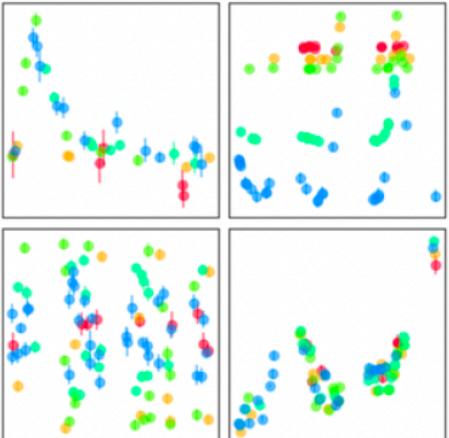
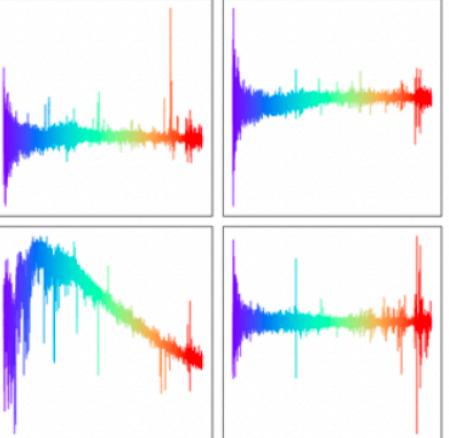
morphologies



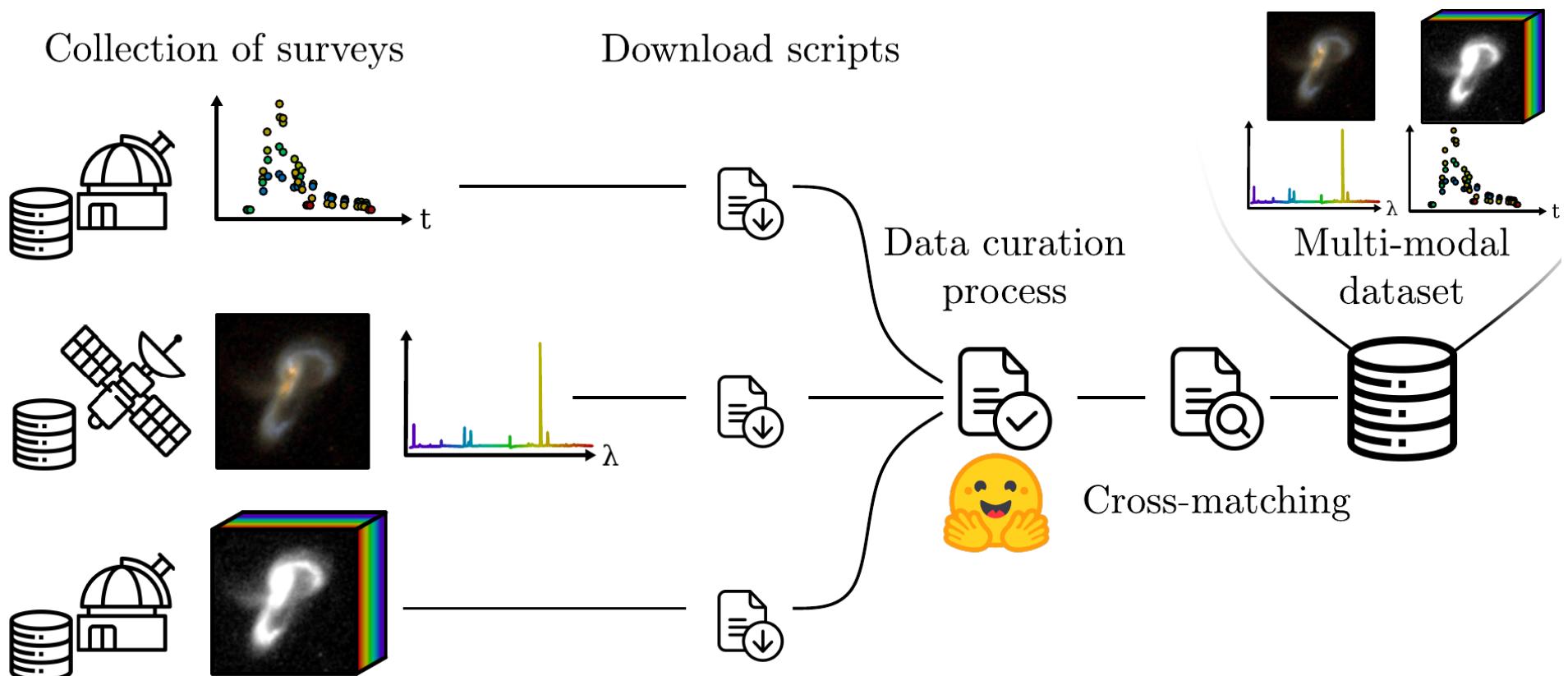
Adding information does not
always result in more informative embedding spaces



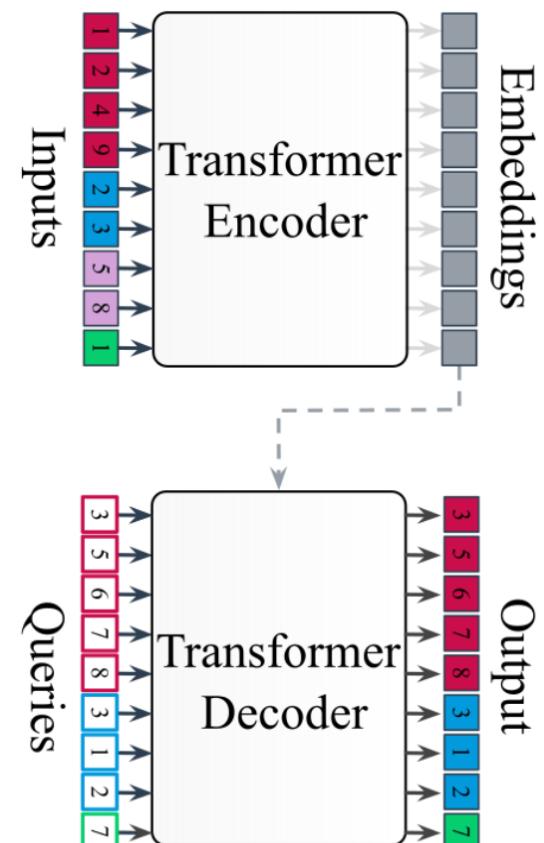
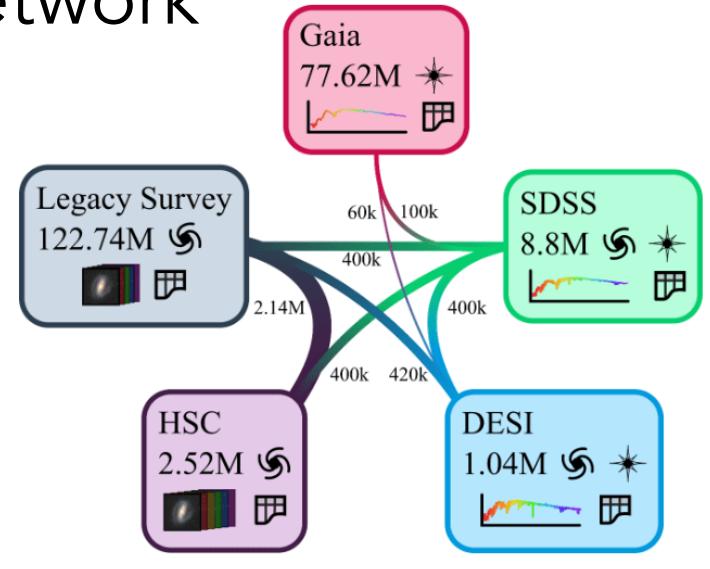
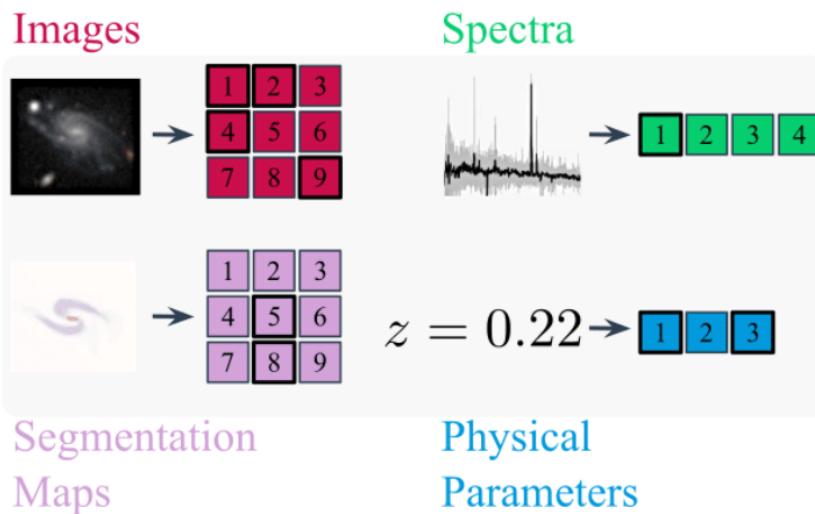
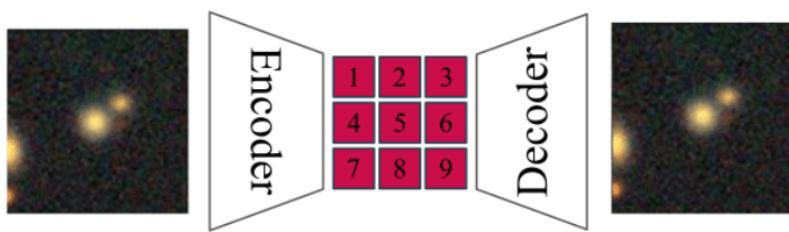
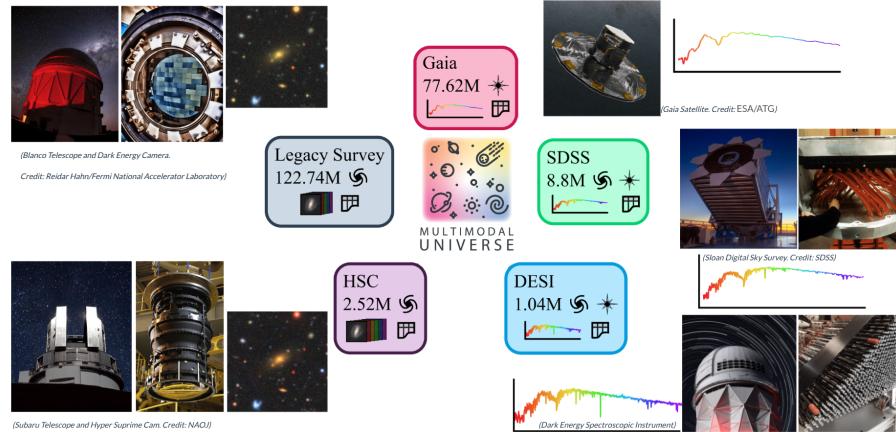
The Multimodal Universe: Enabling Large-Scale Machine Learning with 100 TB of Astronomical Scientific Data

	Images	Time-Series	Spectra
# examples	140M	4.5M	225M
Description	images in a variety of wavelength ranges, including optical and infrared	multivariate time-series of flux + uncertainty in different wavelength ranges	flux as a function of wavelength
Tasks	galaxy classification, physical property estimation	time-series classification, redshift estimation	physical property estimation
Examples			

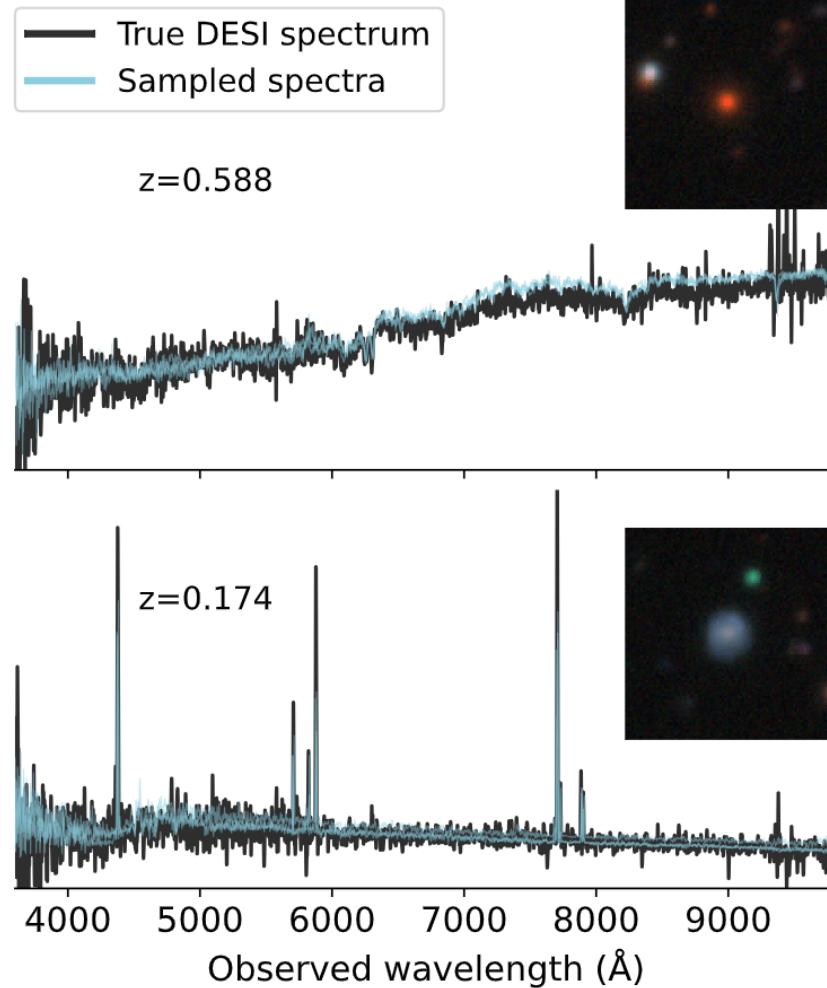
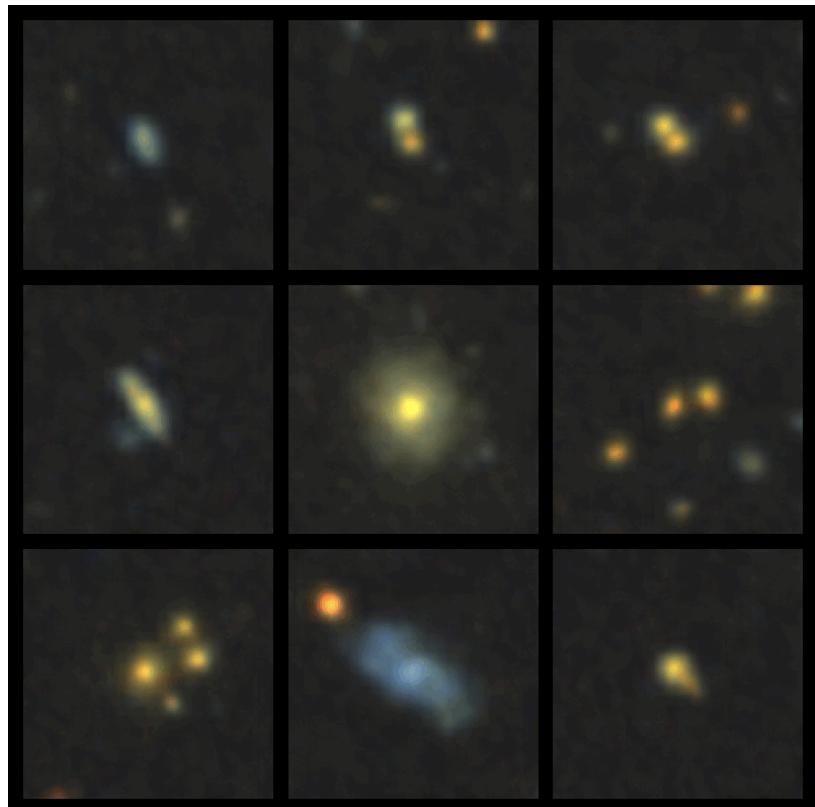
The Multimodal Universe: Enabling Large-Scale Machine Learning with 100 TB of Astronomical Scientific Data



AION-1: Astronomical Omnimodal Network (*polymathic team)



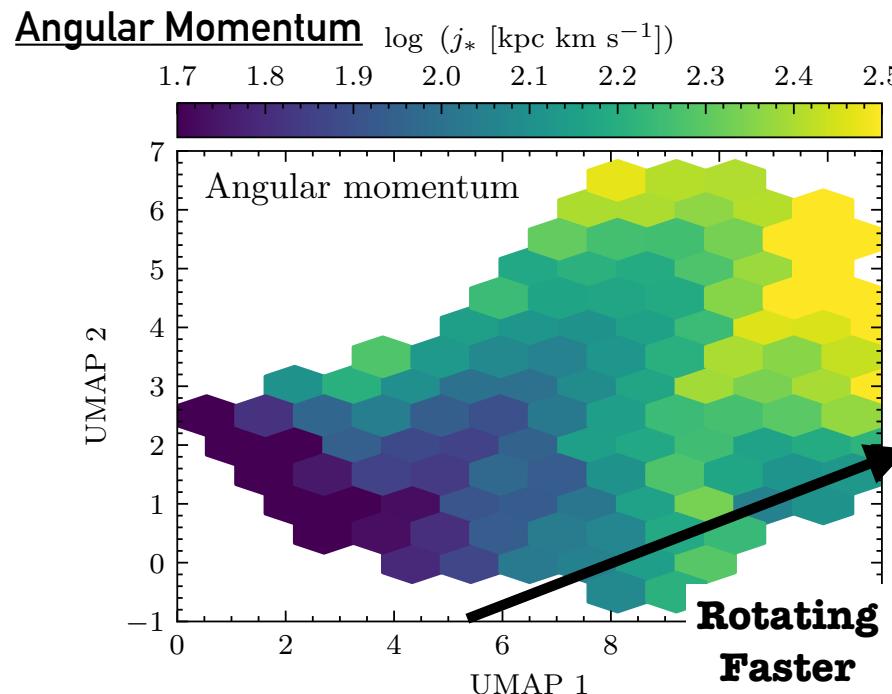
AION-1: Astronomical Omnimodal Network (*polymathic team)



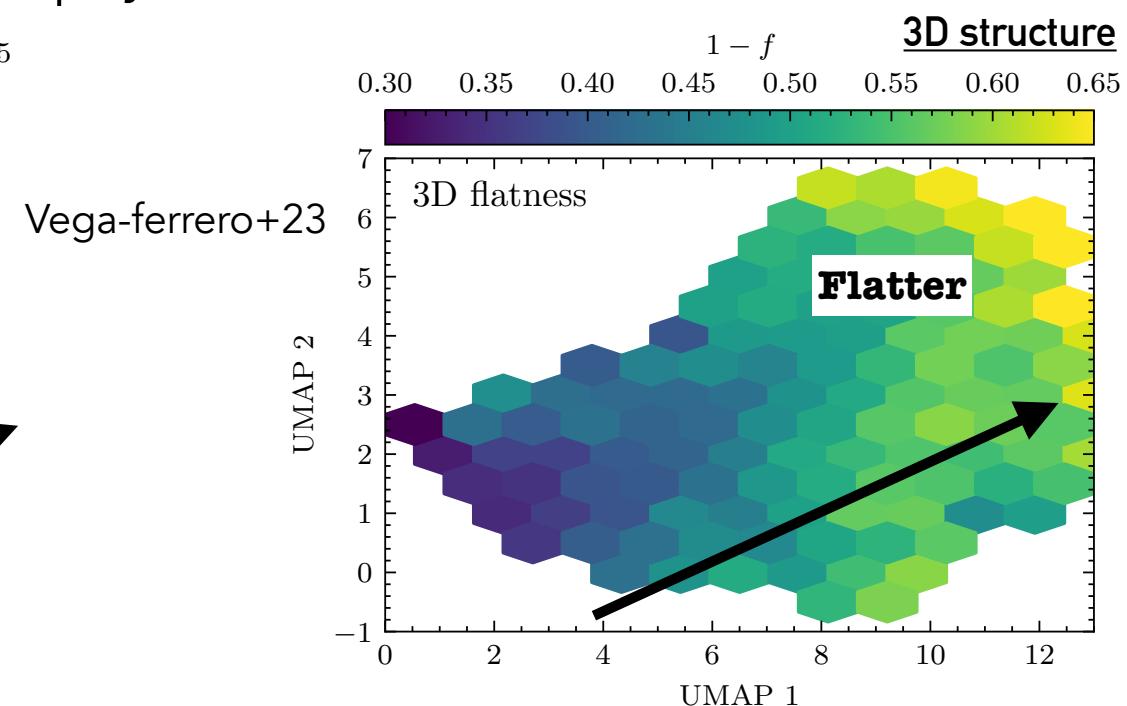
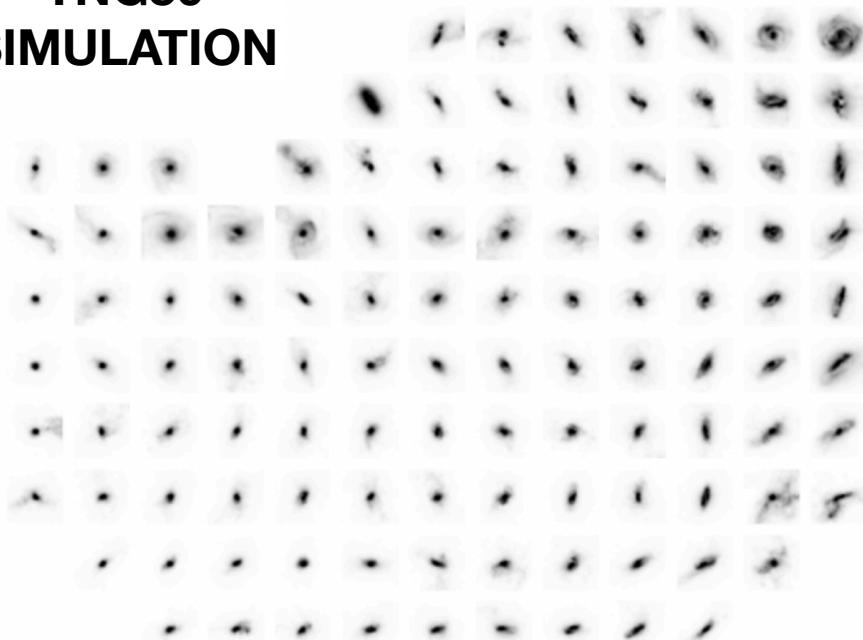
$$p(\mathbf{x}_{HSC} | \mathbf{x}_{DES})$$

$$p(\mathbf{x}_{DES} | \mathbf{x}_{HSC})$$

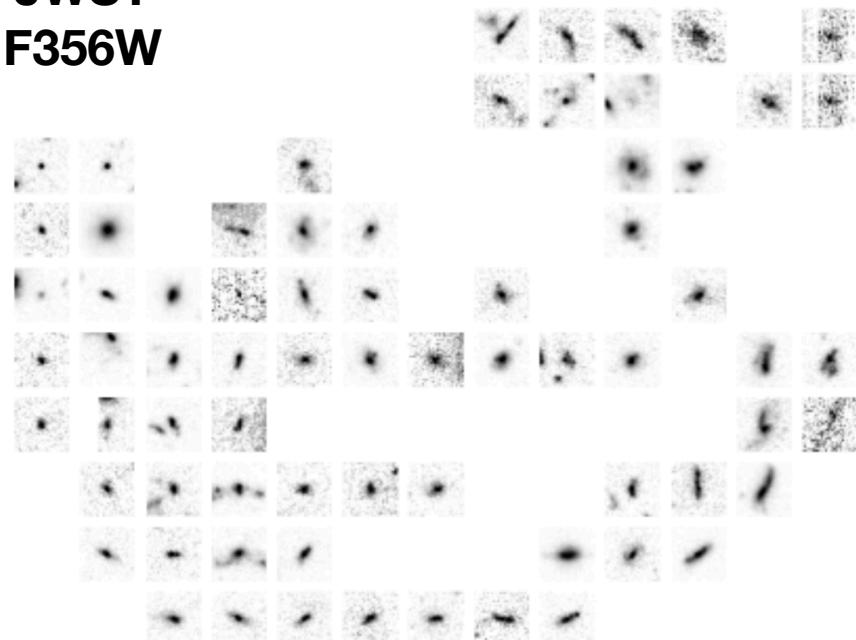
Can we interpret these embedding spaces to learn something about the physics?



**TNG50
SIMULATION**



**JWST
F356W**



projecting simulations and observations in the same embedding space

