# Challenges in AI/ML at the Cosmic Frontier
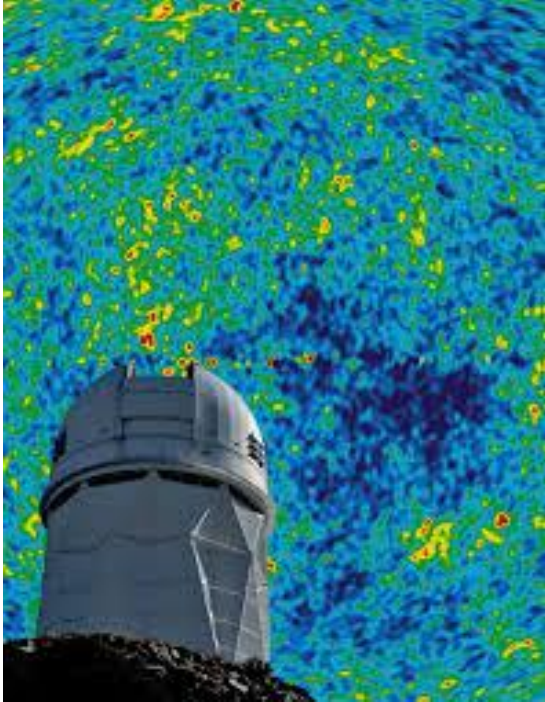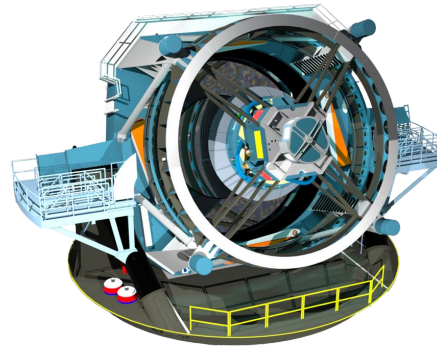
**Simone Ferraro**

**(Lawrence Berkeley National Lab)**

SLAC Summer Institute
Aug 8, 2023

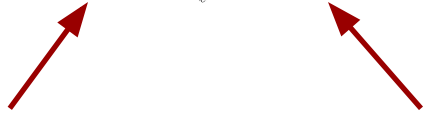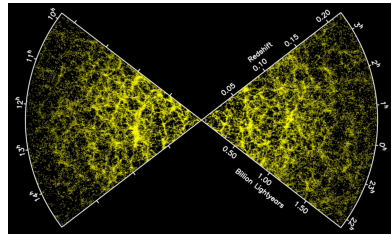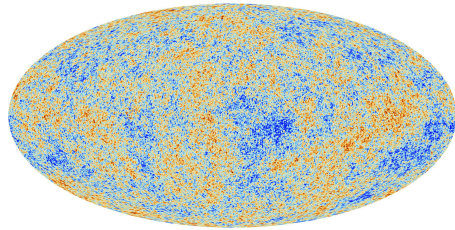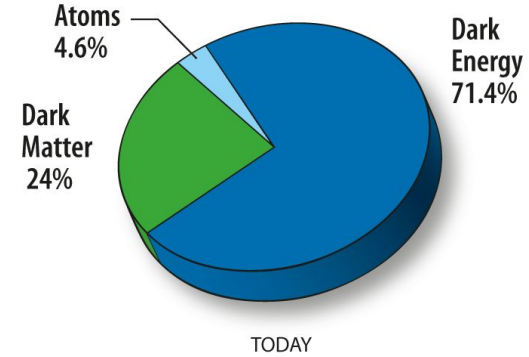# Plan for this talk

- ***Part I:*** What is the Cosmic Frontier and what are the observables?

- ***Part II:*** Are we being efficient (speed)?

- ***Part III:*** Are we extracting the whole information?

- ***Part IV:*** Do we understanding the data?

Simone Ferraro (LBNL)

# Part I:
# What is the Cosmic Frontier and what are the observables?

Simone Ferraro (LBNL)

# A simple yet strange Universe

Planck, BOSS

But the model is based on…
- **Dark Matter** (?)
- **Dark Energy** (?)
- **Inflation** (?)
- **Neutrinos and other light particles** (?)

A major goal of the Cosmic Frontier program is to understand these "ingredients"!

# A simple yet strange Universe
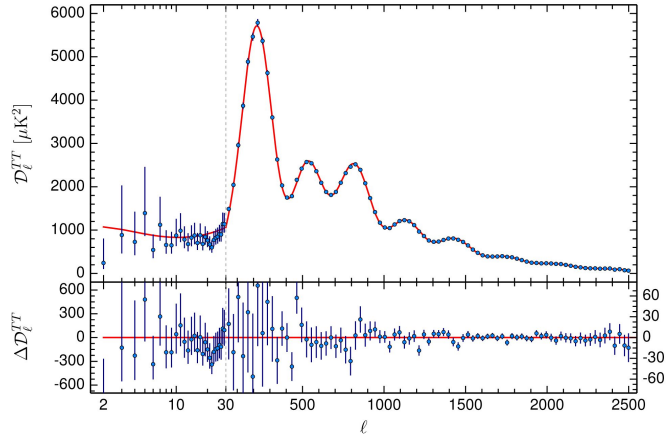
Fit fully characterized by <u>6 numbers*</u>
(with no evidence of needing more):

$$\{\Omega_m, \Omega_b, A_s, n_s, \tau, H_0\}$$

matter

"baryons"

amplitude of
primordial
fluctuations

slope
of primordial
fluctuations

reionization

expansion
rate today

Planck, BOSS

* and a few assumptions such as flat geometry and minimum mass neutrinos

# A brief history of the Universe

**94%** of photons travel from the CMB to us without scattering*

**6%** scatter with matter

On small scales, the Cosmic Microwave Background (CMB) contains a "map" of the entire observable universe

*path slightly deflected by gravitational lensing

Simone Ferraro (Berkeley)

# Cosmic microwave background (CMB)

Planck Satellite (2018)

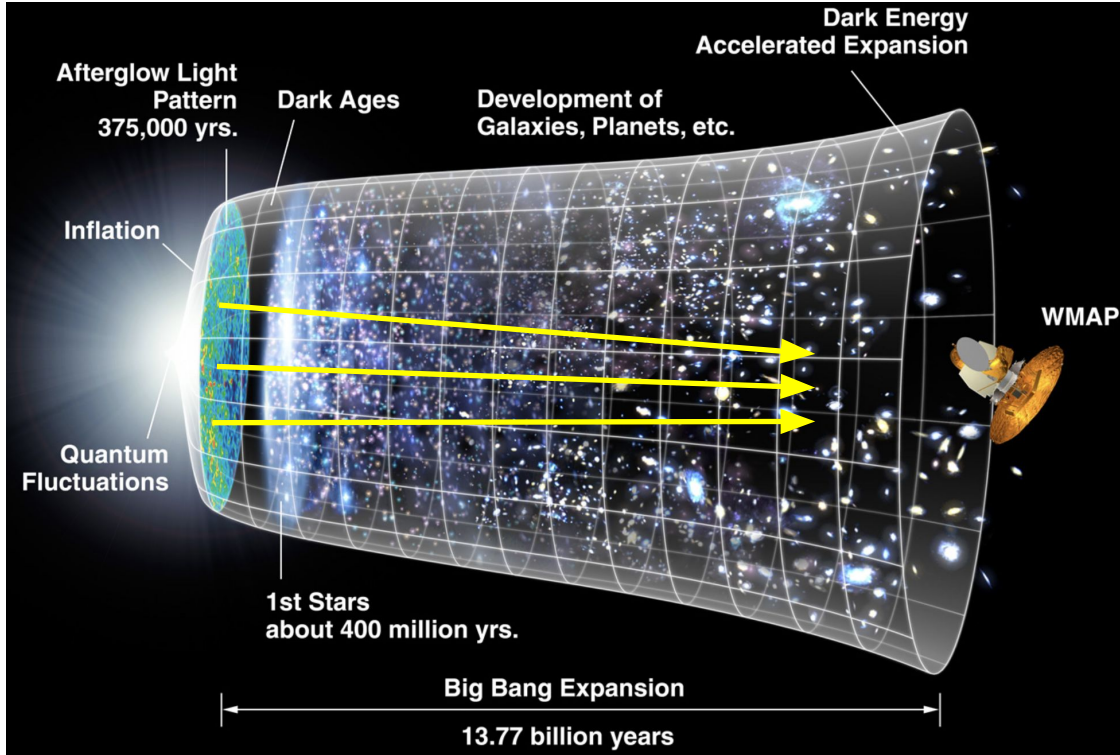**"primary fluctuations"**

- Large scales (< 1 deg) ▫ primordial

- Smaller scales (> 1 deg) ▫ processed by (known) plasma physics + gravity

"CMB (angular) power spectrum"



Simone Ferraro (Berkeley)

# CMB lensing

Paths of CMB photons deflected by matter ⬚ create statistical anisotropy that can be measured
**Can make maps of the projected matter density (including Dark Matter) to the CMB!**





ACT DR6 lensing map

Simone Ferraro (Berkeley)

# Galaxy (weak) lensing

Rubin Observatory
LSST



"cosmic shear power spectrum"



5 redshift bin tomography
$\Omega_s=20000 \deg^2, n_g=50 \text{ arcmin}^{-2}, \sigma_\epsilon=0.22$

$1.6<z_5$
$1.2<z_4<1.6$
$0.8<z_3<1.2$
$0.4<z_2<0.8$
$0<z_1<0.4$

$l(l+1)C_l/2\pi$

$l$

large scales            small scales

Unlensed            Lensed

Without Shape Noise

With Shape Noise

Simone Ferraro (LBNL)            Wikipedia

# Large Scale Structure (LSS)



large scales                small scales

"galaxy power spectrum"

Planck TT
Planck EE
Planck φφ
SDSS DR7 LRG
BOSS DR9 Ly-α forest
DES Y1 cosmic shear

$P_{\mathrm{m}}(k)\ [(h^{-1}\mathrm{Mpc})^3]$

Wavenumber $k\ [h\,\mathrm{Mpc}^{-1}]$

$M_r < -20.44$
$M_r \geq -20.44$
$21.25 \leq \eta < 28.75$

Comoving distance [$h^{-1}$ Mpc]

SDSS

Each dot is a real galaxy!

The small-scale matter power spectrum, $k > k_{fs}$, is reduced in presence of massive neutrinos:

- On larger scales $\nu$s cluster in the same way as cold dark matter
- Free-streaming $\nu$s do not cluster
- The growth rate of CDM and baryon fluctuations is reduced

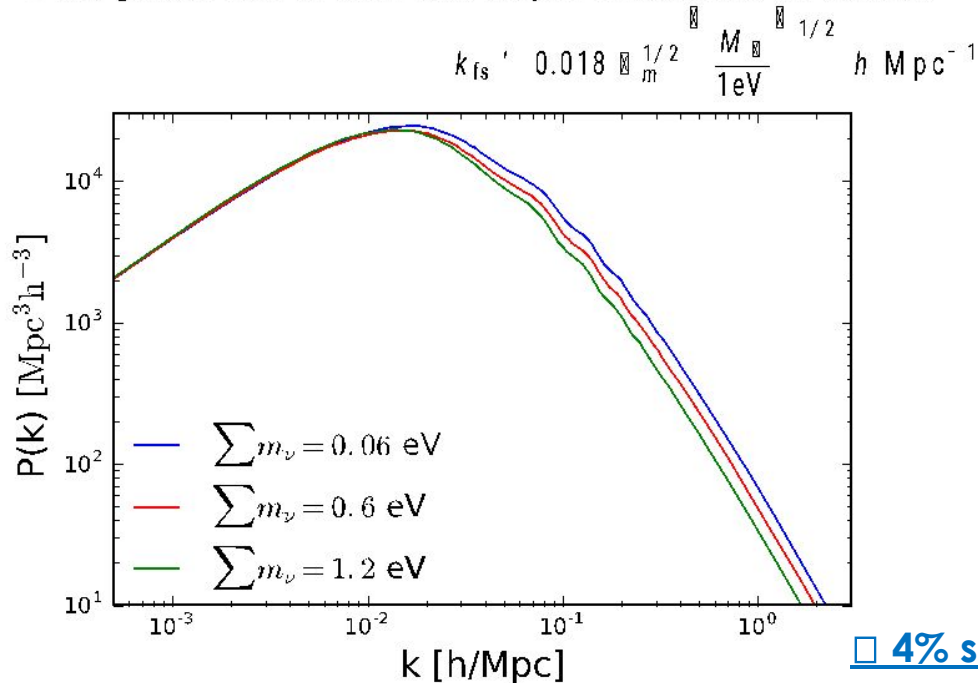$$k_{fs} \simeq 0.018 \, \Omega_m^{1/2} \left( \frac{M_\nu}{1eV} \right)^{1/2} h \; Mpc^{-1}$$



Power suppression at small scales

$$\frac{\Delta P(k)}{P(k)} \simeq -8 f_\nu$$

$$f_\nu = \omega_\nu / \omega_m$$
$$= 0.5\%$$

**4% suppression minimum!**

+ non-linear calculations: additional suppression at large k
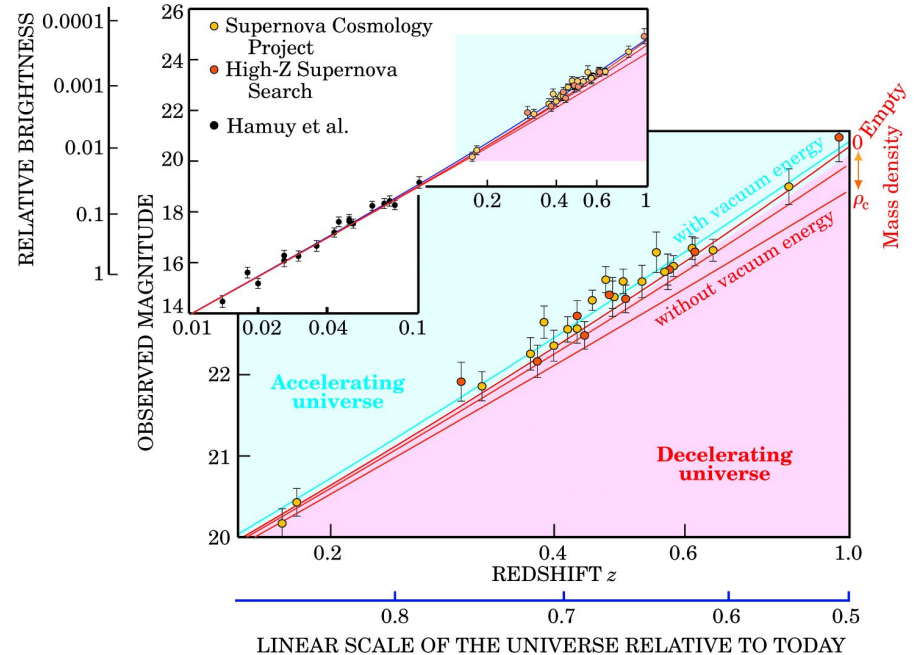see Villaescusa-Navarro et al. 2013

# Transients

**Extragalactic:**
- Supernovae/kilonovae
- Fast Radio Bursts, gamma ray bursts
- Tidal disruption events
- Strong lensing time delays
- …

**Galactic:**
- Asteroids
- Interacting binaries
- Transiting exoplanets
- Microlensing
- Pulsars
- …

**D**ark **E**nergy **S**pectroscopic **I**nstrument: Massively multiplexed spectroscopic survey with 5000 robotic fibers, over ~14,000 sq. deg



Five target classes

**35 million** redshifts

(SDSS x20)

**DESI (2021-2026)**

**2.4 million QSOs**
Lya      z > 2.1
Tracers  1.0 < z < 2.1

**17 million ELGs**
0.6 < z < 1.6

**6 million LRGs**
0.4 < z < 1.0

**10 million Brightest galaxies**
0.0 < z < 0.4

4
2
1
0.7
0.2

Redshift

S. Ferraro

Simone Ferraro (LBNL)

**Dark Energy Spectroscopic Instrument**

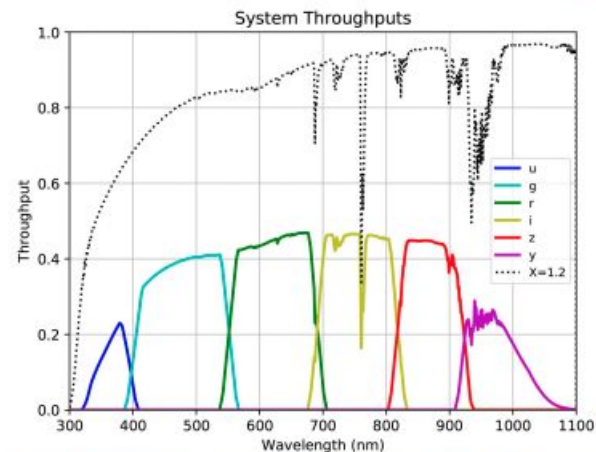Slide 4

# The Vera C. Rubin Observatory

Location: El Peñón, Cerro Pachon, Chile
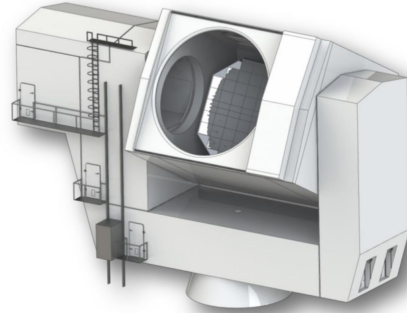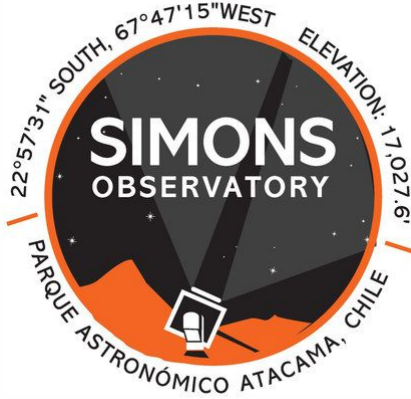  (median seeing 0.67 arcsec)

Specs: 8.4m mirror, 3.2 Gigapixels camera
  9.6 sq. deg. field of view (~40 full moons),
  6 broadband filters (*ugrizy*)

It will perform the 10 year LSST survey of the sky

System Throughputs

Throughput

Wavelength (nm)

u
g
r
i
z
y
X=1.2

Optical elements seen
by thrilled observers

# The CMB landscape – mid 2020s
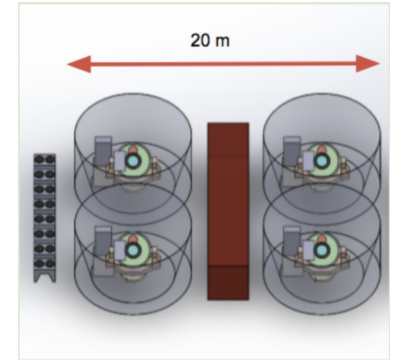
Large Aperture Telescope
one 6 meter in diameter



Small Aperture Telescopes
42 cm refractors

Large frequency coverage (30 – 270 GHz)

- 10 Countries
- 40+ institutions

Fully funded
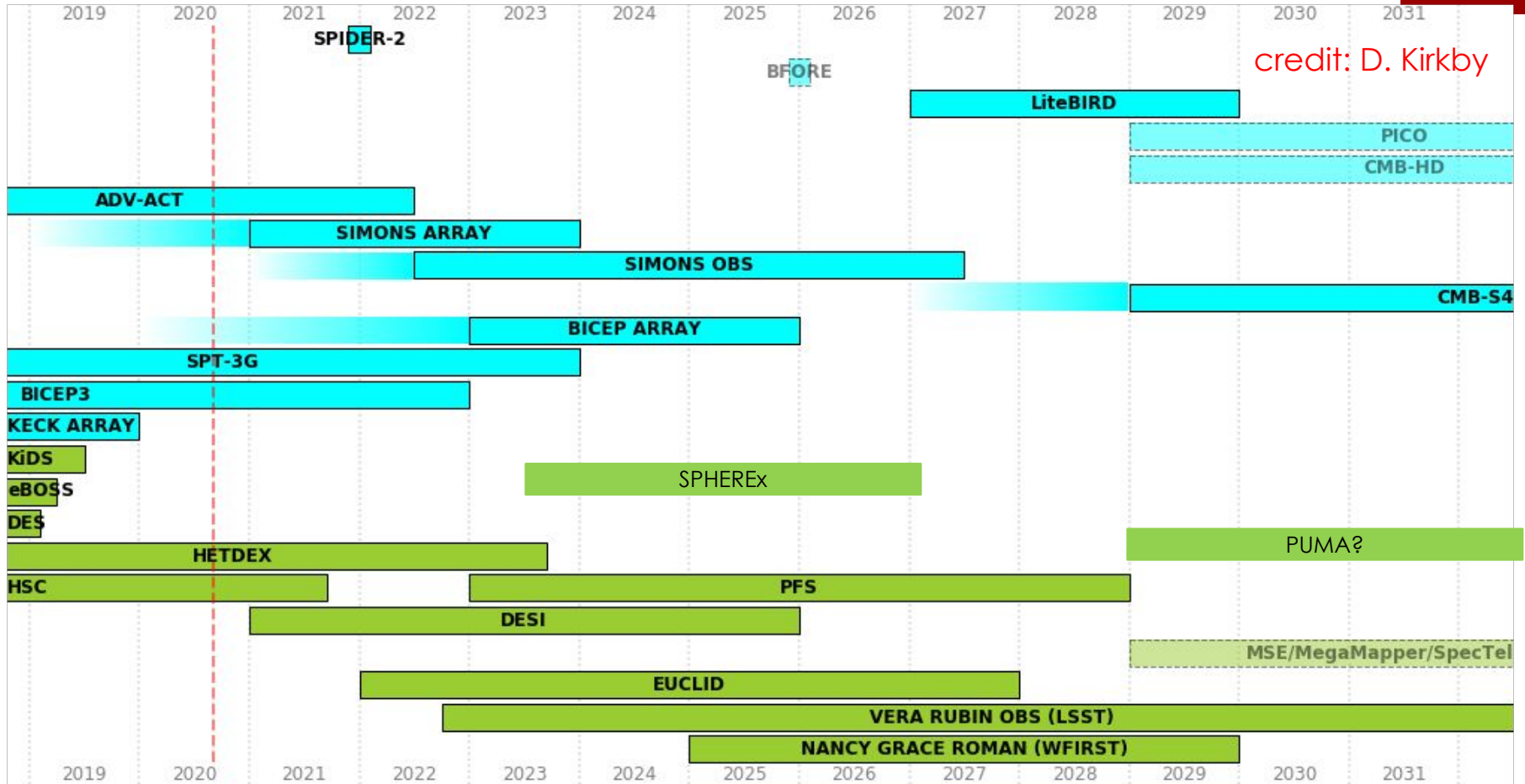6-year program
**First light in 2024!**

# CMB-S4

- **CMB S4:** next generation ground based experiment
- Factor of ~10 increase in sensitivity
- ETA ~late in this decade

- Multi-agency effort (DOE & NSF)

# Looking ahead: the "explosion" of surveys



credit: D. Kirkby

# Dark Matter direct and indirect detection

IF Dark Matter interacts (weakly) with the Standard Model, can look for scattering/recoil (_direct detection_).
Several targets: Xenon, Germanium, etc

Also: "_indirect detection_" in astrophysical systems (Eg. Fermi gamma ray satellite)

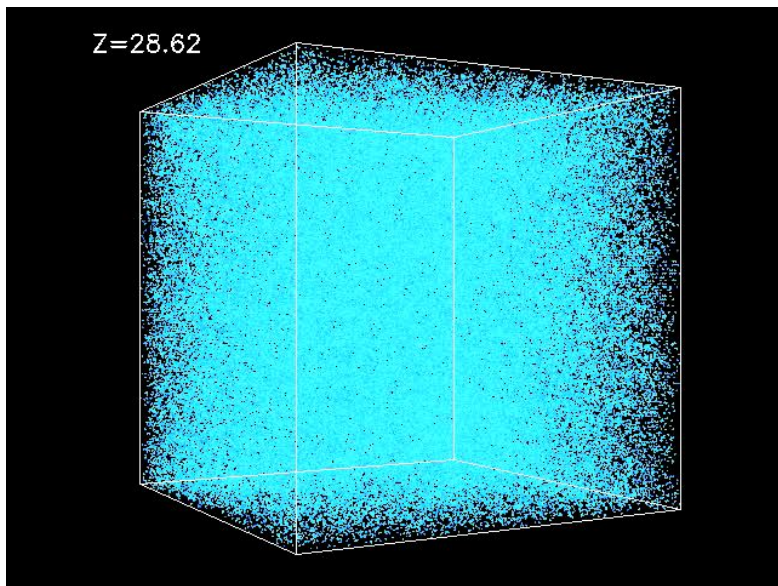❑ **Maria Elena Monzani's lecture on anomaly detection!**



Electrons

Outgoing Particle

Incoming Particle

LZ/SLAC

# Part II:
# are we being efficient (speed)?

# Challenge 1: theoretical model

- <u>Complex & non-linear</u> dependence of theory on cosmological parameters even for power spectrum (2 point function). Often <u>no analytical form</u>, and prediction relies on expensive numerical simulations.



Z=28.62

Calculating

$$\text{theory}(\Omega_m, \Omega_{DE}, A_s, n_s, \tau, \ldots, \{\text{nuisance parameters}\})$$

Can take minutes to hours (or more).
Often too slow for parameter inference!

**SOLUTION: <u>Build emulators</u>!**
Reduce theory calculation to O(ms) per call

☐ **Joe DeRose's lecture on emulators**

# Challenge 2: parameter inference

- <u>High dimensional problem:</u> typically > 100 parameters (dimensions) inference. Slow or impossible.

  - The probability distribution $P(\Theta|\mathbf{d})$ (posterior) for model parameters $\Theta$ given data $\mathbf{d}$ can be related to the probability $P(\mathbf{d}|\Theta)$ (likelihood) of an experiment giving data $\mathbf{d}$ for model parameters $\Theta$ using the Bayes' theorem:

$$P(\Theta|\mathbf{d}) = \frac{P(\mathbf{d}|\Theta)P(\Theta)}{P(\mathbf{d})} \qquad (15)$$

  where $P(\Theta)$ is called the prior and $P(\mathbf{d}) = \sum P(\mathbf{d}|\Theta)P(\Theta)$ is used for the normalization purpose.

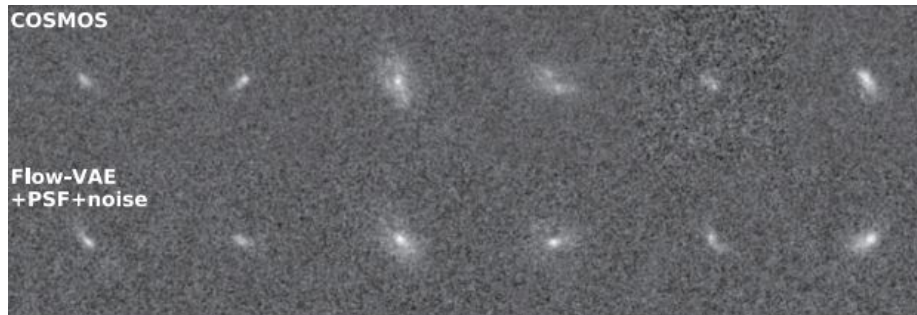Typically you would <u>sample the posterior</u> by **Monte-Carlo**.
In high-dimension (>100 parameters), algorithms involving ***gradients*** are more efficient (eg. Hamiltonian Monte Carlo).

### SOLUTION: <u>Build a differentiable likelihood + differential emulators!</u>
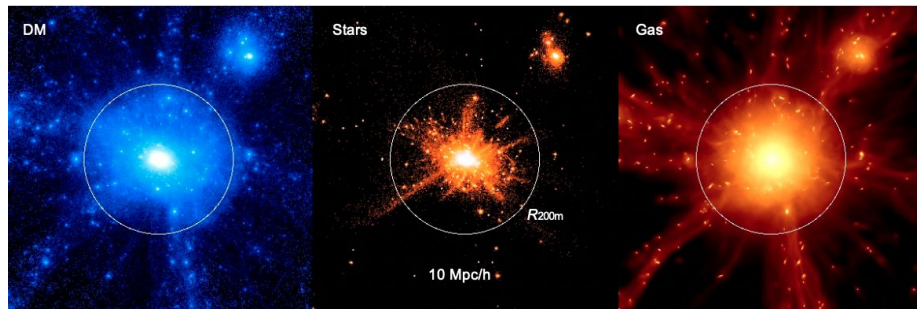### ☐ Joe DeRose's lecture

Simone Ferraro (LBNL)

# Mock data generation

Lensing data



COSMOS

Flow-VAE +PSF+noise

arXiv:2008.03833

"Gas pasting" on Dark-Matter only simulations



DM            Stars            Gas

$R_{200m}$

10 Mpc/h

arXiv:2110.02232
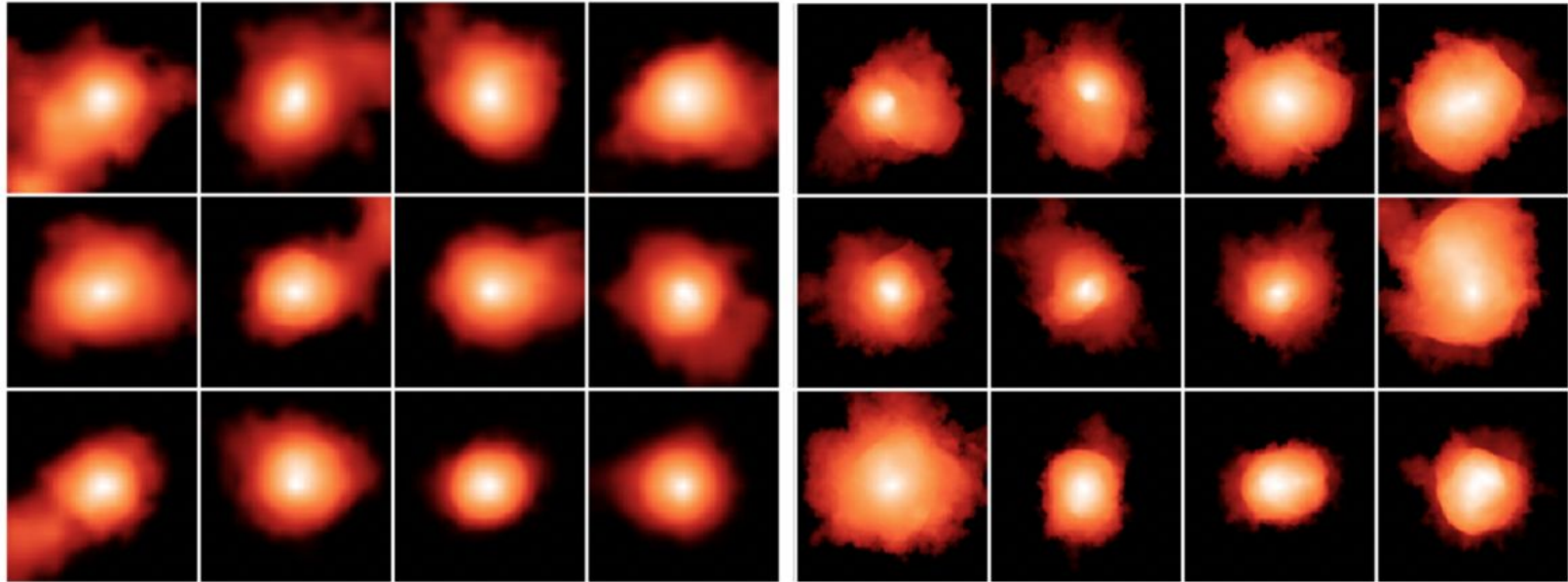
Both cases based on (conditional) Variational Auto-Encoder (VAE)

Similar applications with Generative Adversarial Networks (GANs)

Diffusion models?

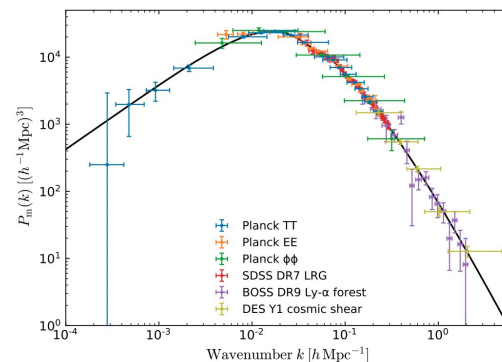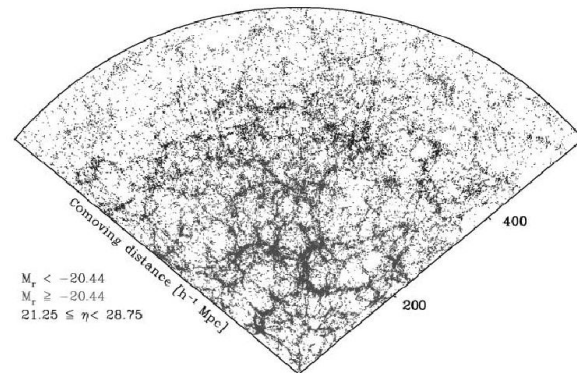 See lectures on generative models applications

Simone Ferraro (LBNL)

Can you tell which one is generated by an (expensive) hydrodynamical simulation and which is generated by CVAE?
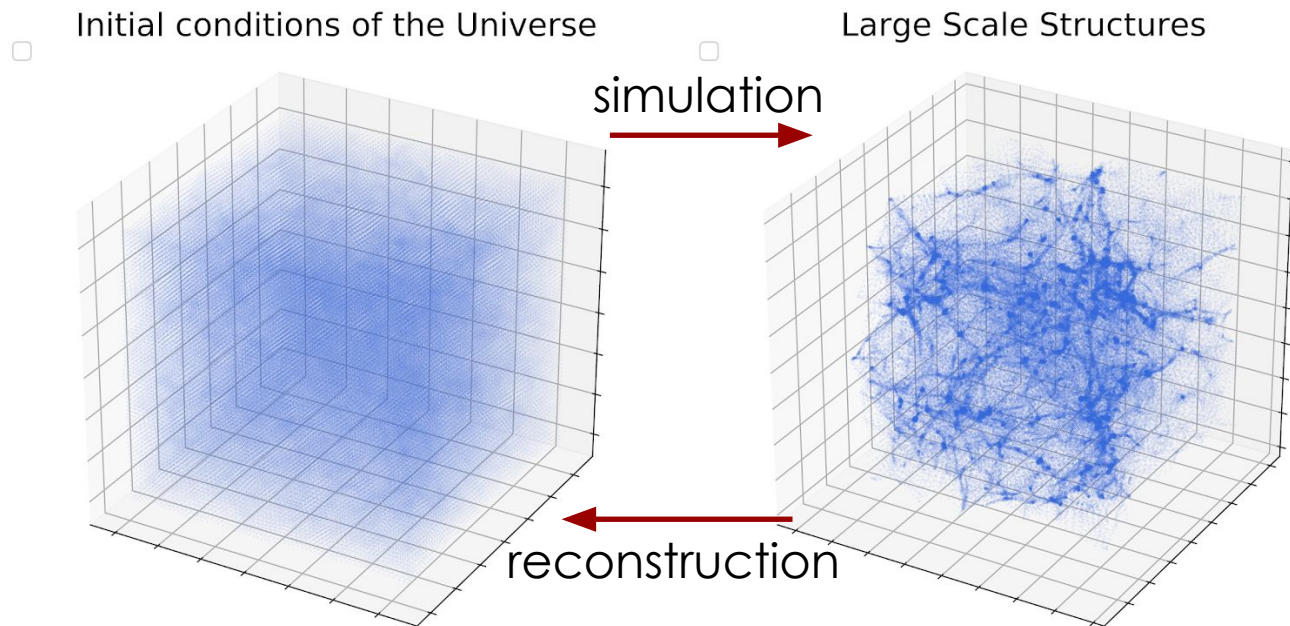
# Part III:
## are we extracting the whole information?

# Beyond the power spectrum

- Power spectrum (2pt function) contains the whole information only for *Gaussian* fields
- Can perturbatively consider 3pt function and higher, but (in general), limited information available.
- But in general, no guidance on what's the most informative statistic…

- Several options available:
  - Field-level inference
  - Compression in a "small" number of summary statistics
  - In both cases, likelihood may not be known analytically ☐ likelihood-free/simulation-based inference)
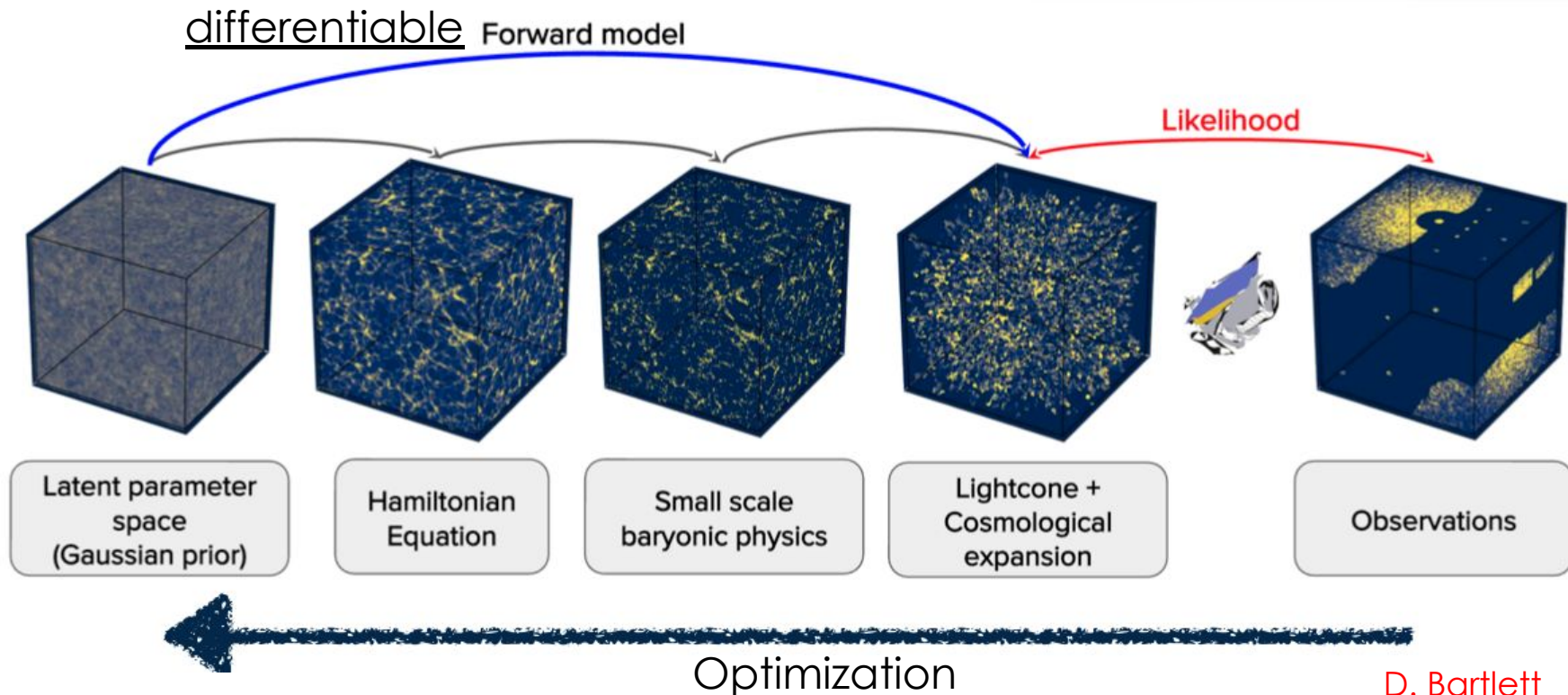
# Field-level inference

Initial conditions of the Universe

simulation

Large Scale Structures

reconstruction

The initial conditions are very close to Gaussian: they contain the whole information. Can we reconstruct them?

# Field-level inference



differentiable Forward model

Likelihood

Latent parameter space (Gaussian prior)

Hamiltonian Equation

Small scale baryonic physics

Lightcone + Cosmological expansion

Observations
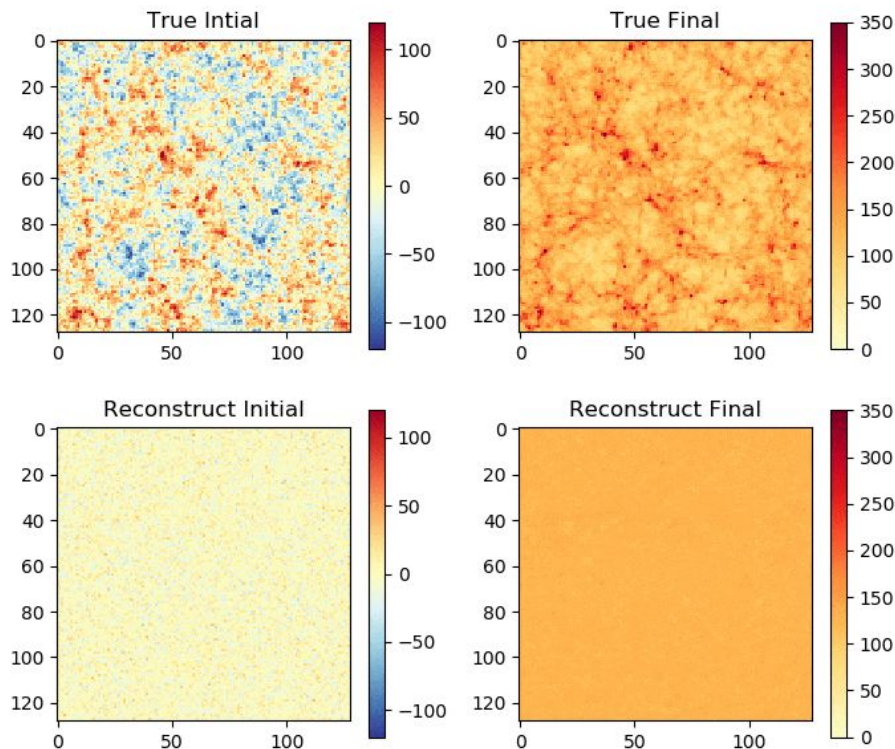
Optimization

D. Bartlett

"solving the inverse problem by optimization"

Simone Ferraro (LBNL)

# Field-level inference

Optimization converges in O(20 steps), even though $N_{dim} = N_{pix} \sim 10^6$ or more!

But: want to marginalize over the initial conditions to extract cosmological parameters. Active area of research and questions remain!

Can use Laplace approximation or MUSE (see arXiv:2112.09354).

Or… full HMC sampling (eg. BORG https://www.aquila-consortium.org/)

☐ **See "Introduction to Differentiable Programming in Jax" (F. Lanusse)**

https://blog.tensorflow.org/2020/03/simulating-universe-in-tensorflow.html (Modi, Lanusse et al)

Simone Ferraro (LBNL)

# More generally: the full problem

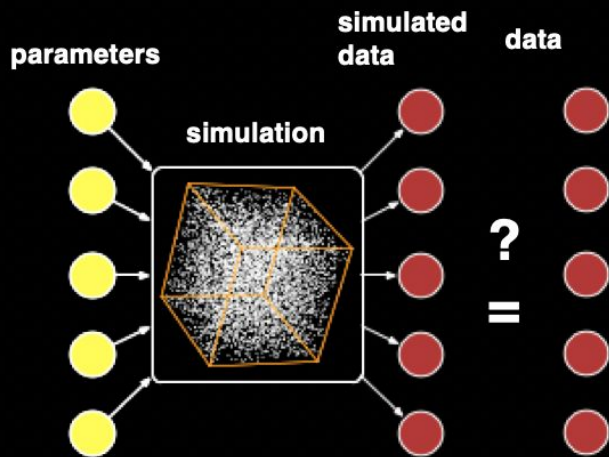Often we need more freedom than a traditional likelihood approach:

- We may not know what the likelihood is (Gaussian approximation is often a bad one!)
- We may summarize, cut, mask the data any way we want
- Observational or instrumental effects are hard to treat analytically but easy to simulate.

Simulating data is often much easier than deriving an accurate likelihood
  ☐ **Simulation-Based Inference** (SBI)

SBI = Inference "engine" when explicit likelihood is intractable or unknown, but simulation is possible.
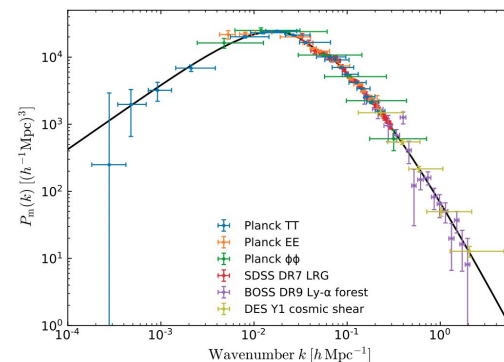
# Simulation based inference introduction



Simulation based inference
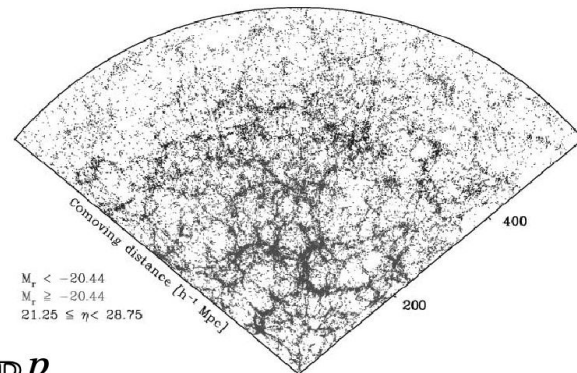
Simplest implementation of "**Approximate Bayesian Computation**" (ABC)
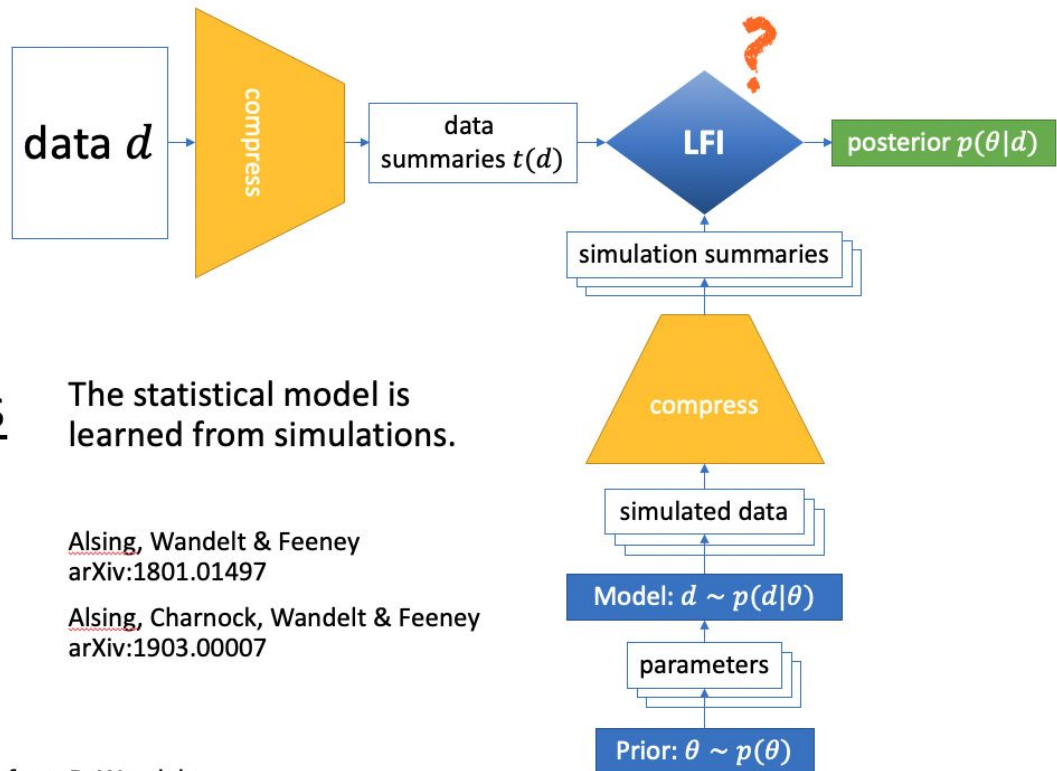
Suffers from severe "Curse of dimensionality"

# Compression beyond the power spectrum

- Beyond the power spectrum, we have little guidance on what's the most informative statistic…

- *Score compression* and *Information Maximizing Neural Networks* (IMNN): $\mathbf{t(d)} : \mathbb{R}^N \rightarrow \mathbb{R}^p$ (p < N) produce a small number of summary statistic that maximize the retained Fisher information.

- See arXiv:1802.03537 for more info.

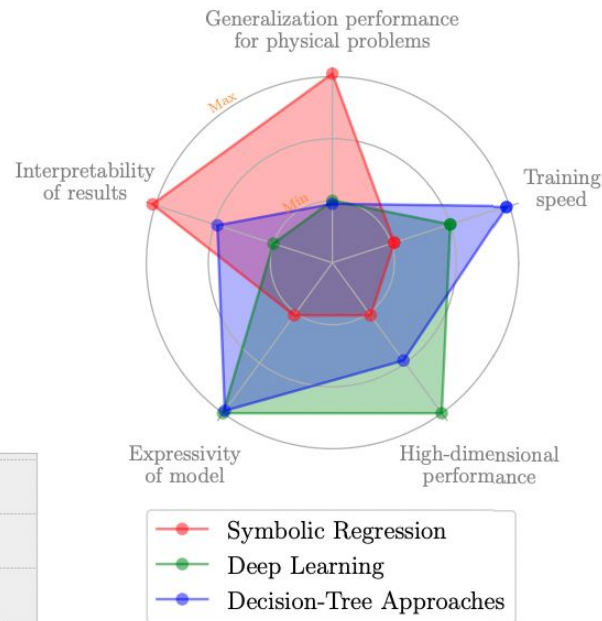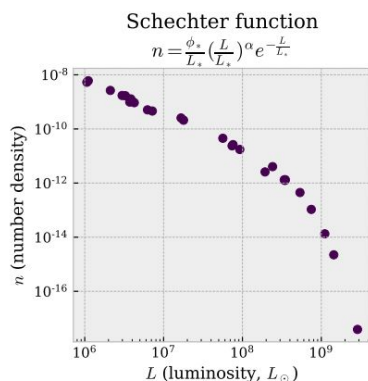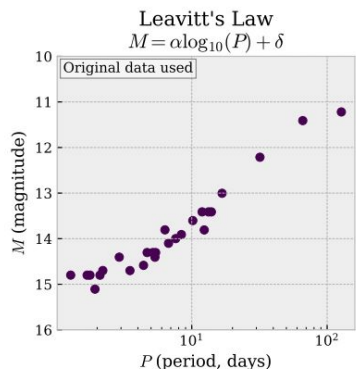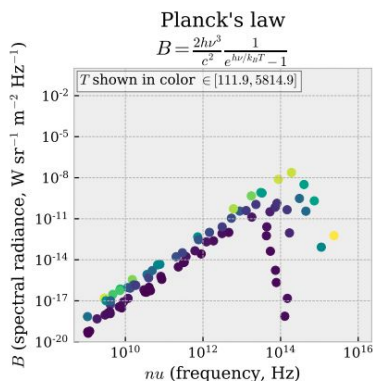# Likelihood-free inference (LFI)



Likelihood-Free
Inference (LFI)
with summaries

The statistical model is
learned from simulations.

Alsing, Wandelt & Feeney
arXiv:1801.01497

Alsing, Charnock, Wandelt & Feeney
arXiv:1903.00007

Slide from B. Wandelt

- Search the space of analytic equations to fit some data
- Often done by hand but efficient algorithms exist!
- Concise
- Interpretable



Planck's law

$$B = \frac{2h\nu^3}{c^2} \frac{1}{e^{h\nu/k_BT}-1}$$

$T$ shown in color $\in [111.9, 5814.9]$

Leavitt's Law

$$M = \alpha \log_{10}(P) + \delta$$

Original data used

Schechter function

$$n = \frac{\phi_*}{L_*} \left(\frac{L}{L_*}\right)^\alpha e^{-\frac{L}{L_*}}$$

arXiv:2305.01582

arXiv:2201.01305

Simone Ferraro (LBNL)

https://github.com/MilesCranmer/PySR
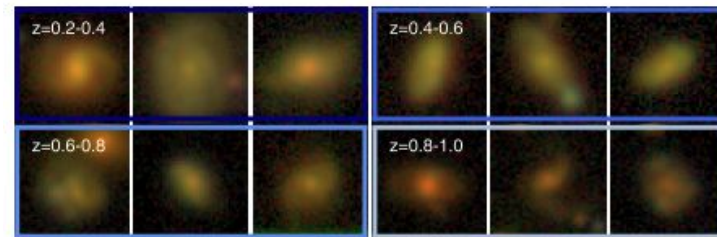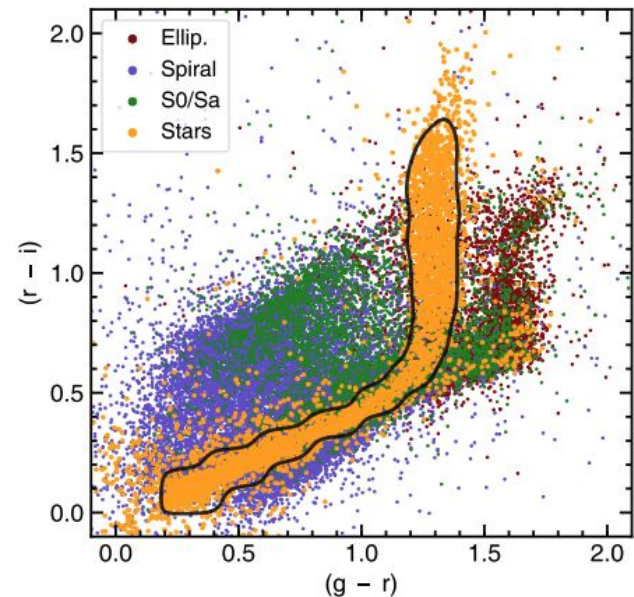
# Part IV:
# Do we understand the data?

# Classification

- Source classification (stars vs galaxies vs quasars etc)
- Transient classification
- Classification of the cosmic web (voids, sheets, filaments etc)
- Photometric redshifts
- …

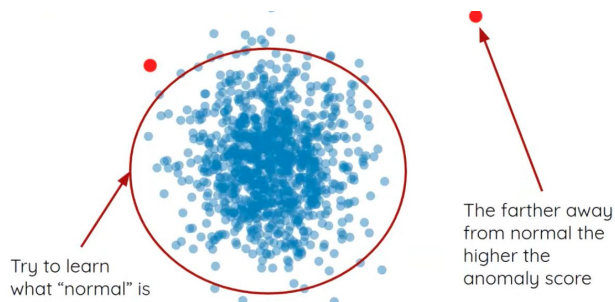Both supervised and unsupervised methods. Need to allow for the existence of unexpected patterns!



(a) *Spiral galaxies.*

Simone Ferraro (LBNL)

arXiv:1909.10537

# Anomaly detection

Several goals:
- Eliminate the influence of outliers or contaminants
- Find new signals when the signal is rare
  - *known unknowns*: supernovae, strong lenses, transients, …
  - *unknown unknowns*: eg. discovery of pulsars



Try to learn what "normal" is

The farther away from normal the higher the anomaly score

credit: M. Lochner

anomaly

vs

real

https://github.com/MichelleLochner/astronomaly

☐ **Maria Elena Monzani's lecture**
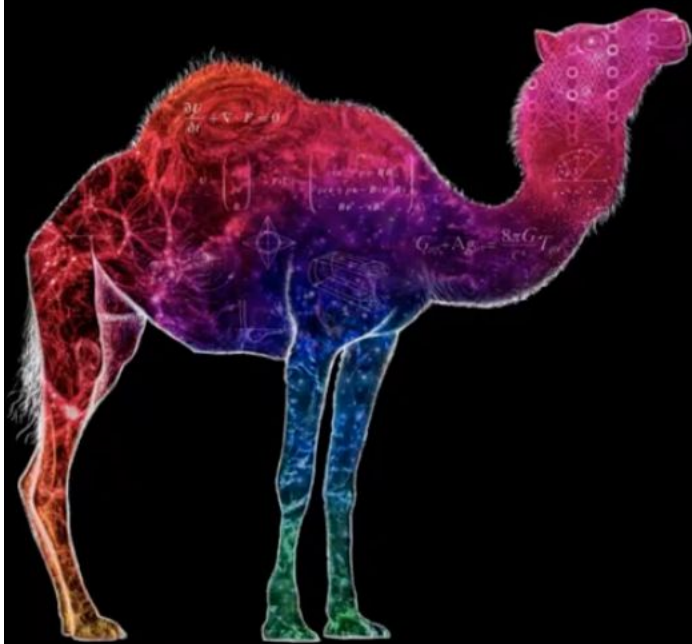
Simone Ferraro (LBNL)
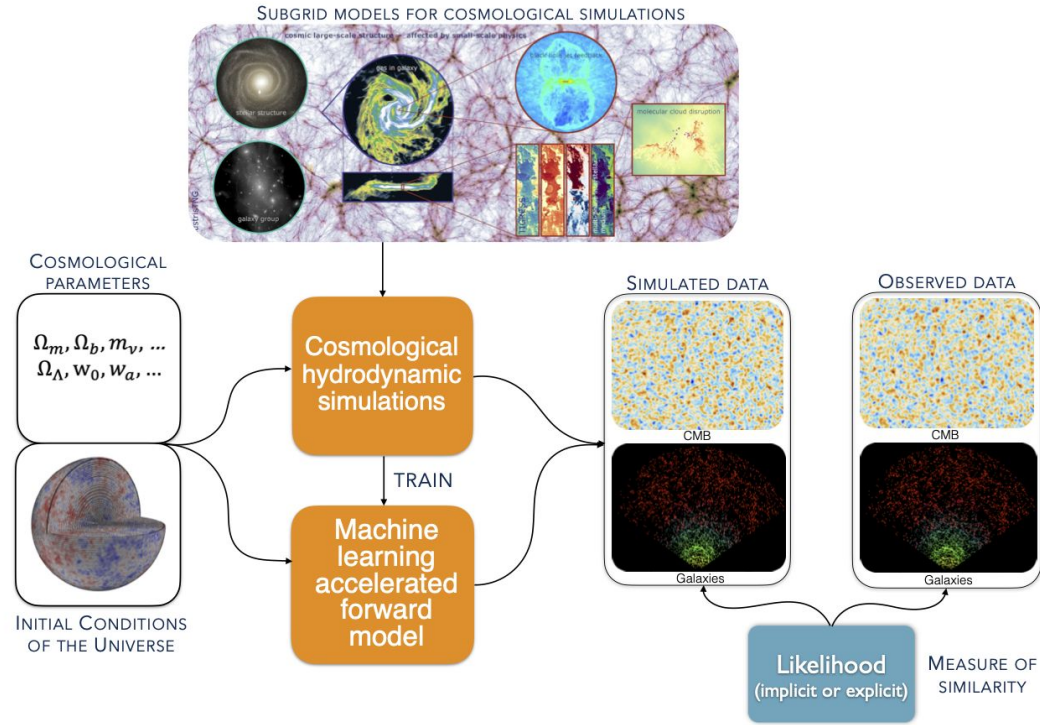
CAMELS

https://www.camel-simulations.org

Cosmology and Astrophysics with MachinE Learning Simulations

- A suite of 4,233 simulations

- 2,049 N-body; Gadget-III

- 2,184 state-of-the-art (magneto-)hydrodynamic sims

- AREPO/IllustrisTNG + GIZMO/SIMBA

- 6 parameters: $\{\Omega_m, \sigma_8, A_{SN1}, A_{SN2}, A_{AGN1}, A_{AGN2}\}$

- More than 100 billion resolution elements over combined volume of ~(400 Mpc/h)$^3$

- More than 2,000 cosmologies & astrophysics models; more than 140,000 snapshots

- Designed for machine learning applications

Simone Ferraro (LBNL)

F. Villaescusa-Navarro

# Learning the Universe

## Simons Collaboration on "Learning the Universe"



https://www.learning-the-universe.org/

# Conclusions

■ ML will help us solve cosmological problems that are intractable today

■ With great power, come great responsibility! Astrophysical systems are complex and often not fully understood. Model misspecification can lead to issues and great care needs to be taken.

■ Finally, for a comprehensive list of ML applications to cosmology, see https://github.com/georgestein/ml-in-cosmology

# **Thanks!**