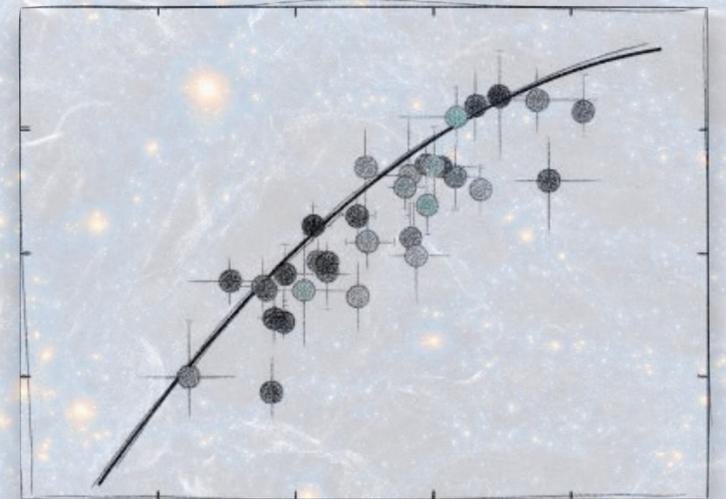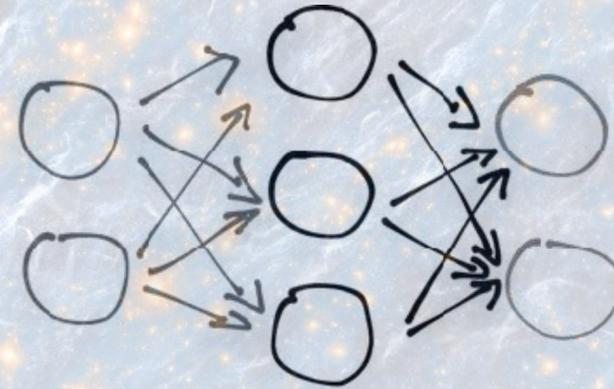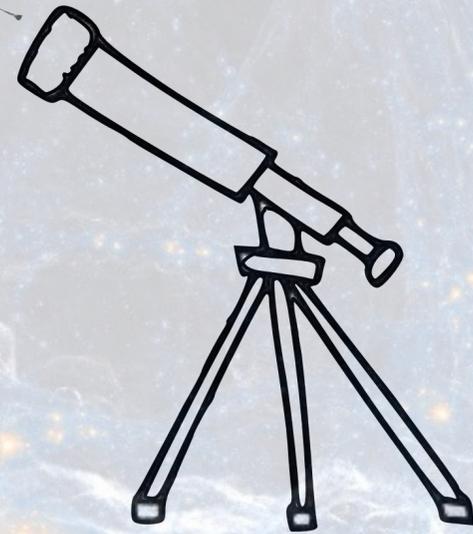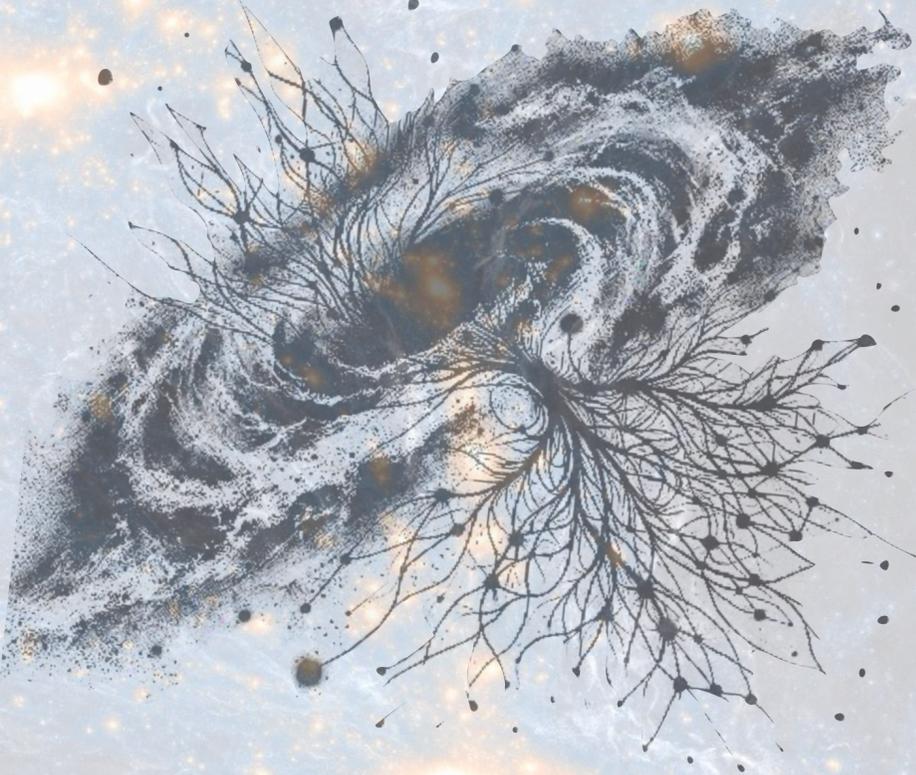# Decoding galaxy properties with machine learning and simulation based inference

Michele Ginolfi — UniFi & INAF

[arXiv:2410.22420; arXiv:2410.16370; arXiv:2502.20448]
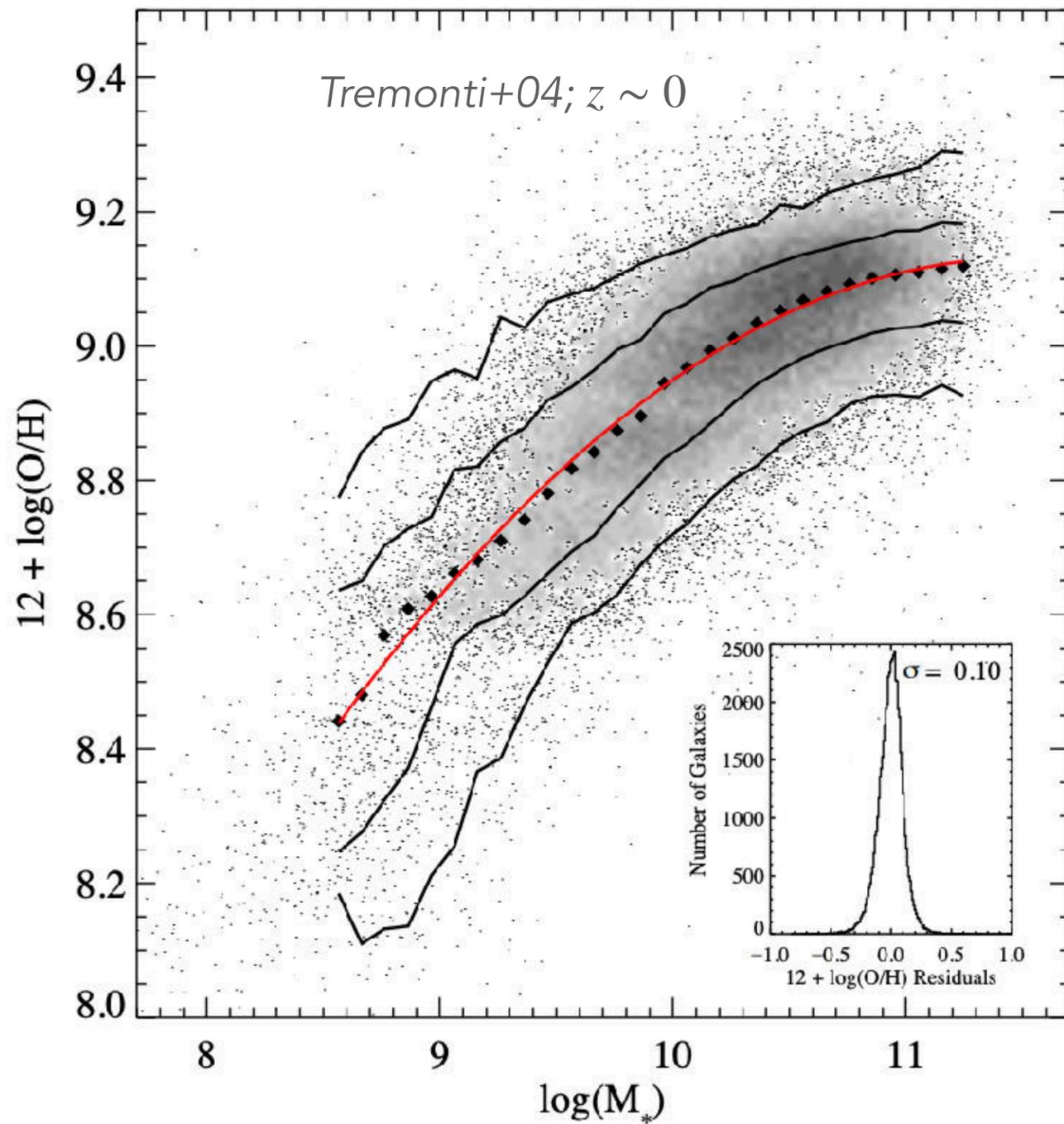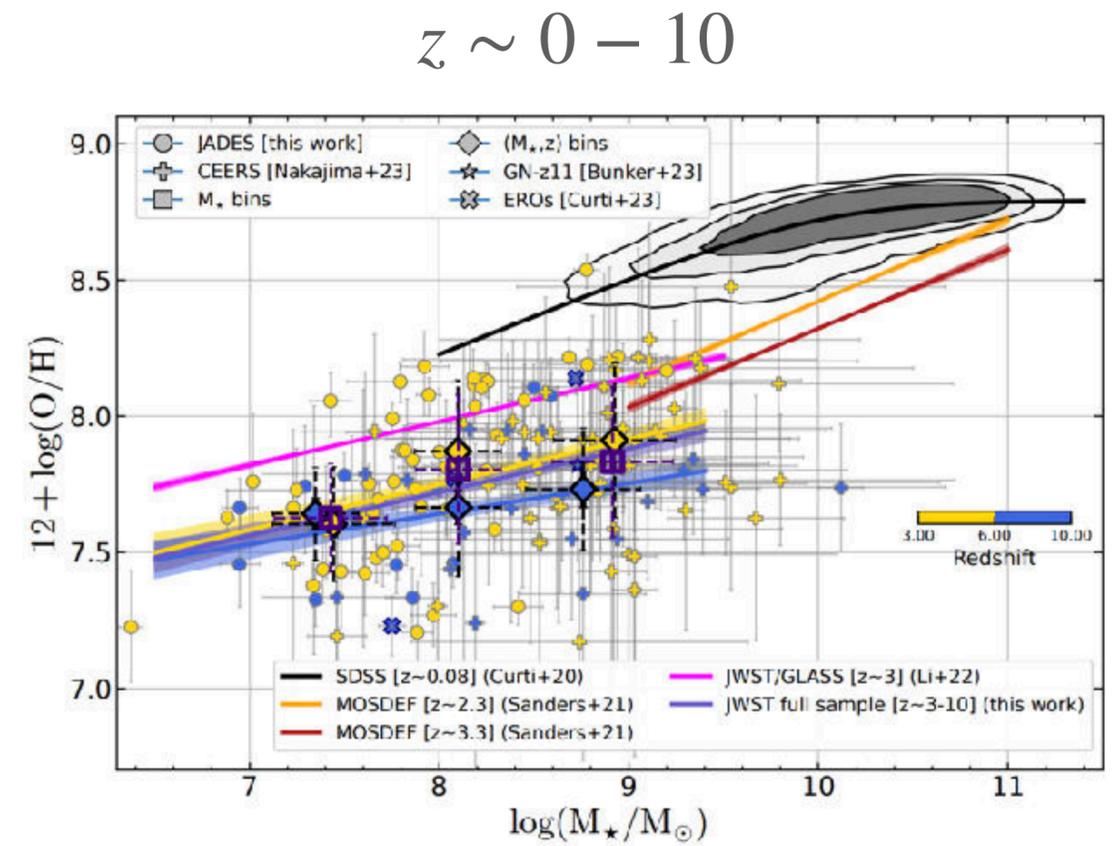
UNIVERSITÀ DEGLI STUDI FIRENZE

INAF
ISTITUTO NAZIONALE DI ASTROFISICA

> 50.000 galaxies from SDSS
(>2.000.000 galaxy spectra)
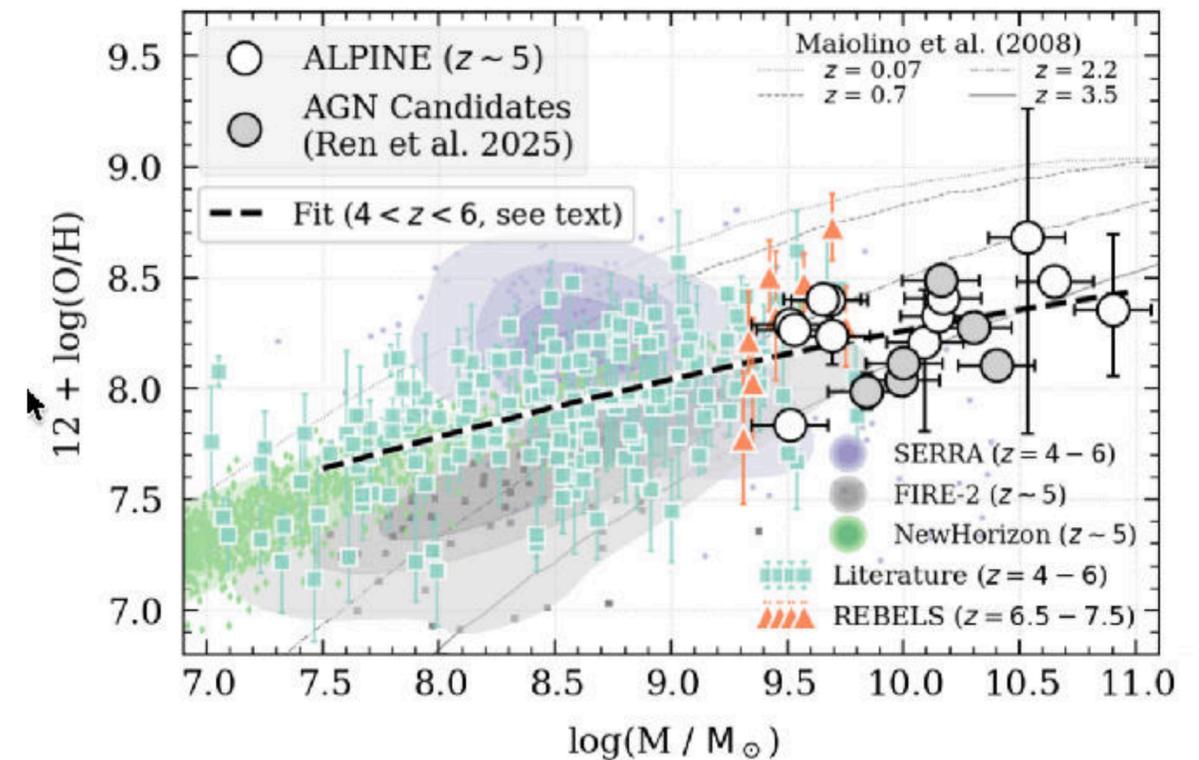
*Tremonti+04; z ∼ 0*

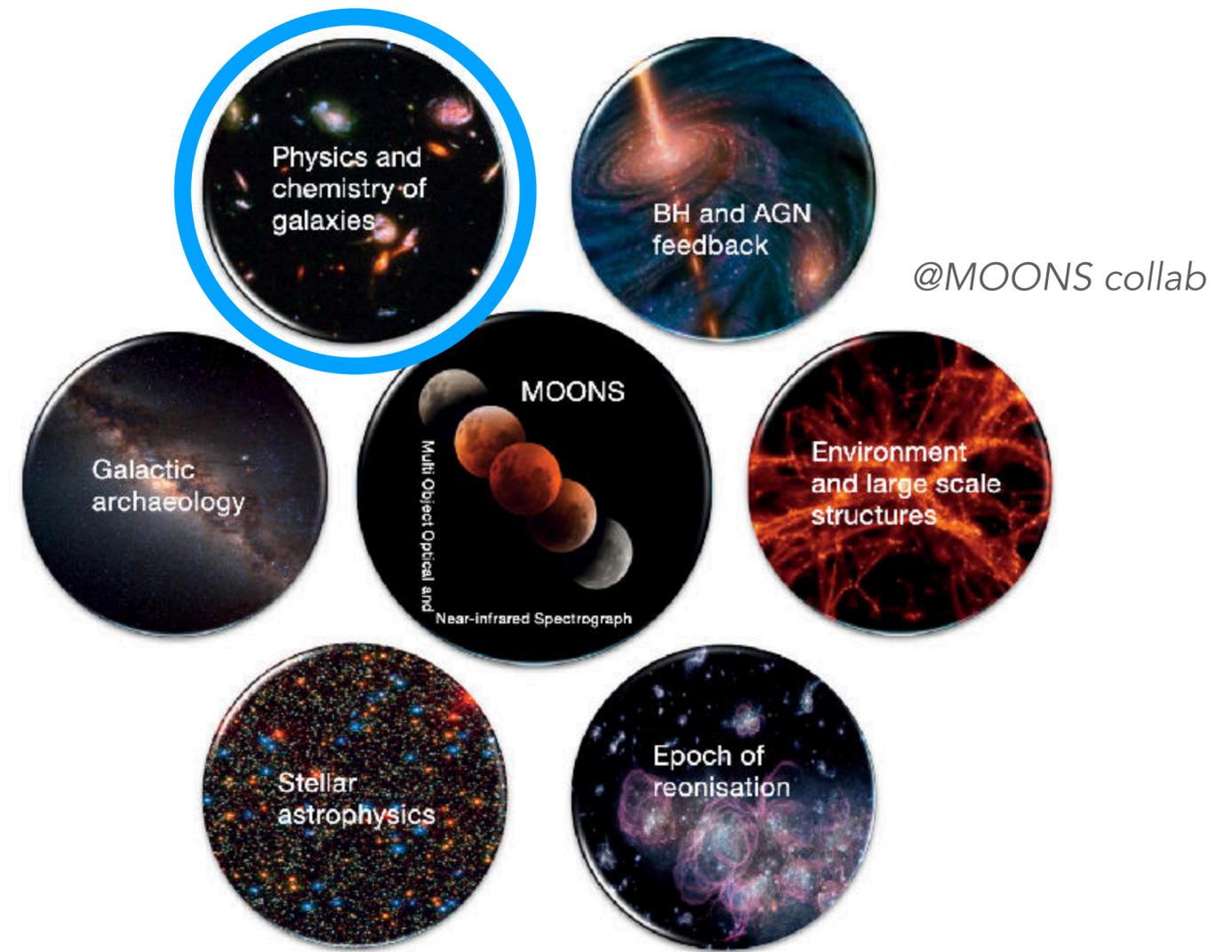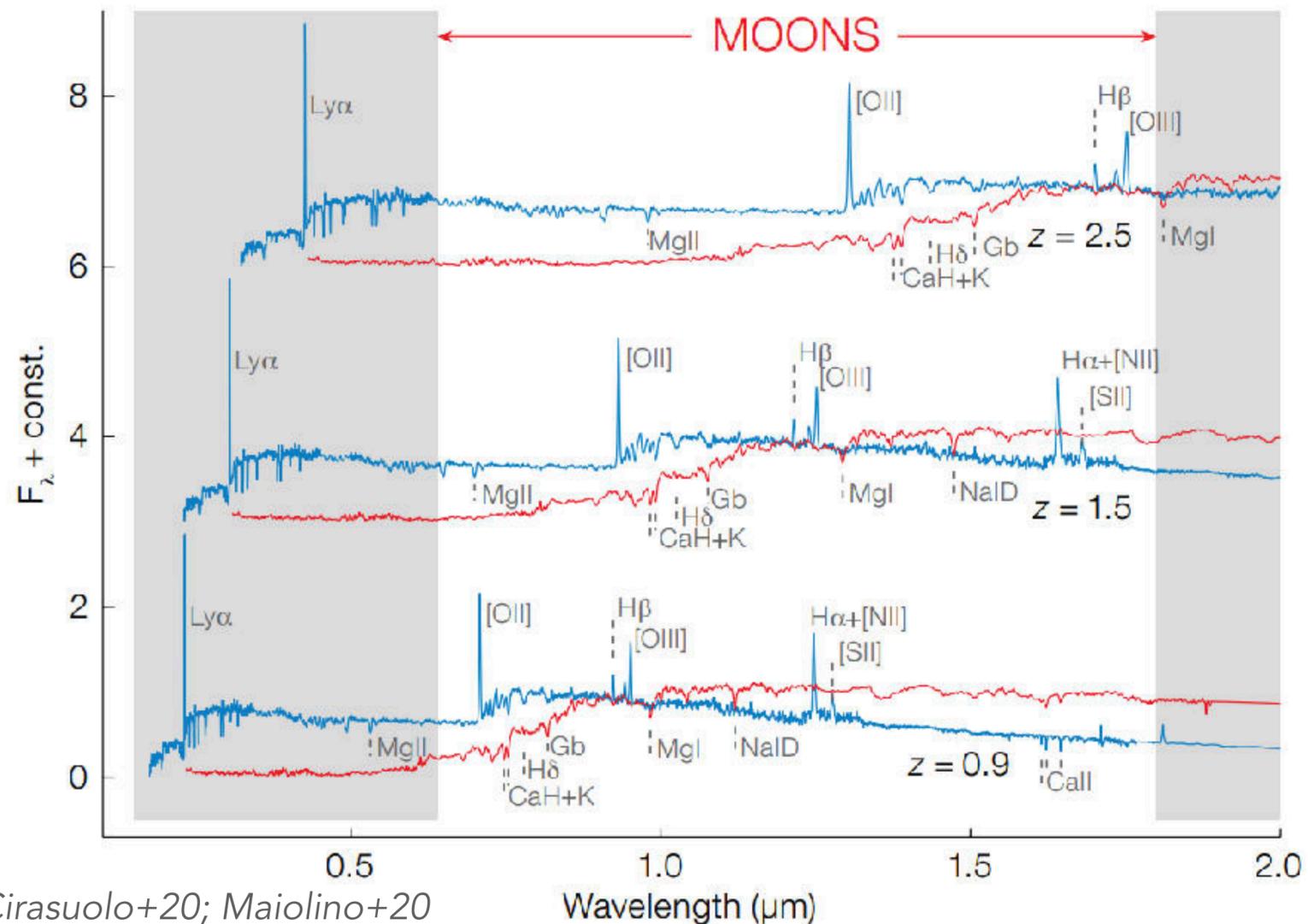$z \sim 0 - 10$

Few years ago

Today

*Curti+24*

*Faisst+25*

Ongoing & upcoming all-sky spectroscopic surveys (DESI, 4MOST, **MOONS**)

**hundreds of millions of spectra** will be acquired over the next half-decade!

# MOONS 🌑🌓🌖🌕 is the new **Multi-Object Optical and Near-infrared Spectrograph**,

soon to be operated @VLT, ESO

- 1000 fibres, over a field of view of $\sim 500^2$ arcmin;
- low- (R~4000–7000) / high-resolution (~19000 in *H*);
- 0.64 – 1.8 μm wavelength range.



*@MOONS collab*



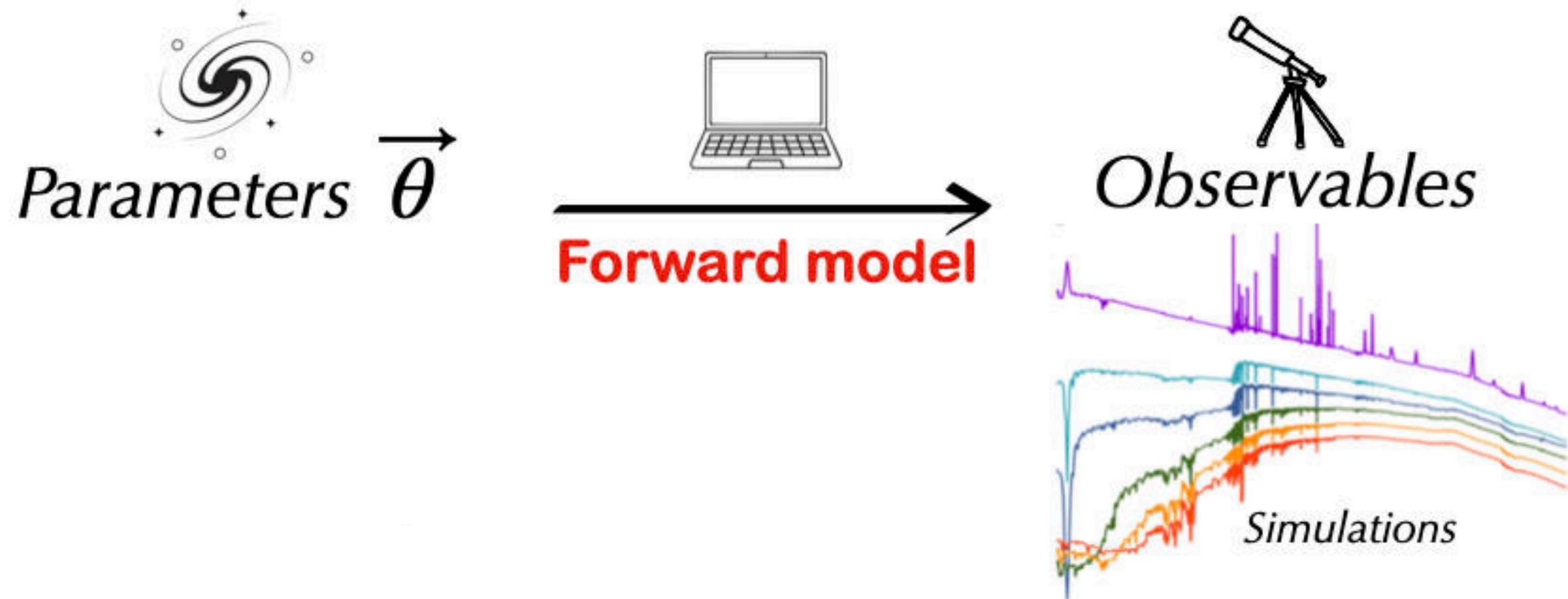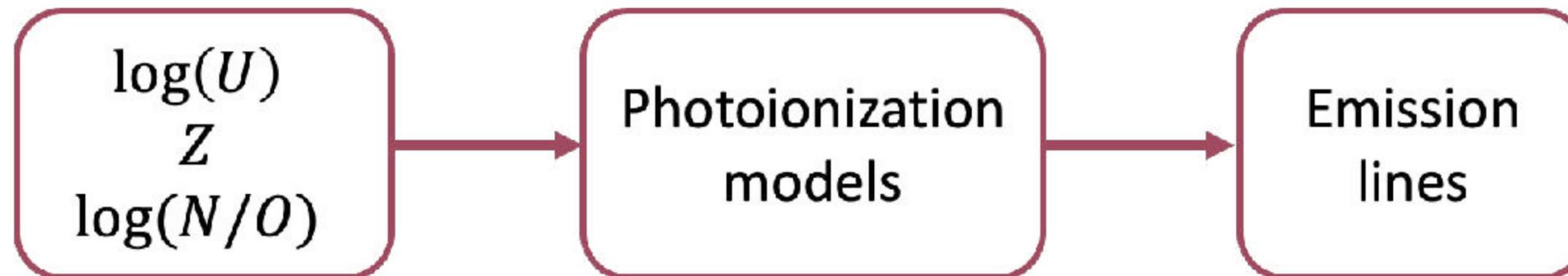*Cirasuolo+20; Maiolino+20*

A data challenge!

- up to about **half a million galaxies** at $0.9 < z < 2.6$
- >12000 elements per spectrum in low-resolution!
- **Standard fitting methods are typically slow and often fail in weak-signal regimes.**

**Mini bibliography** — Taylor+18; Maiolino+20, Cirasuolo+20; Looser+21; Gonzales+22; Cabral+22; Gonzales+23; etc etc
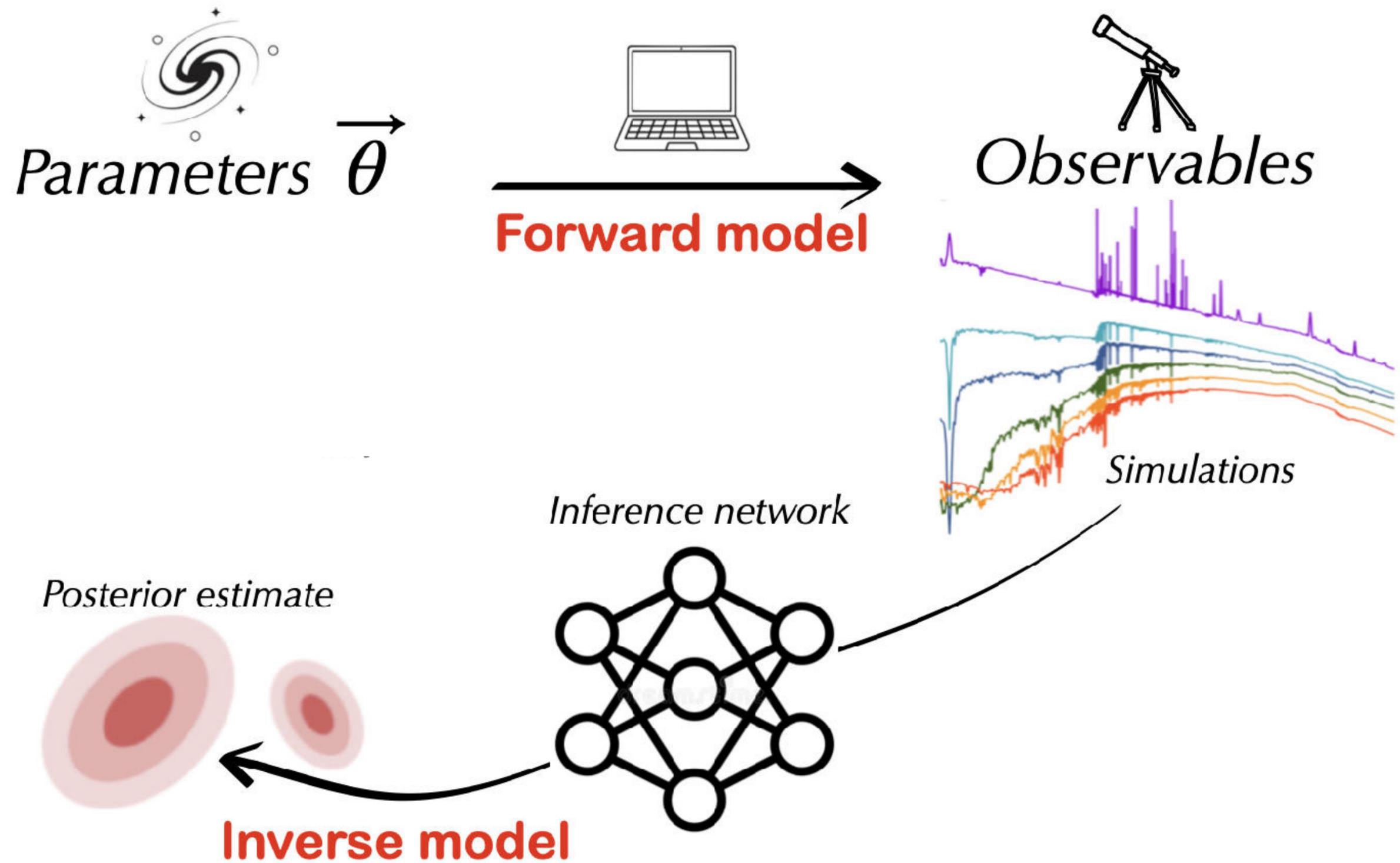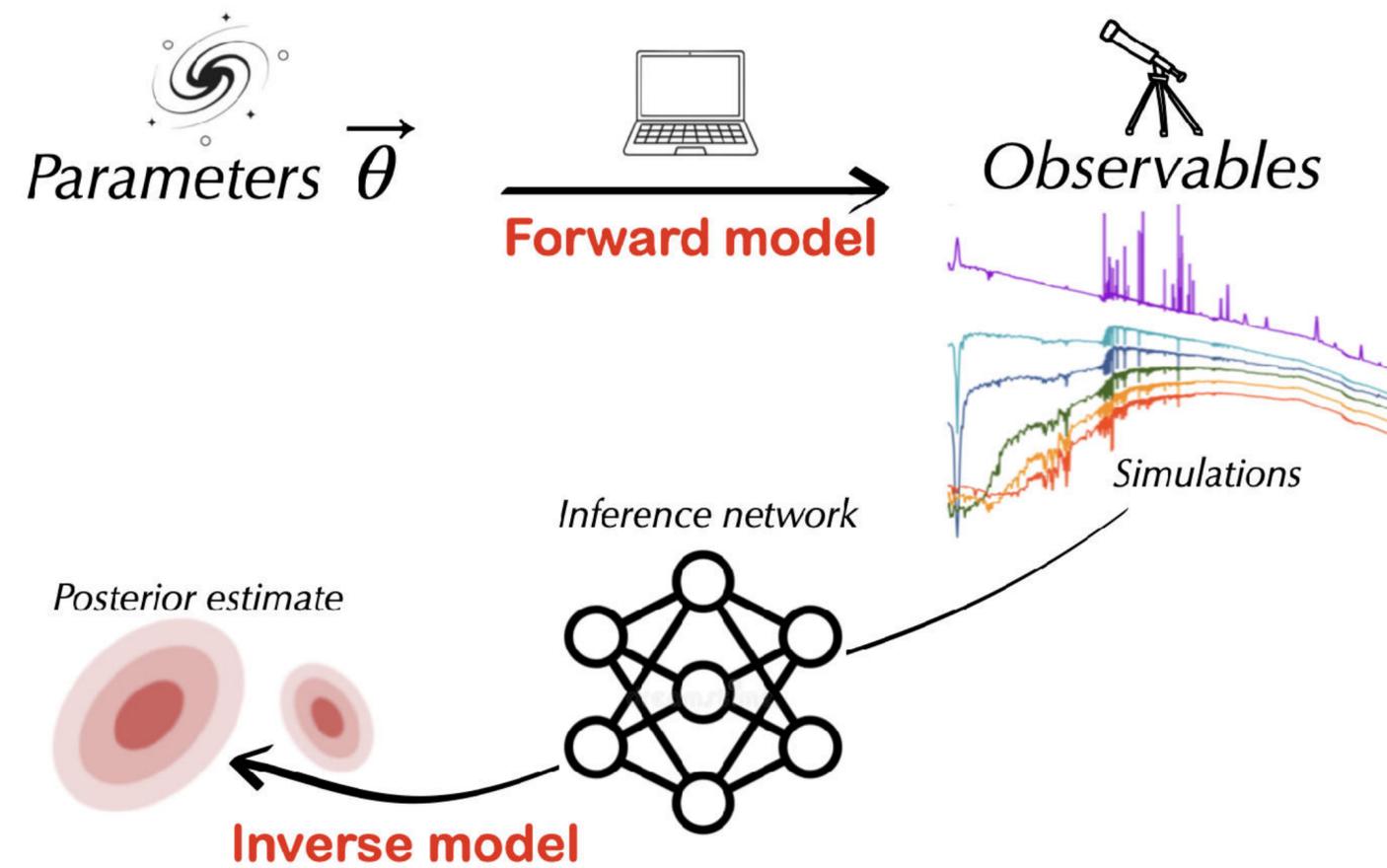
# How do we infer physical properties?



Parameters $\vec{\theta}$

**Forward model**

Observables

Simulations

**e.g.,**

$$\log(U)$$
$$Z$$
$$\log(N/O)$$

Photoionization models
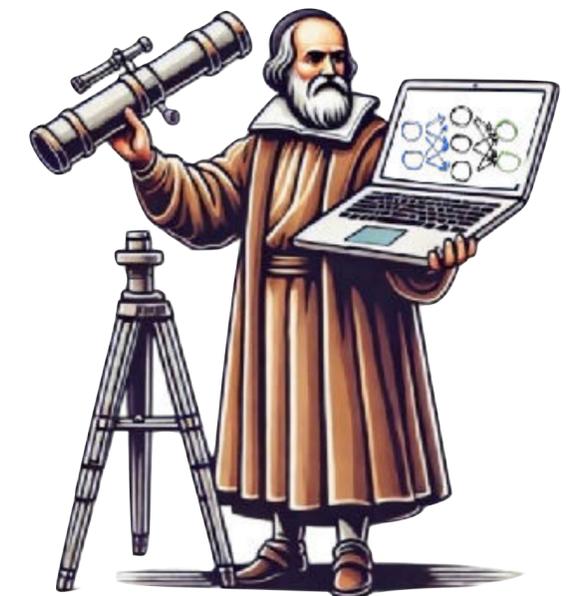
Emission lines

**We need to solve the "inverse problem"**

# Simulation-based inference (SBI) with neural estimators

# Simulation-based inference (SBI) with neural estimators
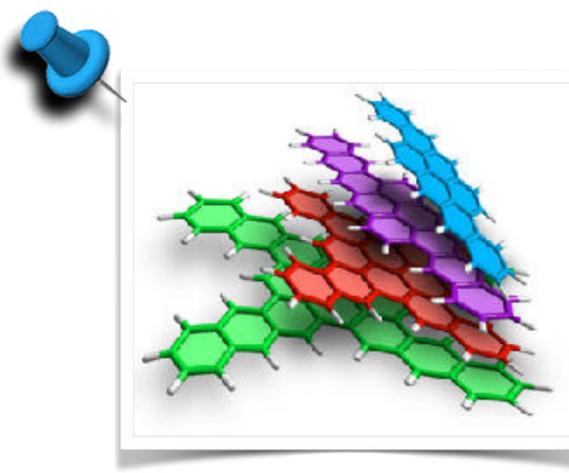


**AI4phys@Florence**

Parameters $\vec{\theta}$

**Forward model**

Observables

Simulations

Inference network
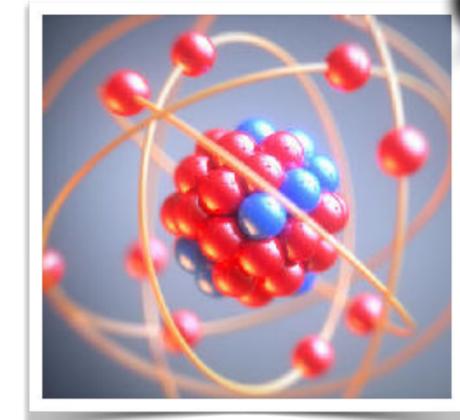
Posterior estimate

**Inverse model**
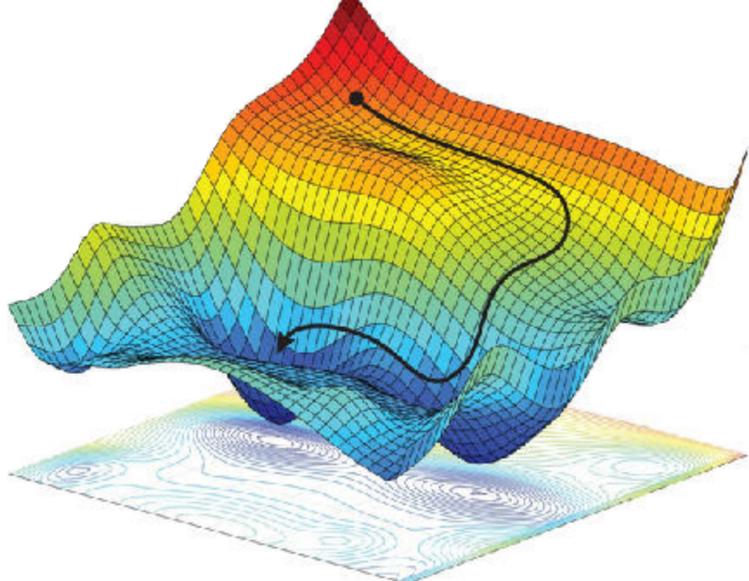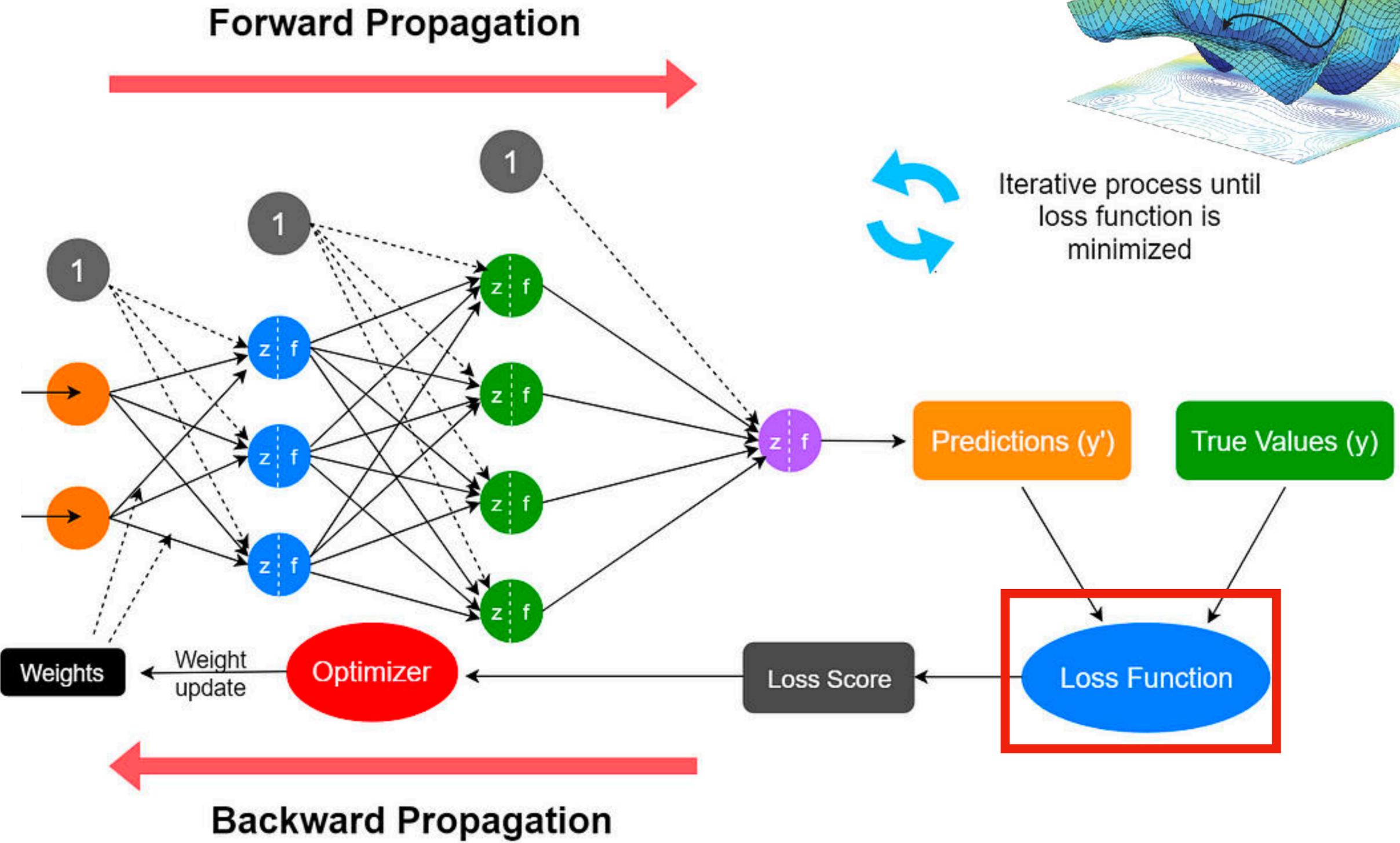
Astro

Medical imaging

Material science

Nuclear physics

**Goal: objective function —** Find parameters that minimize a loss function

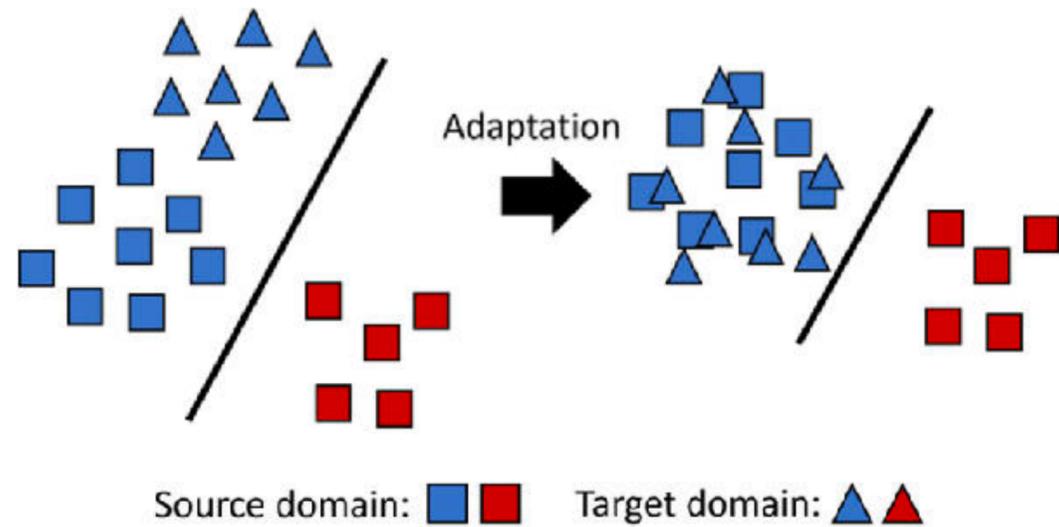e.g, fluxes in different bands for SED fitting

# Three **big** challenges in simulation-based inference with AI

**Domain-invariant** learning

**Interpretable** learning



Source domain: ■ ■  Target domain: ▲ ▲

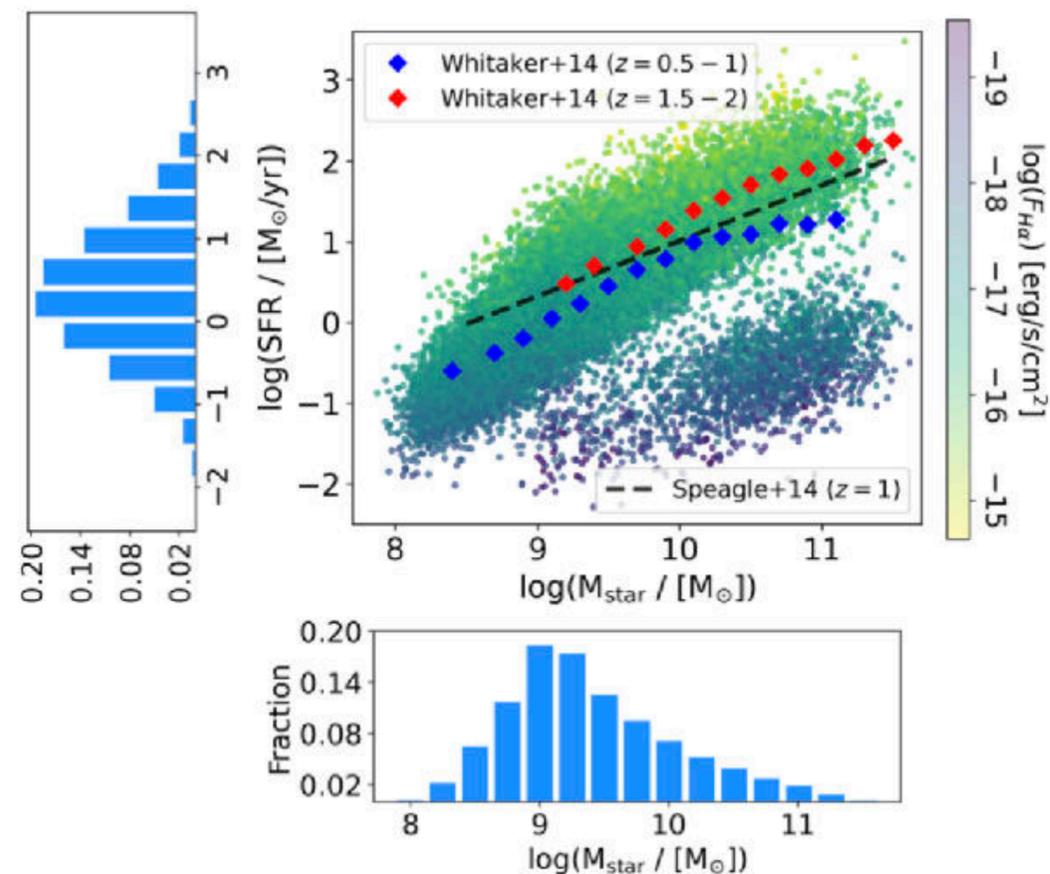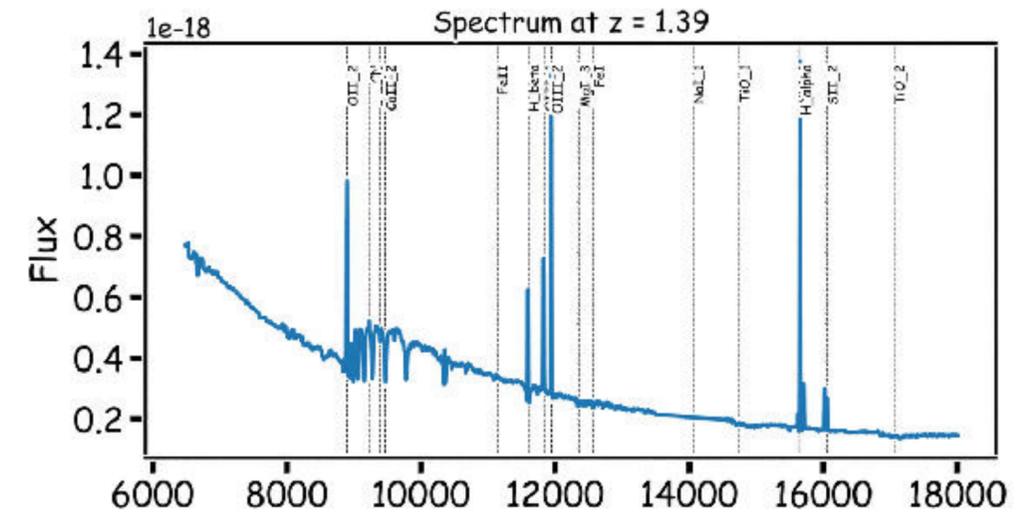**Uncertainty-aware** learning

# Simulated dataset

- **~120.000** spectra
- generated by running **MAMBO** templates through **moons1d**
- moons1d ran with low resolution mode for all 3 channels (RI, YJ, H), **0.64-1.8 µm**
- a seeing of 0.8'' and airmass of 1.2
- $t_{exp} = 2 - 4 - 8 \ h$





Model

Moons1d — MOONS simulator for 1D data

sky condition, obs strategy, etc



Simulated spectrum

# Dataset

- $t_{exp} = 2, 4, 8\ h$

- 0.64 – 1.8 μm

- 12.217 channels

## Target physics

- redshift, $z$

- stellar mass, $M_{\text{star}}$

- star formation rate, $SFR$

Simulated spectra

# Deep learning MOONS spectra 🌑🌒

## The case of redshift

*Classical scheme: a regression problem*



Spectrum at z = 1.00

$x_1$
$x_2$
$x_3$
⋮
$x_n$

$z = 1.00$

# Deep learning MOONS spectra 🌑🌒🌓🌔

Why not discretise continuous redshift values into finely spaced bins? 💡

## Switch to a classification task      We adopt $dz = 0.003$



$x_1$
$x_2$
$x_3$
$\vdots$
$x_n$

$z = 1.00$

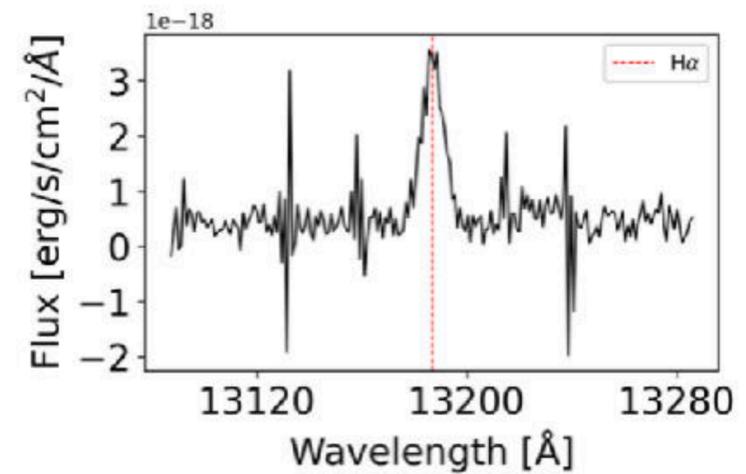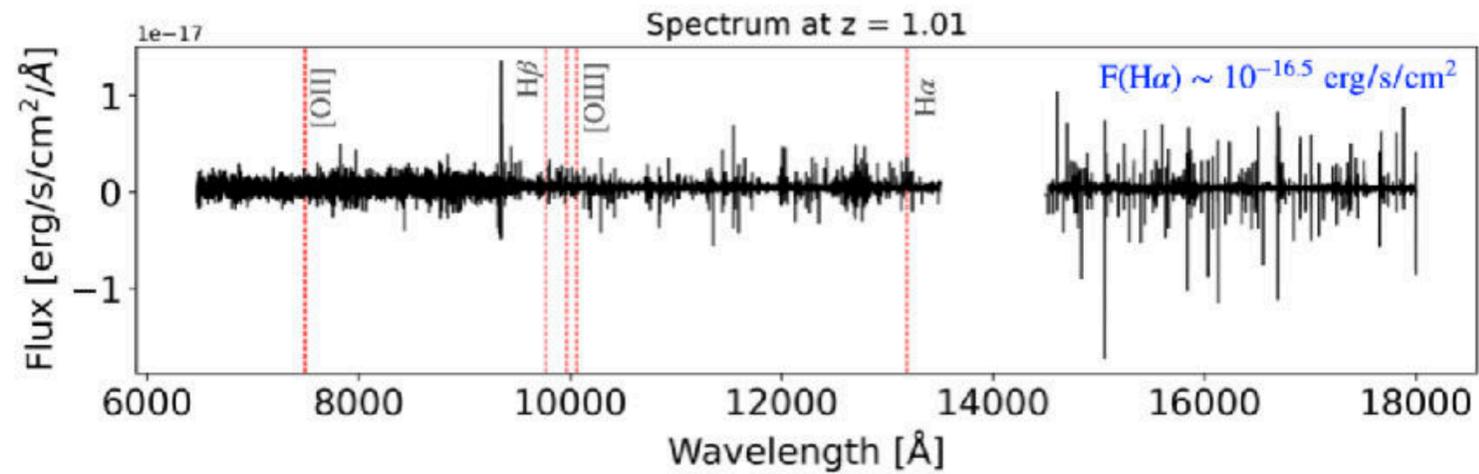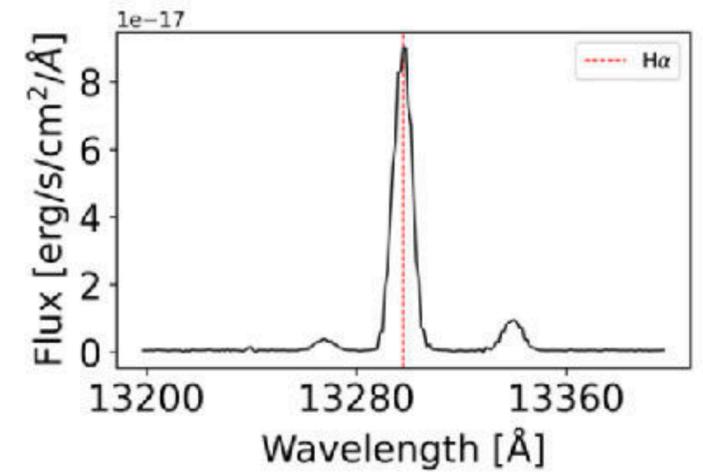$P(y=y_1|x)$   $z = 0.997$  $(P = 0.1)$

$P(y=y_2|x)$   $z = 1.000$  $(P = 0.95)$

$P(y=y_3|x)$   $z = 1.003$  $(P = 0.05)$

$P(y=y_m|x)$   $z = 2$  $(P = 0)$

**Mini bibliography** — Carrasco Kind & Brunner 2013;
Stivaktakis et al. 2018; Stewart et al. 2022; Pankaj et al. 2022

# Switch to a classification task

We can also account for the *quality* of spectra

(for instance using the rms of the spectra)

Finely-spaced input

| $z = 1$ | $z = 1.1$ | $z = 1.2$ | $z = 1.3$ | $\cdots$ | $z = 3$ |

| $z = 1$ | $z = 1.1$ | $z = 1.2$ | $z = 1.3$ | $\cdots$ | $z = 3$ |



Trained with a cross-entropy loss and
a softmax activation in the last layer

Finely-spaced output

| $z = 1$ | $z = 1.1$ | $z = 1.2$ | $z = 1.3$ | $\cdots$ | $z = 3$ |

| $z = 1$ | $z = 1.1$ | $z = 1.2$ | $z = 1.3$ | $\cdots$ | $z = 3$ |

# Learning through multi-task training 🤹



**Mini bibliography** — Caruana 1997; Ruder 2017; Crawshaw 2020, Hervella+24

# M-TOPnet (Multi-Task network Outputting Probabilities)

# M-TOPnet (Multi-Task network Outputting Probabilities)

# Predictions on the test set 🔮🪄

- **~18000** spectra
- Same distributions as training set

# A physically-motivated test sub-set 🔭

MOONRISE-like selection: $m_H < 25$; $z < 2.6$; $t_{exp} = 8\ h$ for passive galaxies; $t_{exp} = 2\ h$ for star-forming galaxies



$\mu = 0$
$\sigma = 0.001$

$\mu = 0.02$
$\sigma = 0.16$

$\mu = -0.006$
$\sigma = 0.27$

$t_{exp} = 2\ h$

$t_{exp} = 8\ h$

# Metrics: let's help design observational strategies 🧐

# Let's have a look at the predictions 🧙



Good prediction



*redshift pdf*

# Let's have a look at the predictions 🧙

Bad predictions

# Good vs bad predictions 🎯



$$|\Delta z| < 0.01$$

Passive

No bright emission lines

No ongoing SF

No ongoing SF

Well Classified
Poorly Classified

# Analysis and screening of the output redshift PDFs 🤔



They look different, right?

Intuitively, one might expect that galaxy spectra with unsuccessfully predicted redshifts (for instance, due to the reasons discussed before) would have broader or more dispersed PDFs with multiple peaks.

Can this intuition be quantified objectively, potentially enabling further a posteriori screening of the output?

# Analysis and screening of the output redshift PDFs 🤔

Well Classified
Poorly Classified

$t_{exp} = 2\ h$

$t_{exp} = 8\ h$

Entropy

Entropy

accuracy before filtering: 0.77;
**accuracy after filtering: 0.97**
Filtered-out spectra ~27%

accuracy before filtering: 0.9;
**accuracy after filtering: 0.99** 🤯
Filtered-out spectra ~12%

# Information encoded in the last embedding layers



~100 D space

⬇

UMAP dimensionality reduction

⬇

2-D projection

*The model encodes information that can be described through classes and quantities that "make sense" and that the model has not seen in training — explainable learning.*

**Mini bibliography** — Portillo et al. 2020; Pat et al. 2020; Liang et al. 2023; Stoppa et al. 2023; Sarmiento et al. 2021; Melchior et al. 2023; Huertas-Company & Lanusse 2023

# Domain Gap



Neural networks learn by optimising their performance on the training set.
**They tend to converge to domain-specific solutions.**

# A domain adversarial neural net to classify nebulae in the ISM



Belfiore, Ginolfi+ 25

ML Classification Map
SNR+DIG  HII  PNe

a)

PNe
SNRs

$H\alpha$ S/N

HII region
PNe
SNR + DIG

- 24 MUSE pointings
- 8 billion pixel (spatial + spectral)
- Needs to work well on other IFU data too!

Bracci+25, PhD @florence

Normalised, continuum-subtracted spectrum

Convolutional block

Residual Block — 16 | Avg. Pooling 1D | Residual Block — 32 | Avg. Pooling 1D | Residual Block — 64 | Avg. Pooling 1D

$$x \rightarrow \text{Conv1D} \rightarrow \text{Conv1D} \rightarrow \oplus \rightarrow x + F(x)$$

$F(x)$

Flatten

Spectral lines-based representation

MLP block

Dense — 64 | Dropout | Dense — 128 | Dropout

Line location

Normalised, 10-points continuum vector

Dense — 32

Dense layer

BatchNorm + Concatenate

"Shared info" layer

MLP block

Dense — 64 | Dropout | Dense — 128 | Dropout

$z$

MLP block

Dense — 64 | Dropout | Dense — 32 | Dropout

$M_{\text{star}}$

MLP block

Dense — 64 | Dropout | Dense — 32 | Dropout

$SFR$

Our plan 👀
**a tool ready-to-use with any data**

$X_{\text{Source}}$

$X_{\text{Target}}$

Domain-Invariant Feature Extractor

$f$

Classifier

Label Classification

Discriminator

Source/ Target Classification

# Degenerate problems → Bayesian approach (e.g., MCMC sampling)



$$P(\theta|data) \;=\; \frac{P(data|\theta)P(\theta)}{\int P(data|\theta)P(\theta)\,\mathrm{d}\theta}$$

# Simulation-based inference with conditional flows



Yang +19

**Normalising flows** —> they learn invertible transformation from simple to complex distribution

**Mini bibliography for simulation based inference in astro using conditional flows** — List et al. 2021; Mishra-Sharma & Cranmer 2021; Hahn et al. 2023a,b; Lemos et al. 2023; Bhardwaj et al. 2023; Alvey et al. 2023; Aubin et al. 2023; Hahn et al. 2024; Massara et al. 2024; Candebat et al. 2024; Angeloudi et al. 2024; Iglesias-Navarro et al. 2024; Barret & Dupourqué 2024 etc etc

# Simulation-based inference with conditional flows



$f_1(\mathbf{z}_0)$    $f_i(\mathbf{z}_{i-1})$    $f_{i+1}(\mathbf{z}_i)$

$\mathbf{z}_0$   $\mathbf{z}_1$   ...   $\mathbf{z}_{i-1}$   $\mathbf{z}_i$   ...   $\mathbf{z}_K = \mathbf{x}$

$\mathbf{z}_0 \sim p_0(\mathbf{z}_0)$    $\mathbf{z}_i \sim p_i(\mathbf{z}_i)$    $\mathbf{z}_K \sim p_K(\mathbf{z}_K)$

Yang +19

Normalizing Flow

Initial density    $f(\theta)$

Transformation conditioned to an observable

**Normalising flows** —> they learn invertible transformation from simple to complex distribution

**Mini bibliography for simulation based inference in astro using conditional flows** — List et al. 2021; Mishra-Sharma & Cranmer 2021; Hahn et al. 2023a,b; Lemos et al. 2023; Bhardwaj et al. 2023; Alvey et al. 2023; Aubin et al. 2023; Hahn et al. 2024; Massara et al. 2024; Candebat et al. 2024; Angeloudi et al. 2024; Iglesias-Navarro et al. 2024; Barret & Dupourqué 2024 etc etc

# Cloudy$_{\mathrm{SBI}}^{-1}$

Observations $x^{obs}$

PHANGS
MUSE

Joint posteriors $p(\theta, x)$

$\theta_1$

Probability

$\theta_2$

Prior $p(\theta)$

$\theta$

$\log(U)$
$Z$
$\log(N/O)$
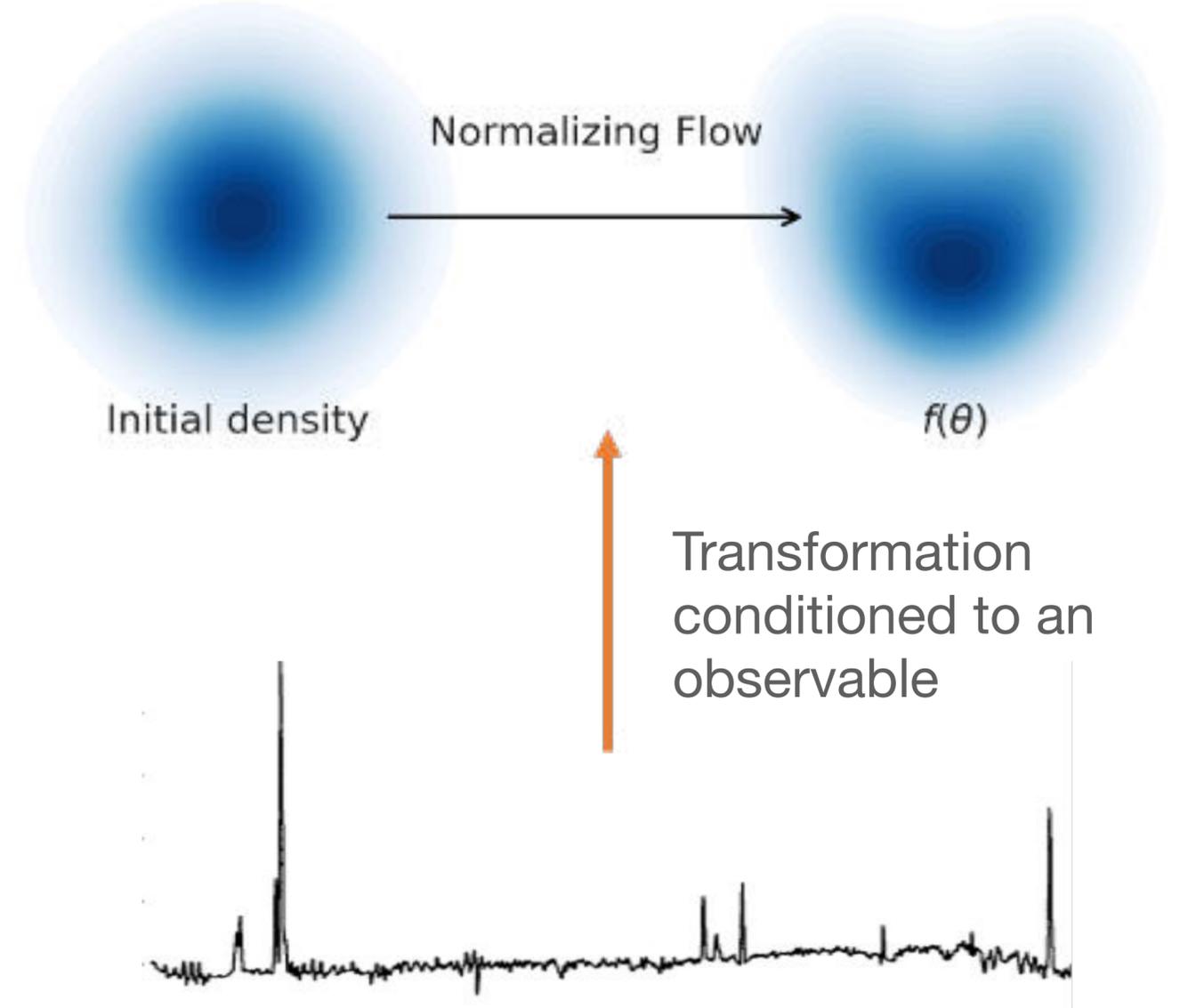
Model x = $s(\theta)$

$x$

Simulated nebular emission

Flux

$\lambda$

Density estimator
(normalising flows)

**Cloudy** photoionization models

Posterior Comparison for Instance 2866

$-2.651^{+0.041}_{-0.040}$

$8.150^{+0.043}_{-0.041}$

$-1.393^{+0.038}_{-0.037}$

logU

Met

log(N/O)

MCMC
SBI NoiseNet
True
NoiseNet median
NoiseNet 68% CI

Bracci+, in prep

Advantages: accurate & <u>amortized</u>

A proposed framework:

## SBI w. domain-invariant neural flows

- Use virtual data from state-of-the-art simulators
- Works with any type of data & ensures that summary statistics are domain-invariant
- Handles uncertainties effectively
- Currently being tested…

# Take away messages

- The upcoming volume and complexity of spectral data, especially around cosmic noon, call for the help of deep learning.

- The integration of simulation-based inference (SBI) with machine learning has emerged as a transformative approach across scientific disciplines, and it's starting to impact galaxy spectroscopy.

- It is crucial to address the domain gap between simulations and real data; domain adaptation techniques can help reduce this gap.

- Methods such as conditional neural flows can make models uncertainty-aware, effectively mimicking Bayesian approaches.