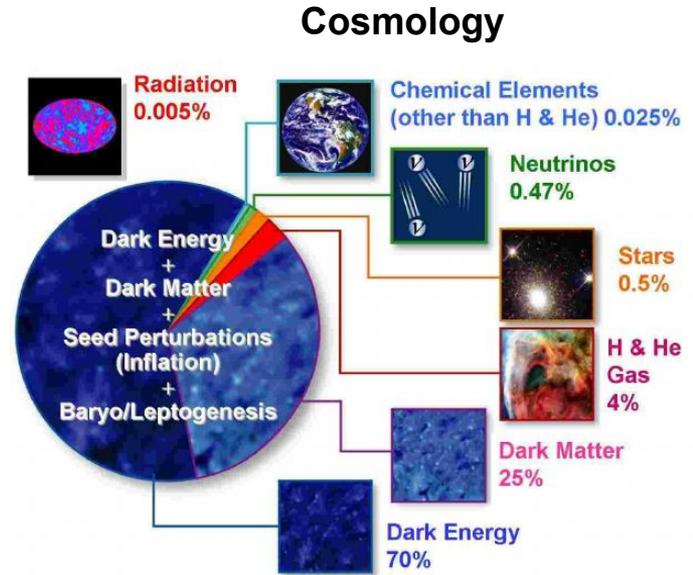
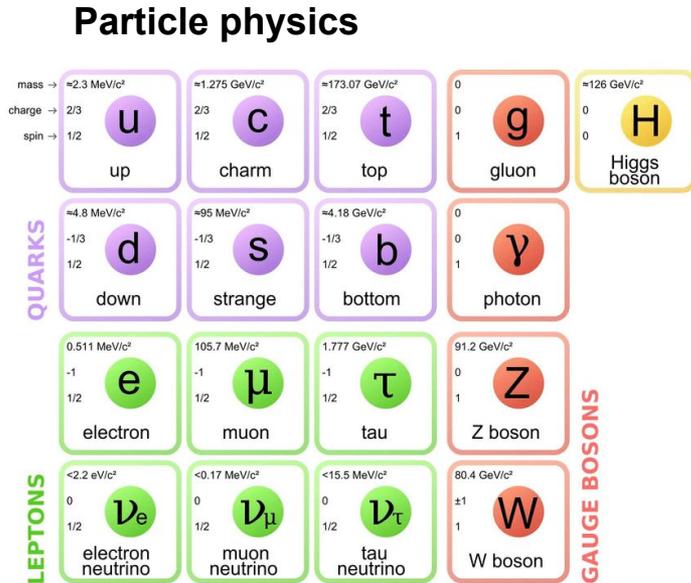


Enabling precision dark matter searches with Truncated Marginal Neural Ratio Estimation

Christoph Weniger

James Alvey (UvA), Uddipta Bhardwaj (UvA), Alex Cole (UvA), Adam Coogan (U. Montreal), Androniki Dimitriou (U. Valencia), Elias Dubbeldam (UvA), Mathis Gerdes (UvA), Kosio Karchev (SISSA), Ben Miller (UvA), Noemi Anau Montel (UvA), Roberto Trotta (SISSA)
Gilles Louppe (U. Liège), Anchal Saxena (Groningen), Patrick Forré (UvA), Samaya Nissanke (UvA), Maxwell Cai, Meiert Grootes, Francesco Nattino (eScience)

Standard models and open questions

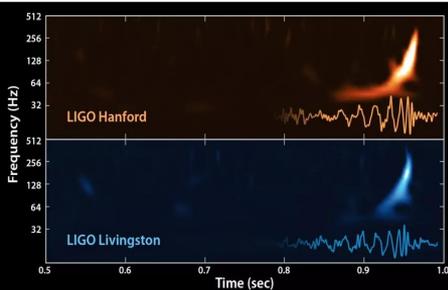


Many open questions: What is dark matter? Is dark energy dynamic? Where is the anti-matter? What caused seed-perturbations? How do black holes grow and merge? How do neutron stars develop? How did the first stars form? What stabilized the electroweak scale? Is there grand unification? How do galaxies form and grow?

Astrophysical searches for breaks in the standard models

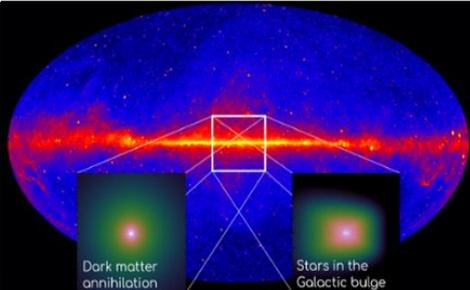
Strategy: Search for deviations from standard-model predictions in order to get answers on some of the open questions.

Image: LIGO



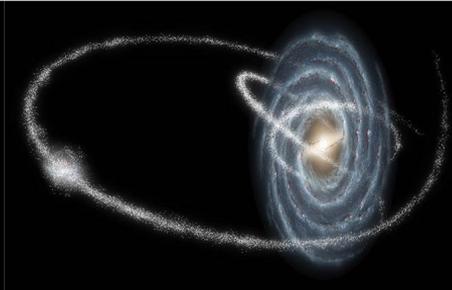
GRAVITATIONAL WAVES

Image: Fermi LAT/CW

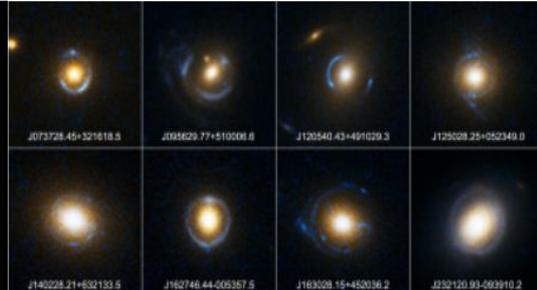


GAMMA RAYS

Image: NASA/JPL-Caltech/R Hurt (SSC/Caltech)



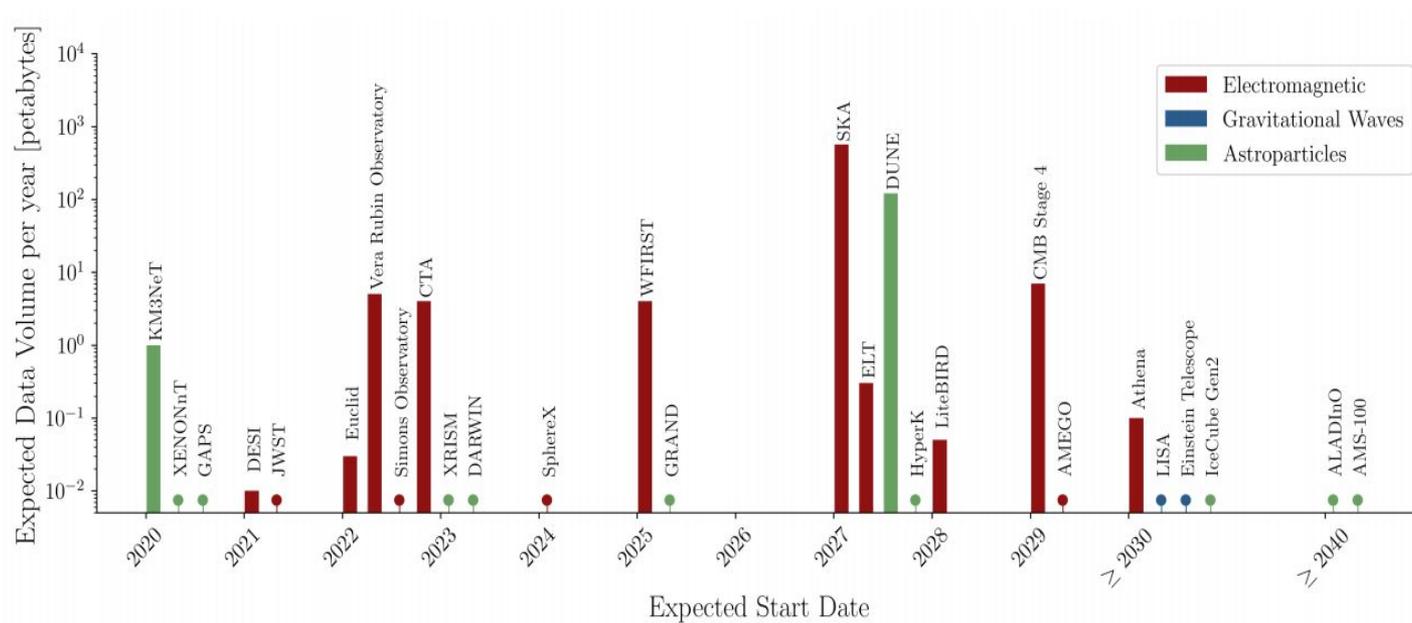
STELLAR STREAMS



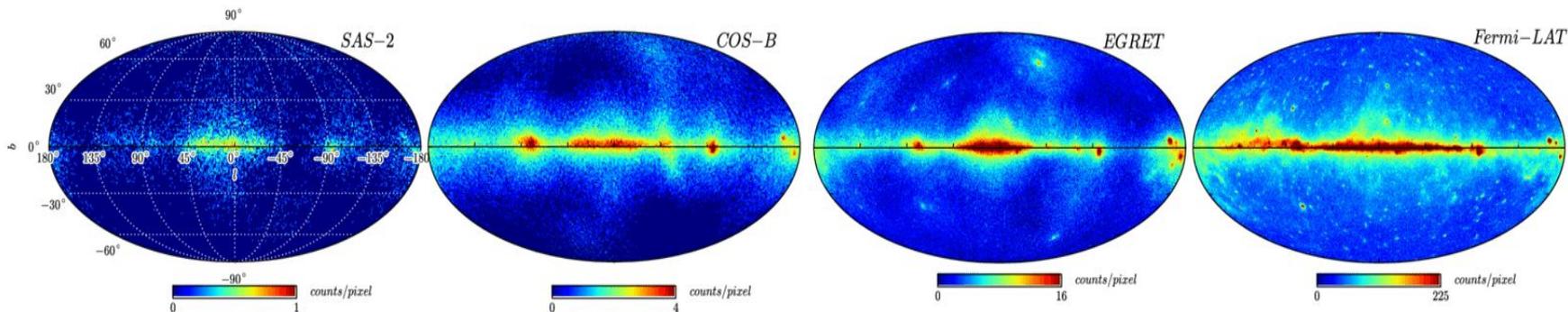
GRAVITATIONAL LENSING

Challenge 1) High-dimensional data

Thousands of petabytes of upcoming data
Preprocessed data can have millions of dimensions



Challenge 2) High-dimensional models

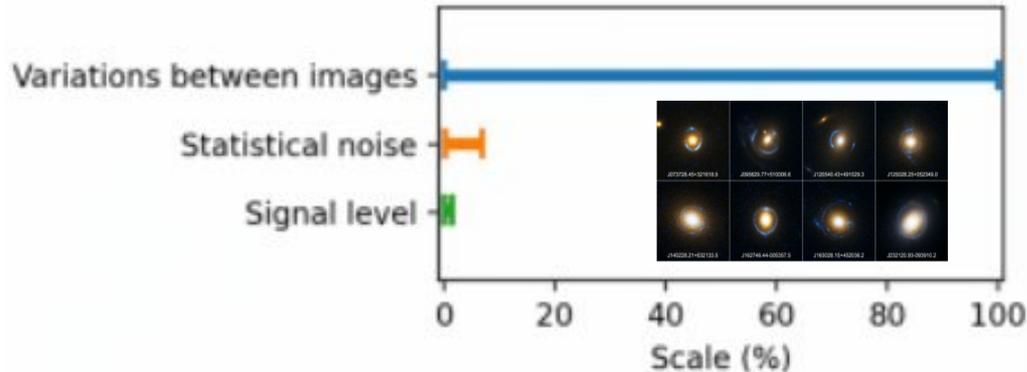
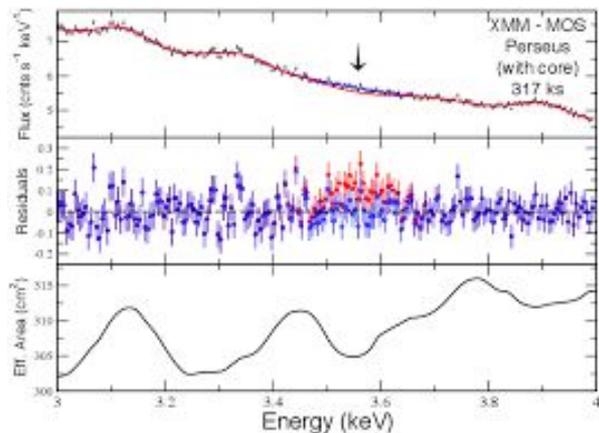


More data → More details → More accurate models
 ↘ **More model parameters**
 ↘ **Slower simulations**

Endless statistical analysis challenges

- Hierarchical models (source populations)
- Trans-dimensional models (number of sources)
- Label switching problem (instance detection)
- Parameter degeneracies (distances unclear)
- Non-parametric components (gas maps)
- Inference of fields (diffusion zone)
- Millions of parameters & minutes - hours to evaluate...

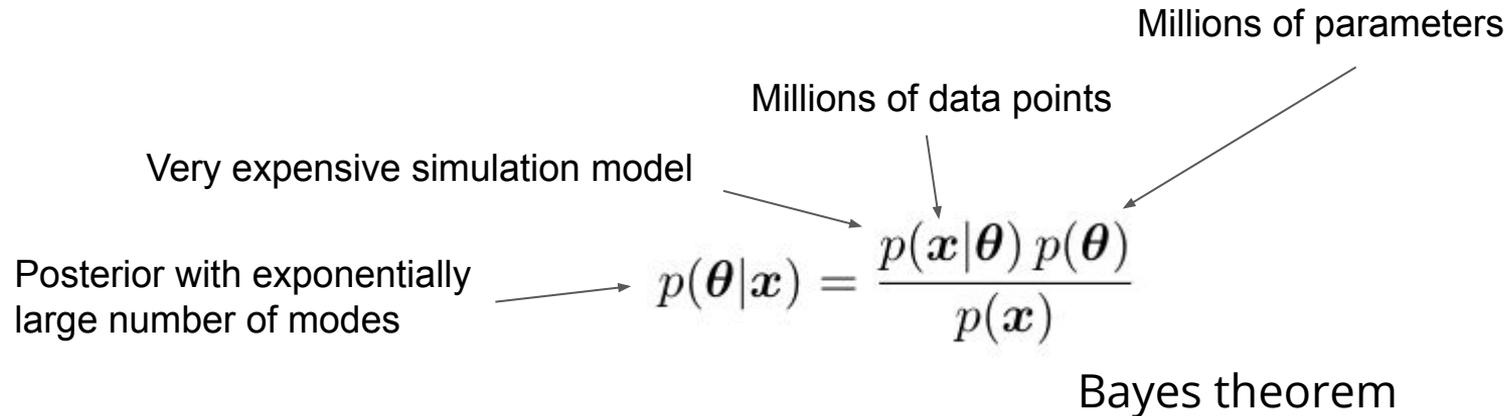
Challenge 3) Signals \ll Background



Any mismodeling of backgrounds, or uncharacterized model uncertainty, or unjustified simplification, does backfire.

New physics searches = Avoiding to shoot yourself in the foot

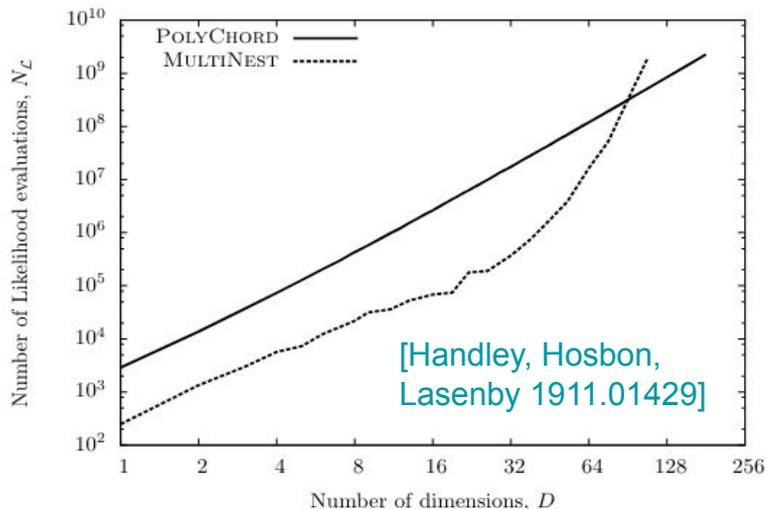
The inverse problem = the impossible problem?



Solving the inverse problem

de facto standard: Markov Chain Monte Carlo

$$p(\boldsymbol{\theta}|\mathbf{x}) = \frac{\overset{\text{Likelihood}}{p(\mathbf{x}|\boldsymbol{\theta})} \overset{\text{Prior}}{p(\boldsymbol{\theta})}}{\underset{\text{Posterior}}{p(\mathbf{x})} \underset{\text{Evidence}}{p(\mathbf{x})}}$$

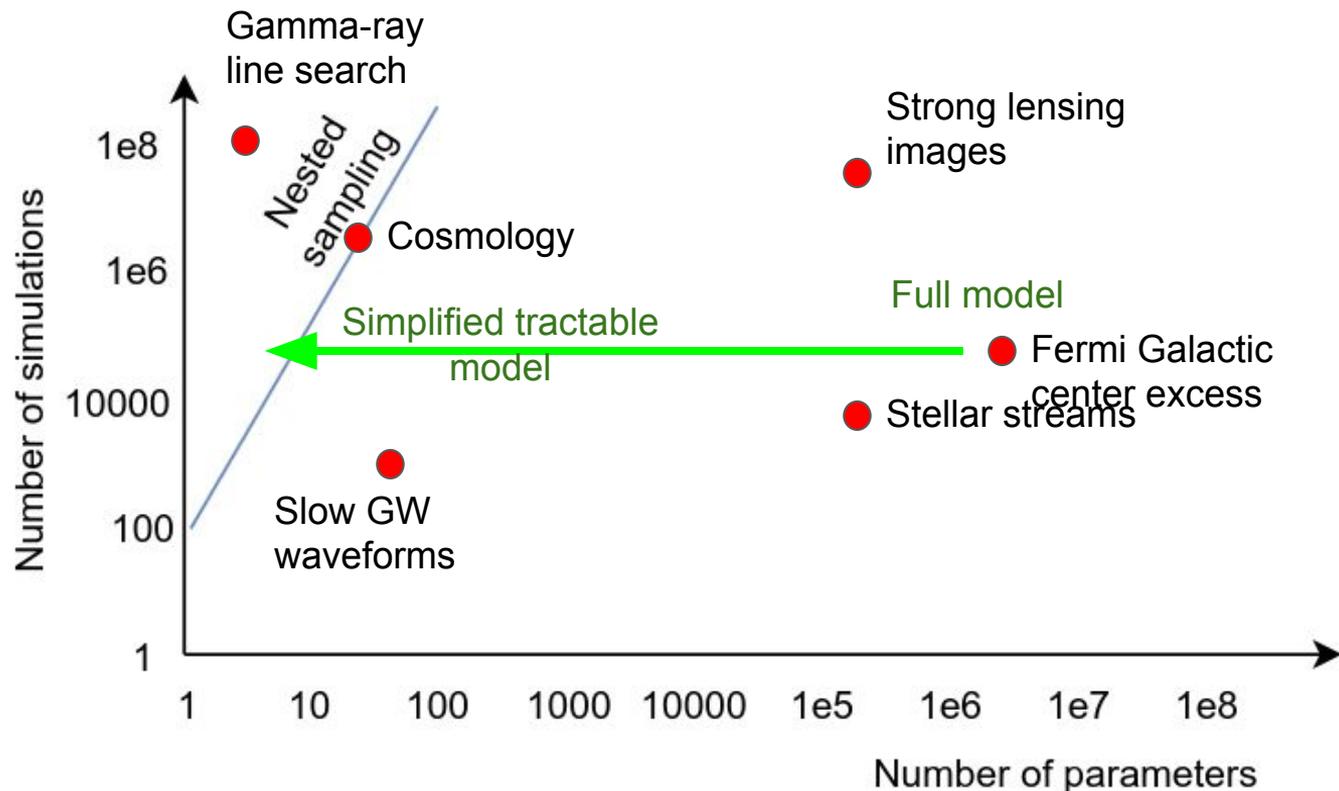


- Step 1: Samples from joined posterior
$$\boldsymbol{\theta} \sim p(\boldsymbol{\theta}|\mathbf{x}) \quad \boldsymbol{\theta} \in \mathbb{R}^D$$
$$D: \text{Number of parameters}$$
- Step 2: Marginalization to parameters of interest

$$\boldsymbol{\theta} \equiv (\theta_1, \theta_2, \dots, \theta_D)^T \rightarrow (\theta_i, \theta_j)^T \in \mathbb{R}^2$$

Typical likelihood-based inference algorithms (MH, HMC, VI, ...) require a small enough number of parameters.

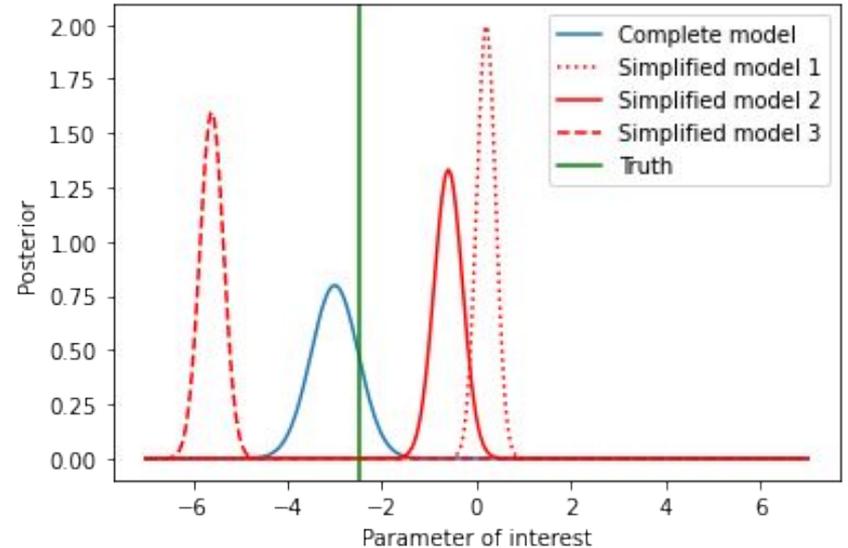
Likelihood-based inference enforces low-dim models



The price of model simplification

Exaggerated (?) illustration of the potential dangers of model simplification

- Biases
- Overly optimistic uncertainties

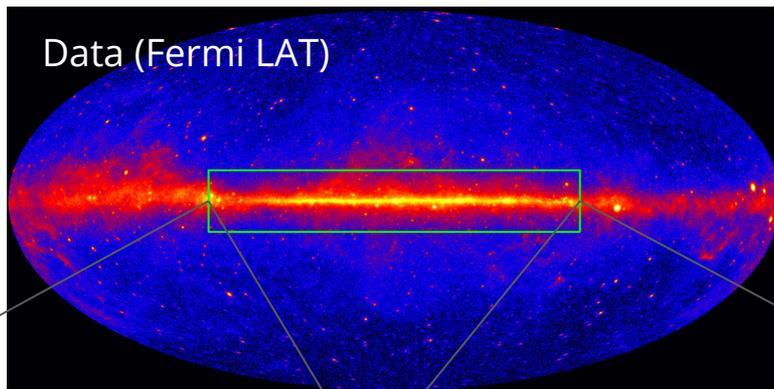


The benefits of high-dimensional models

Almost all existing analysis of Fermi LAT data have these kind of residuals.

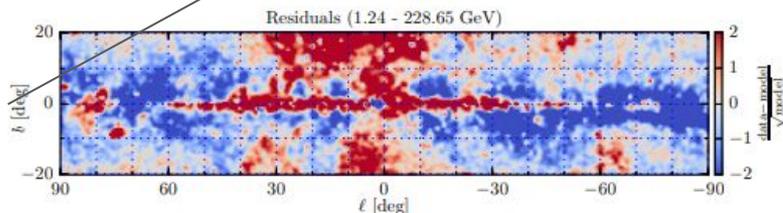
There is no shortage in anomalies in astrophysical data...

Consequences: Large modeling errors because of **simplistic low-dim models**



We pulled this off with **gradient-based optimization**. **Very hard to use** in practice, only a handful of examples in the literature.

Residuals
Low-dim model (10 dims)



Residuals
High-dim model (10,000 dims)



Neural simulation-based inference (SBI)

Very active young research field

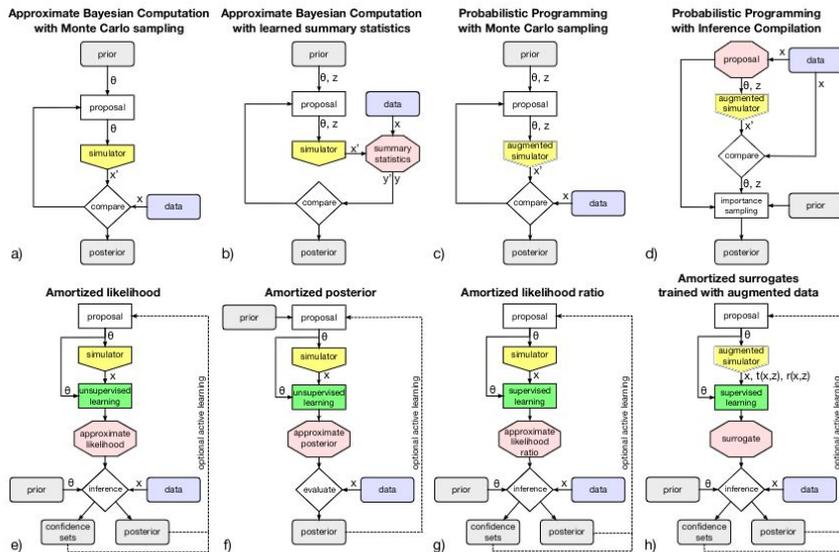


Fig. 3. Overview of different approaches to simulation-based inference.

General goal: obtain neural network approximator for one of the following:

- Posterior* $p(\theta|x)$
- Likelihood* $p(x|\theta)$

- Ratios of posteriors and priors = ratios of likelihood and evidence

$$r(x, \theta) = \frac{p(x|\theta)}{p(x)} = \frac{p(x, \theta)}{p(x)p(\theta)} = \frac{p(\theta|x)}{p(\theta)}$$

- Various variations of the above quantities...

[Cranmer, Brehmer, Louppe 1911.01429]

* require normalization of densities

High-dimensional inference can be simple

Examples
 $x, z \sim p(x|z)p(z)$

Image, x



“Simulated images”

Parameter vector, z

1, 3, 2, **1**, 5, 4, 3, 1, 6, 7, 9, ...
 6, 2, 5, **8**, 6, 8, 4, 3, 2 1, 3, 4, ...
 2, 3, 4, **3**, 1, 7, 8, 9, 5, 3, 2, ...
 4, 2, 1, **4**, 6, 8, 6, 4, 3, 2, 4, ...
 1, 3, 2, **9**, 5, 4, 3, 1, 6, 7, 9, ...
 6, 2, 5, **8**, 6, 8, 4, 3, 2 1, 3, 4, ...
 2, 3, 4, **1**, 1, 7, 8, 9, 5, 3, 2, ...
 4, 2, 1, **2**, 6, 8, 6, 4, 3, 2, 4, ...
 1, 3, 2, **4**, 5, 4, 3, 1, 6, 7, 9, ...
 6, 2, 5, **4**, 6, 8, 4, 3, 2 1, 3, 4, ...
 ?, ?, ?, **8**, ?, ?, ?, ?, ?, ?, ?, ?, ...

Red:
Parameter of interest

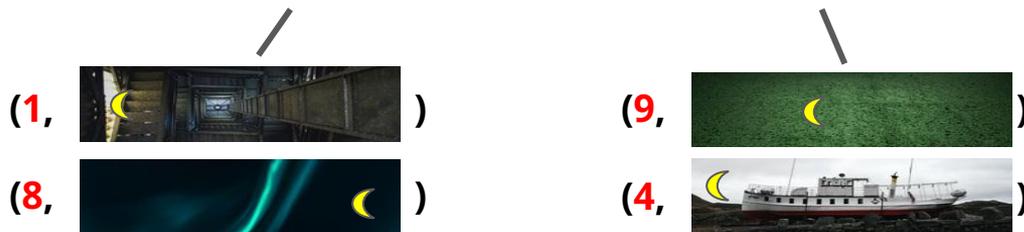
Black:
Nuisance parameters
 (parametrizing *all* possible background images)

Observed data



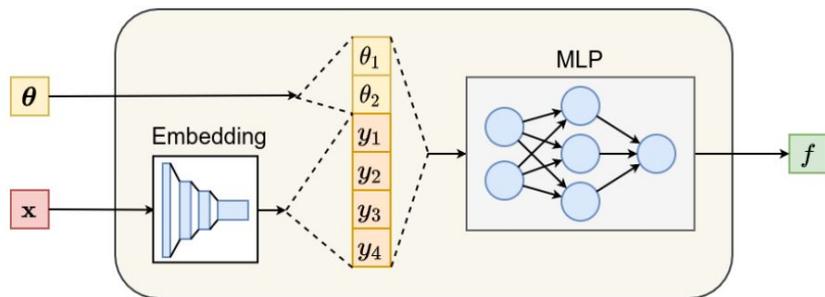
Neural ratio estimation (NRE) in a nutshell

Strategy: Learning to distinguish between **matching (parameter, data) pairs** and **random pairs**.



Loss function: Binary cross entropy

$$\ell[f_\phi]_{\text{NRE}} = - \int d\mathbf{x} d\boldsymbol{\theta} [p(\mathbf{x}, \boldsymbol{\theta}) \ln \sigma(f_\phi(\mathbf{x}, \boldsymbol{\theta})) + p(\mathbf{x})p(\boldsymbol{\theta}) \ln (1 - \sigma(f_\phi(\mathbf{x}, \boldsymbol{\theta})))]$$

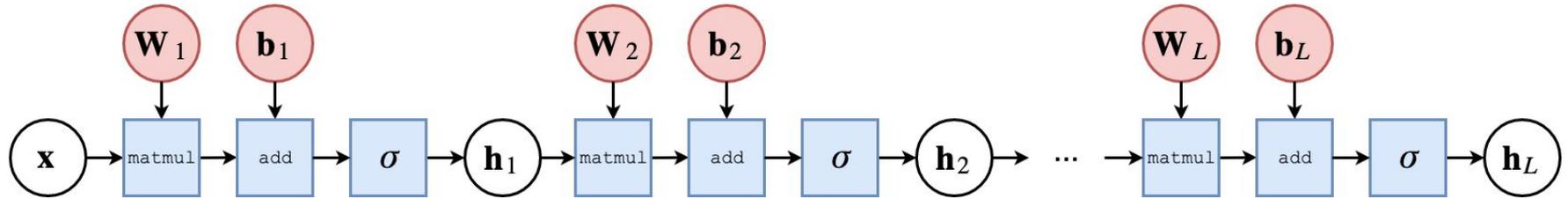


Minimizing network approximates posteriors

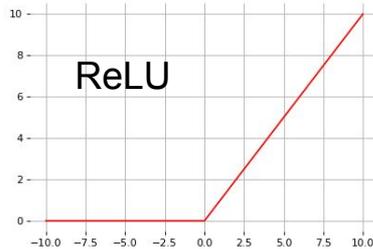
$$f_\phi(\boldsymbol{\theta}, \mathbf{x}) \approx \ln \frac{p(\mathbf{x}, \boldsymbol{\theta})}{p(\mathbf{x})p(\boldsymbol{\theta})} = \ln \frac{p(\boldsymbol{\theta}|\mathbf{x})}{p(\boldsymbol{\theta})}$$

What is a Multi-Layer Perceptron?

Multi-layer perceptron = MLP = dense network



Activation function



Training would optimize transformation matrices W and biases b .

Three conjectures for scalable SBI

1) Marginal posterior rather than joint posteriors

- A “universal” approach must scale to millions of parameters, and outrageously complex posteriors (transdimensional models, label switching, strong correlations, ...)

$$p(z_1, z_2, \dots, z_{1000000} | \mathbf{x})$$

Joined: In general intractable
(any approach)

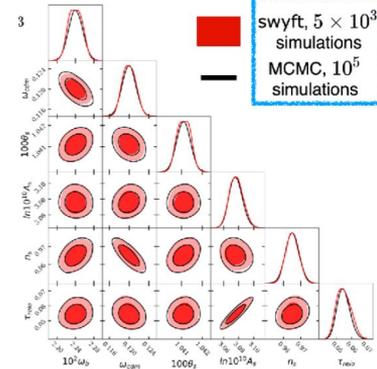
$$p(z_1 | \mathbf{x}), p(z_2, z_3 | \mathbf{x}), p(\max(\mathbf{z}) | \mathbf{x}), \dots$$

Marginals: Often tractable
(NRE, forward-KL based approaches, ...)

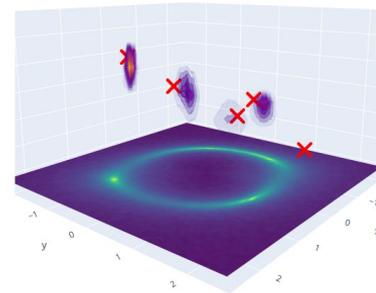
- Scientifically, we are usually only interested in marginal posteriors anyway
 - Parameter regression: 1-dim marginals
 - Parameter correlations: 2-dim marginals
 - Bayesian model comparison: ratios of marginals
 - Object identification: density functions
 - ...

[for discussions see e.g. Alsing+ 1903.01473, Jeffrey+ 2011.05991, Miller+ 2011.13951]

- Caveats: Goodness-of-fit tests, posterior predictive distribution, requires upfront intuition about what matters



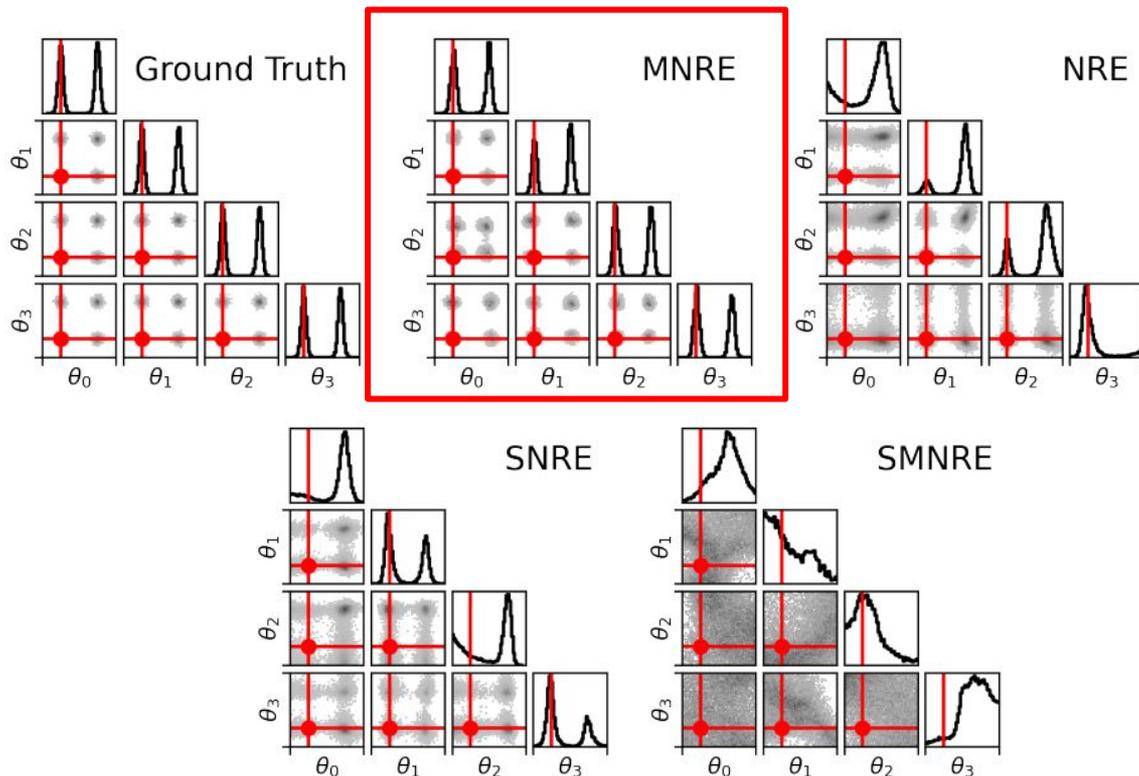
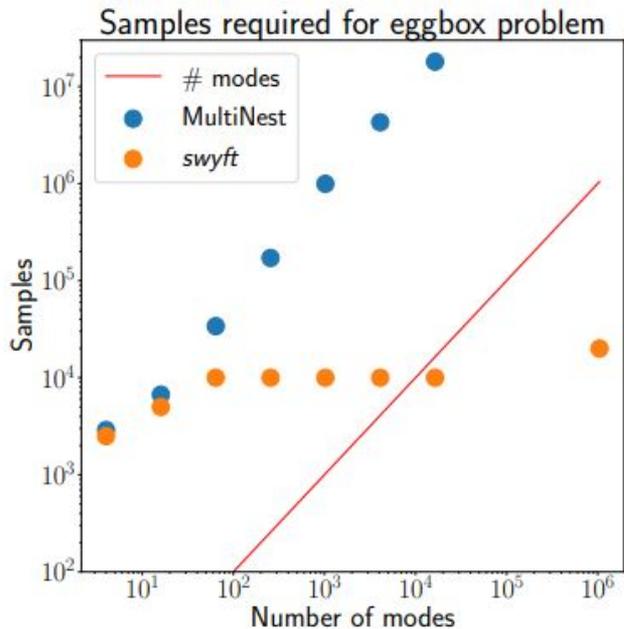
1-dim and 2-dim marginals for corner plots



Density functions for object detection

1) Marginal posterior rather than joint posteriors

$$\# \text{modes} = 2^{\# \text{nparams}}$$



Estimating marginals breaks naive scaling laws.

2) Truncated priors as sequential proposals

- Sequential techniques are based on targeted training data

$$\mathbf{x}, \mathbf{z} \sim p(\mathbf{x}|\mathbf{z})\tilde{p}(\mathbf{z})$$

[Durkan+ 2002.03712 for a discussion]

$$\tilde{p}(\mathbf{z}) \approx p(\mathbf{z}|\mathbf{x}_o)$$

- This is fine if the goal is to locally train, e.g., the likelihood (which is prior independent)

$$p(\mathbf{x}|\mathbf{z})$$

[Alsing+ 1903.00007 as example (pydelfi)]

- But:** *Marginal* likelihoods/posteriors will be affected by the proposal distribution

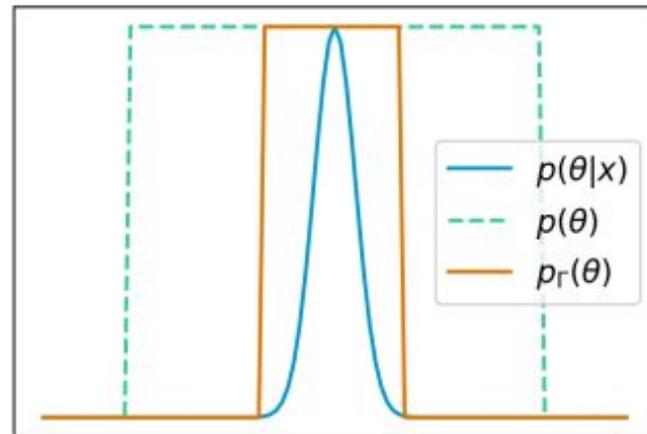
$$p(\mathbf{x}|z_1) = \int dz_2 \dots dz_N p(\mathbf{x}|\mathbf{z})\tilde{p}(z_2, \dots, z_N)$$

[see e.g. Alsing+ 1903.01473 for a possible summary-statistics related solution]

- To alleviate this we proposed to use a *truncation scheme*

$$\tilde{p}(\mathbf{z}) = \mathbb{I}(\mathbf{z} \in \Gamma)p(\mathbf{z})$$

[Miller+ 2011.13951, 2107.01214 - swyft & TMNRE]



3) Likelihood-to-evidence ratios rather than densities

- Ratio estimation \rightarrow Binary classification = Battle-proven simplicity

$$f_\phi(\boldsymbol{\theta}, \mathbf{x}) \approx \ln \frac{p(\mathbf{x}, \boldsymbol{\theta})}{p(\mathbf{x})p(\boldsymbol{\theta})} = \ln \frac{p(\boldsymbol{\theta}|\mathbf{x})}{p(\boldsymbol{\theta})}$$

[Hermans+ 1903.04057]

[see Cranmer+ 1911.01429 for discussion of many alternatives]

- Usually remains conservative (works well in a truncation scheme)

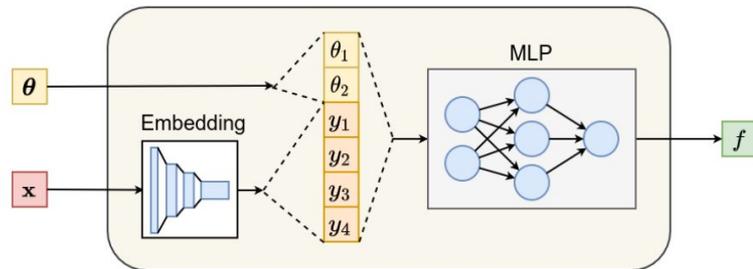
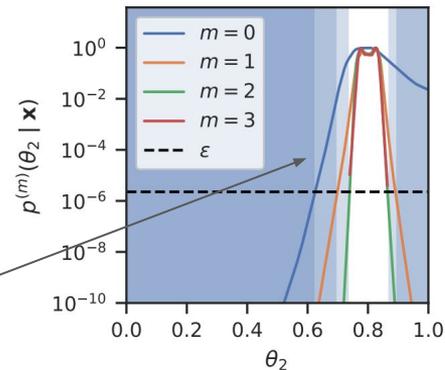
[but see Hermans+ 2110.06581]

- Ratio estimation automatically generates information maximizing data compression

$$\ell[\hat{\rho}_\phi] = -2\mathbb{E}_{p(\mathbf{x})} [\text{JSD}(p(\boldsymbol{\theta}|\mathbf{s}(\mathbf{x}))||p(\boldsymbol{\theta}))] + \text{const}$$

[see Alsing+ 1903.00007 for related discussions in context of likelihood estimation]

- When focusing on low-dim marginals, sampling is simple (no MCMC or flow-based models required).

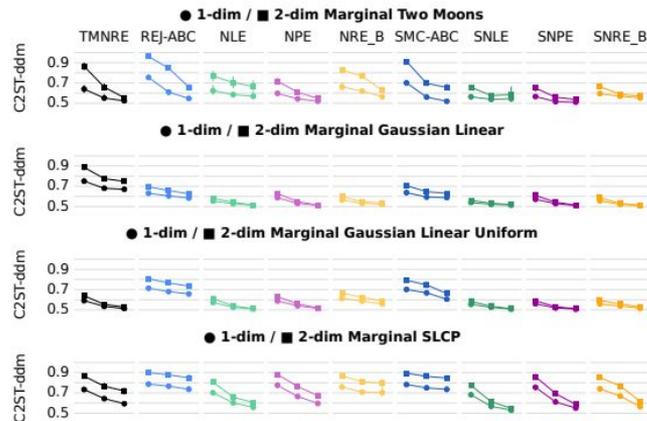


1+2+3 = Truncated Marginal Neural Ratio Estimation

NeurIPS | 2021

Thirty-fifth Conference on Neural Information
Processing Systems

Competitive performance on
standard tasks, but more
scalable.

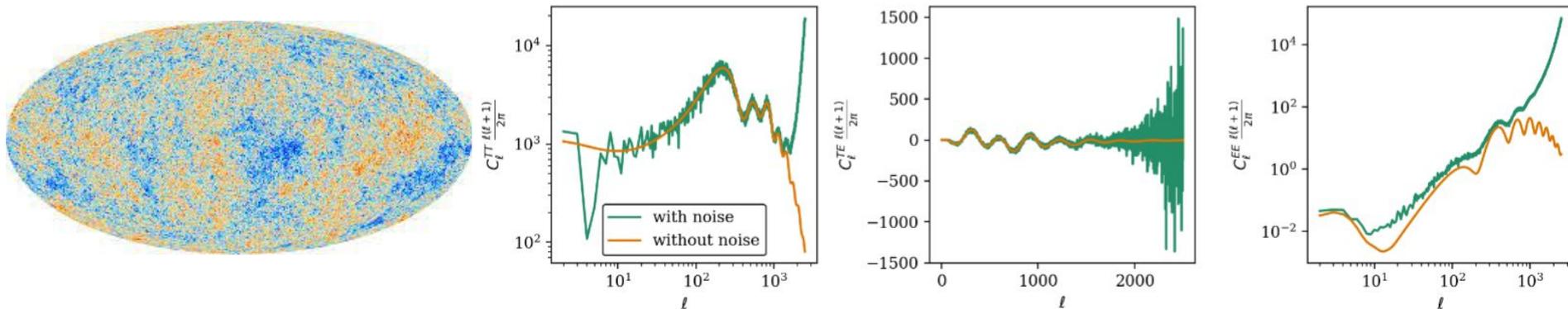


Combination of various properties of existing algorithms

Property / Method	Likelihood-based	ABC	NRE	NPE	SNRE	SNPE	TMNRE
Targeted inference	✓	•	✗	✗	✓	✓	✓
Simulator efficient <i>direct</i> marginals	✗	✓	•	•	✗	✗	✓
(Local) amortization	✗	✗	✓	✓	✗	✗	✓

Example 1: Cosmic microwave background

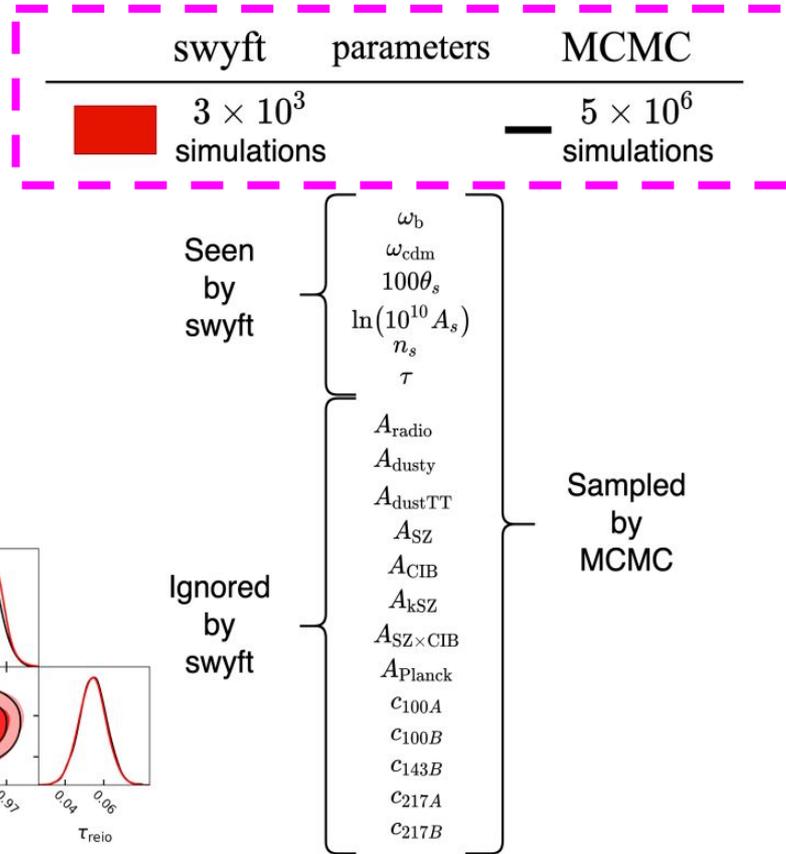
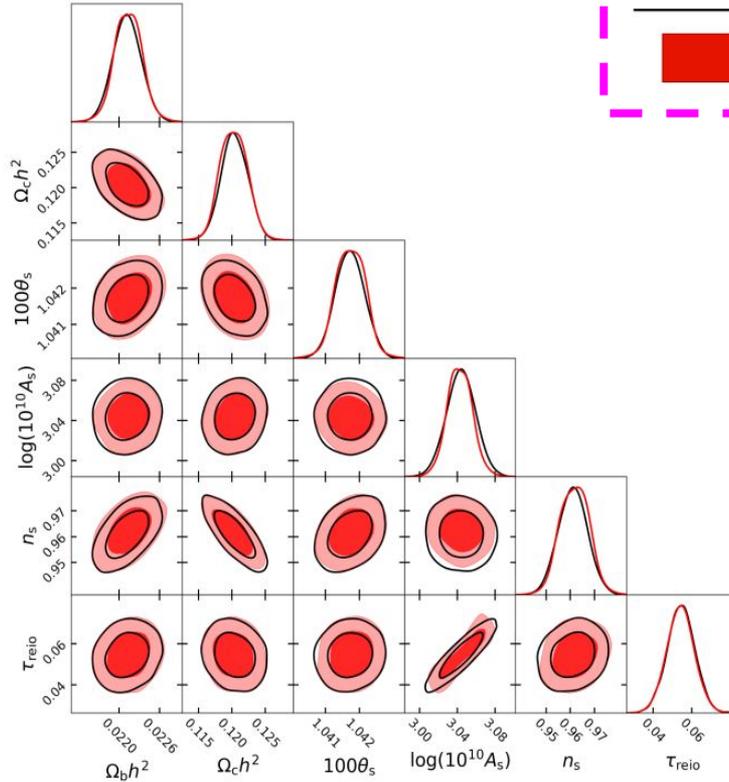
Planck cosmology



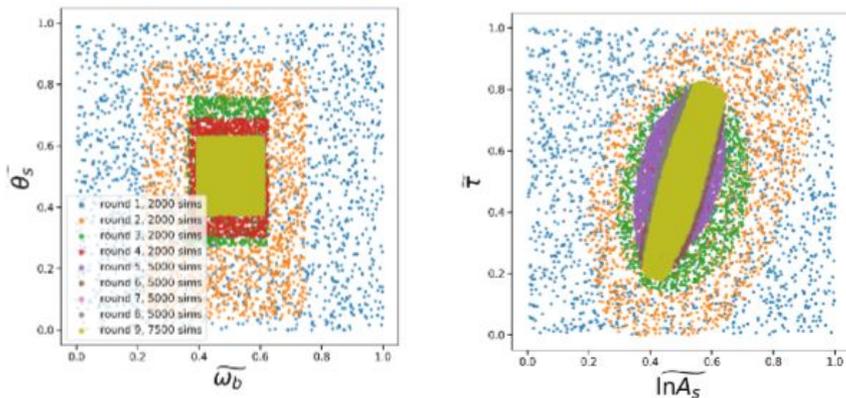
Noise = instrument contribution + cosmic variance

- TT, TE, EE angular power spectrum of CMB with Planck-like noise (Di Valentino+ 2016)
- 6 cosmology parameter to infer, using tight priors (+- 5 sigma Fisher estimate)
- HiLLiPoP likelihood: Planck likelihood, 13 varying nuisance parameters [Couchot et al. '16]
- Comparison with MCMC is feasible and straightforward
- We use a linear embedding network to go from 7500 \rightarrow 10 features

Cosmology with ~ 1000 times less simulations



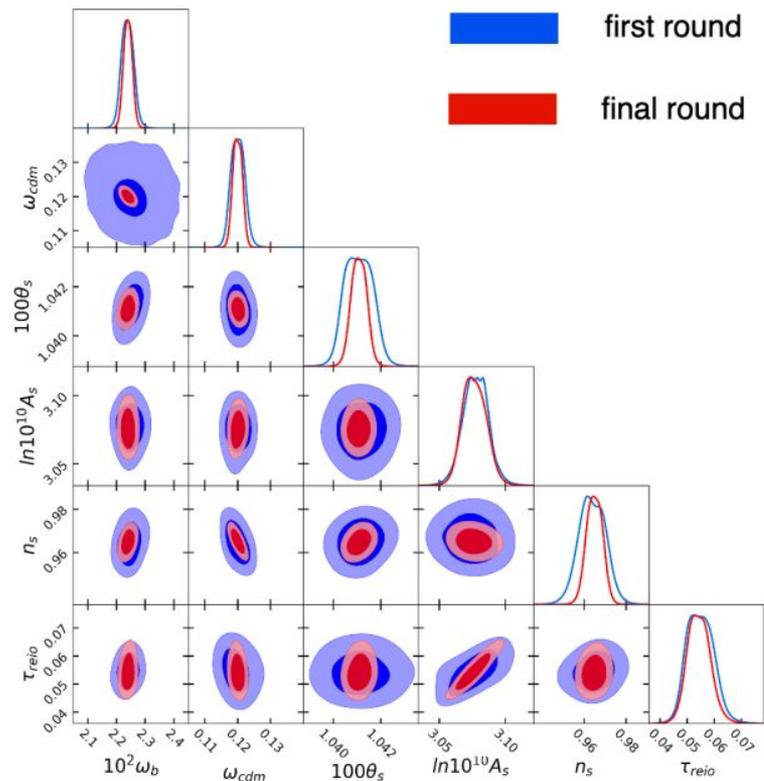
Simulation efficiency through truncation



- Demonstration of prior that is “too big” by a factor of 5 for the cosmological parameters
- Truncation effectively identifies region with 20000 extra sims.

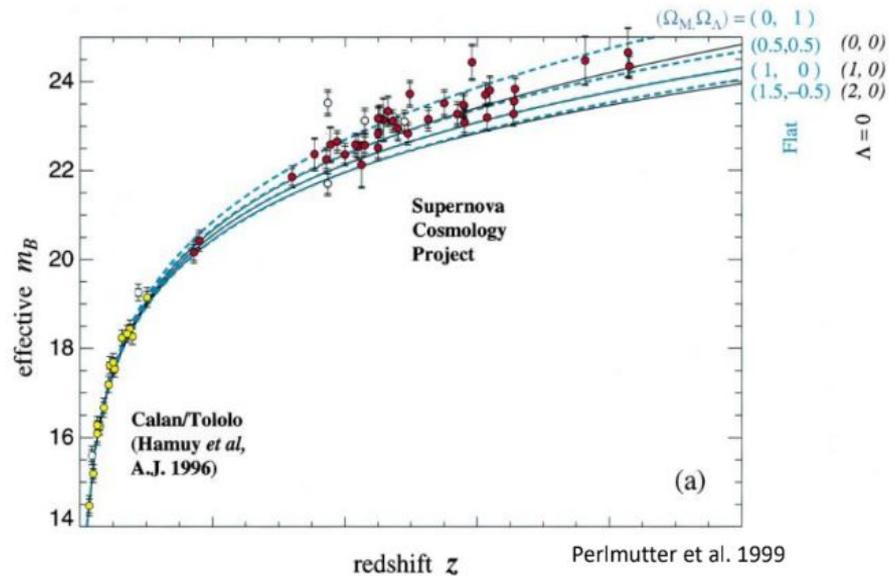
Structure of ratio estimator

- Input: Vector (7500)
- Embedding: Linear (7500 \rightarrow 10)
- Marginals: MLP (19 1-dim, 15 2-dim)

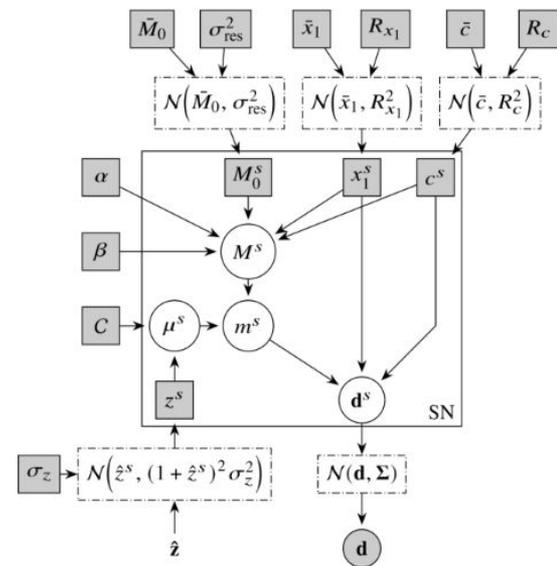


Example 2: Supernova cosmology

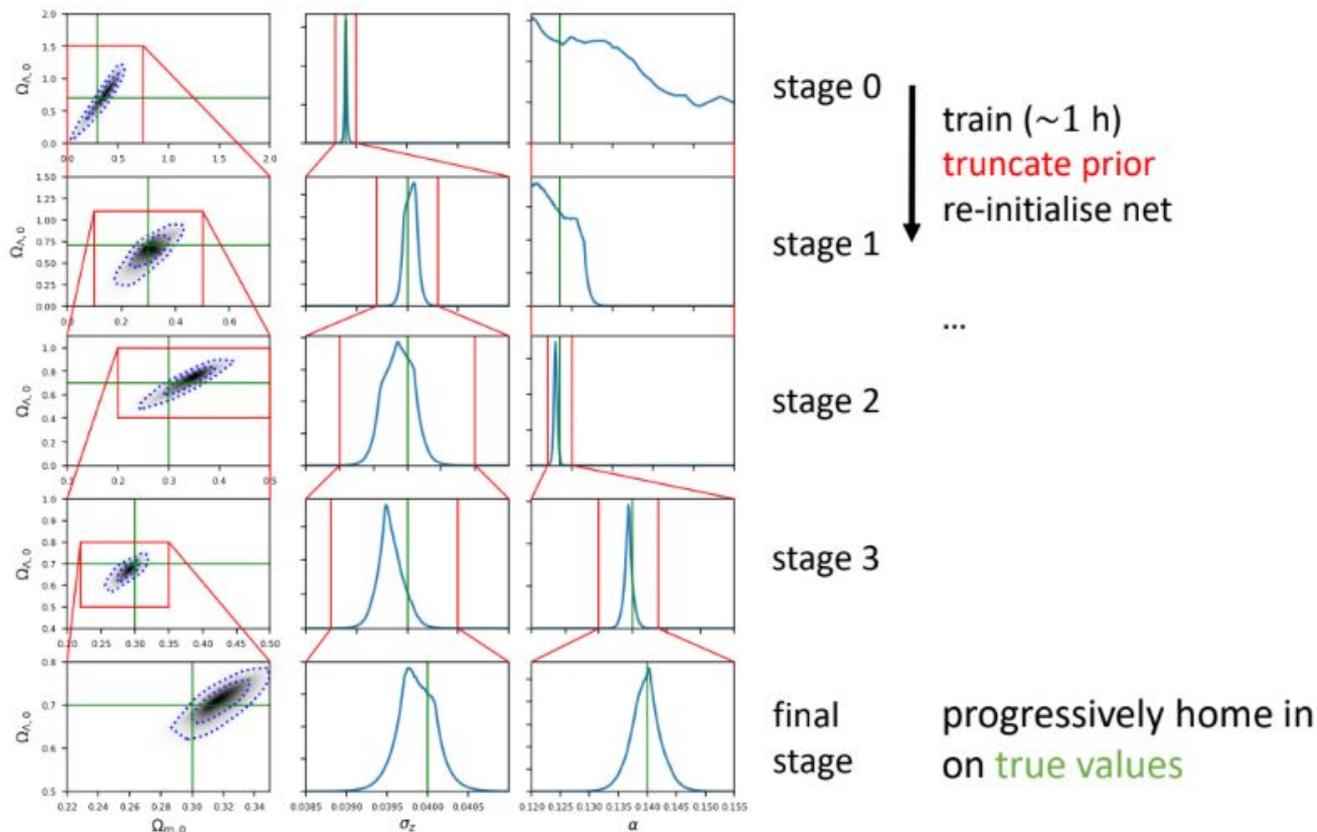
Supernova cosmology



$$m = M + \mu(z, \mathcal{C}) + \text{"noise"}$$



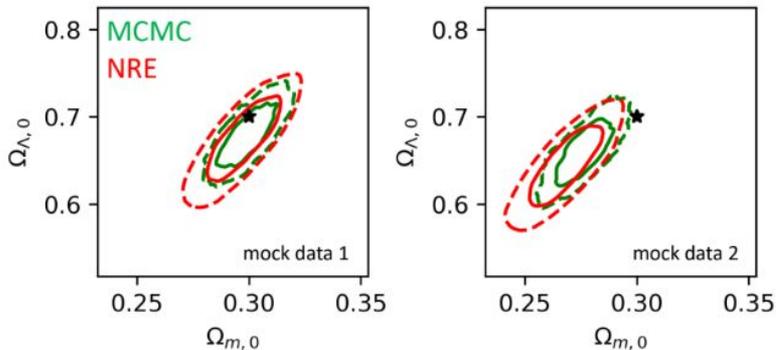
Simulation efficiency through truncation



Ongoing work with Kosio Karchev and Roberto Trotta

(Marginal) measurements for 100000 parameters

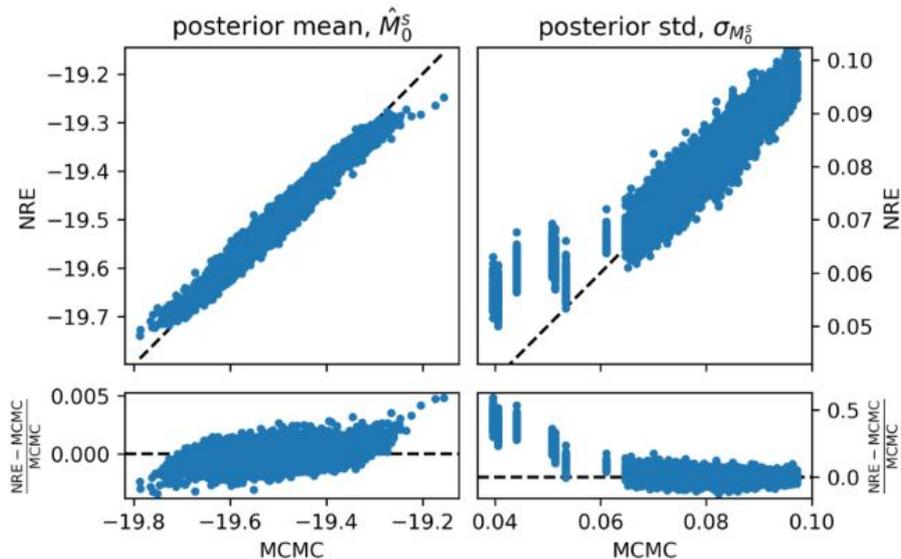
100 000 supernovae



- “MCMC” results were obtained using pre-marginalized likelihoods (only possible under assumptions of Gaussianity and SN independence).
- Instead, NRE marginalizes automatically and *assumption-free*.

Ongoing work with Kosio Karchev and Roberto Trotta

MALFOI: marginal likelihood-free object-by-object inference

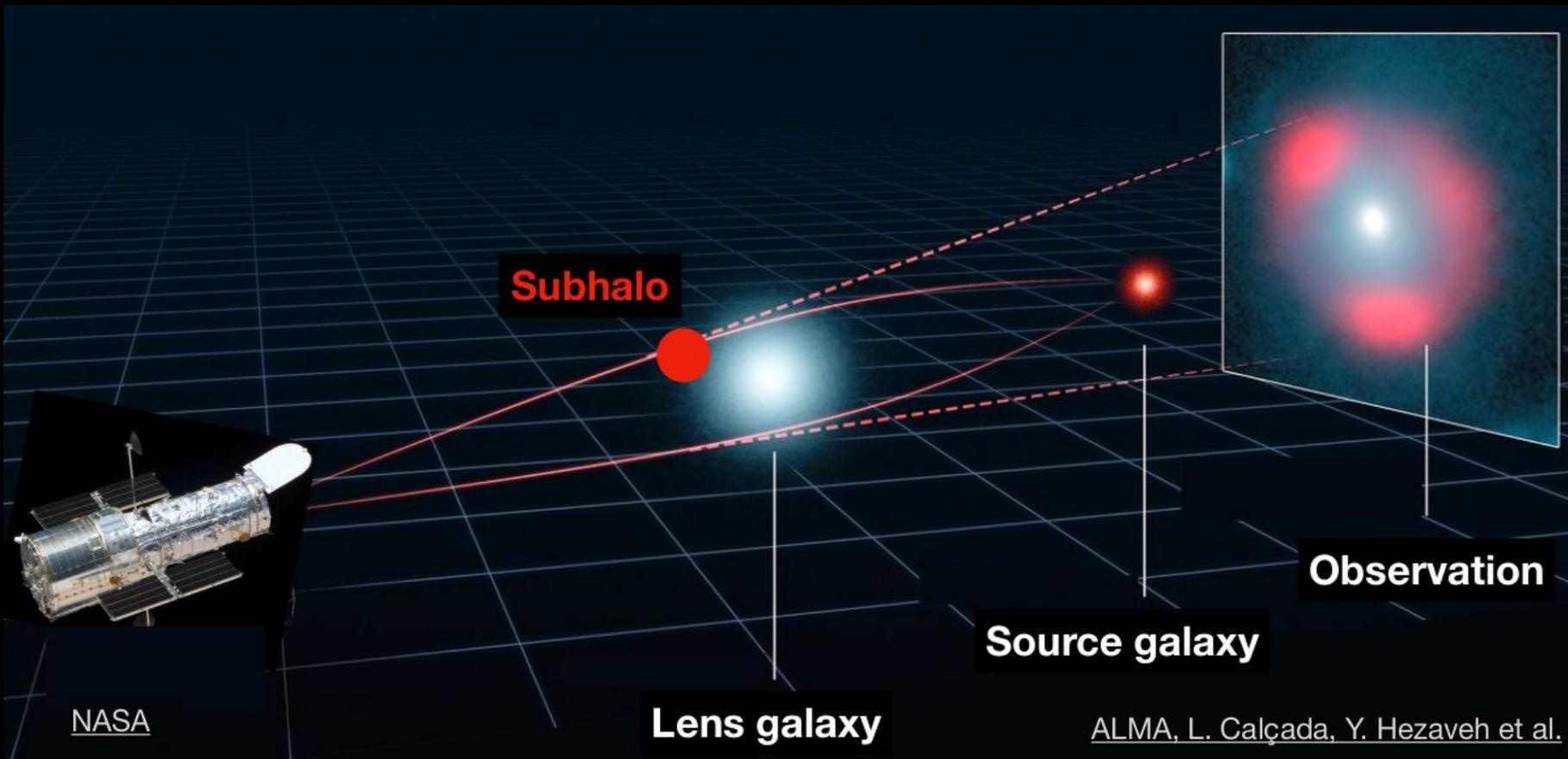


Structure of ratio estimator

- Input: 100.000 Spectra (100000, 3)
- Embedding: Linear (300000 \rightarrow 256)
- Marginals: MLP (100009 1-dim, 1 2-dim)

Example 3: Strong lensing

Strong galaxy-galaxy lensing



Warm vs cold dark matter

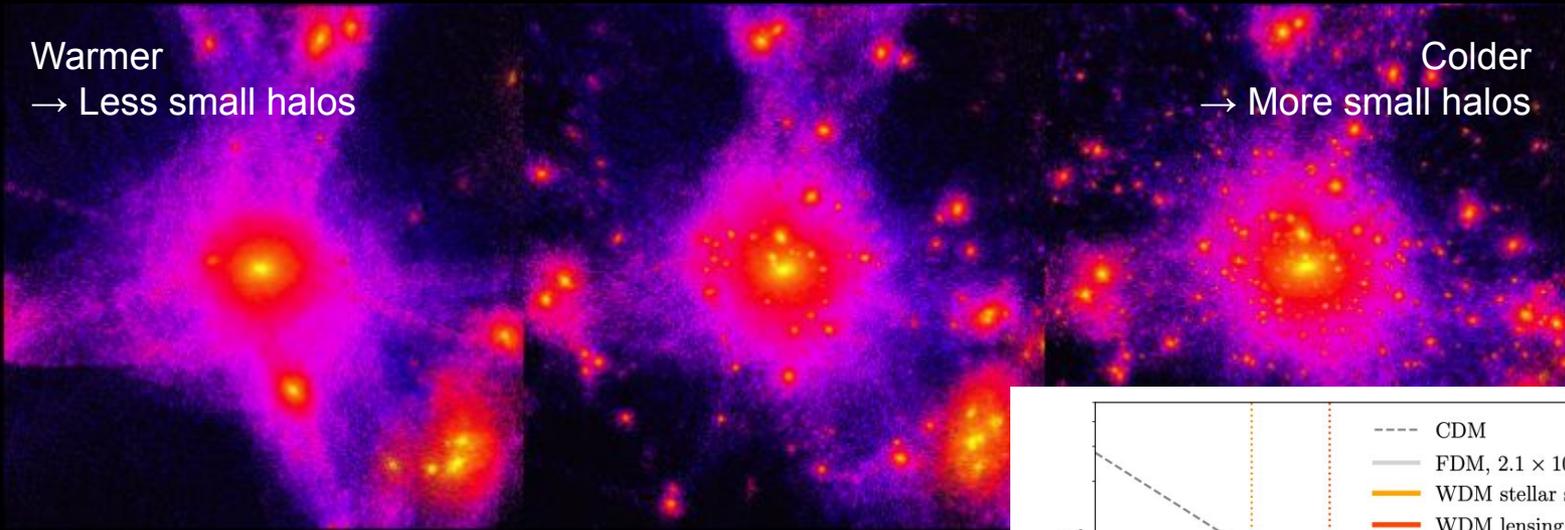
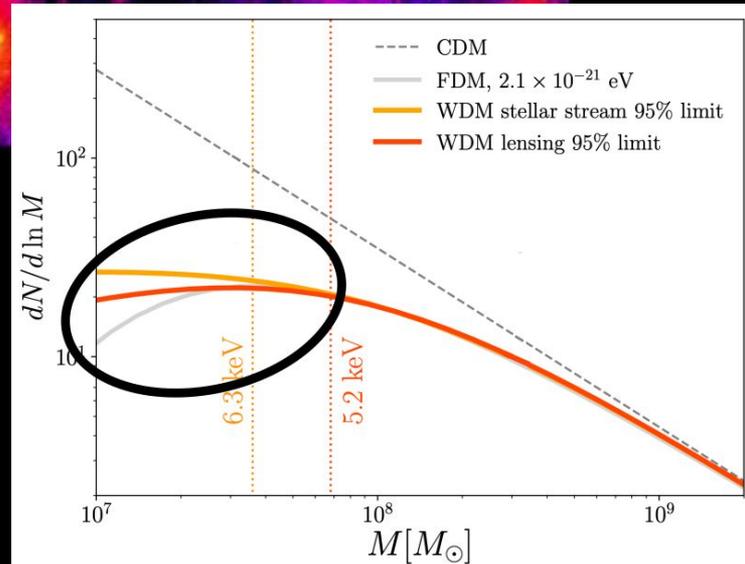


Image credit: ITC @ University of Zurich

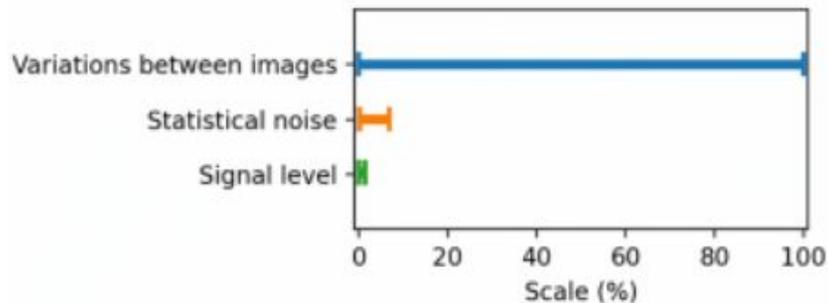
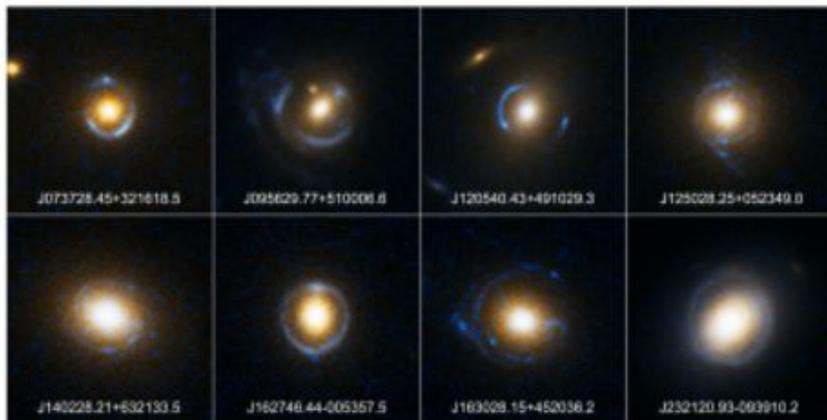


Strong lensing animated

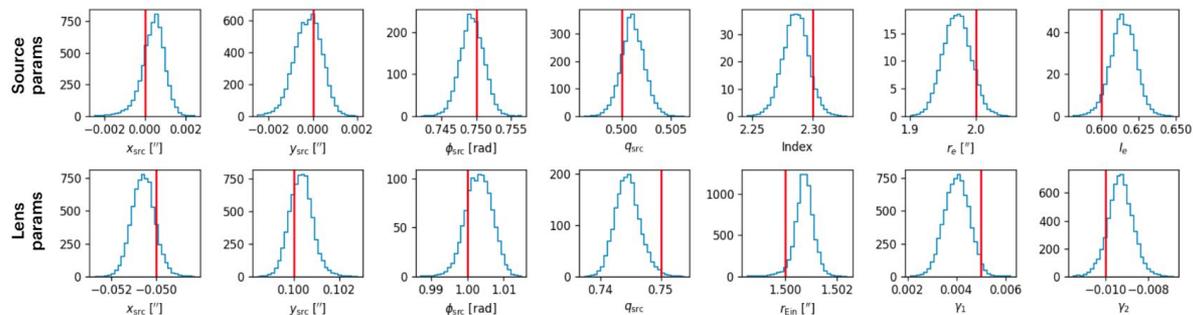
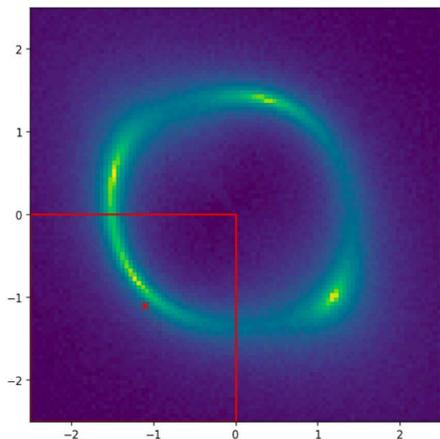
<https://adam-coogan.github.io/lensing-multisub/>

Inference challenges

- **Signal is small** compared to noise and variations between images
- **Marginalization** over numerous source, lens and halo parameters
- Joint posterior has $\sim N_{\text{sub}}!$ **modes**; likelihood can be **intractable**



A) Single subhalo, simple source model

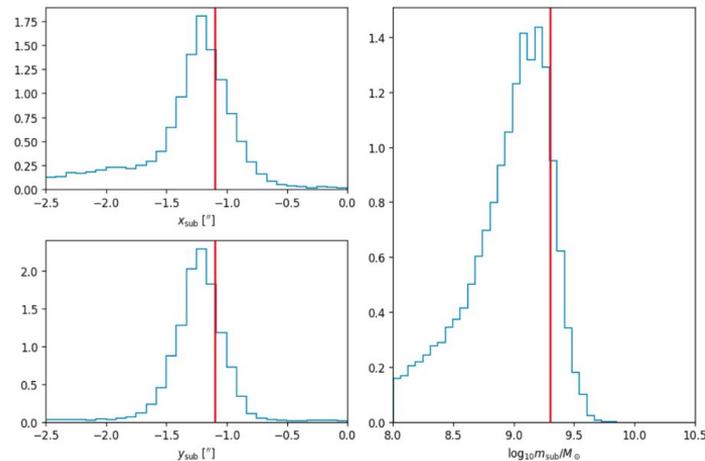


Structure of ratio estimator

- Input: Images (typically 200x200)
- Embedding: CNN
- Marginals: MLP (17 1-dim)

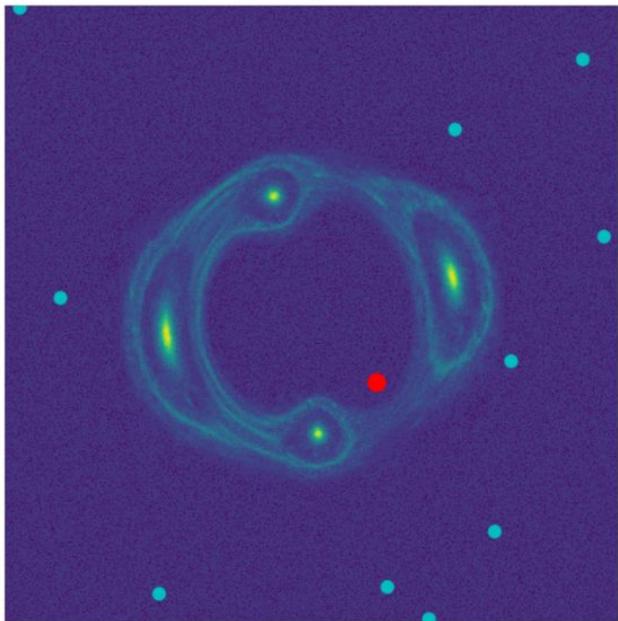
Ongoing work led by Adam Coogan

Slide credit: Noemi Anau Montel



B) Multiple subhalos, complex source model

Training data

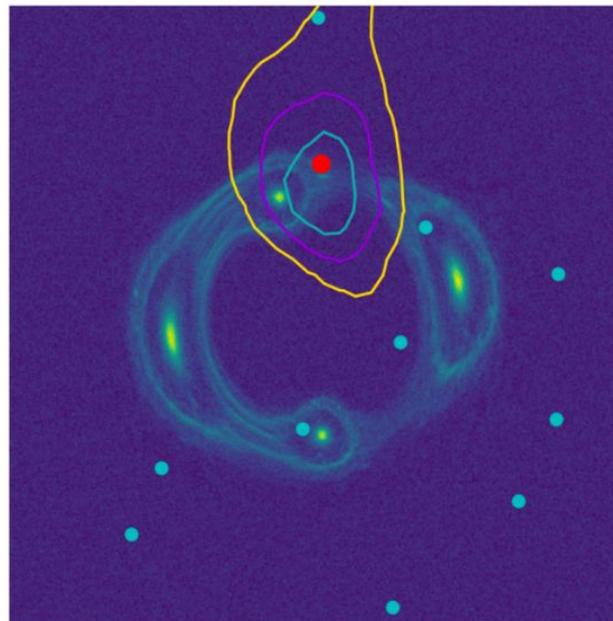


• = $5 \times 10^9 M_{\odot}$, • = $10^8 - 10^9 M_{\odot}$

Marginalized over source, lens and halo population

Ongoing work led by Adam Coogan

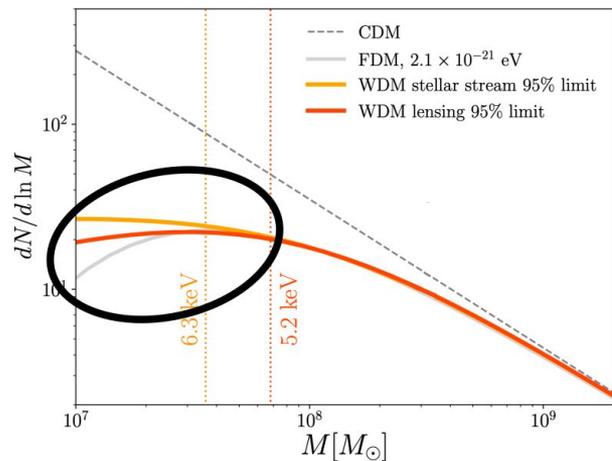
Inference



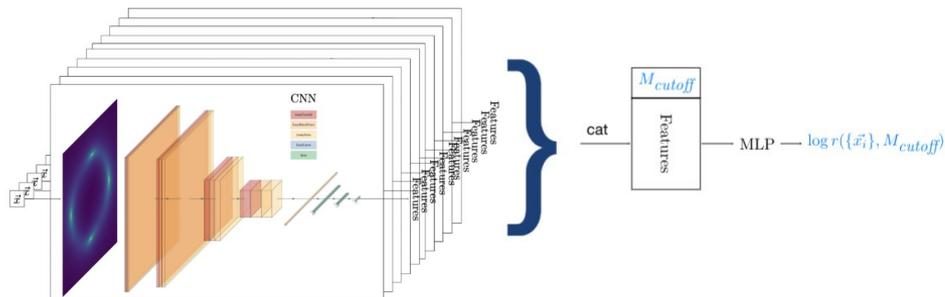
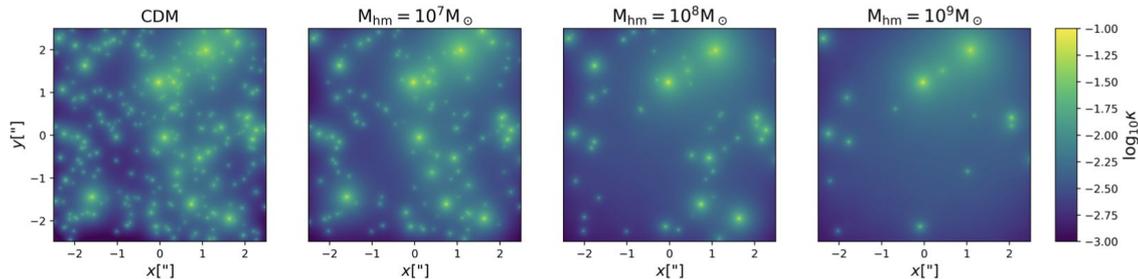
Structure of ratio estimator

- Input: Images (typically 200x200)
- Embedding: CNN
- Marginals: MLP (2-dim)

C) Subhalo mass function cutoff - Combined analysis



Combining observations to reduce subhalo shot noise



Structure of ratio estimator

- Input: 10 Images ($10 \times 100 \times 100$)
- Embedding: Stack of CNNs
- Marginals: MLP (1-dim)

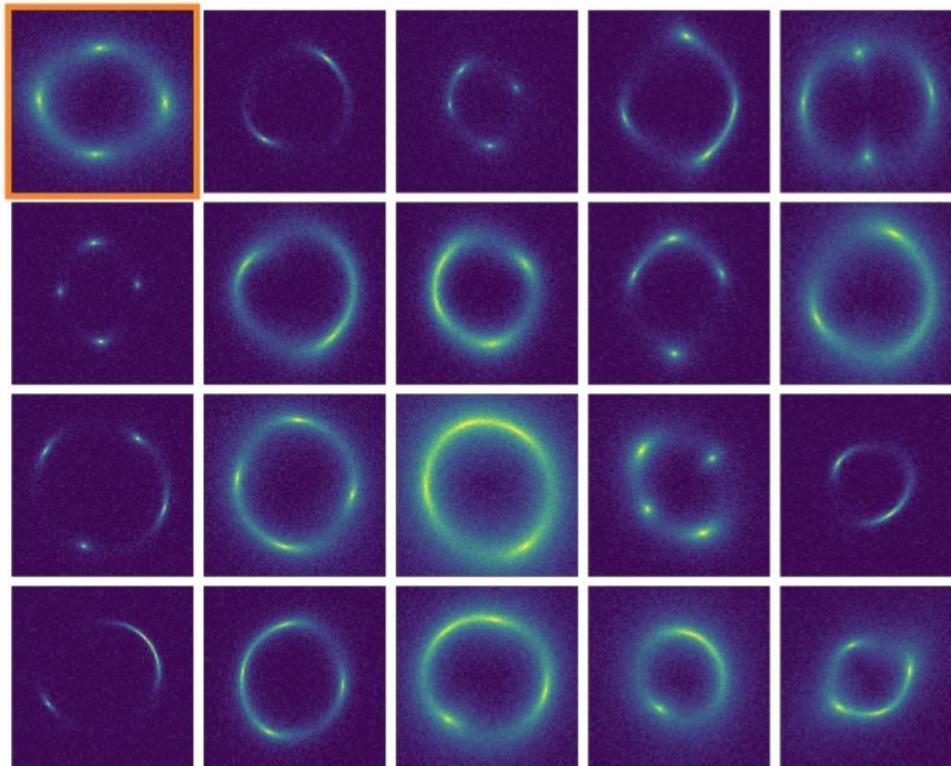
Anau Montel+ 2205.09126

Related work: He+ 2010.13221 (similar in spirit, using ABC), Wagner-Carena+ 2203.00690 (constraining subhalo mass function normalization) 38

C) Subhalo mass function cutoff - Combined analysis

Subset of 20 target images in our analysis.

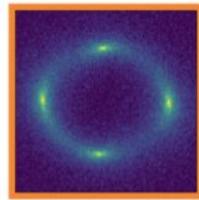
Let's focus on this one,
and call it "A"



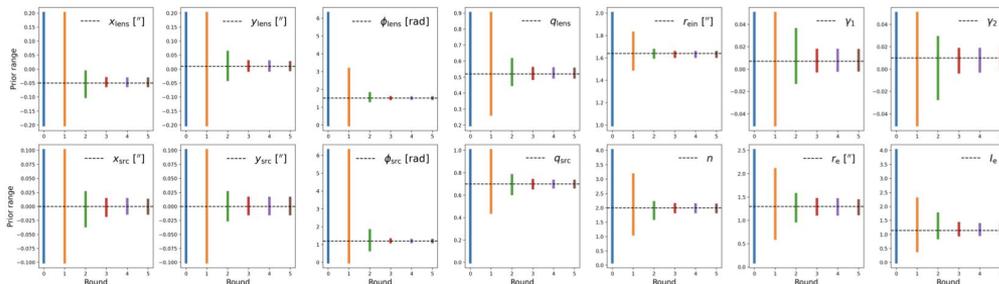
Anau Montel+ 2205.09126

C) Subhalo mass function cutoff - Prior truncation

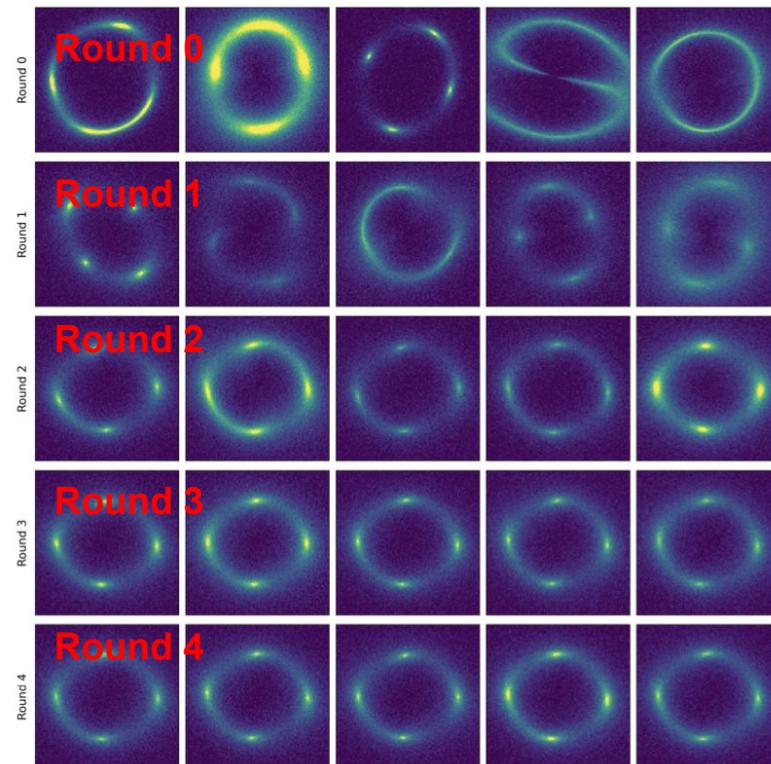
Target image "A"



Constrained prior ranges (round 0 - round 4) for all 14 main lens & source parameters



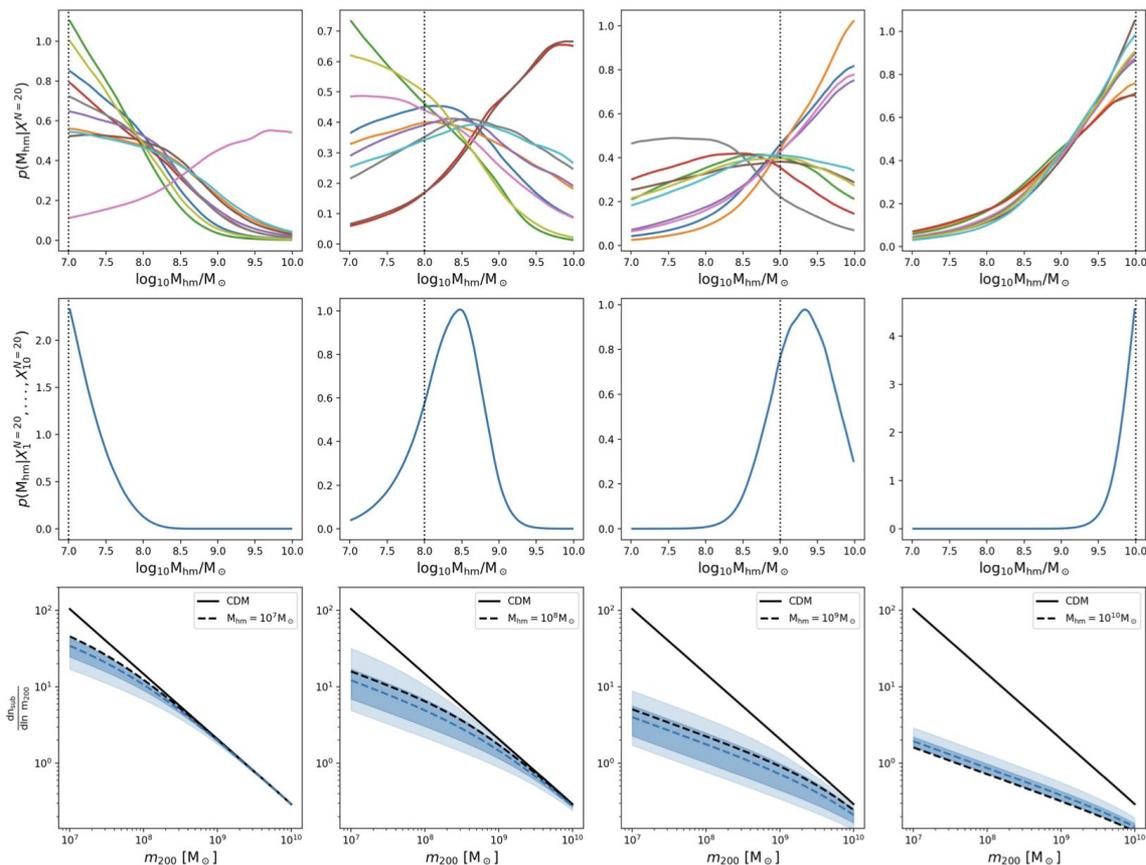
Training data for target image "A"



C) Subhalo mass function cutoff - Results

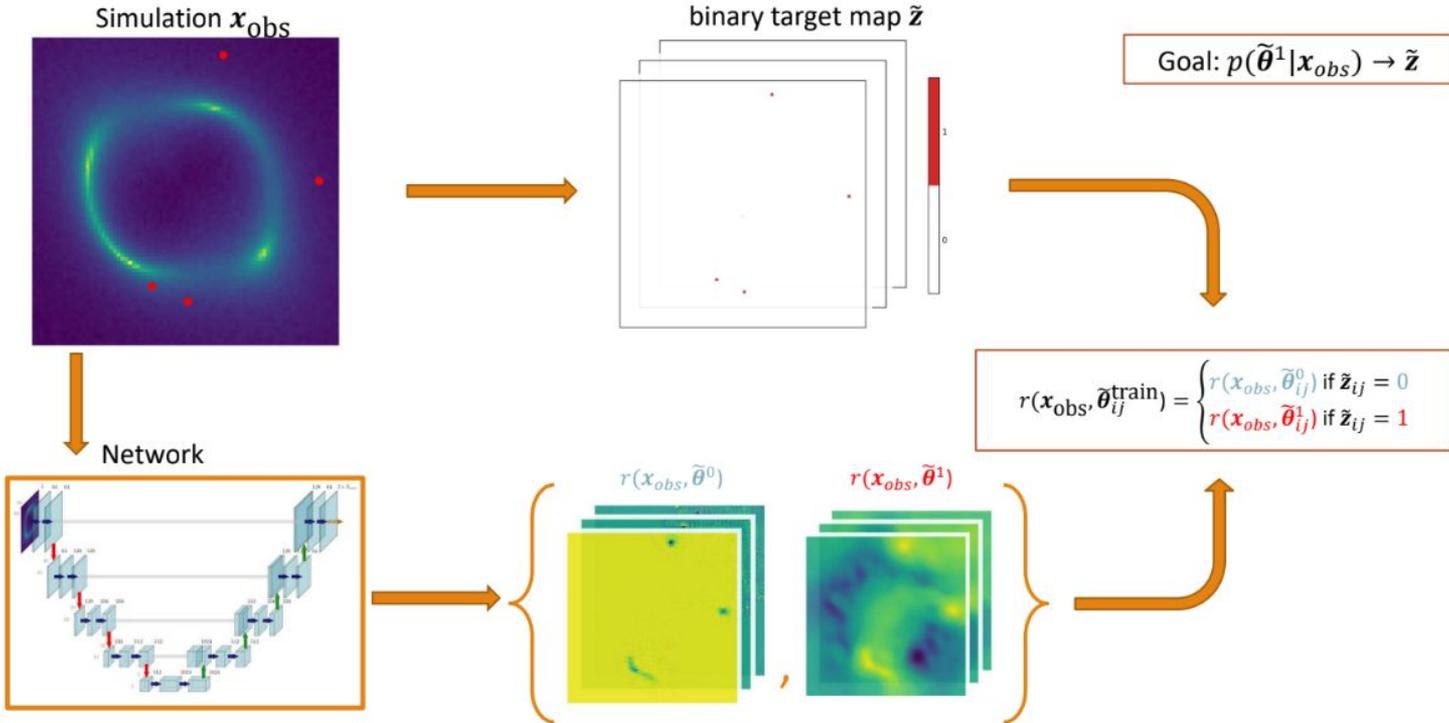
Generating targeted training data for 100 images and combining their constraining power gives tight constraints on the subhalo mass function.

$$p(M_{\text{cutoff}} | \{\vec{x}_i\}_{i=1, \dots, 10})$$



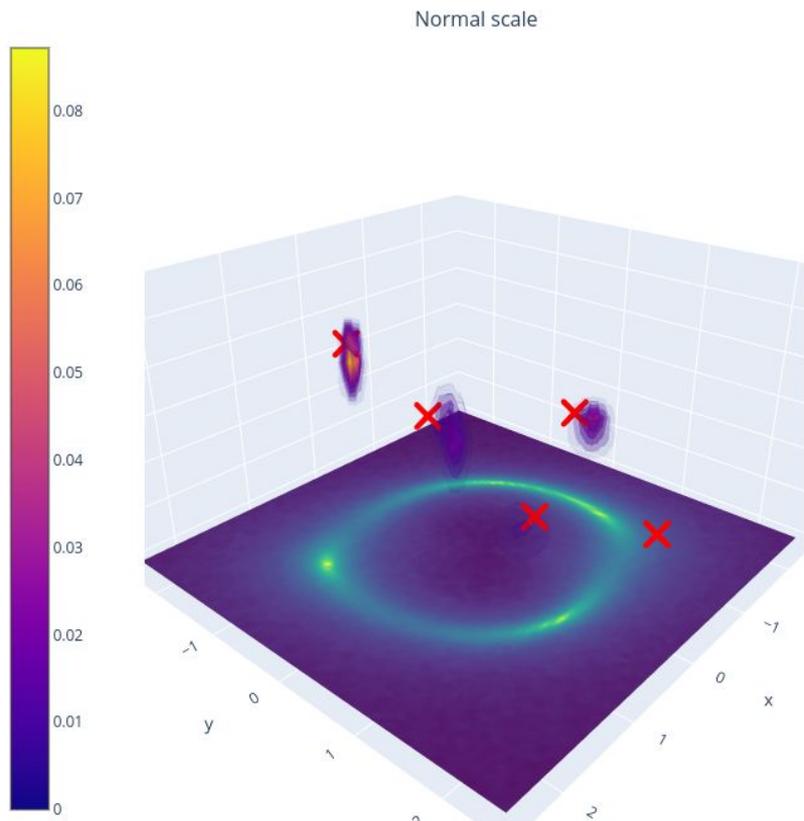
D) Halo detection - Probabilistic image segmentation

In the presence of multiple subhalos, we can also estimate the subhalo density function (which can be understood as marginal of the more complex joined subhalo distribution).



D) Halo detection - Probabilistic image segmentation

Subhalo posteriors. Transparency decreases with posterior value.



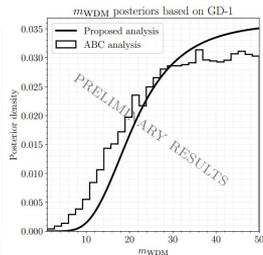
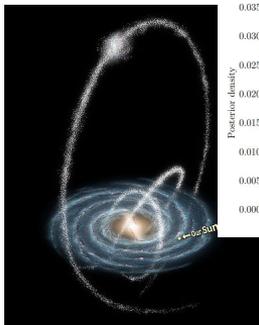
Structure of ratio estimator

- Input: Image (typically 100x100)
- Embedding: U-Net
- Marginals: Binary marginals 100x100x10 (ten mass bins)

<https://dm-lensing-parislfi.github.io/>

TMNRE/SWYFT appear to be broadly applicable

Stellar streams



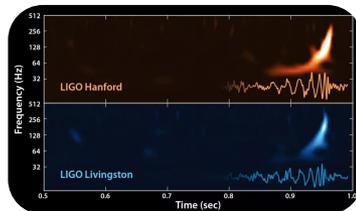
Hermans et al., 2020
James Alvey, Mathis Gerdes, in progress

21 cm cosmology
LHC pheno fits
Fermi data

...

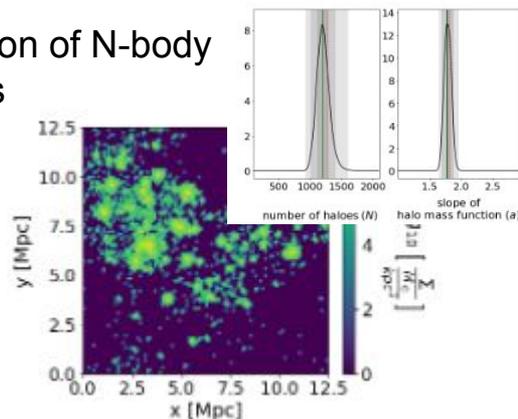
TMNRE/SWYFT

Gravitational waves



Delaunoy+ 2020
Uddipta Bhardwaj+, in progress

Interpretation of N-body simulations



Androniki Dimitriou+, soon

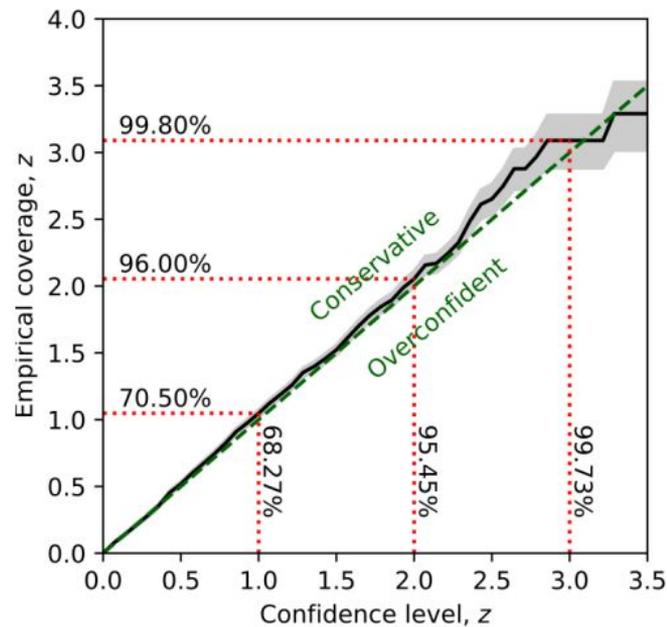
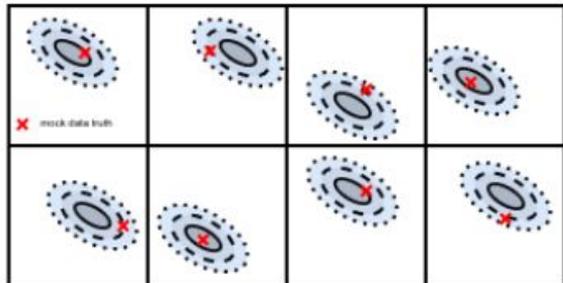
How can one trust results?

Credibility of inference results can be tested

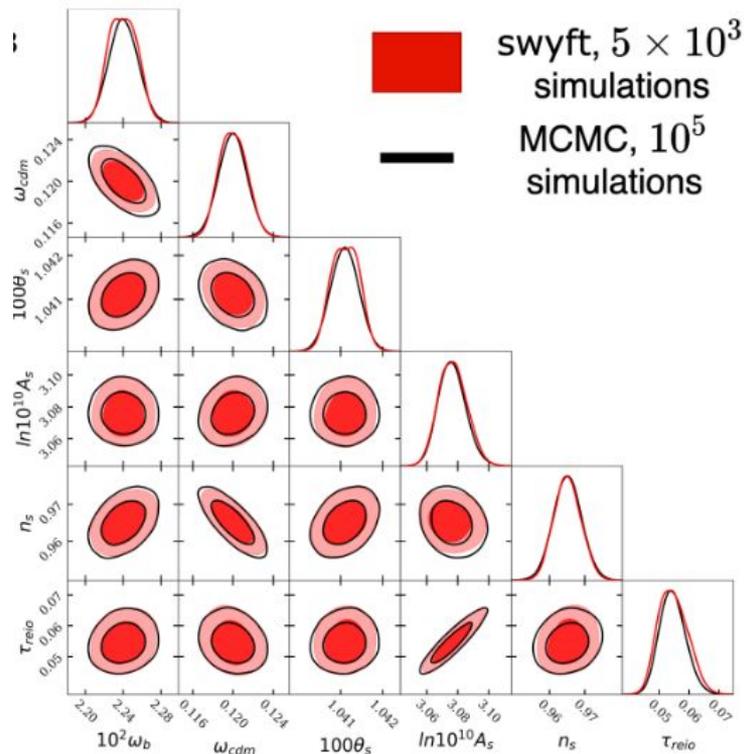
Let $\Theta_{p(\vartheta|\mathbf{x})}(1 - \alpha)$ denote the $1 - \alpha$ highest posterior density region

Expected coverage of the 68% and 95%

$$1 - \hat{\alpha} = \mathbb{E}_{p(\vartheta, \mathbf{x})} [\mathbb{1} [\vartheta \in \Theta_{\hat{p}(\vartheta|\mathbf{x})}(1 - \alpha)]]$$

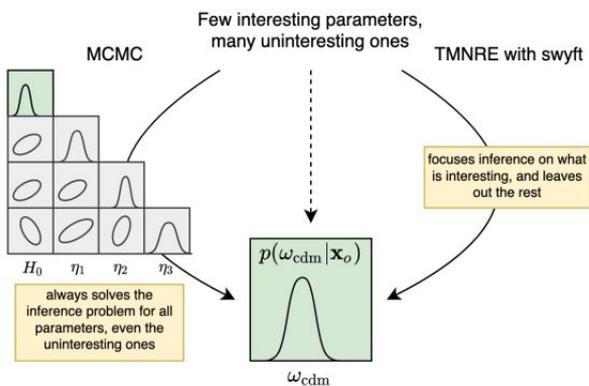


Coverage tests!

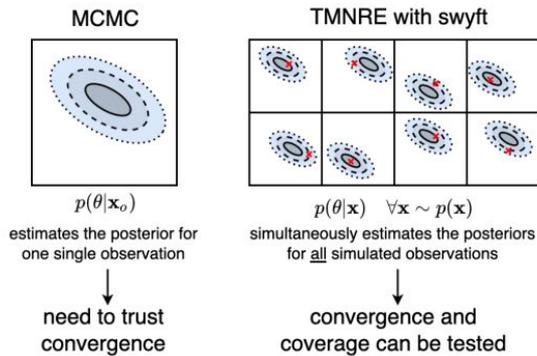


Open source package SWYFT

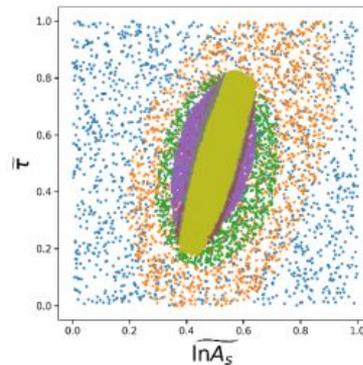
Estimating marginals of interest



Coverage tests



Truncation schemes



Check it out on: <https://github.com/undark-lab/swyft>
(under heavy development)



Conclusions

Conclusions

- Simulation-based inference (SBI) has the potential to deal with **ultra-high dimensional inference problems**.
- We discussed a few components that we found very useful in practice, and which are part of **TMNRE**
 - **Neural ratio estimation** offers flexibility and simplicity
 - Focus on **marginal posteriors** rather than the joint
 - **Prior truncation**
- We demonstrated that this framework is promising in tackling a wide range of astrophysical / cosmological data analysis problems. Domain knowledge enters the analysis in terms of network architectures.
 - **CMB Cosmology**
 - **SN Cosmology**
 - **Strong lensing image analysis**
- We provide a **software implementation for TMNRE (“swyft”)**, which we currently use for a much wider range of dark-matter-related analysis problems.