



Simulation-based inference **[Deserved] hype, tutorial, and future outlook and challenges**

Becky Nevin and DeepSkies Lab

December, 2023

COFI AI/ML Winter School

About me

The Deepskies Lab SBI team



Jason Poh



Sreevani Jarugula



Amanda Pagul



Humna Aman



Egor Danilov



Andrea Roncoli



Gourav Khullar



Becky Nevin



Natalie Malgon



Kabelo Tsiane



Anvi Padhi



Marcos Tamargo



Moonzarin Reza



Neil Kumar



Brian Nord

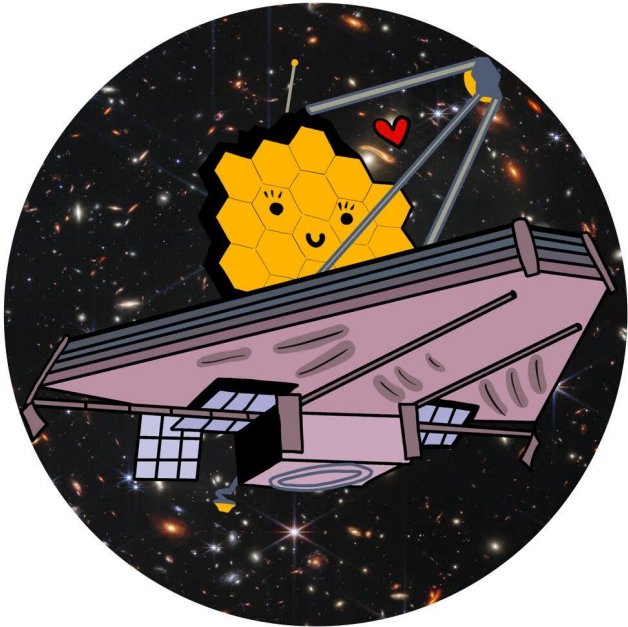


Aleksandra Ćiprijanović



Yuanyuan Zhang

It's Saturday, prizes!



Iteratively updating our presentation prior



Sreevani Jarugula



Jason Poh



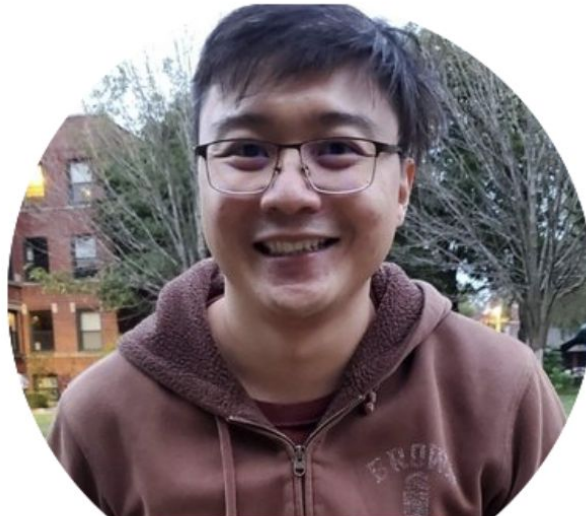
Becky Nevin

Iteratively updating our presentation prior

This presentation is Bayesian



Sreevani Jarugula



Jason Poh



Becky Nevin

What do you know about simulation-based inference?

- Heard of it?
- Talked about it in a class?
- Done research with it?

Also sometimes called
“likelihood-free inference” or
“approximate bayesian
computation”

What do you know about simulation-based inference?

- Heard of it?
- Talked about it in a class?
- Done research with it?

Also sometimes called
“likelihood-free inference” or
“approximate bayesian
computation”

- Have you heard of likelihood-based inference? Ever run MCMC?

Visual outline

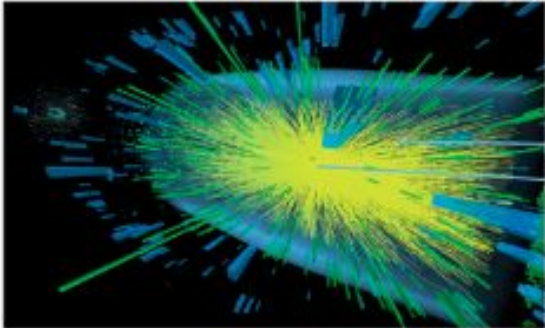
1. Simulation-based versus likelihood-based inference
2. Advantages of SBI relative to LBI
3. Tutorial
4. Applications / future outlook / challenges
5. Deepskies Lab SBI projects

In physics, [and science in general], we love our simulations

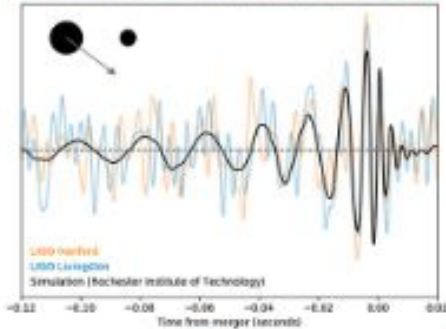
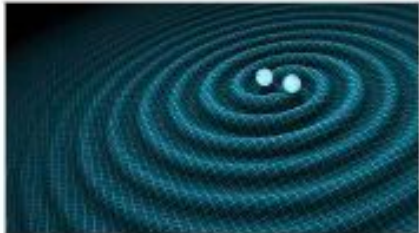
Smallest imaginable scales

astrophysical objects

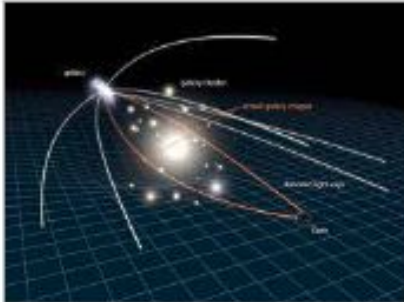
the Universe



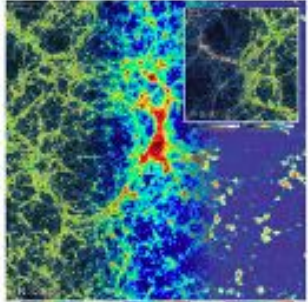
Particle physics



Gravitational wave detectors

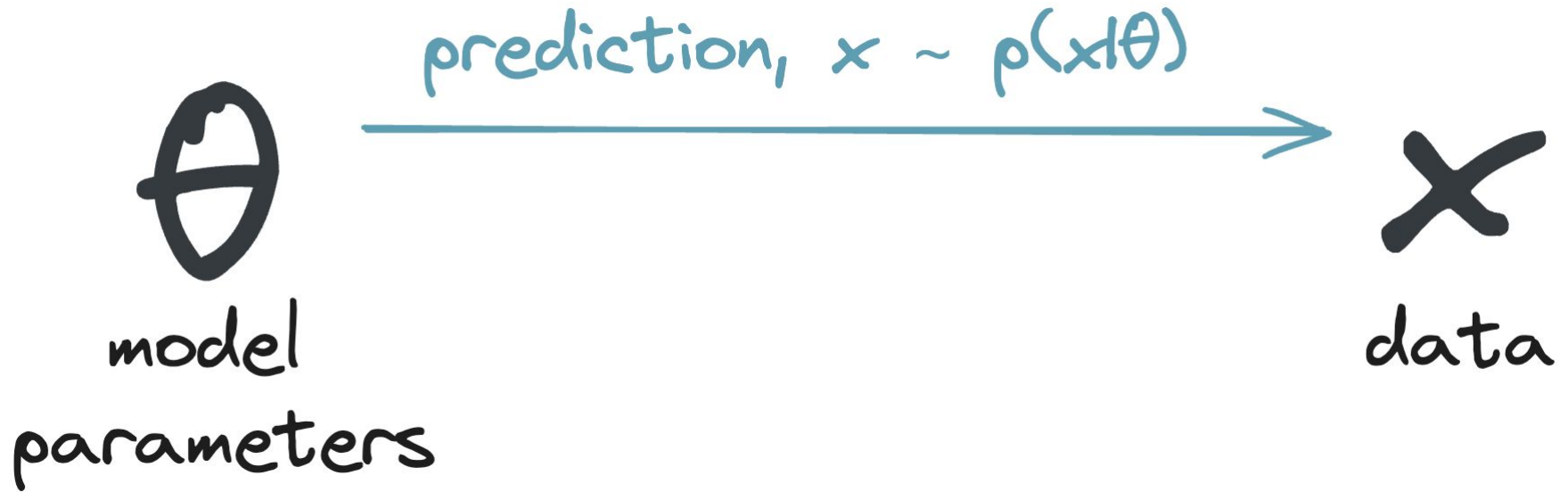


Strong gravitational lens

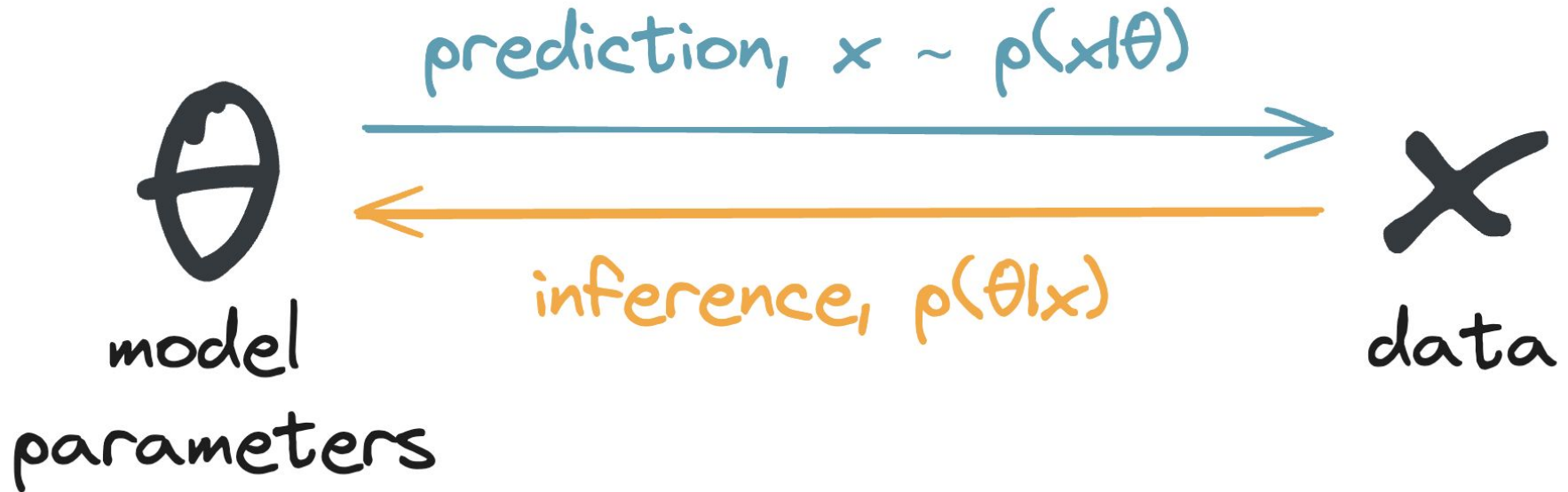


Cosmological simulations

These simulators are giving us a picture of the universe (x)



But what we really want is to be able to say thing statistically about the likely θ s



Inference

Given an observation x , we want to infer the *distribution* of probable parameters θ of the model

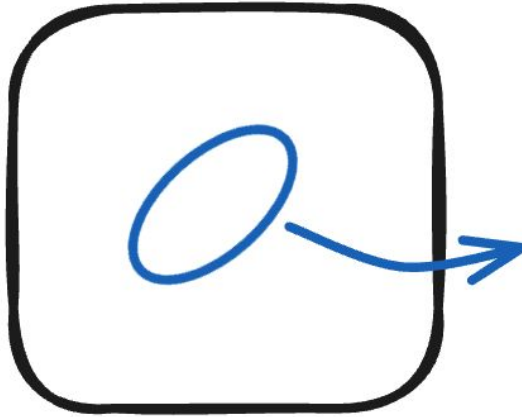
Inference

Given an observation x , we want to infer the *distribution* of probable parameters θ of the model, $p(\theta|x)$

important astrophysical
parameter #1



model
parameters



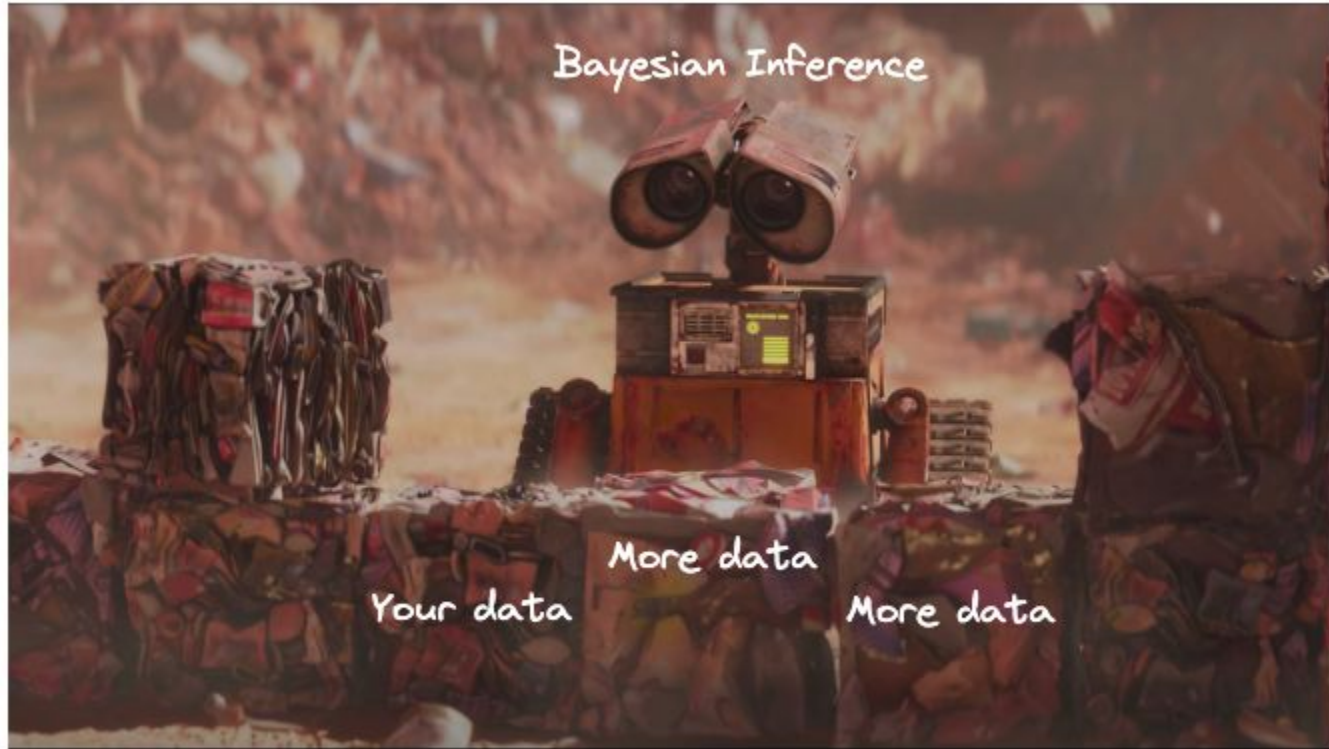
important astrophysical
parameter #2

estimate of posterior
from a posterior
inference method

(Classic) Bayesian Inference (Bayes' Rule)



(Classic) Bayesian Inference (Bayes' Rule)



Bayesian Inference (Bayes' Rule), likelihood-based

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

Bayesian Inference (Bayes' Rule)

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

I'm going to try to convince you why this is often hard/impractical/intractable/impossible to do practically in physics

Bayesian Inference (Bayes' Rule)

$$\text{Posterior} \\ p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

Bayesian Inference (Bayes' Rule)

Posterior

$$p(\theta|x) =$$

θ #1

Marginal
Distribution
 $\int p(\theta_2|x) d\theta_1$



θ #2



Bayesian Inference (Bayes' Rule)

$$p(\theta|x) = \frac{\overset{\text{Likelihood}}{p(x|\theta)} p(\theta)}{p(x)}$$

Bayesian Inference (Bayes' Rule)

$$p(\theta|x) = \frac{\overset{\text{Likelihood}}{p(x|\theta)} p(\theta)}{p(x)} = \int p(x, z|\theta) dz$$

Bayesian Inference (Bayes' Rule)



$$\text{Likelihood } p(x|\theta) p(\theta) = \int p(x, z|\theta) dz$$

$$p(x)$$

Bayesian Inference (Bayes' Rule)

$$p(\theta|x) = \frac{p(x|\theta) \overset{\text{Prior}}{p(\theta)}}{p(x)}$$

Bayesian Inference (Bayes' Rule)

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{\text{"Evidence"} p(x)}$$

Bayesian Inference (Bayes' Rule)

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{\text{"Evidence"} p(x)}$$
$$= \int p(x|\theta)p(\theta)d\theta$$

Likelihood-based inference versus simulation-based

Likelihood-based inference: Bayes' rule (analytically do some math), MCMC sampling approaches

Simulation-based inference: ????



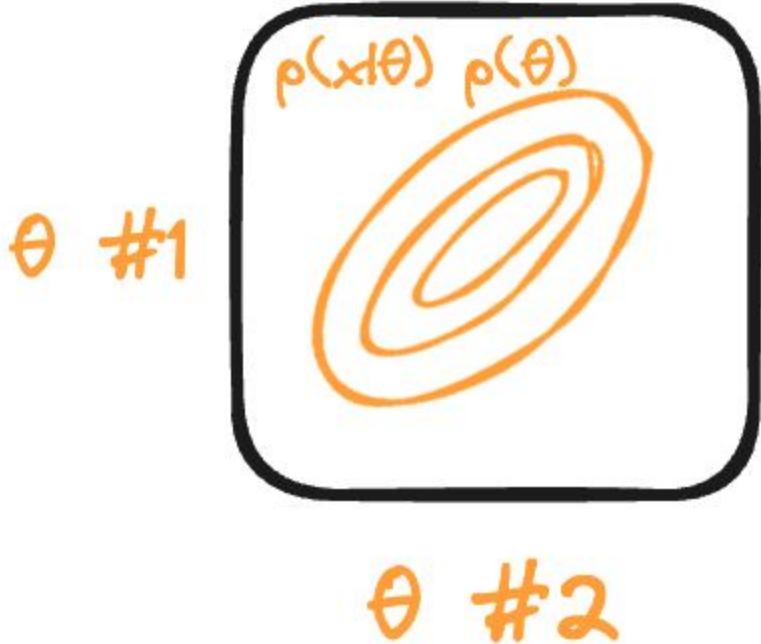
Bayesian Inference (Bayes' Rule)

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{\text{"Evidence"} p(x)}$$
$$= \int p(x|\theta)p(\theta)d\theta$$

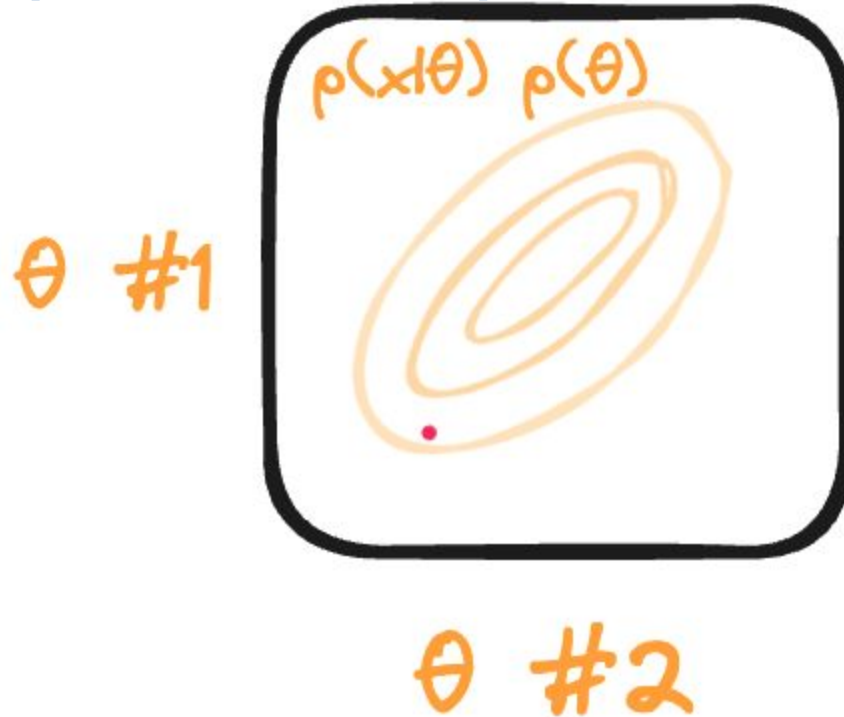
Likelihood-based inference

$$p(\theta|x) \stackrel{\alpha}{\neq} \frac{p(x|\theta)p(\theta)}{\cancel{p(x)}}$$

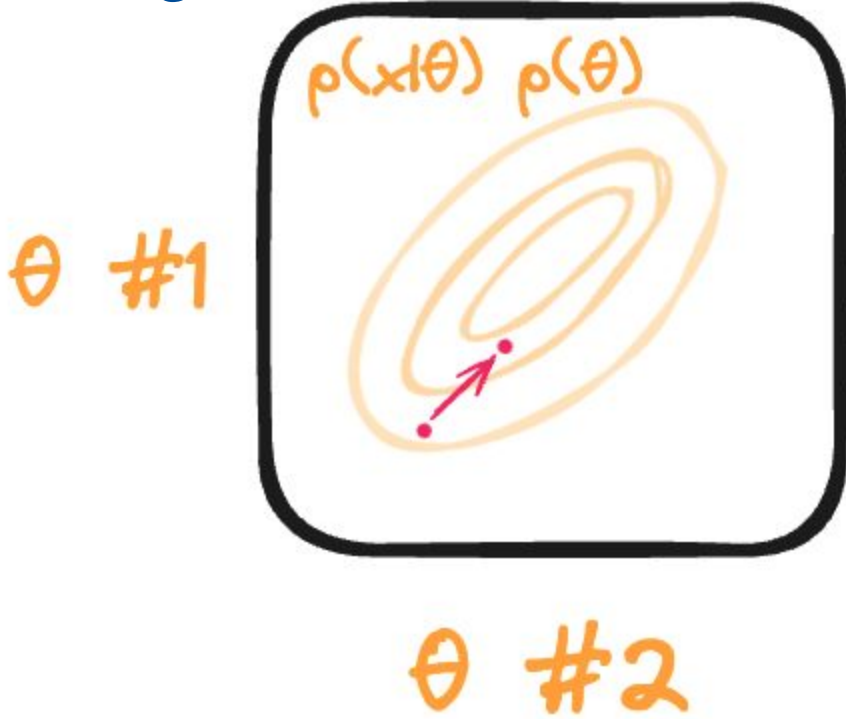
The goal of MCMC is to sample from the approximate posterior manifold



MCMC works iteratively by evaluating the likelihood and prior at a given (θ_1, θ_2) value

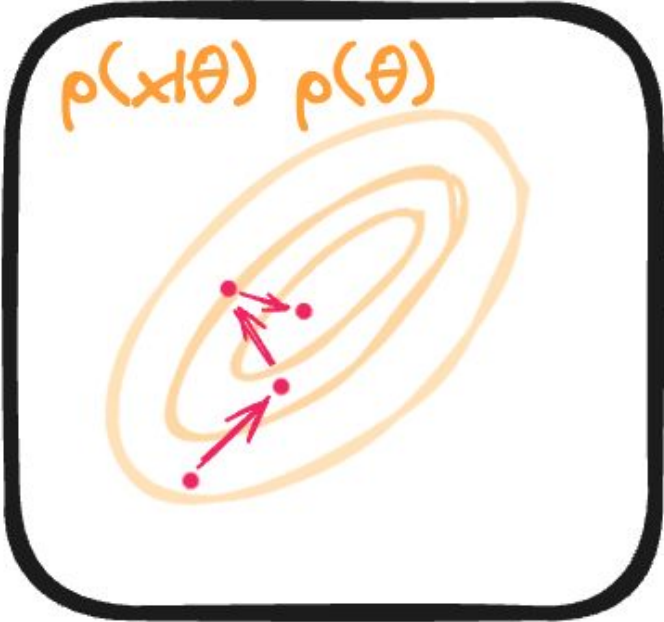


Then it takes a step; this can use a gradient or not depending on the algo



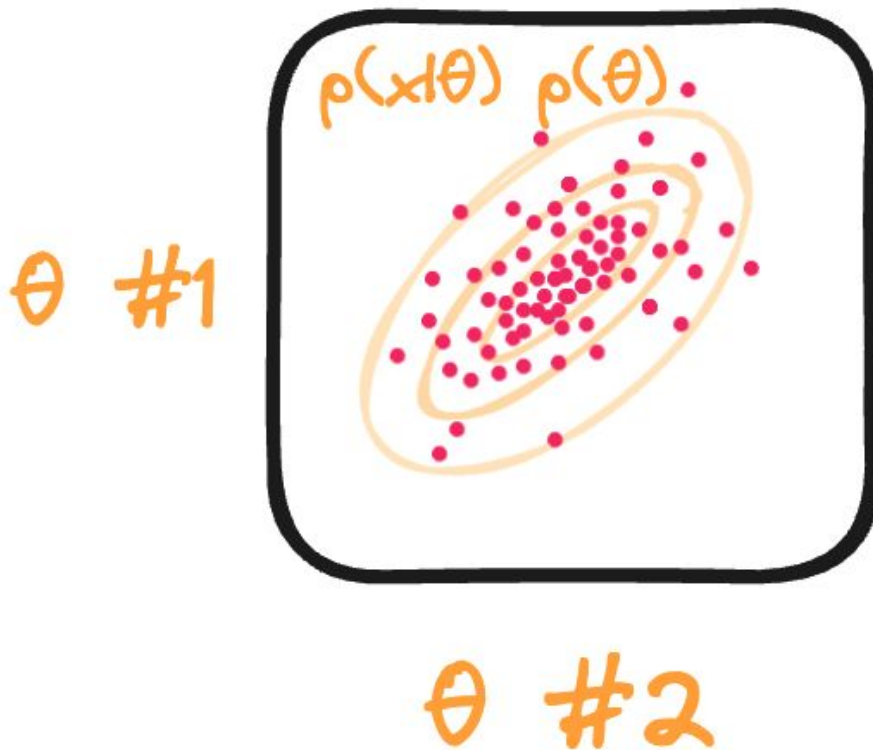
By taking a series of steps, it will walk towards areas of high probability space

θ #1



θ #2

The idea is to do this enough times with enough walkers that you build up a picture of the approximate posterior

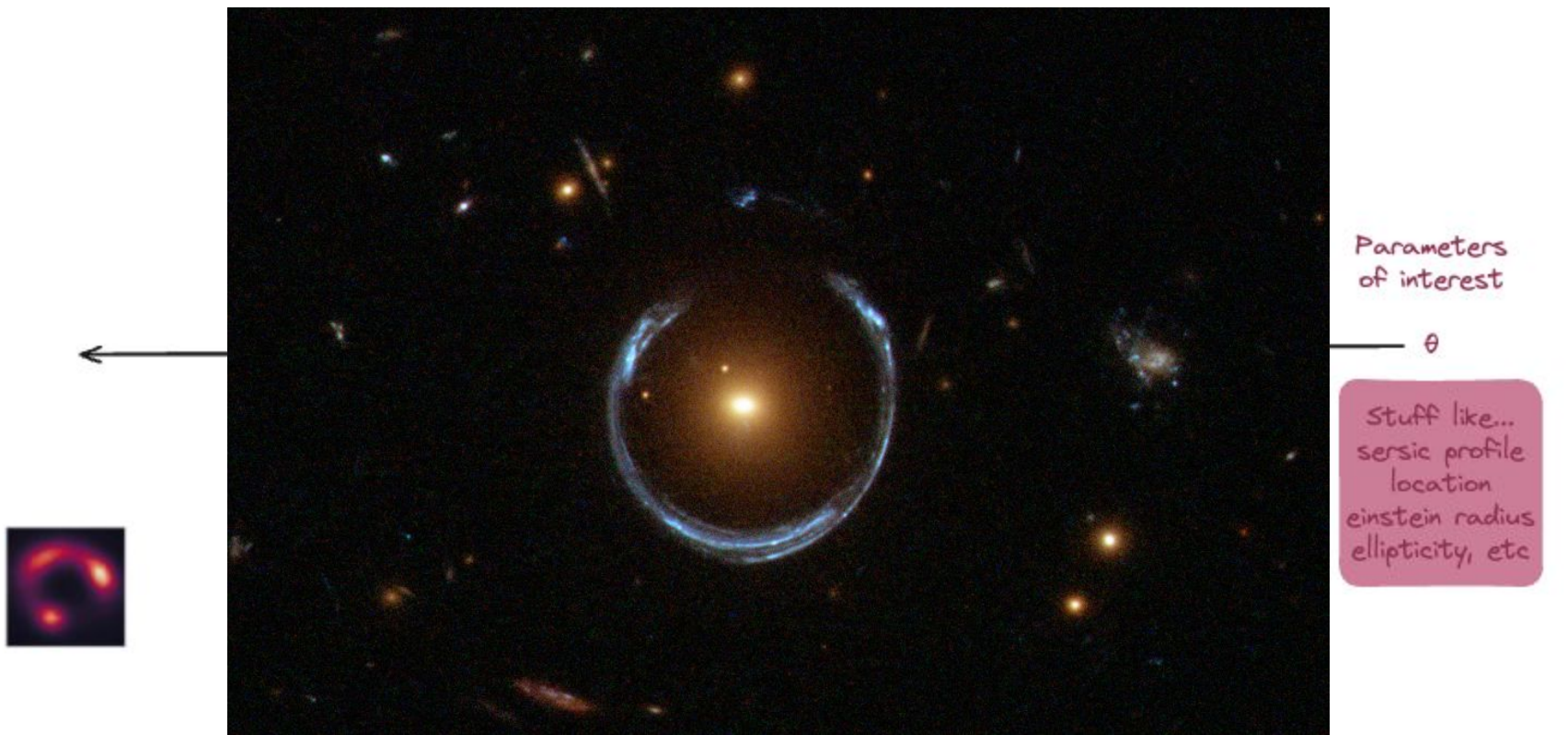


But back to the elephant in the room - what does calculating this likelihood entail?



$$\text{Likelihood} = \int p(x|\theta, z) dz$$
$$\frac{p(x|\theta) p(\theta)}{p(x)}$$

An example simulation from Deepskies lab; generating images of gravitational strong lenses



An example simulation from Deepskies lab; generating images of gravitational strong lenses

deepskies/
deeplenstronomy



[Morgan+2021](#)
Birrer+2021
Birrer & Amara 2018

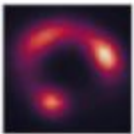
A pipeline for versatile strong lens sample simulations

9 Contributors 28 Issues 6 Discussions 25 Stars 7 Forks

X

Parameters of interest

θ



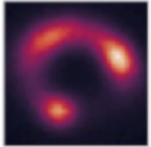
$$x \sim p(x|\theta)$$

stuff like...
sersic profile
location
einstein radius
ellipticity, etc

Latent parameters

Observables

x



Cosmic
evolution

Parameters
of interest

cosmo ← θ

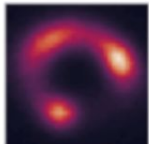
cosmological
parameters,
coupling of
dark matter
and galaxies

Stuff like...
seraic profile
location
einstein radius
ellipticity, etc

Latent parameters

Observables

x



Ray tracing and
thin lens approximation

Cosmic
evolution

Parameters
of interest

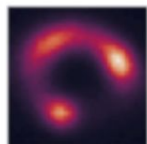
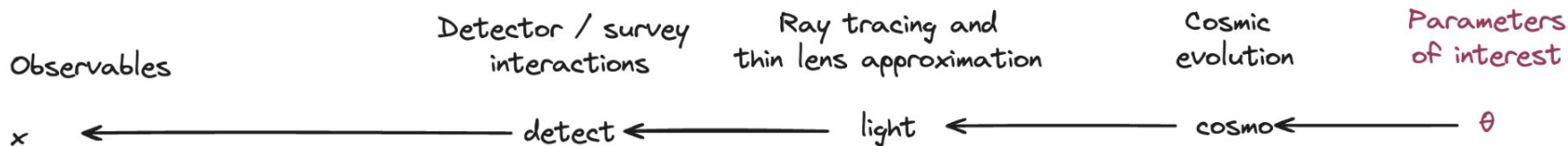
light ← cosmo ← θ

deterministic
physics
equations,
integrating
along the line
of sight

cosmological
parameters,
coupling of
dark matter
and galaxies

Stuff like...
seraic profile
location
einstein radius
ellipticity, etc

Latent parameters



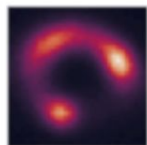
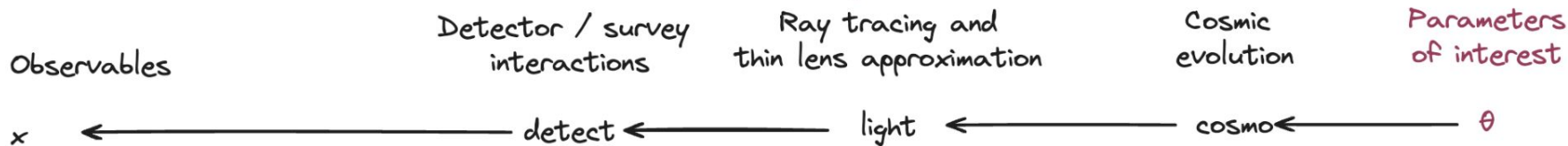
Poisson noise,
detector
properties,
various
stochastic
processes

deterministic
physics
equations,
integrating
along the line
of sight

cosmological
parameters,
coupling of
dark matter
and galaxies

Stuff like...
seraic profile
location
einstein radius
ellipticity, etc

$$x \sim p(x|\theta)$$



Poisson noise,
detector
properties,
various
stochastic
processes

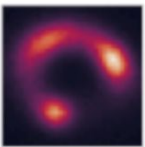
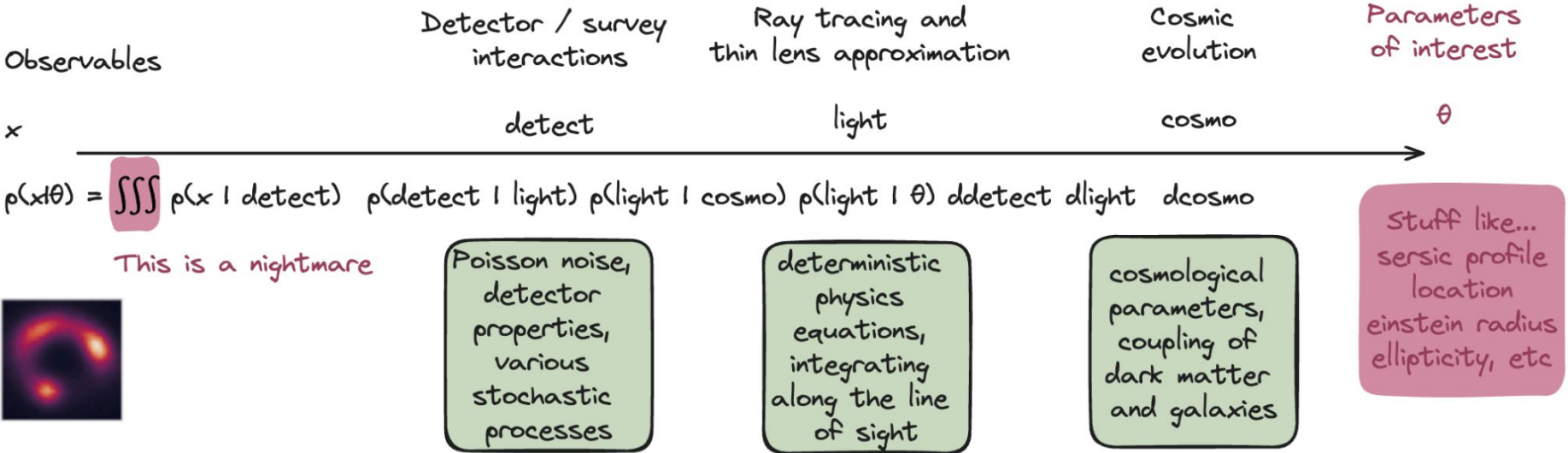
deterministic
physics
equations,
integrating
along the line
of sight

cosmological
parameters,
coupling of
dark matter
and galaxies

Stuff like...
seraic profile
location
einstein radius
ellipticity, etc

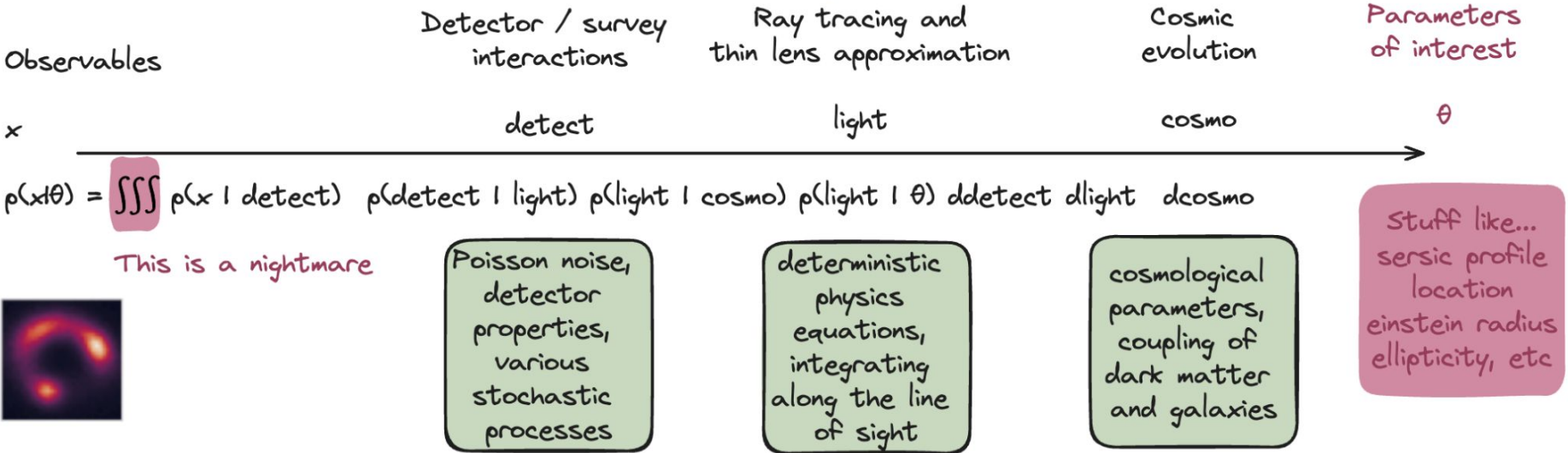
Likelihood-based inference is challenging if the simulation is very high dimensional with lots of params

Latent parameters



Likelihood-based inference is challenging if the simulation is very high dimensional with lots of params

Latent parameters



Idea credit: Kyle Cranmer

Disadvantages of MCMC / Likelihood-based inference

1. You need to evaluate the likelihood
2. You have to do this for every piece of data

Disadvantages of MCMC / Likelihood-based inference

1. You need to evaluate the likelihood
2. You have to do this for every piece of data (trash)

Big surveys are like:



But what if I told you:

Is not always necessary to evaluate the likelihood



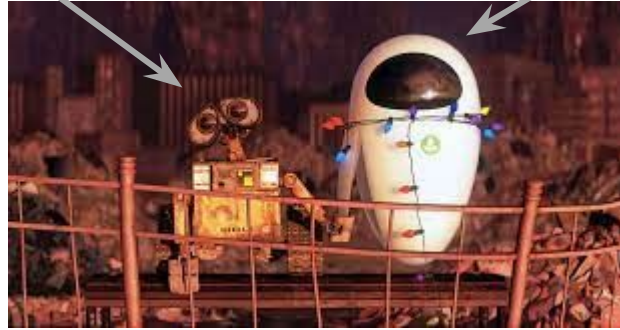
Simulation-based inference does not require that you write out or evaluate a likelihood

A word on “likelihood-free” inference

- Can still be floating around in there, “implicit likelihood”
- There might still be a likelihood you can write out

Approximate bayesian computation (ABC) is an example of SBI

Marriage of Bayesian inference and computation



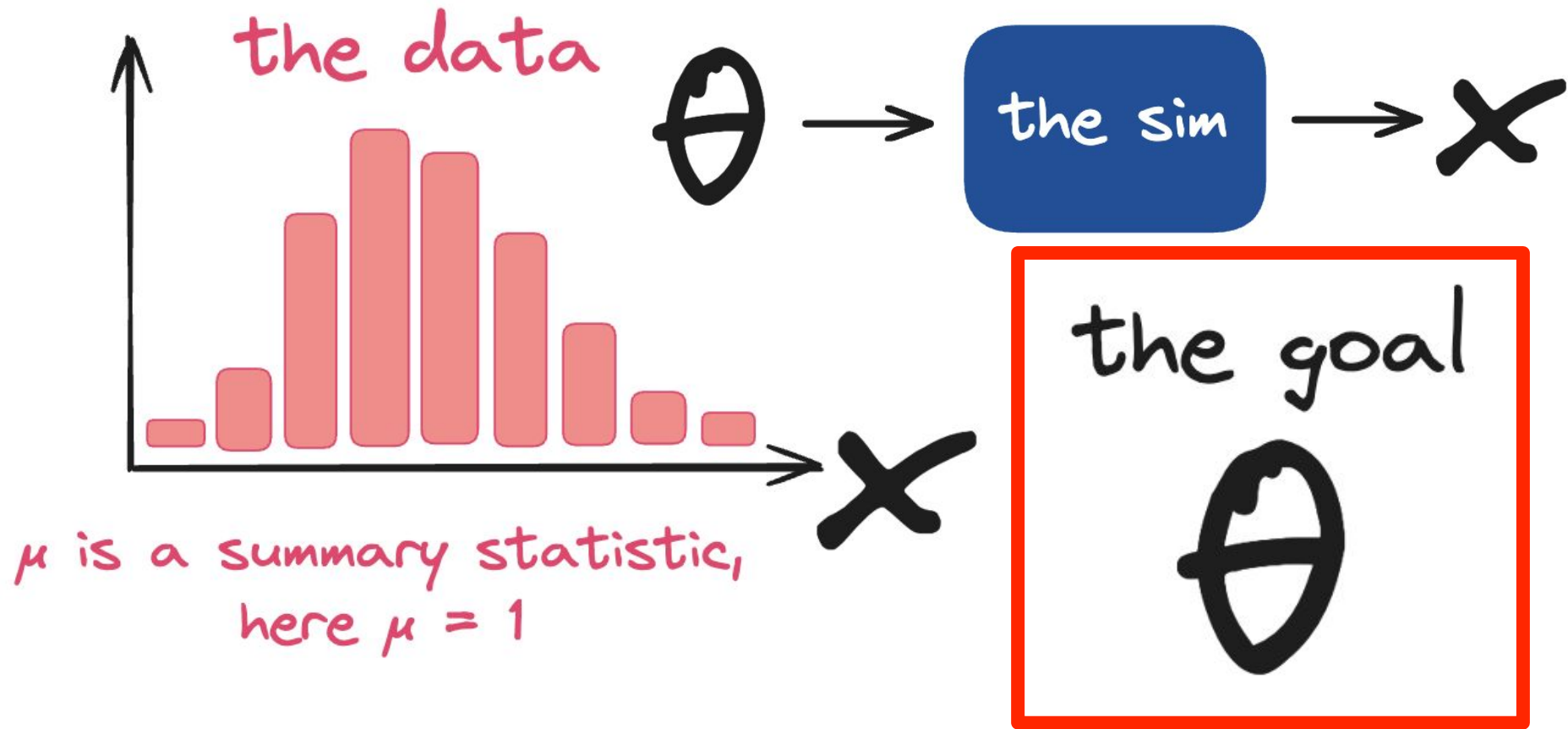
If you have access to an efficient simulation...

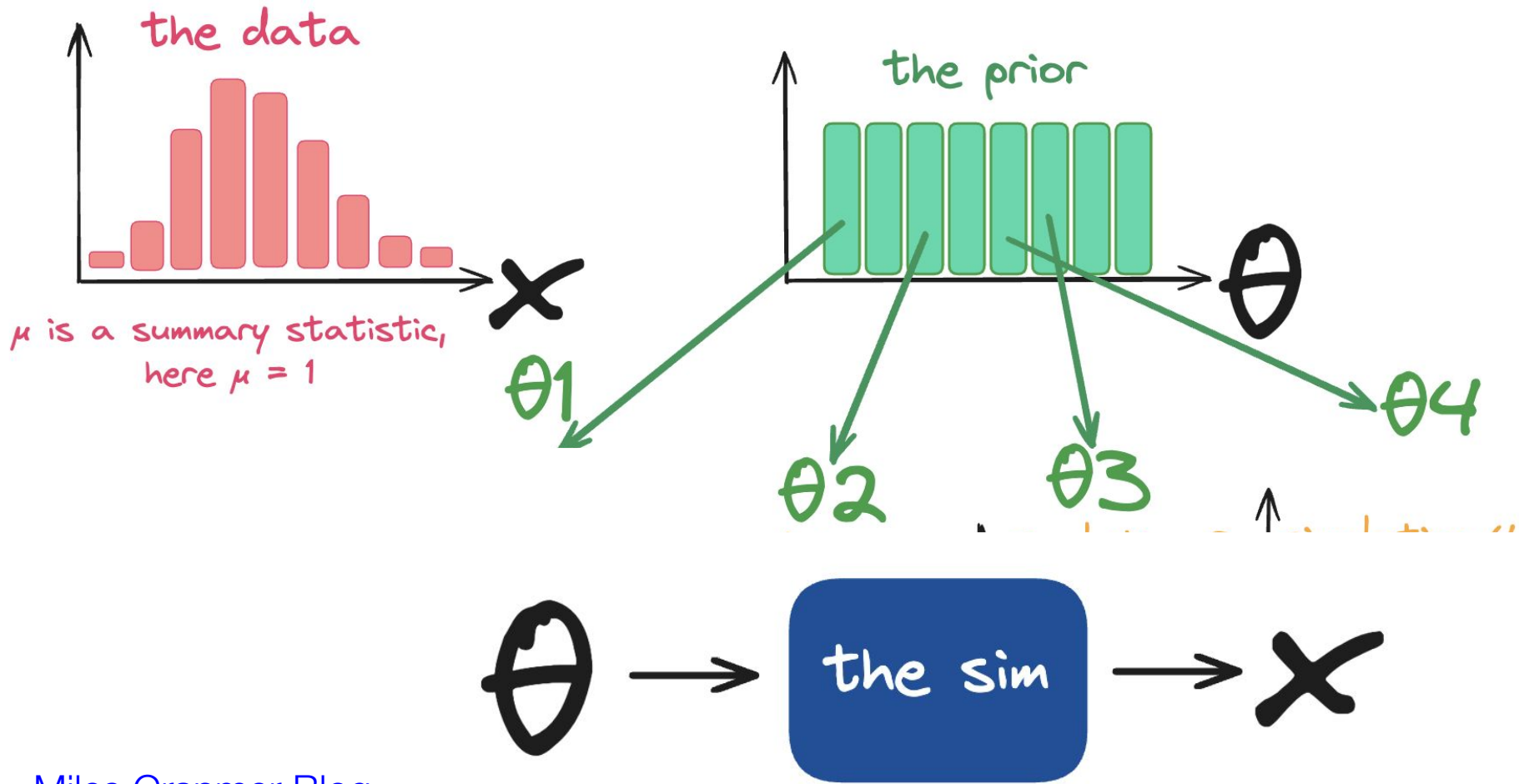
the data

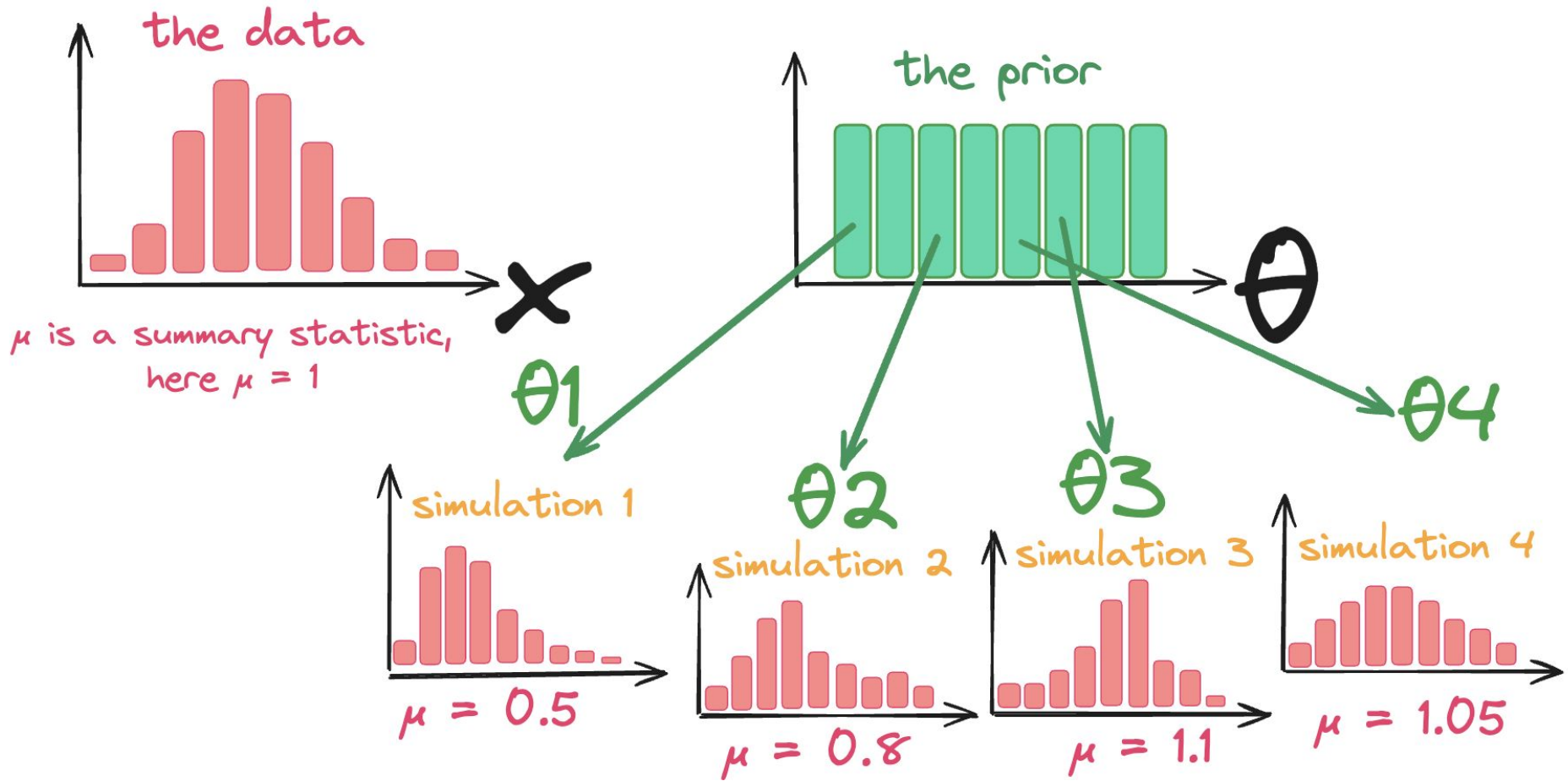


μ is a summary statistic,
here $\mu = 1$







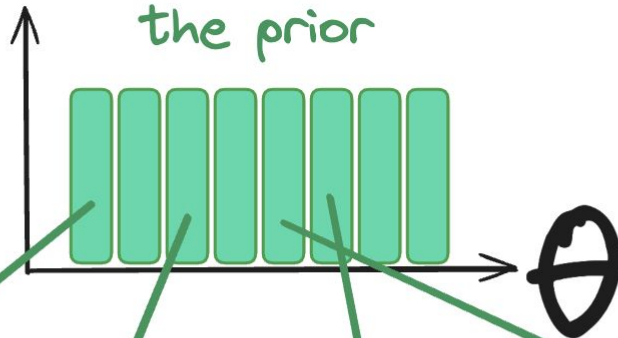


the data



μ is a summary statistic,
here $\mu = 1$

the prior



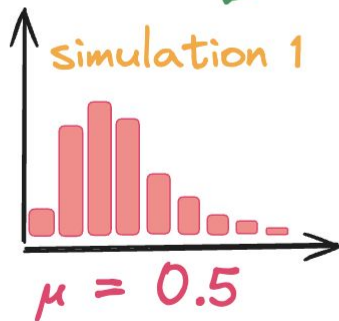
θ_1

θ_2

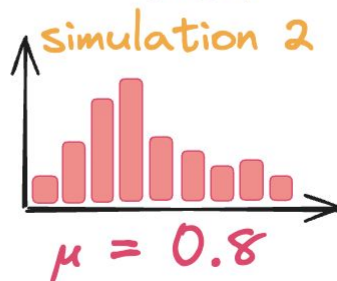
θ_3

θ_4

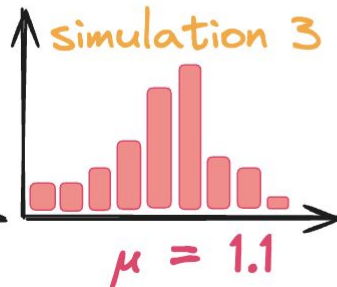
now compare the μ values with
that of the data,
if $\rho(\mu_{\text{data}}, \mu_{\text{sim}}) < \text{epsilon}$,
keep that simulation,
here $\text{epsilon} = 0.1$



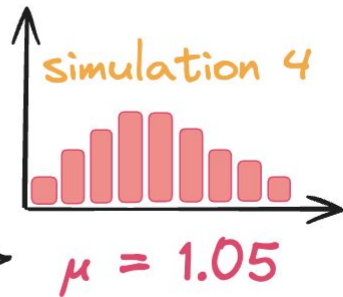
✗



✗

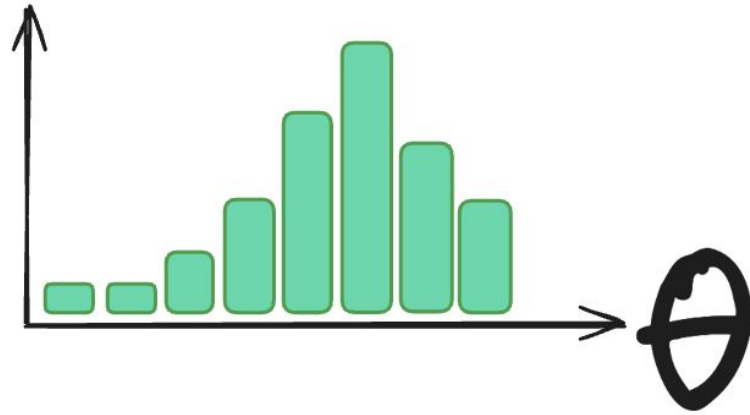


✓



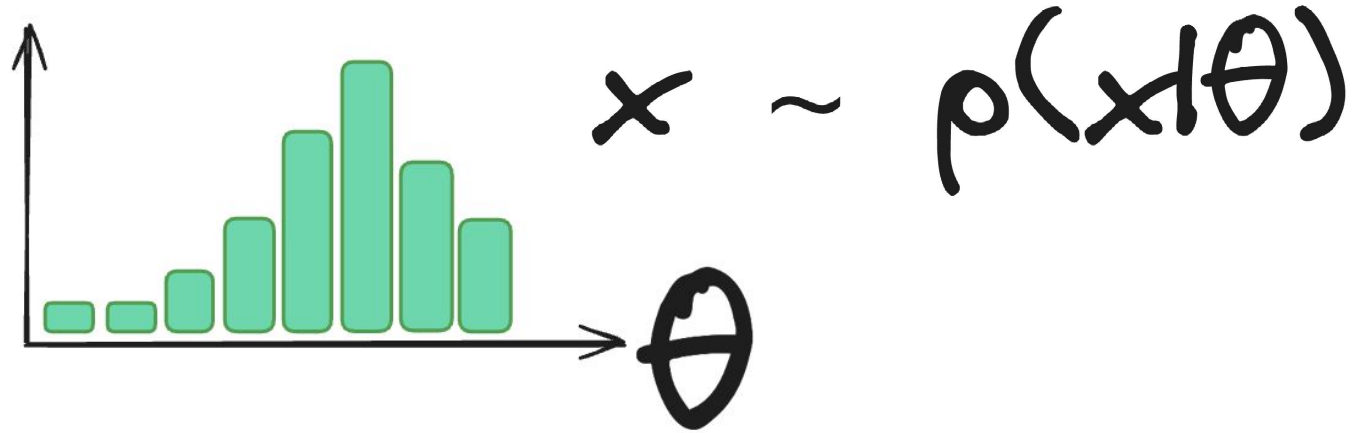
✓

combine the θ values from the simulators that pass the test
use these to construct the posterior distribution as a fxn of θ



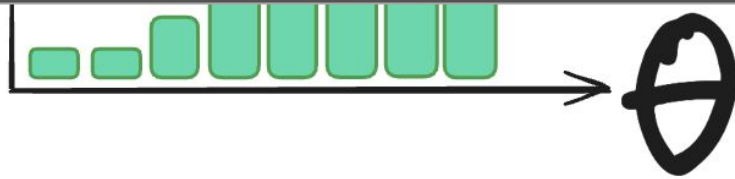
ABC involves the likelihood implicitly via sampling

combine the θ values from the simulators that pass the test
use these to construct the posterior distribution as a fxn of θ



combine the θ values from the simulators that pass the test
use of θ

Surprise! ABC is an early form of simulation-based inference! It uses a simulation and not an explicit likelihood to approximate the posterior landscape.



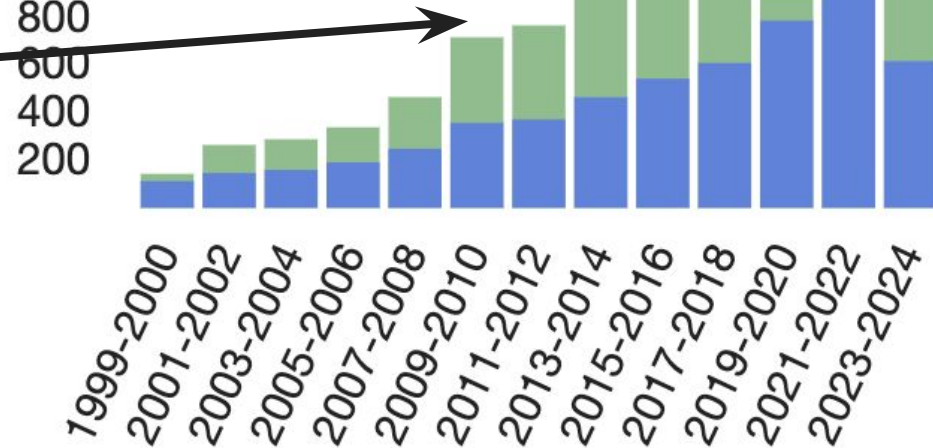
Simulation-based inference has undergone a recent glow-up

First normalizing flows article I could find:

Agnelli+ 2010 “Clustering and Classification through Normalizing Flows in Feature Space”

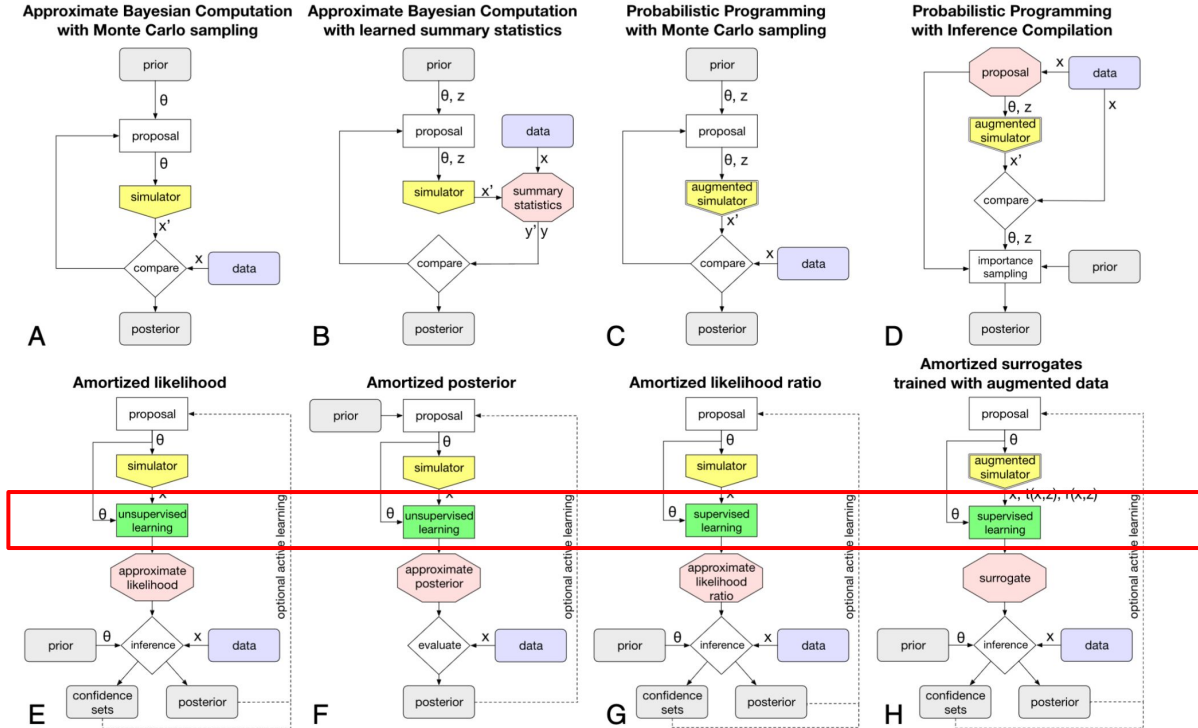
Rezende & Mohamad 2015

2.6k
2.4k
2.2k
2k
1.8k
1.6k
1.4k
1.2k
1k
800
600
400
200



The SBI landscape circa 2020

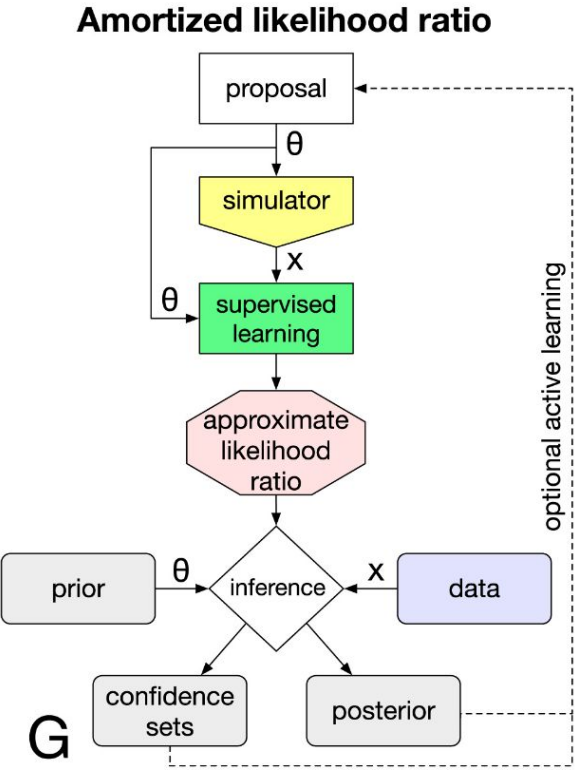
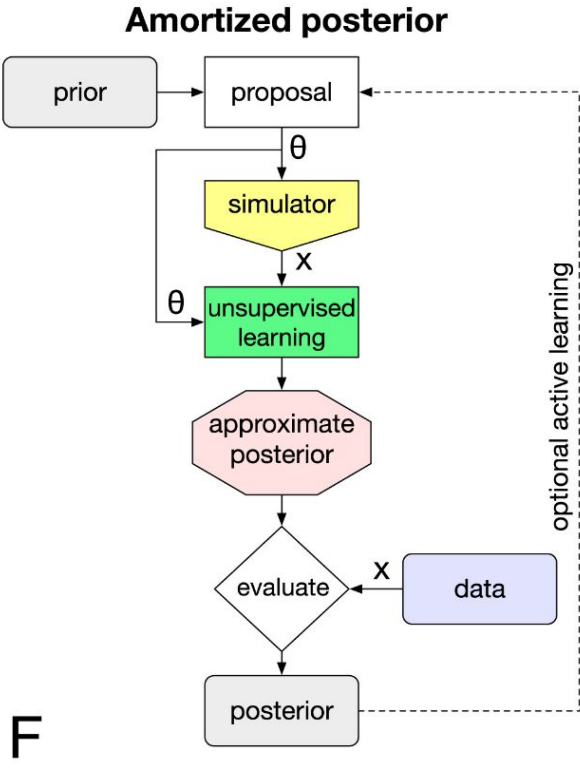
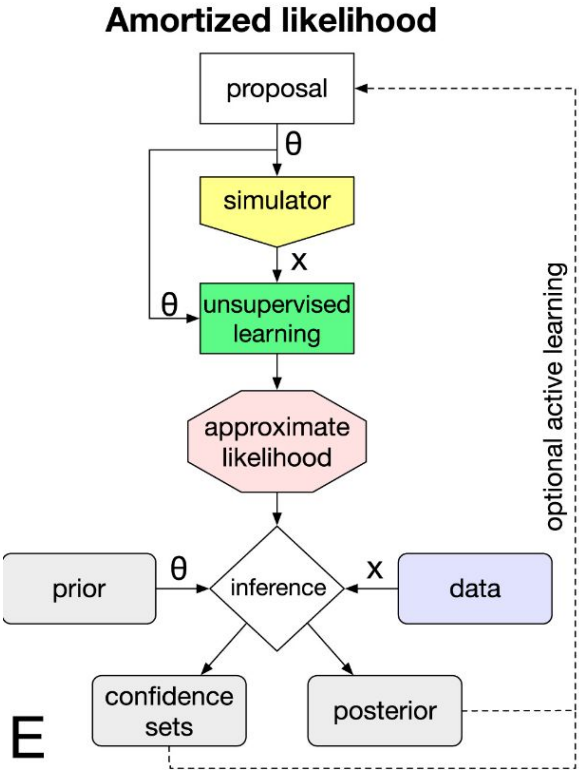
Credit: [Cranmer et al \(2020\)](#)
“The Frontier of SBI”



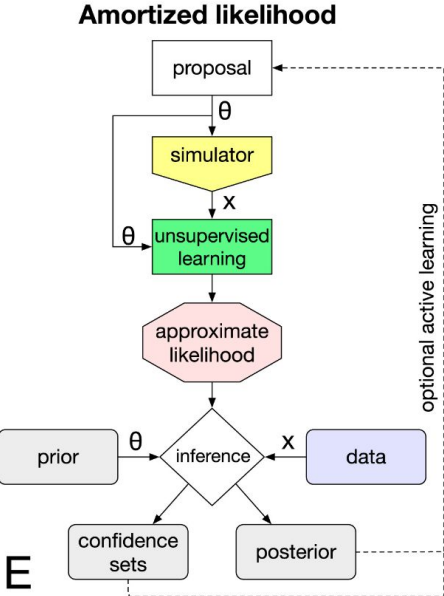
Machine Learning

Fig. 1. (A–H) Overview of different approaches to simulation-based inference.

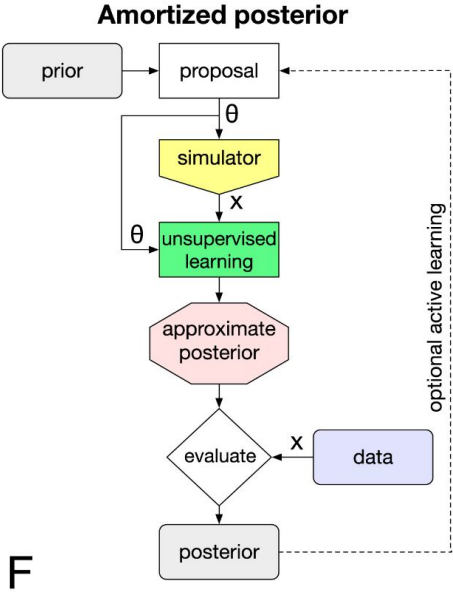
A HUGE advantage of modern SBI methods is that they are amortized



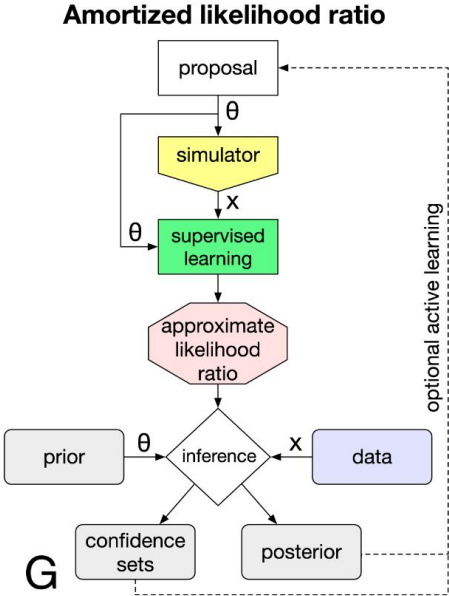
A HUGE advantage of modern SBI methods is that they are amortized



Likelihood



Posterior



Likelihood ratio

Targets of SBI

- Likelihood ratio
- Likelihood / evidence = posterior / prior
- Likelihood (Neural likelihood estimation), estimating the probability density of the data conditional on parameters
- Neural posterior estimation
- Dimensionality of data and parameter vector

Advantages of SBI

- No likelihood required, skips integrating over latent parameters
- You don't have to re-run the inference for every data point
- Very advantageous in cases of large latent dimensionality

Additional resources

- [Talk by Stephen Green about gravitational wave parameter estimation](#)
- [Talk by Kyle Cranmer](#) also about SBI
- https://www.youtube.com/watch?v=315xKcYX-1w&ab_channel=TheTWIMLAIPodcastwithSamCharrington → an ML podcast with George Papamakarios about masked autoregressive flows
- https://www.youtube.com/watch?v=7q4ueFiJjAY&ab_channel=KapilSachdeva → MADE paper review
- This awesome blog about flow models:
<https://lilianweng.github.io/posts/2018-10-13-flow-models/>

Details of the neural density estimation

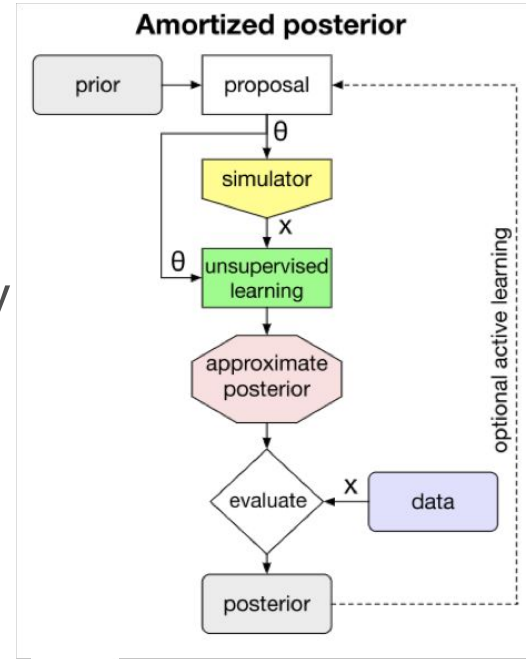
- Review of normalizing flows
- Masked autoregressive flows

Modern SBI methods, neural posterior estimation (NPE)

Modern SBI methods leverage **machine learning** techniques to overcome the issues faced by traditional SBI methods.

The general workflow:

1. Pay upfront simulation cost of generating many pairs of $\{\theta_i|x_i\}$.
2. Use ML model (neural density estimators) to model the posterior distribution from the simulation data.
3. Evaluate new observations x_{new} using the trained model.



Credit: [Cranmer et al \(2020\)](#)

Sequential slide

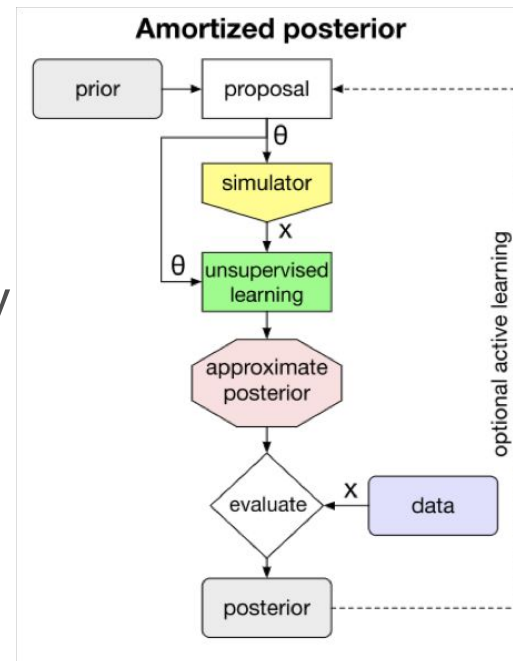
- When the posterior is way more concentrated than the prior you don't really need to explore the likelihood everywhere in parameter space
- Motivates active learning aka sequential methods such as:
 - Sequential Neural Likelihood Estimation [SNLE]
 - Sequential Neural Posterior Estimation [SNPE]
 - Sequential Neural Ratio Estimation [SNRE]
- NOT amortized, partially amortized

Modern SBI methods

Modern SBI methods leverage **machine learning** techniques to overcome the issues faced by traditional SBI methods.

The general workflow:

1. Pay upfront simulation cost of generating many pairs of $\{\theta_i|x_i\}$.
2. Use ML model (neural density estimators) to model the posterior distribution from the simulation data.
3. Evaluate new observations x_{new} using the trained model.

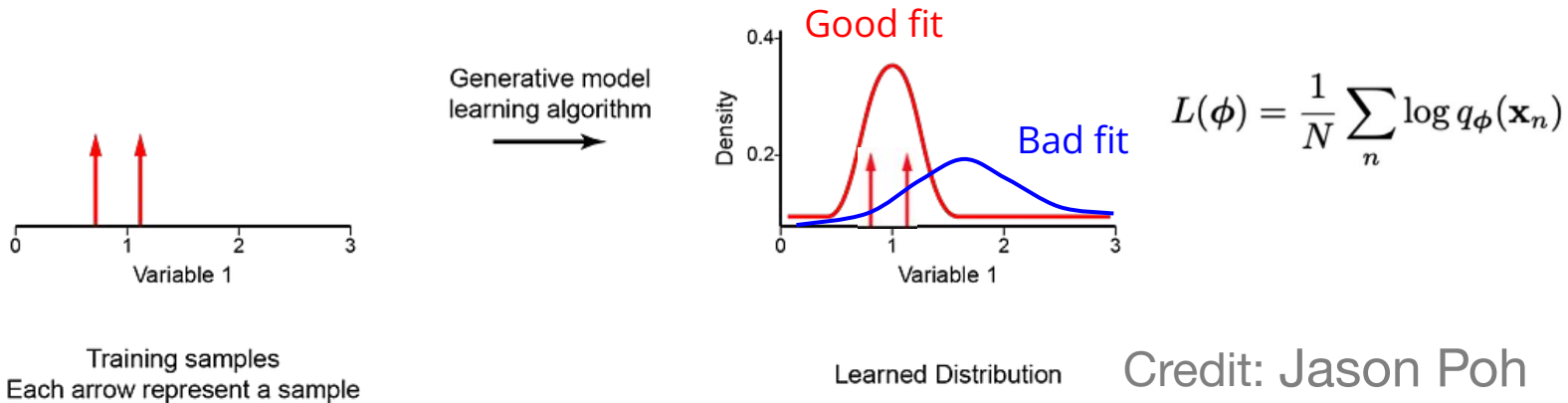


Credit: [Cranmer et al \(2020\)](#)

Neural density estimators

Problem statement: We have some data $\{x_1, \dots, x_n\}$ generated from a process $p(x)$. We want to fit a parametric model, $q_\phi(x)$ to approximate $p(x)$.

- We want to fit $q_\phi(x)$ to the data such that $q_\phi(x)$ is **most likely** to generate $\{x_1, \dots, x_n\}$.



Credit: Jason Poh

Neural density estimators

There are non-neural density estimators, think Gaussians (parametric) or kernel density estimators (non-parametric)

But **neural** network based density estimators have more freedom

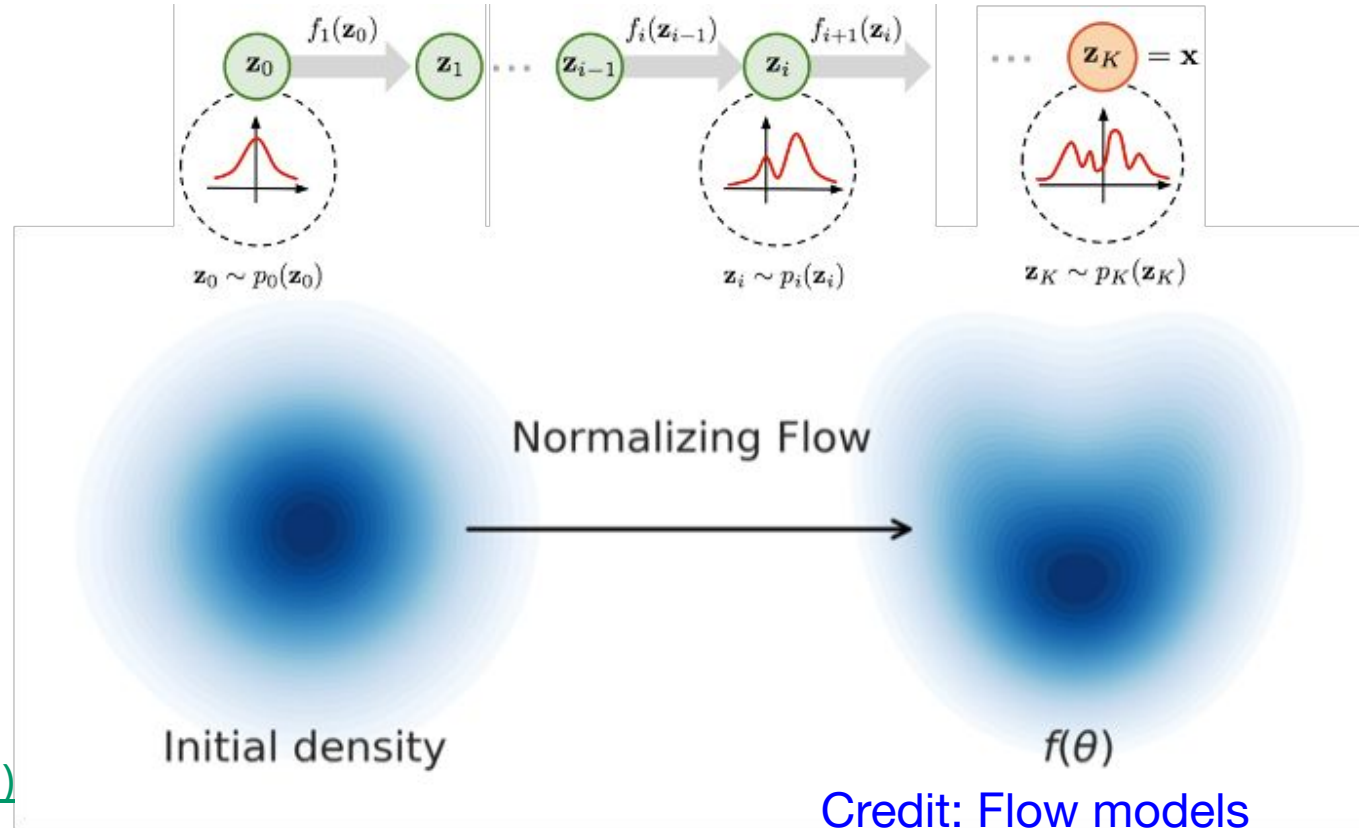


Fast ϵ -free Inference of Simulation Models with Bayesian Conditional Density Estimation

George Papamakarios
School of Informatics
University of Edinburgh
g.papamakarios@ed.ac.uk

Iain Murray
School of Informatics
University of Edinburgh
i.murray@ed.ac.uk

Building blocks of modern NDEs



[A 1 minute description](#)

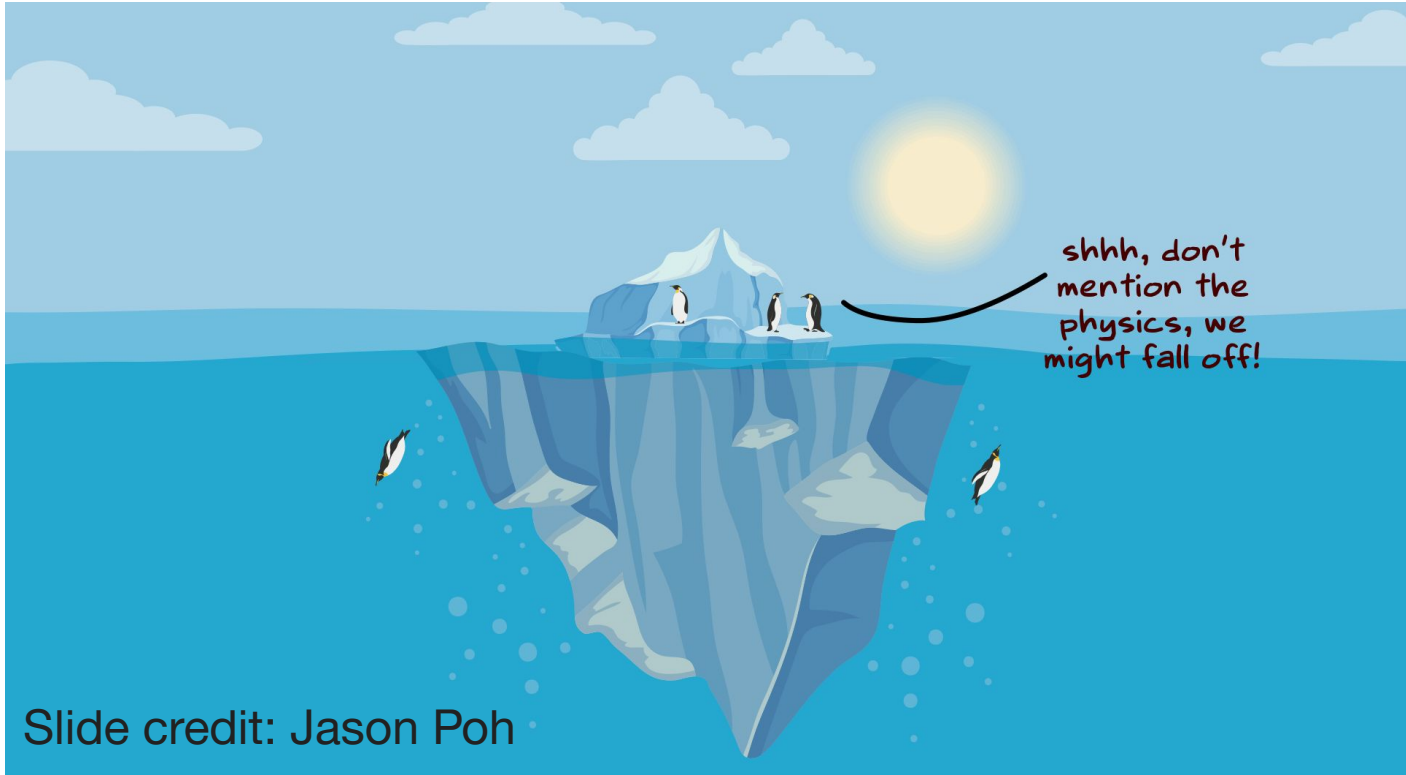
[Video Tutorial \(12 mins\)](#)

[Credit: Flow models](#)

Normalizing flows (and the other methods within this family) can transform any distribution into a simpler distribution in an invertible way (bijective)



Masked Autoregressive Flows



Autoregressive models model densities conditionally

- A class of models where a target joint distribution is factorized over n-dimensional probability conditionals, and those conditionals are modeled in turn.
- For example, if we are interested in the probability density of an image in an autoregressive model, the joint probability of an image could be the combination of the probability of all its pixels (as was used in PixelCNN)

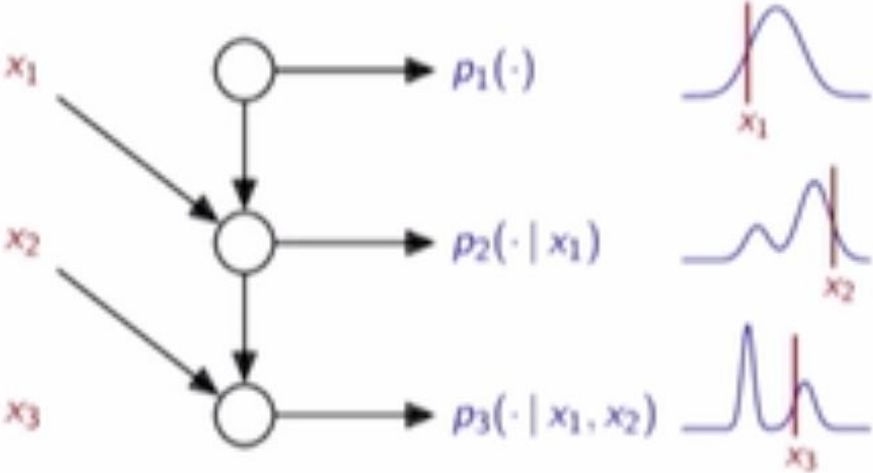
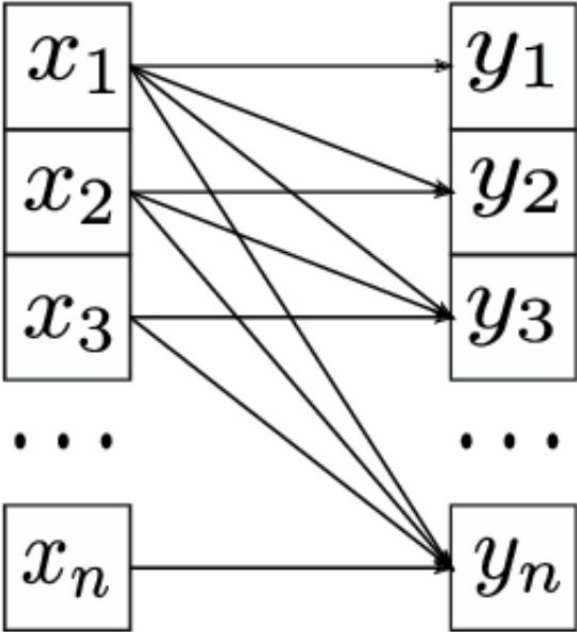
$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_n)$$
$$p(\mathbf{x}) = p(x_1)p(x_2) \dots p(x_n)$$
$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i)$$
$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, \dots, x_{i-1})$$

Likelihood of image \mathbf{x} Probability of i'th pixel value given all previous pixels

[Credits: PixelCNN](#)

Slide credit: Jason Poh

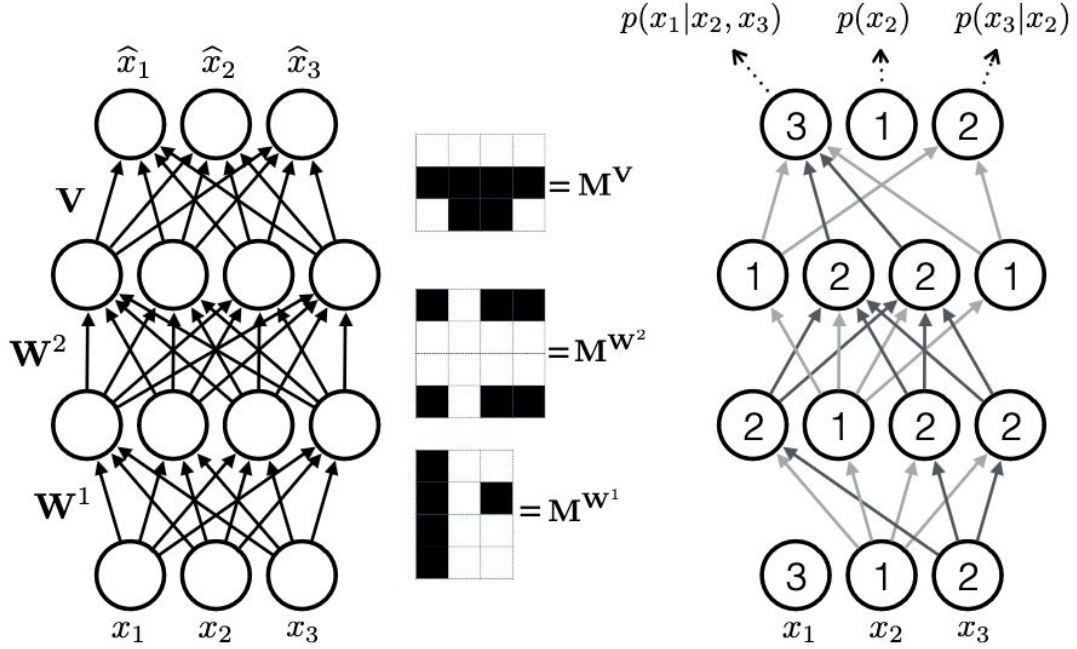
An autoregressive network uses the previous inputs to estimate the density



$$\rho(\mathbf{x}) = \rho_1(x_1) \rho_2(x_2 | x_1) \rho_3(x_3 | x_1, x_2)$$

Masked Autoencoders for Density Estimation (MADE)

$$\hat{x}_d = p(x_d | \mathbf{x}_{<d})$$



[Video Tutorial](#)
[Germain+2015](#) MADE paper
[Papamakarios+2018](#) MAF paper

Autoencoder × **Masks** → **MADE**

Slide credit: Jason Poh

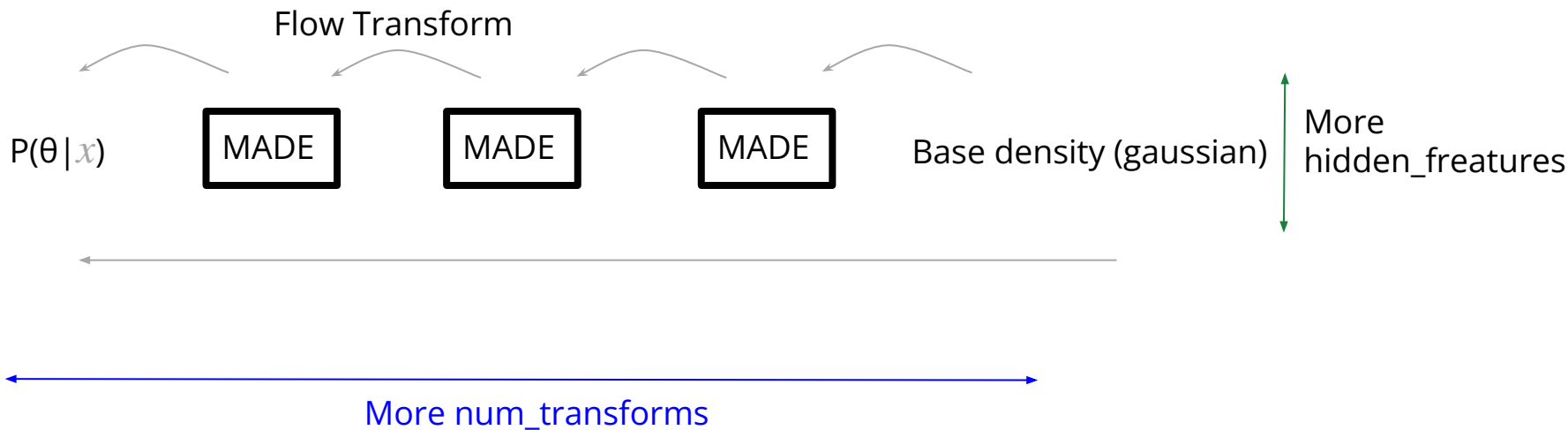
$$p(\mathbf{x}) = \pi_u(f^{-1}(\mathbf{x})) \left| \det \left(\frac{\partial f^{-1}}{\partial \mathbf{x}} \right) \right|$$

$$u_i = (x_i - \mu_i) \exp(-\alpha_i) \quad \text{where} \quad \mu_i = f_{\mu_i}(\mathbf{x}_{1:i-1}) \quad \text{and} \quad \alpha_i = f_{\alpha_i}(\mathbf{x}_{1:i-1}).$$

$$\left| \det \left(\frac{\partial f^{-1}}{\partial \mathbf{x}} \right) \right| = \exp \left(-\sum_i \alpha_i \right) \quad \text{where} \quad \alpha_i = f_{\alpha_i}(\mathbf{x}_{1:i-1}).$$

Masked Autoregressive Flows are stacks of MADEs

NPE is constructed from stacks of MADE that are flow transformed.



Slide credit: Jason Poh

Embedding network to learn summary statistics

Simulator
Input

parameters

Simulator
output

images

Embedding
network input

images

Embedding
network output

Summary stats

Slide credit: Jason Poh

Overall trains to learn embedding network and NPE

Training an NPE + embedding



Simulator
Input

Simulator
output

Embedding
network input

Embedding
network output

NPE input

NPE output

parameters

images

images

Summary stats

[Summary stats] $p(\text{param} | \text{data})$
+
[parameter set]

Slide credit: Jason Poh

Masked Autoregressive Flows

There are mainly two families of neural density estimators that are both flexible and tractable:

- Normalizing flows
 - Transform a base density into target density by invertible transformation with tractable jacobian.
- Autoregressive models
 - Decomposes target density as a product of conditionals and models each conditional in turn.

This method is **both** - it is a normalizing flow of autoregressive models (Masked Autoencoders for Density Estimation (MADEs to be exact))

Neural Density Estimators

- An estimate of the exact posterior $p(\theta|x)$ can be learnt from neural density estimators $q_\phi(\theta|x)$
- Maximize the average log probability $\frac{1}{N} \sum_n \log q_\phi(\theta_n | x_n)$

I think I'd like to get more into NPE because that's what we'll do in the tutorial. I'm planning on talking about normalizing flows and I'll probably mention autoregressive flows, but not focus on them. I understand that a normalizing flow can estimate the likelihood, since we're sampling from $x \sim p(x|\theta)$, but I'm struggling with how to explain that it does this also for the posterior. I understand that it's accounting for the prior via drawing from it, but how is it accounting for the evidence?

Take a break



sbi : A toolkit for simulation-based inference

Alvaro Tejero-Cantero^{e, 1}, Jan Boelts^{e, 1}, Michael Deistler^{e, 1},
Jan-Matthis Lueckmann^{e, 1}, Conor Durkan^{e, 2}, Pedro J. Gonçalves^{1, 3},
David S. Greenberg^{1, 4}, and Jakob H. Macke^{1, 5, 6}

<https://joss.theoj.org/papers/10.21105/joss.02505>

<https://github.com/sbi-dev/sbi>

```
from sbi.inference import infer
# import your simulator, define your prior over the parameters
parameter_posterior = infer(simulator, prior, method='SNPE', num_simulations=100)
```

Tutorial

- Linefit example:

<https://colab.research.google.com/drive/1CRBQqSim3KZV6s5hcwz-zMVmbJzTeaWz?usp=sharing>

- Lenstronomy example:

<https://colab.research.google.com/drive/1NpwdTy98lfo-vPul5Rt-HTFTalgP3cZH?usp=sharing>

Flavors/algorithms of SNPE

SNPE-A: only allows Gaussian mixtures and mixture density models for modeling density

SNPE-B: more flexible but has some technical issues that are documented on github, having to do with the proposal distribution

SNPE-C: current state-of-the-art

Challenges / the future of SBI / State of the art / things that should keep you up at night



Challenges / the future of SBI / State of the art / things that should keep you up at night

General papers: <https://simulation-based-inference.org/papers/>

Applying to different areas for dramatic speed-up in inference:
[Khullar+2022](#)

Challenges / the future of SBI / State of the art / things that should keep you up at night

General papers: <https://simulation-based-inference.org/papers/>

Applying to different areas for dramatic speed-up in inference: [Khullar+2022](#)

Differentiable simulators: [Zhegal+2022](#)

Calibration / a “crisis” in SBI?: [Hermans+2021](#)

SBI with many variables: [Poh+2022](#)

Guarantees on error?

Hierarchical

Challenges / the future of SBI / State of the art / things that should keep you up at night

General papers: <https://simulation-based-inference.org/papers/>

Applying to different areas for dramatic speed-up in inference: [Khullar+2022](#)

Differentiable simulators: [Zhegal+2022](#)

Calibration / a “crisis” in SBI?: [Hermans+2021](#)

SBI with many variables: [Poh+2022](#)

Guarantees on error?

Hierarchical

Graph neural nets / other architectures

Understanding really complicated simulations?

Domain adaptation!

The Deepskies Lab SBI team



Jason Poh



Sreevani Jarugula



Amanda Pagul



Brian Nord



Gourav Khullar



Egor Danilov



Humna Aman



Anvi Padhi



Becky Nevin



Andrea Roncoli



Aleksandra Ćiprijanović



Marcos Tamargo



Natalie Malgon



Kabelo Tsiane



Yuanyuan Zhang



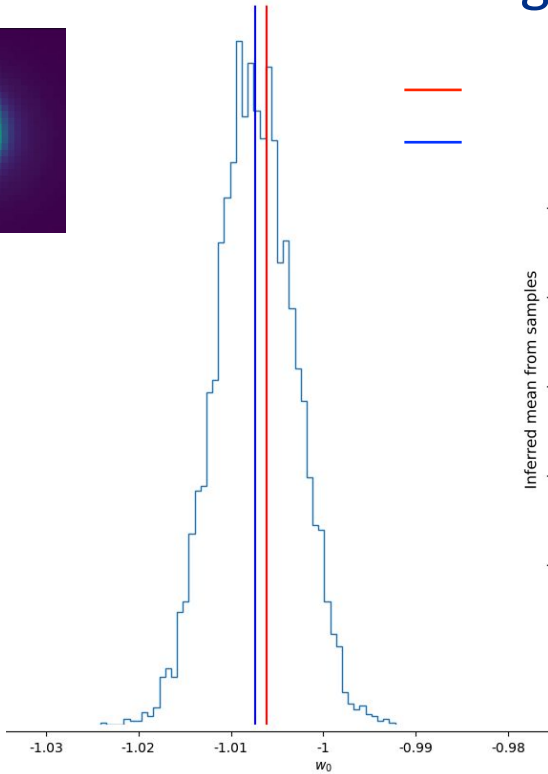
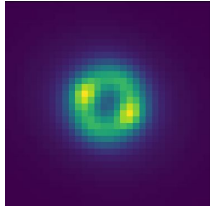
Moonzarin Reza



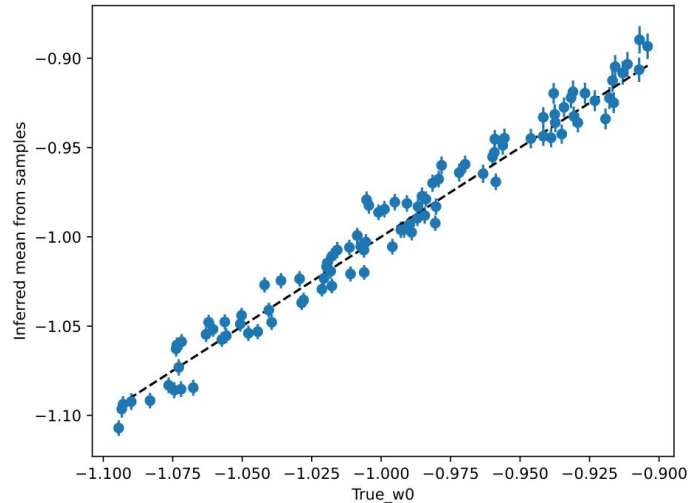
Neil Kumar

Dark Energy equation of state parameter

Inference from a single lens



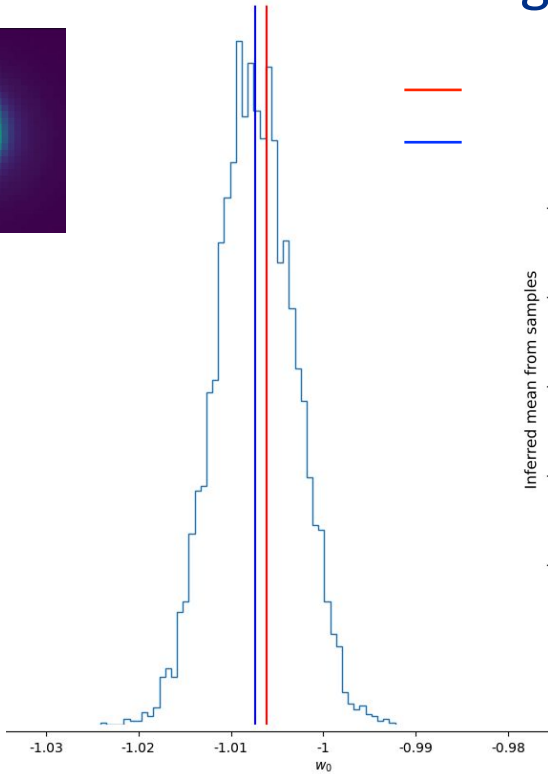
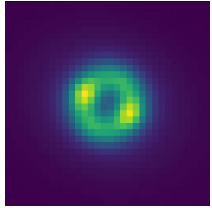
Inference from 100 lens



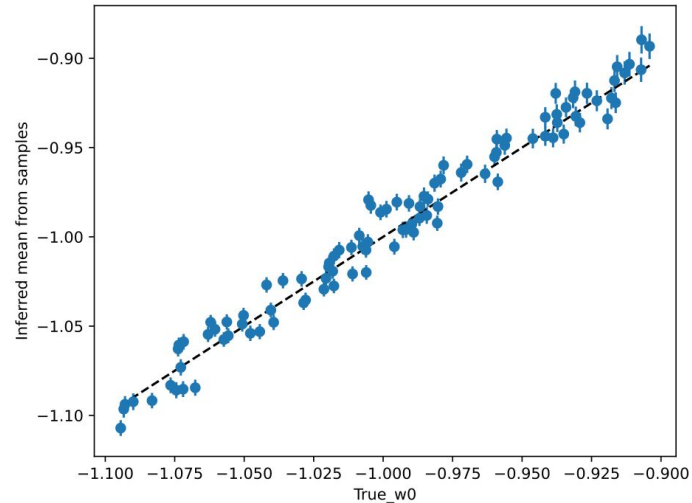
Sreevani Jarugula

Dark Energy equation of state parameter

Inference from a single lens



Inference from 100 lens



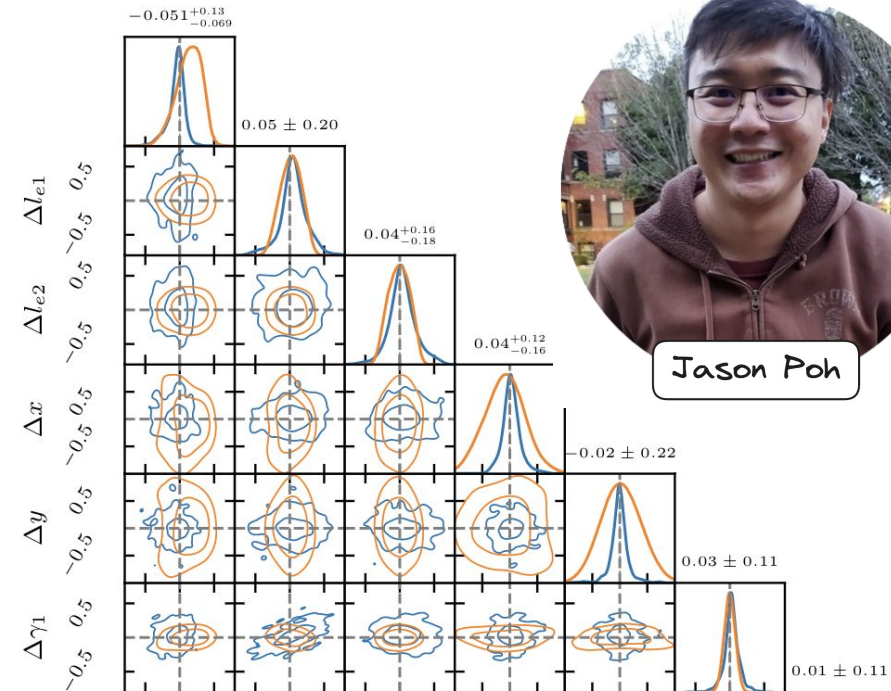
Sreevani Jarugula

Next steps:

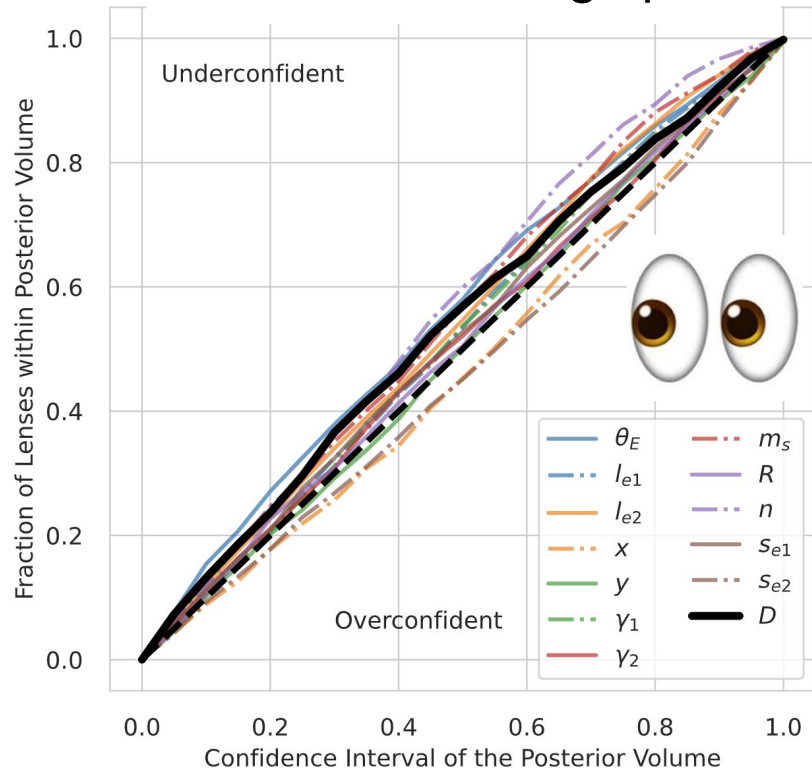
- astrophysics + cosmology
- Joint inference

SBI to infer astrophysics parameters from strong lensing

Jason Poh+ and DeepSkies

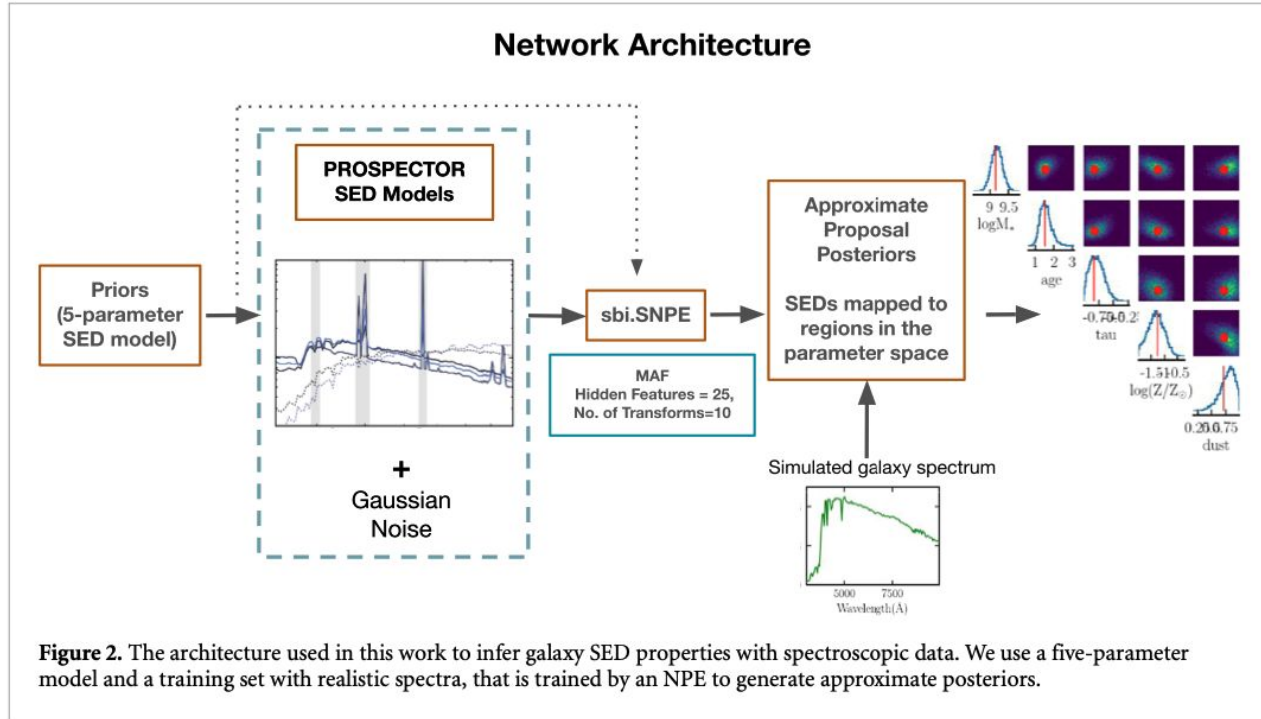


Posterior coverage plot!



SBI to infer galaxy properties from spectra

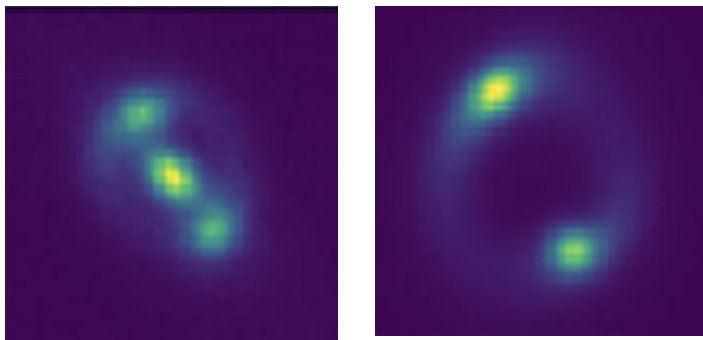
Khullar et al. 2022 and DeepSkies



Domain adaptation to improve the performance of SBI on real observational lenses with noise!

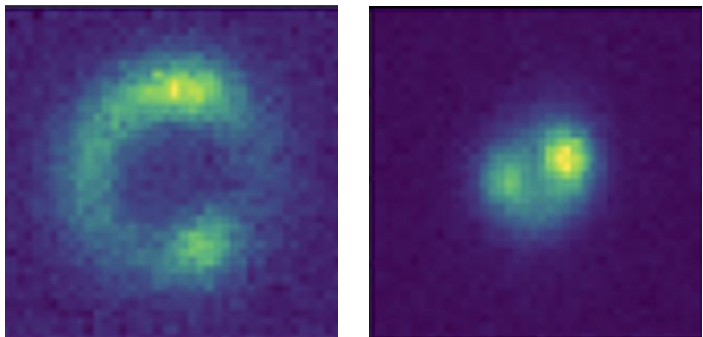
Source domain:

low noise
simulated
lenses



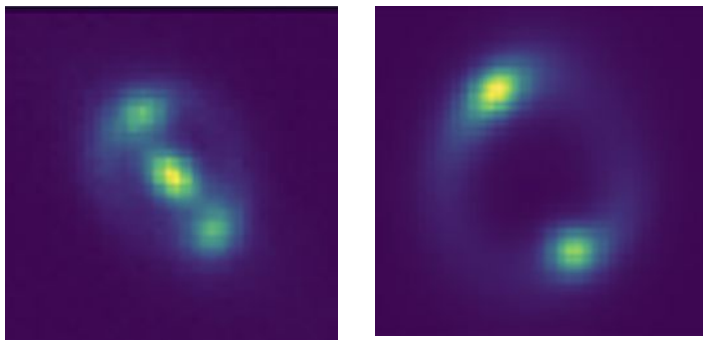
Target domain:

DES noise
simulated
lenses

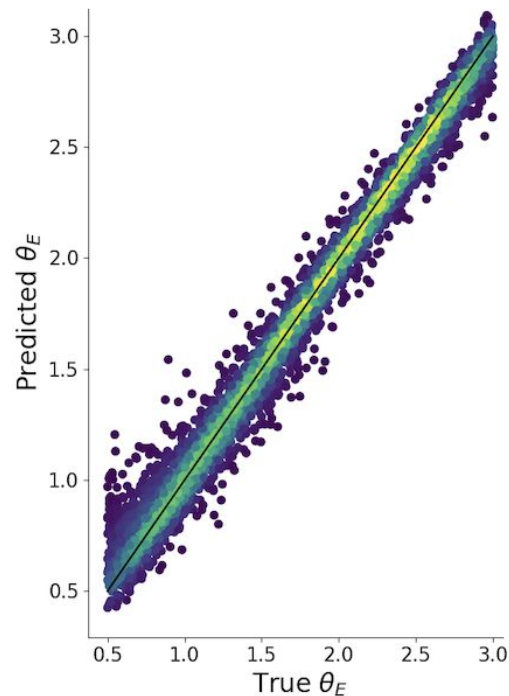
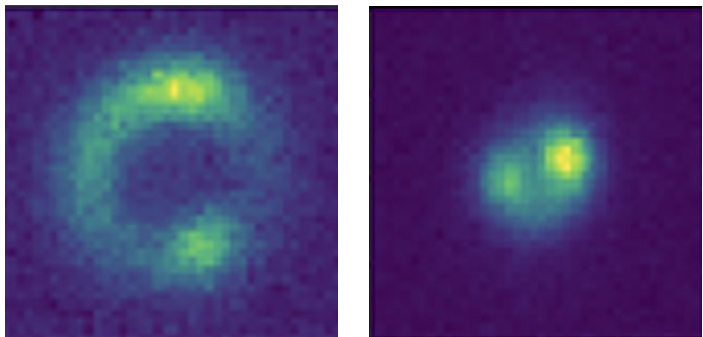


Domain adaptation to improve the performance of SBI on real observational lenses with noise!

Source domain:
low noise
simulated
lenses



Target domain:
DES noise
simulated
lenses

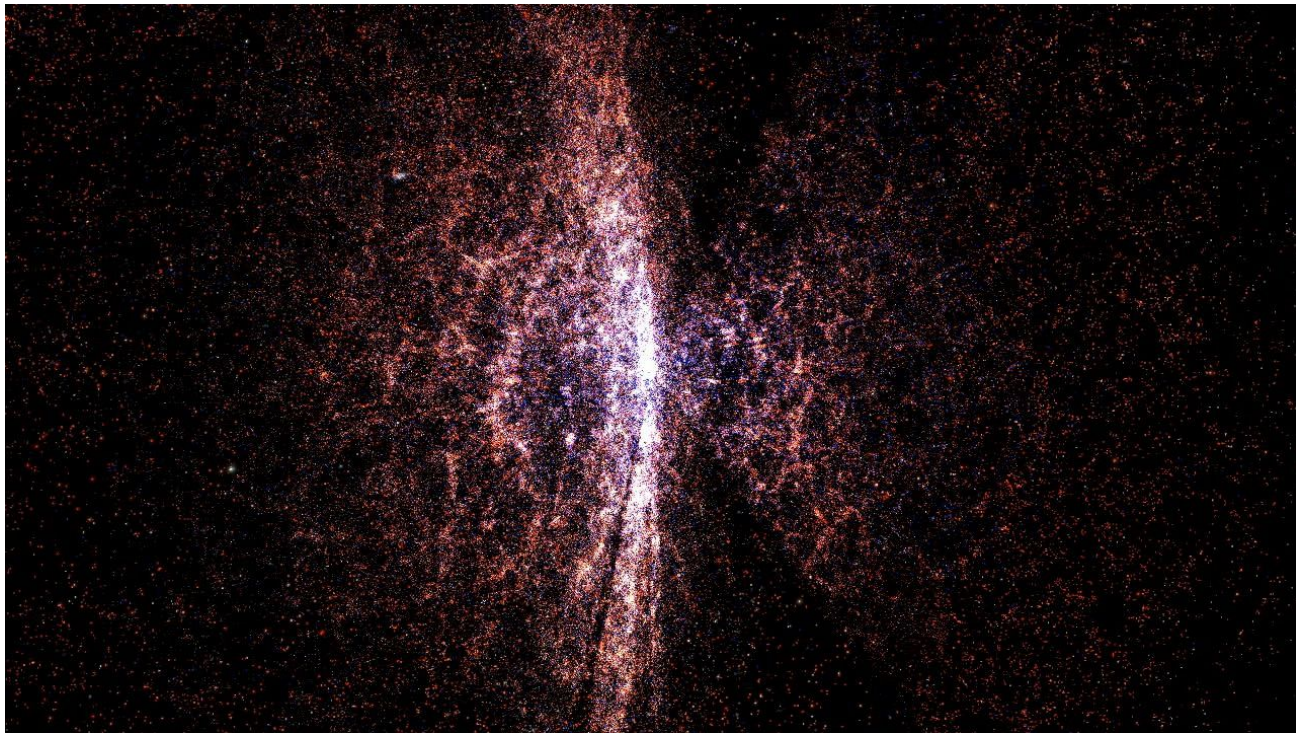


SBI for alternate networks, like graph neural networks



Anvi Padhi

Goal: Infer key cosmological parameters σ_8 and Ω_{matter}



If you take anything from this let it be:

- Simulation-based Inference is an alternative to likelihood-based inference
 - Does not require the computation of a likelihood
 - Machine learning methods are amortized, i.e. can evaluate posterior from new data without retraining the model
 - Computationally more efficient than MCMC based methods



If you take anything from this let it be:

- Simulation-based Inference is an alternative to likelihood-based inference
- Lots of options for what type of probability you want to target (likelihood, likelihood ratio, posterior)
- The future is exciting, lots of research here into expanding into new data types, looking at uncertainty guarantees, lots of opportunities, especially for those interested in explainability, ie the HEP community

Join the Deepskies Lab!



<https://deepskieslab.com/>

[Google form to join.](#)

The end



If you take anything from this let it be:

- Simulation-based Inference is an alternative to likelihood-based inference
- Lots of options for what type of probability you want to target (likelihood, likelihood ratio, posterior)
- The future is exciting, lots of research here into expanding into new data types, looking at uncertainty guarantees, lots of opportunities, especially for those interested in explainability, ie the HEP community



EXTRAS

Notes from Jason's SBI session

- Where to not use SBI - if your simulator is not good, doesn't capture actual physics?
 - Activation energy of it being worth it → if MCMC is quick enough don't need to use it, it has additional upfront cost,
- How complex model can be before you switch?
 - Sam: curse of dimensionality, SBI can do 20 or so parameters, also embedding network
 - Complexity of likelihood versus complexity of parameters

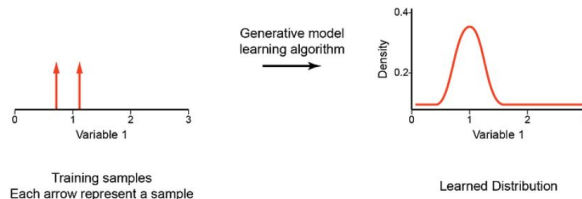
Density Estimators (Thanks to Jason!)

- Unsupervised method of getting structure from data:
 - Given \mathbf{x} , what is $p(\mathbf{x})$?
 - Given $p(\mathbf{x})$, what is \mathbf{x} ?
- Gaussian Density Model

Model: $q_{\phi}(\mathbf{x}) = \frac{1}{|\det(2\pi\Sigma)|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$ where $\phi = \{\boldsymbol{\mu}, \Sigma\}$.

Training: Parametric density models are typically estimated by maximum likelihood. Given a set of training datapoints $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ that have been independently and identically generated by a process with density $p(\mathbf{x})$, we seek a setting of the model's parameters ϕ that maximize the average log likelihood on the training data:

$$L(\phi) = \frac{1}{N} \sum_n \log q_{\phi}(\mathbf{x}_n). \quad (2.13)$$



From [Papamakarios \(2019\)](#)

- Neural Networks to parameterize density model

$$L(\phi) = \frac{1}{N} \sum_n \log q_\phi(\mathbf{x}_n) = \frac{1}{N} \sum_n f_\phi(\mathbf{x}_n).$$

- The parameters of the Neural network are updated through gradient descent

$$\nabla_\phi \hat{L}(\phi) = \frac{1}{M} \sum_m \nabla_\phi f_\phi(\mathbf{x}_{n_m}).$$

- SBI has 4 built-in density estimators:
 - Masked Autoregressive Flow (MAF)
 - Neural Spline Flow (NSF)
 - Masked Autoencoder for Distribution Estimation (MADE)
 - Mixture Density Network (MDN)