

Hands-on Simulation Based Inference

David Prelogović

Setup

- Go to <https://github.com/dprelogo/SBI-tutorial>
- Follow installation instructions

What is Inference?

Bayesian inference, obviously



Bayesian inference

- How to properly learn from the new data?
- Update prior knowledge about the parameters $P(\theta)$ by including data \mathcal{D}

$$\text{\color{red}posterior } P(\theta|\mathcal{D}) = \frac{\text{\color{red}likelihood} \quad \text{\color{red}prior}}{\text{\color{red}evidence}} \frac{P(\mathcal{D}|\theta) \cdot P(\theta)}{P(\mathcal{D})}$$

Bayesian inference

- How to properly learn from the new data?
- Update prior knowledge about the parameters $P(\theta)$ by including data \mathcal{D}

Good “physically-motivated” model

$$posterior \ P(\theta|\mathcal{D}) = \frac{\underset{likelihood}{P(\mathcal{D}|\theta)} \cdot \underset{prior}{P(\theta)}}{\underset{evidence}{P(\mathcal{D})}}$$

= good prior

Bayesian inference

- How to properly learn from the new data?
 - Update prior knowledge about the parameters $P(\theta)$ by including data \mathcal{D}
- Good “physically-motivated” model
=

$$posterior \ P(\theta|\mathcal{D}) = \frac{\underset{likelihood}{P(\mathcal{D}|\theta)} \cdot \underset{prior}{P(\theta)}}{\underset{evidence}{P(\mathcal{D})}}$$

good prior
mostly ignored
(except for model selection)

Bayesian inference

- How to properly learn from the new data?
- Update prior knowledge about the parameters $P(\theta)$ by including data \mathcal{D}

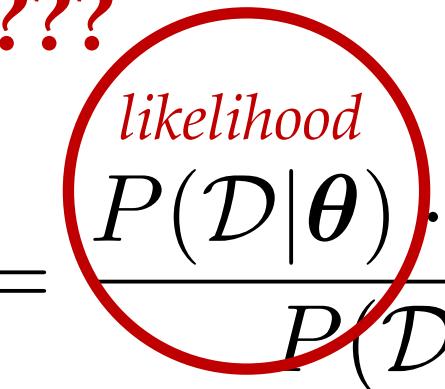
$$posterior \ P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta) \cdot P(\theta)}{P(\mathcal{D})}$$

Good “physically-motivated” model
=
good prior

mostly ignored
(except for model selection)

???

likelihood *prior*
evidence



Bayesian inference

- How to properly learn from the new data?
- Update prior knowledge about the parameters $P(\theta)$ by including data \mathcal{D}

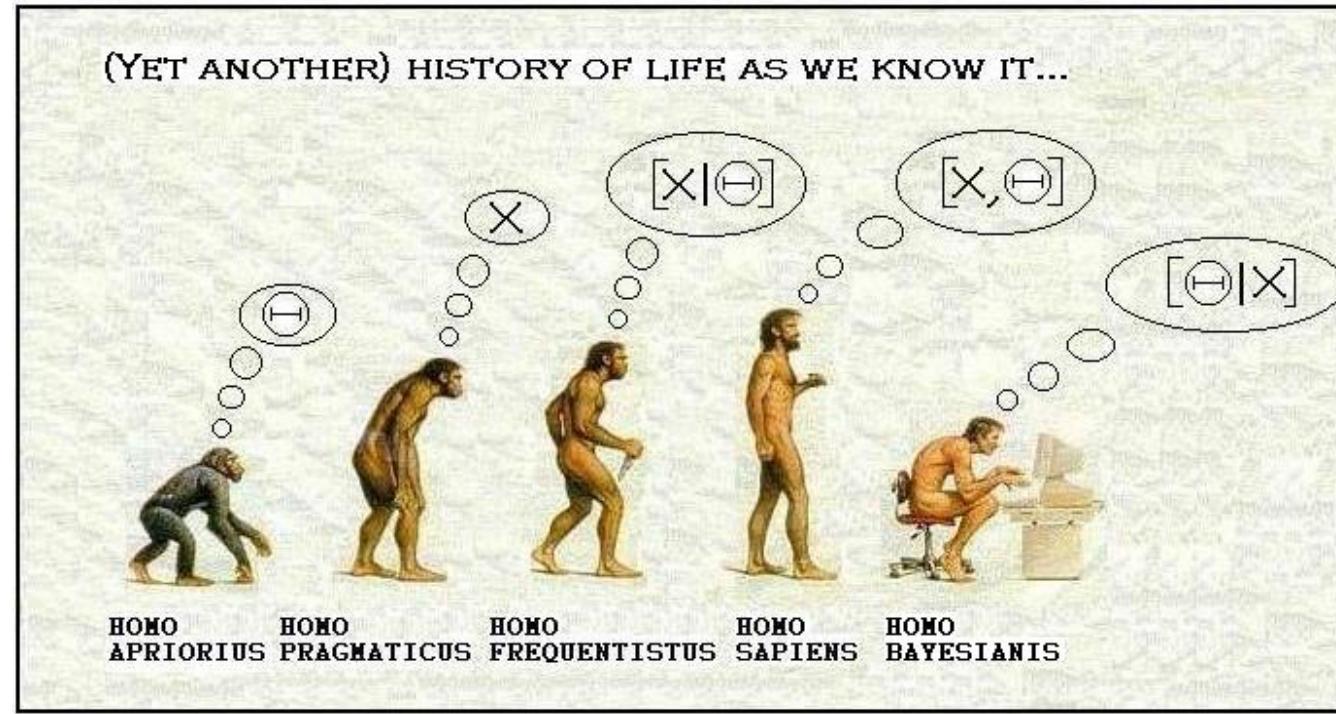
$$posterior \ P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta) \cdot P(\theta)}{P(\mathcal{D})}$$

??? **Good “physically-motivated” model**
likelihood *prior*
evidence **good prior**

mostly ignored
(except for model selection)

- Very often: data compressed to a set of summaries $\mathcal{D} \xrightarrow{f} \mathcal{S}$

What is Classical inference?



Classical inference

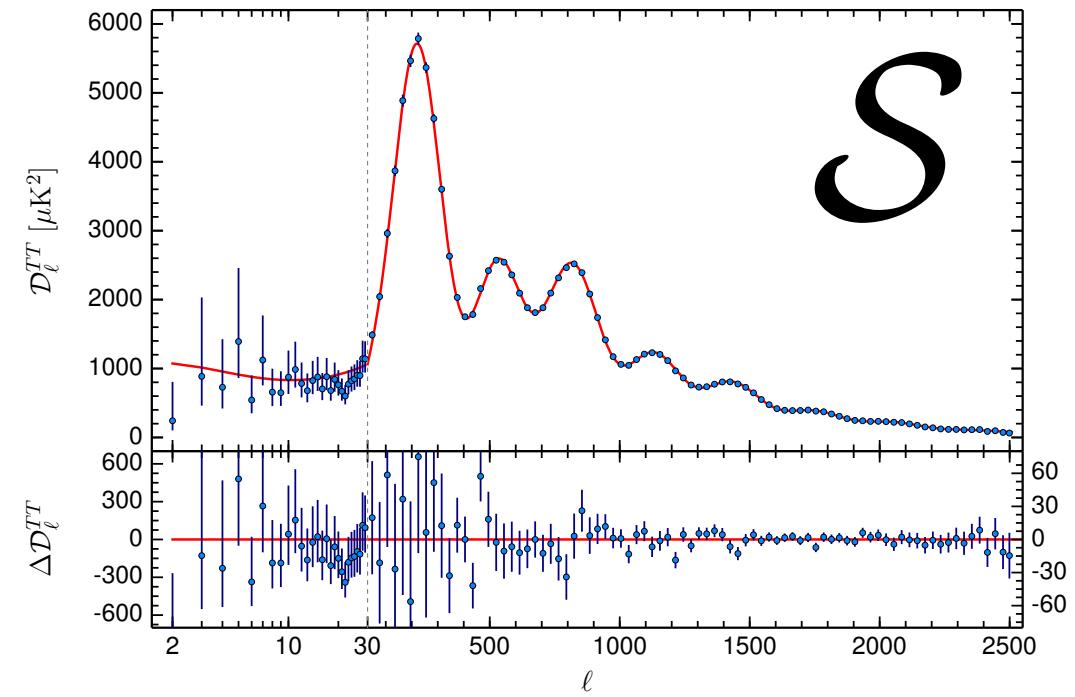
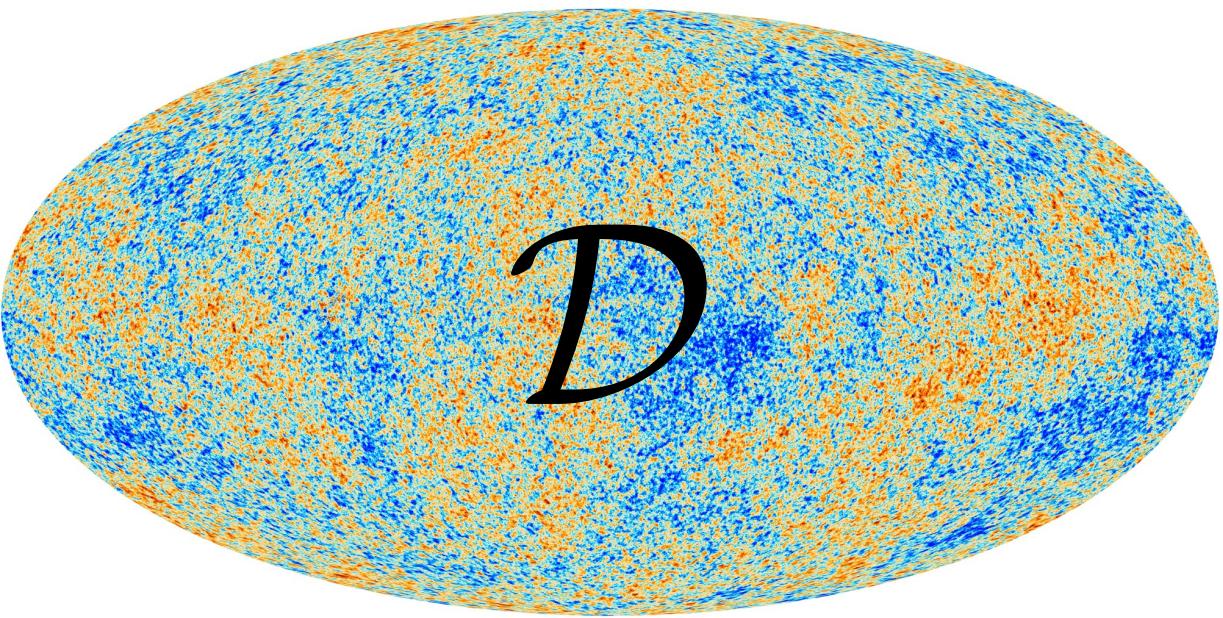
- Knowing your likelihood
- Eg.

$$P(\mathcal{D}|\boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma(\boldsymbol{\theta})|}} e^{-\frac{1}{2}(\mathbf{d}-\boldsymbol{\mu}(\boldsymbol{\theta}))^T \Sigma(\boldsymbol{\theta})^{-1} (\mathbf{d}-\boldsymbol{\mu}(\boldsymbol{\theta}))}$$

$$P(\mathcal{D}|\boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2} \sqrt{\prod_i \sigma_i(\boldsymbol{\theta})^2}} e^{-\frac{1}{2} \sum_i (x_i - \mu_i(\boldsymbol{\theta}))^2 / \sigma_i(\boldsymbol{\theta})^2}$$

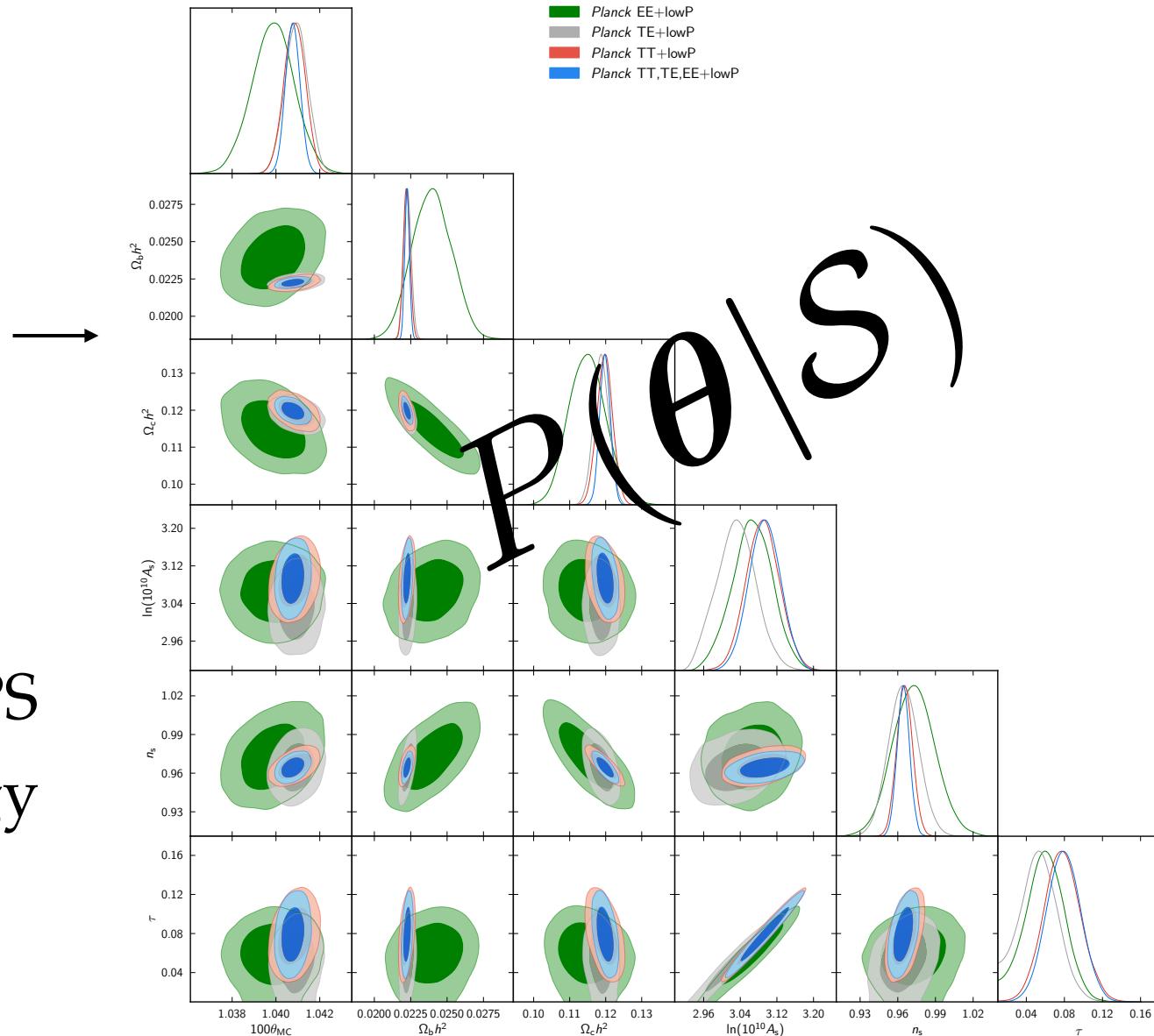
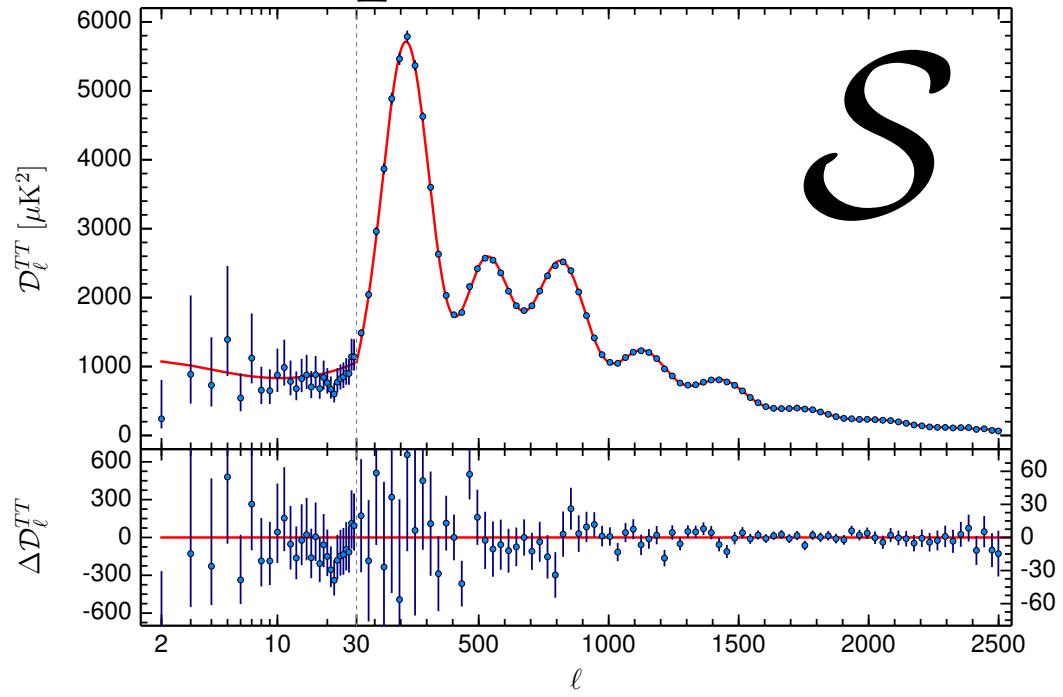
-> i.e. χ^2

Example: CMB



- Full sky map compressed to PS

Example: CMB

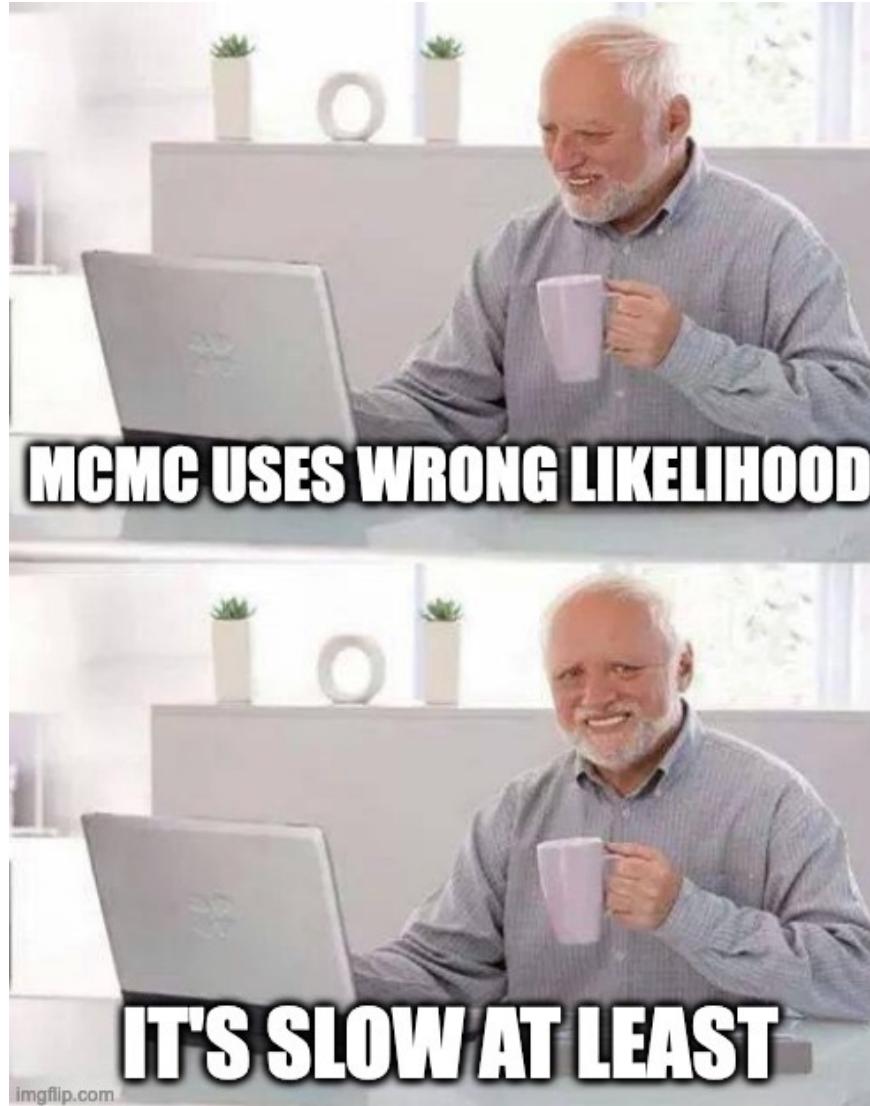


- Full sky map compressed to PS
- From it we infer the cosmology
 - Known compression, known likelihood

How do we construct $P(\theta|\mathcal{S})$

- <https://chi-feng.github.io/mcmc-demo/app.html>

Why isn't classical inference enough?



Why isn't classical inference enough?

- Likelihood is not analytic or not tractable?
- Parameter space is huge?
- Maybe there's only few important parameters, a lot of nuisance ones?
- New data comes in, should I run expensive inference from scratch?

SBI

Classical inference

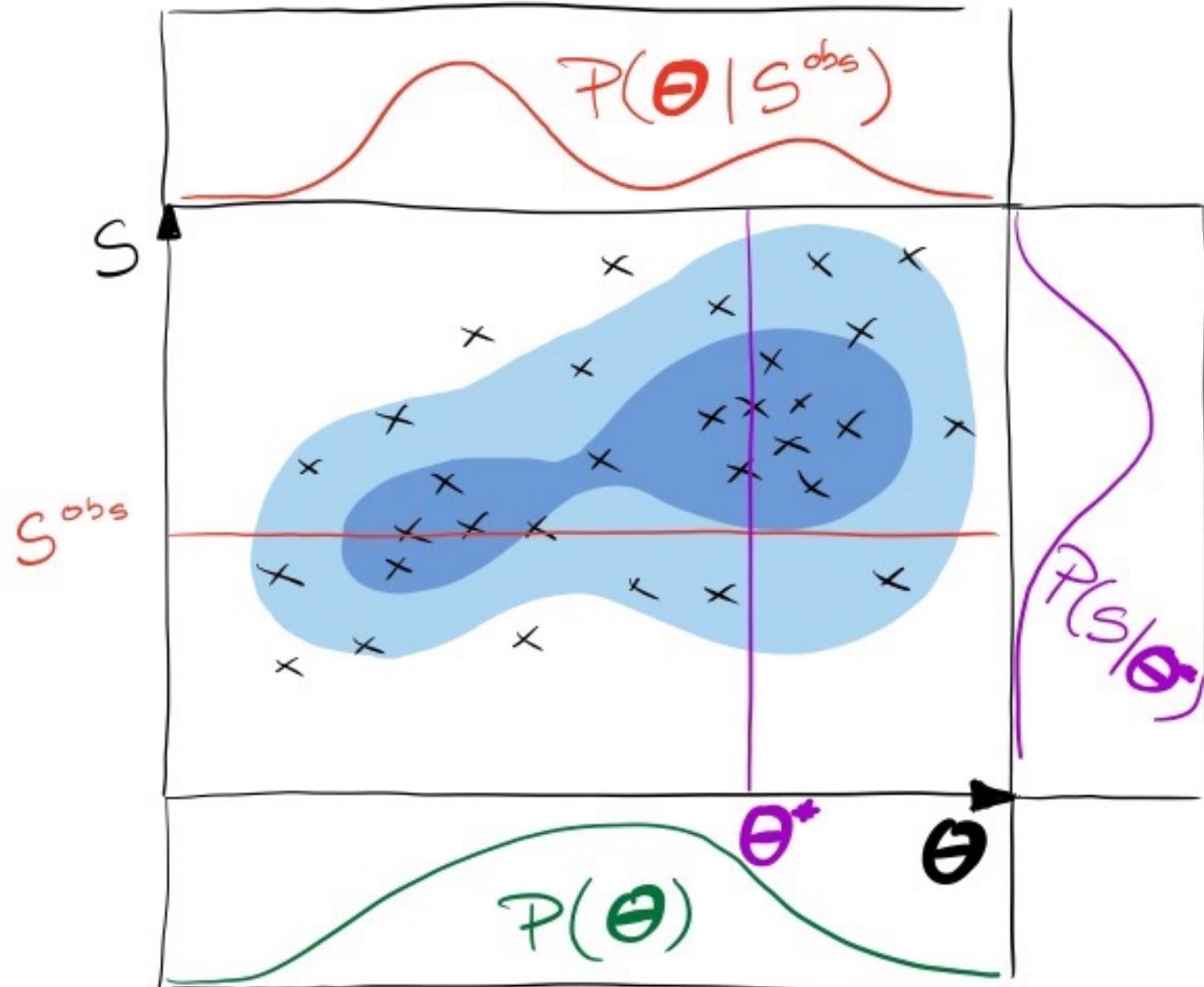


imgflip.com

What is
Simulation Based Inference?

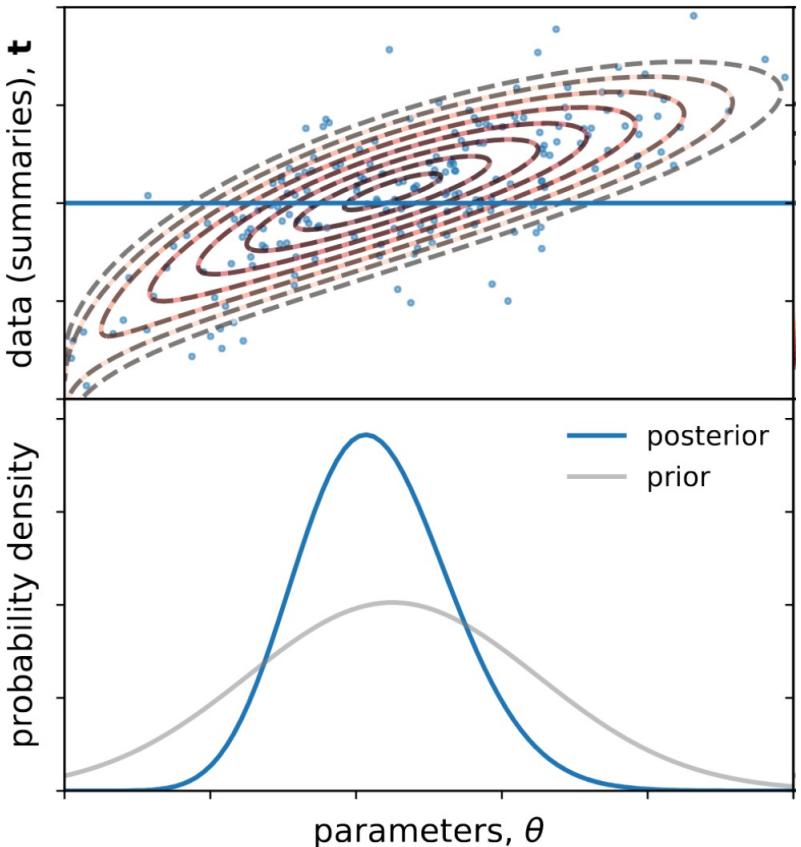
Simulation-Based Inference (SBI)

$$P(\theta|\mathcal{D}) = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}} P(\mathcal{D}|\theta) \cdot P(\theta)$$



Credit: Tom Charnock

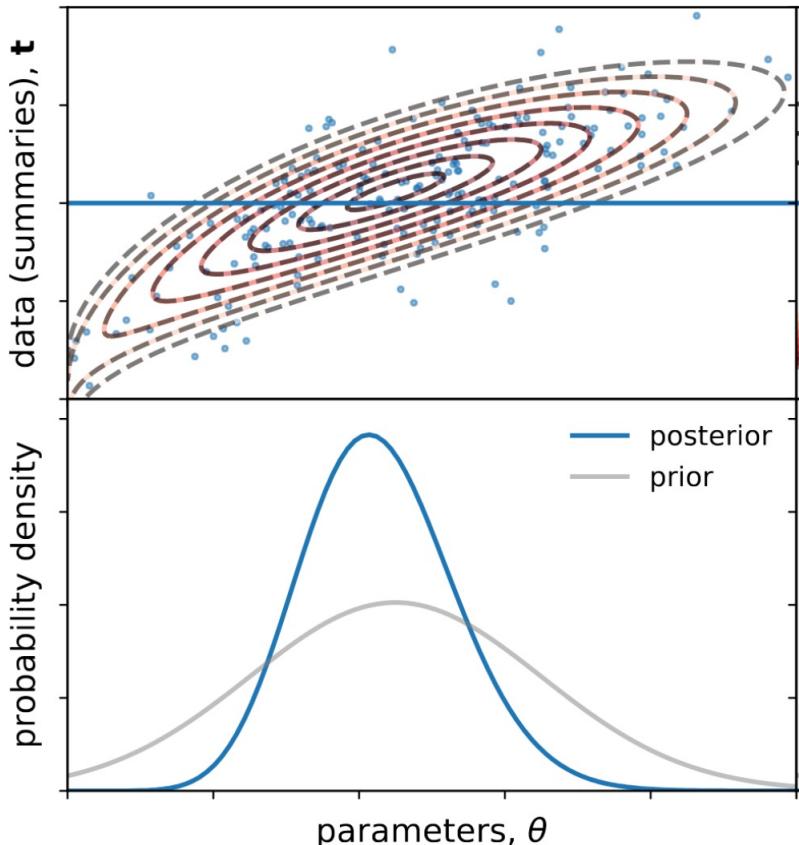
Simulation-Based Inference (SBI)



- Pull from the prior $\tilde{\theta} \sim p(\theta)$
- Pull from the likelihood $\tilde{t} \sim p(t|\tilde{\theta})$
 - Use simulator to simulate the data
 $\tilde{d} = \text{simulator}(\tilde{\theta})$
 - Optionally compress the data to summary
 $\tilde{t} = \text{compressor}(\tilde{d})$
- Repeat many times to construct a set
$$\{(\tilde{t}_1, \tilde{\theta}_1), (\tilde{t}_2, \tilde{\theta}_2), \dots, (\tilde{t}_N, \tilde{\theta}_N)\}$$
- **It is assumed the simulator is perfect, i.e. it pulls from the actual likelihood**

$$p(d|\tilde{\theta})$$

Approximate Bayesian Computation (ABC)



- Define an “ ε -ball” in the data-space

$$B_\epsilon(\mathbf{x}_o) = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon\}$$

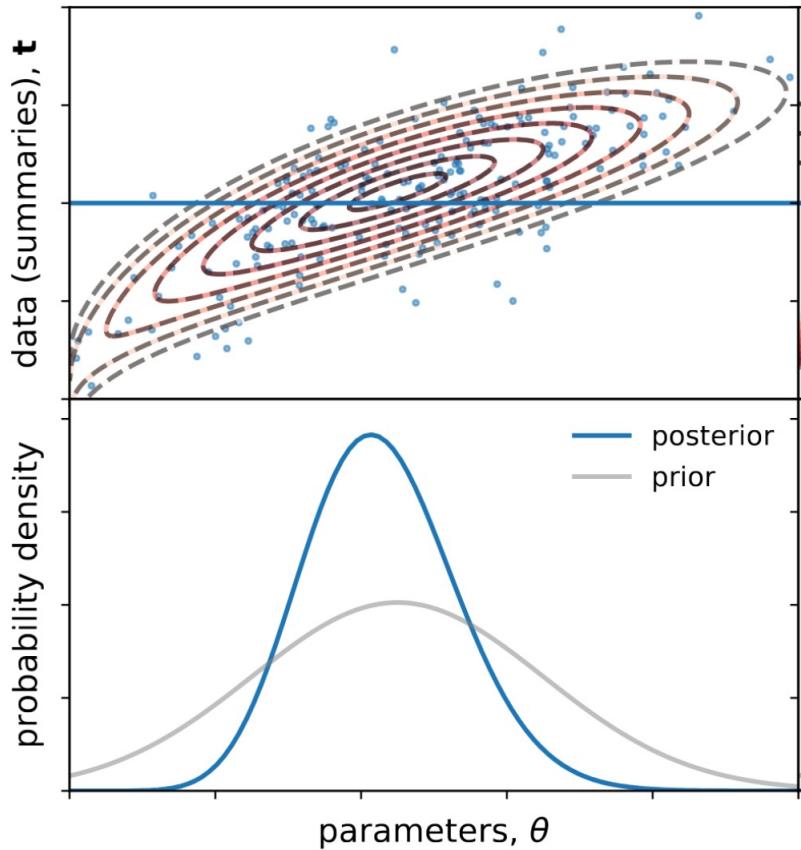
- Likelihood is then approximated by

$$p(\mathbf{x}_o | \boldsymbol{\theta}) \approx \frac{\Pr(\|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon | \boldsymbol{\theta})}{|B_\epsilon(\mathbf{x}_o)|}$$

- Inverting with Bayes rule gives:

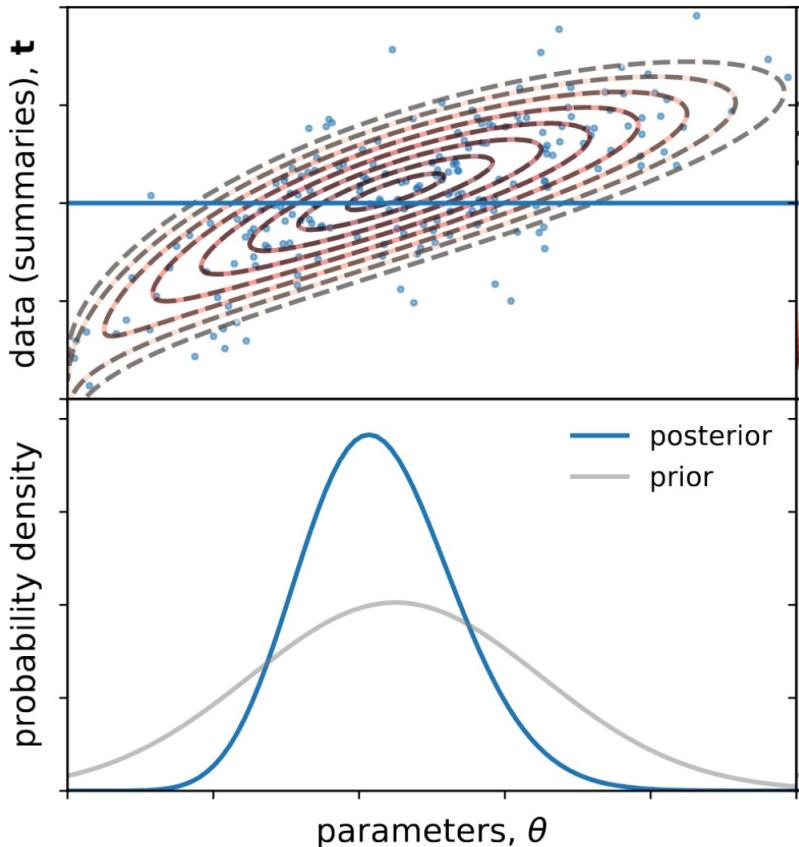
$$p(\boldsymbol{\theta} | \mathbf{x} = \mathbf{x}_o) \approx \frac{\Pr(\|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int \Pr(\|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon | \boldsymbol{\theta}') p(\boldsymbol{\theta}') d\boldsymbol{\theta}'} = p(\boldsymbol{\theta} | \|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon)$$

Approximate Bayesian Computation (ABC)

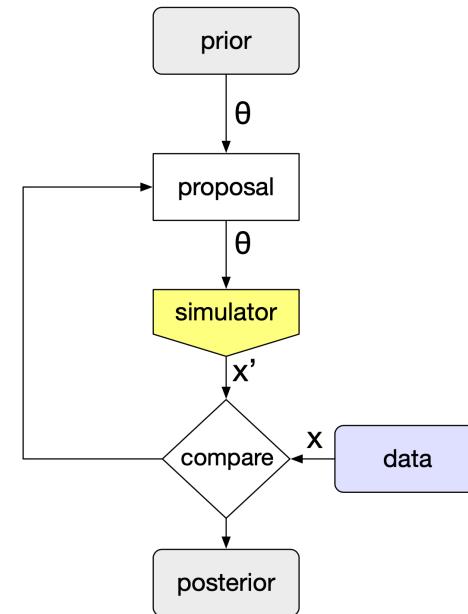


Algorithm	Rejection ABC
repeat	
Sample $\theta \sim p(\theta)$	
Simulate $\mathbf{x} \sim p(\mathbf{x} \theta)$	
until $\ \mathbf{x} - \mathbf{x}_o\ \leq \epsilon$	
return θ	

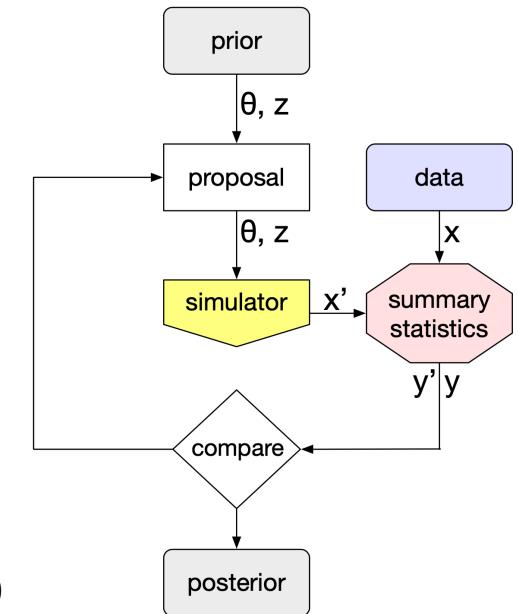
Approximate Bayesian Computation (ABC)



Approximate Bayesian Computation
with Monte Carlo sampling



Approximate Bayesian Computation
with learned summary statistics

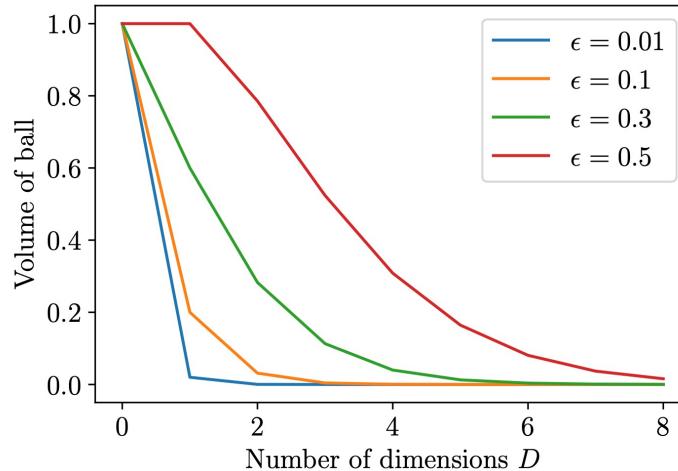


- For new data the whole procedure has to be repeated
- Curse of dimensionality

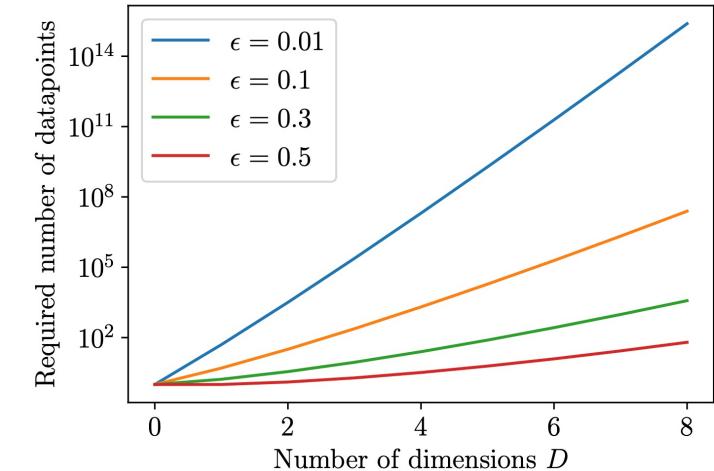
ABC & curse of dimensionality

- For an ϵ -cube around the data point:

$$\Pr(\|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon | \boldsymbol{\theta}) \approx p(\mathbf{x}_o | \boldsymbol{\theta}) |B_\epsilon(\mathbf{x}_o)| = p(\mathbf{x}_o | \boldsymbol{\theta}) (2\epsilon)^D$$



(a) Volume of a D -dimensional ball of radius ϵ .



(b) Expected number of datapoints we need to generate until one falls in the ball.

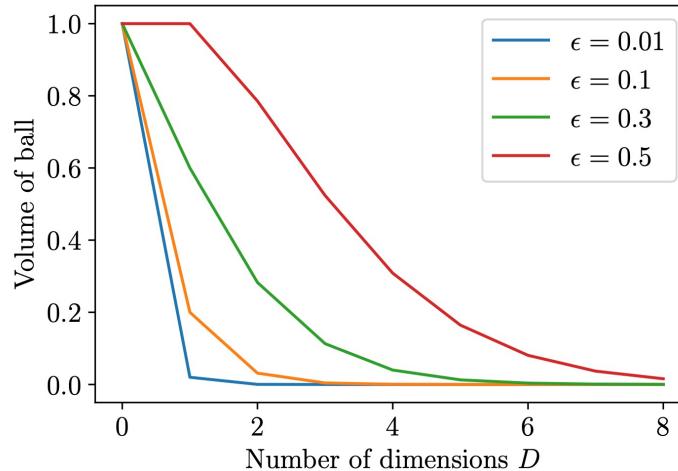
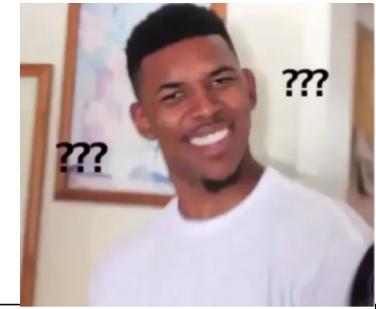
Solutions:

1. Summary statistics
2. High ϵ = lower-quality posterior
3. Smarter samplers (MCMC ABC , Sequential MC ABC)

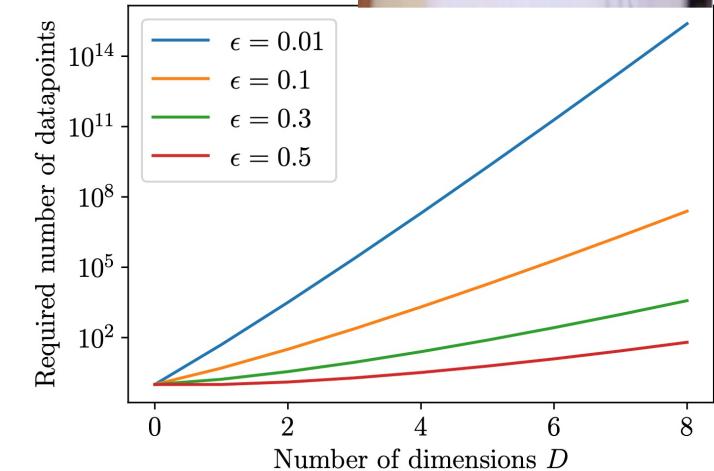
ABC & curse of dimensionality

- For an ϵ -cube around the data point:

$$\Pr(\|\mathbf{x} - \mathbf{x}_o\| \leq \epsilon | \boldsymbol{\theta}) \approx p(\mathbf{x}_o | \boldsymbol{\theta}) |B_\epsilon(\mathbf{x}_o)| = p(\mathbf{x}_o | \boldsymbol{\theta}) (2\epsilon)^D$$



(a) Volume of a D -dimensional ball of radius ϵ .

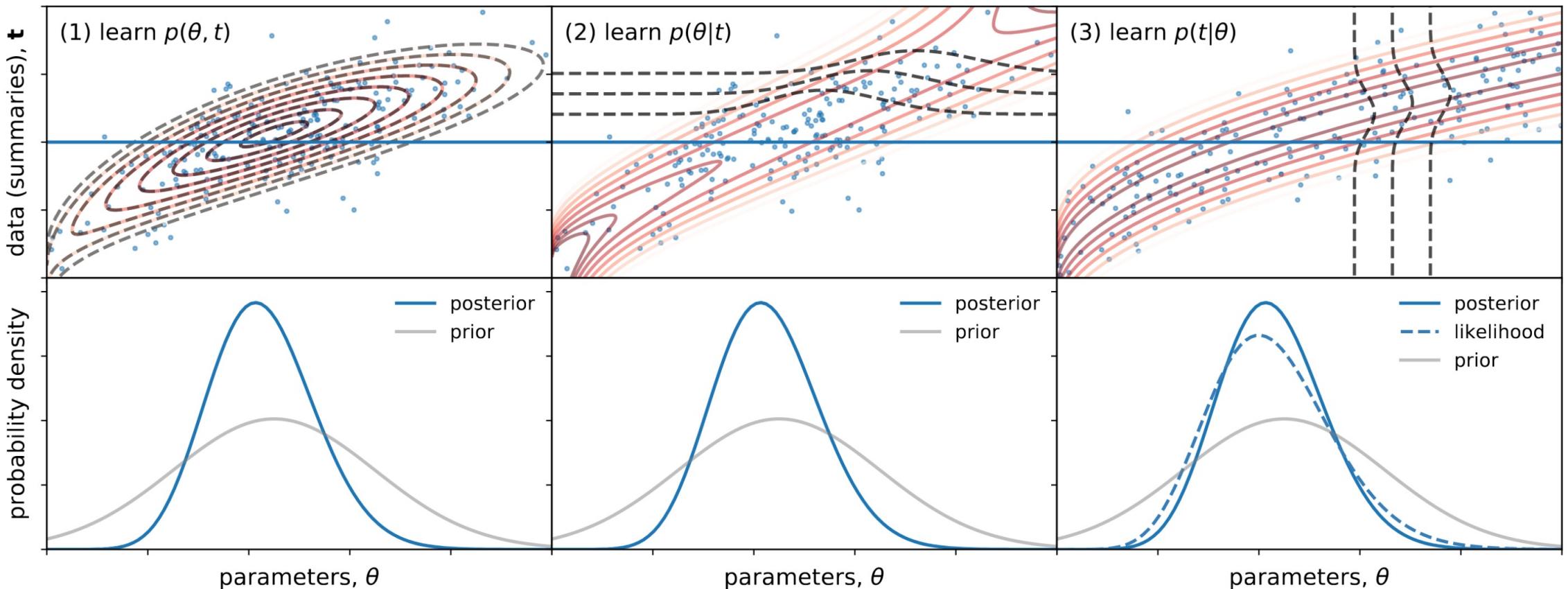


(b) Expected number of datapoints we need to generate until one falls in the ball.

Solutions:

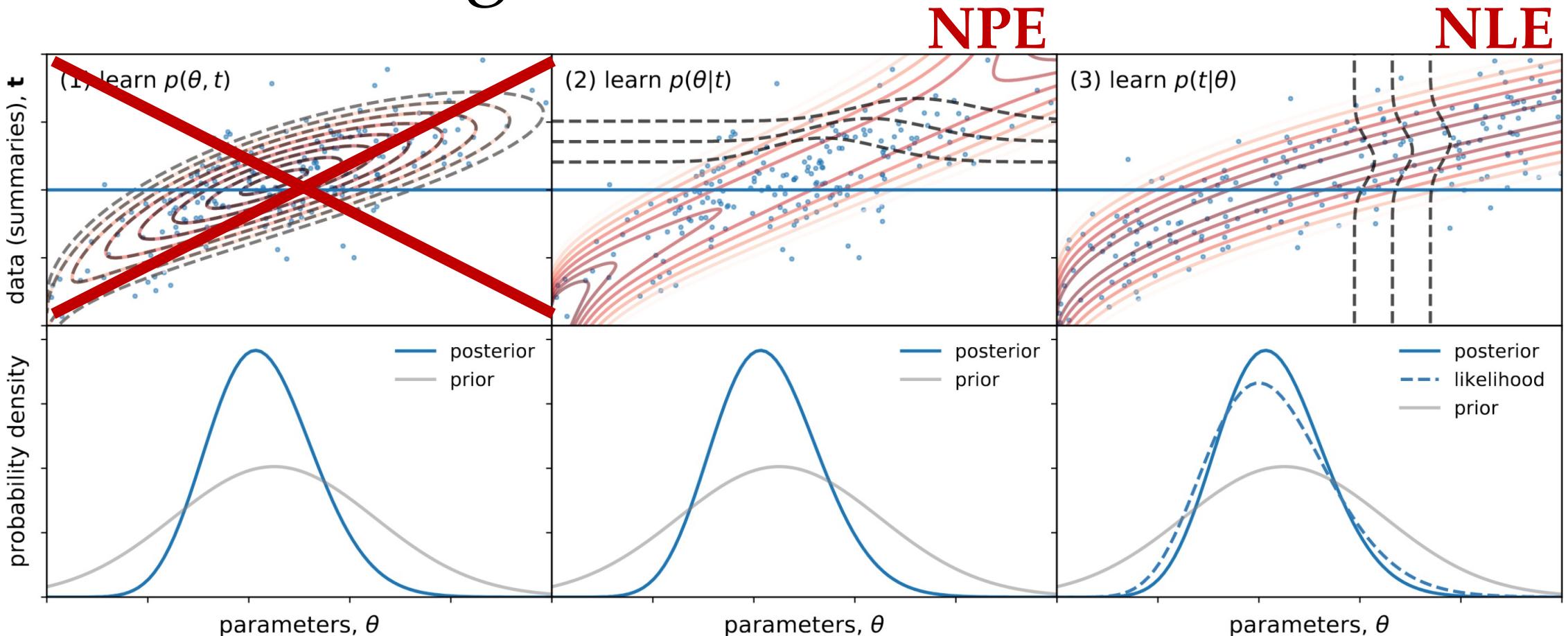
1. Summary statistics
2. High ϵ = lower-quality posterior
3. Smarter samplers (MCMC ABC , Sequential MC ABC)

SBI - Learning the distribution



- Or: learn likelihood-evidence ratio $p(t|\theta)/p(t)$
→ *(Neural) density estimators*

SBI - Learning the distribution



- Or: learn likelihood-evidence ratio $p(t|\theta)/p(t)$ **NRE**
→ *(Neural) density estimators*

NPE vs NLE



NPE vs NLE

- NPE includes prior implicitly through the data
- NLE requires explicit prior to get the posterior
- Depending on the
 - dimensionality of the data vs. parameter space
 - Complexity of the likelihood vs. posteriorone or the other could work better



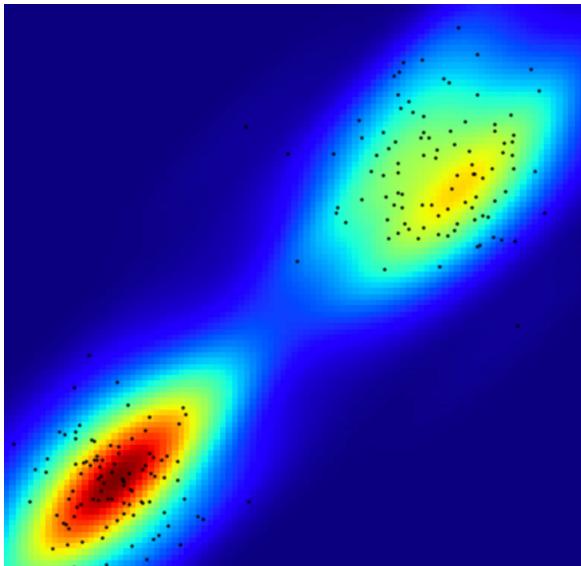
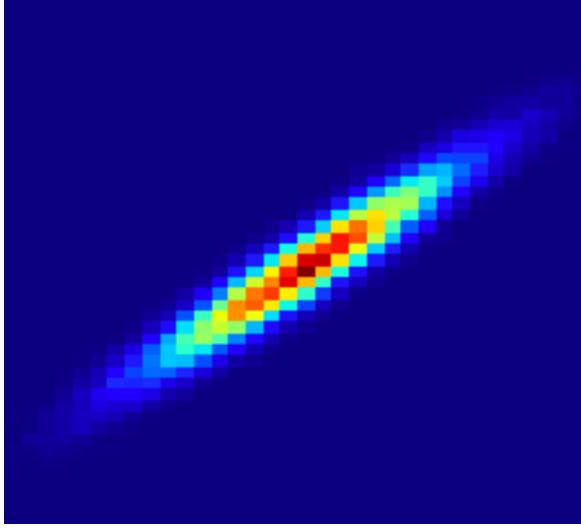
Density estimators (non-parametric)

- Histograms
- Kernel density estimators
 - can be parametric if the width is tuned

$$q_\epsilon(\mathbf{x}) = \frac{1}{N} \sum_n k_\epsilon(\mathbf{x} - \mathbf{x}_n)$$

$$k_1(\mathbf{u}) = \frac{1}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2} \|\mathbf{u}\|^2\right)$$

- Both will in general be cursed by dimensionality

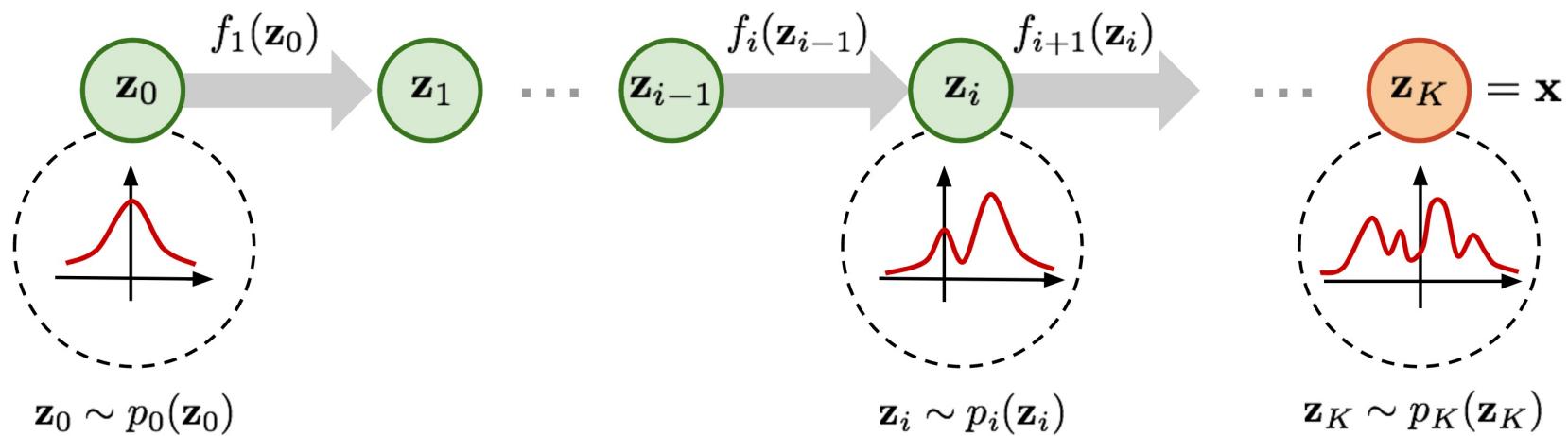


Density estimators (parametric)

- Gaussian $q_{\phi}(\mathbf{x}) = \frac{1}{|\det(2\pi\boldsymbol{\Sigma})|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$ where $\phi = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$
 - “MLE” for ϕ : $\boldsymbol{\mu}^* = \frac{1}{N} \sum_n \mathbf{x}_n$ and $\boldsymbol{\Sigma}^* = \frac{1}{N} \sum_n (\mathbf{x}_n - \boldsymbol{\mu}^*)(\mathbf{x}_n - \boldsymbol{\mu}^*)^T$
- Mixture models $q_{\phi}(\mathbf{x}) = \sum_k \alpha_k q_{\phi_k}^{(k)}(\mathbf{x})$ where $\sum_k \alpha_k = 1$ and $\alpha_k \geq 0$
 - Gaussian mixture model: $q_{\phi_k}^{(k)} = \mathcal{N}(\mathbf{x} | \phi_k = (\boldsymbol{\mu}, \boldsymbol{\Sigma}))$

Density estimators (parametric)

- Normalizing flows
 - Deforming a known distribution to fit the unknown, target distribution



Credit: Lilian Weng

Density estimators (parametric)

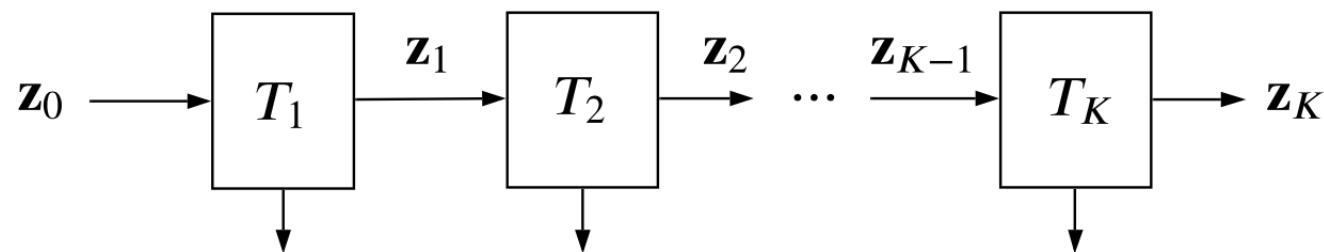
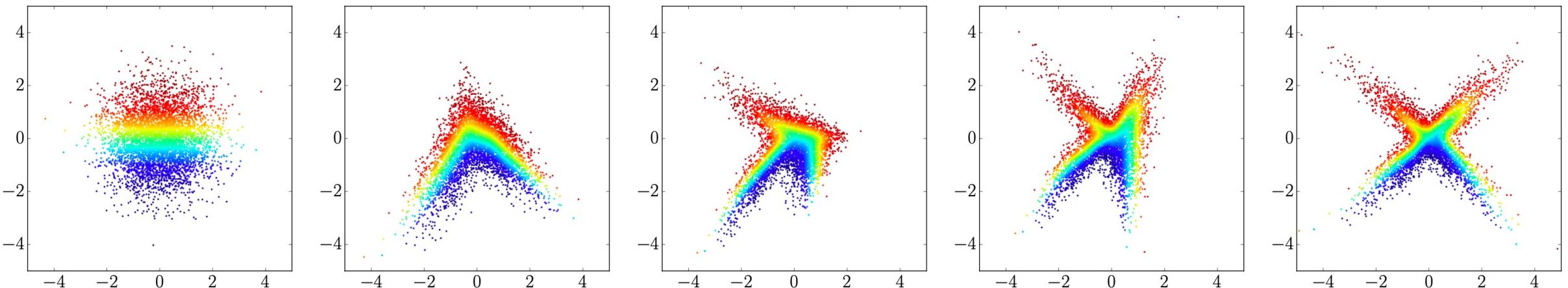
- Normalizing flows
 - Deforming a known distribution to fit the unknown, target distribution

$$\mathbf{z} \sim \pi(\mathbf{z}), \mathbf{x} = f(\mathbf{z}), \mathbf{z} = f^{-1}(\mathbf{x})$$

$$p(\mathbf{x}) = \pi(\mathbf{z}) \left| \det \frac{d\mathbf{z}}{d\mathbf{x}} \right| = \pi(f^{-1}(\mathbf{x})) \left| \det \frac{df^{-1}}{d\mathbf{x}} \right|$$

Density estimators (parametric)

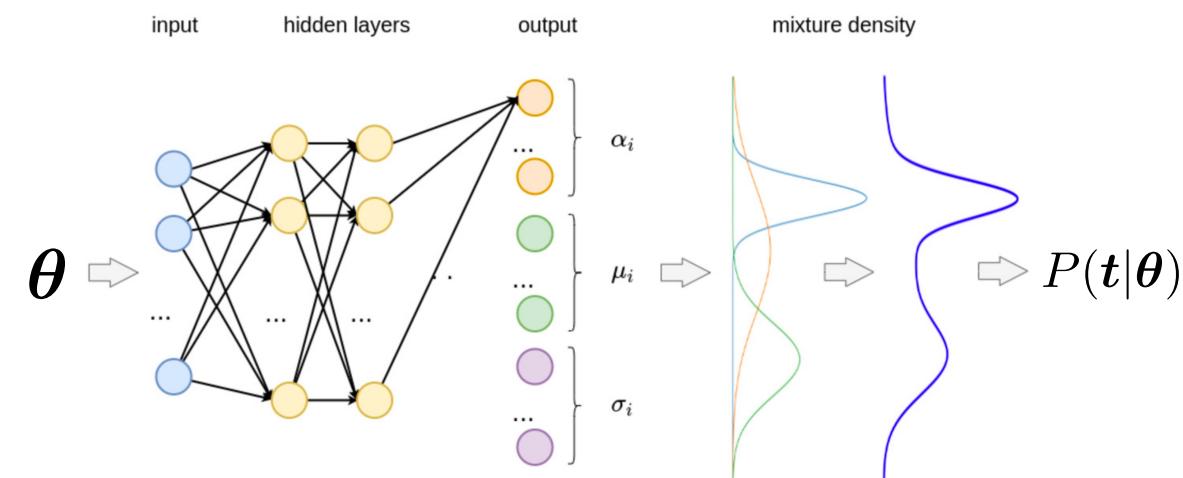
- Normalizing flows
 - Deforming a known distribution to fit the unknown, target distribution



$$\log |\det J_{T_1}(\mathbf{z}_0)| + \log |\det J_{T_2}(\mathbf{z}_1)| + \dots + \log |\det J_{T_K}(\mathbf{z}_{K-1})| = \log |\det J_T(\mathbf{z}_0)|$$

Conditioning the density

- Example #1: KDE -> [example](#)
- Example #2 Gaussian / Gaussian mixture
- Example #3: Normalizing flow
 - Conditioning each transformation on parameters θ



Neural Likelihood Estimation

Neural Likelihood Estimation

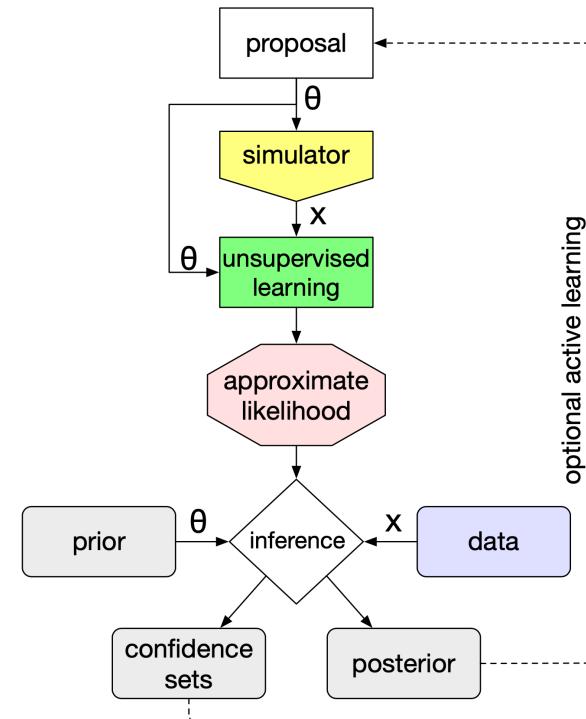
- Take a constructed set $\{(\tilde{t}_1, \tilde{\theta}_1), (\tilde{t}_2, \tilde{\theta}_2), \dots, (\tilde{t}_N, \tilde{\theta}_N)\}$ pulled from the total distribution $p(t, \theta) = p(t|\theta) \cdot p(\theta)$

- Pick a (conditional) neural density of your choice

- Minimize KL divergence

$$\begin{aligned} D_{\text{KL}}(P\|Q) &= \int_{-\infty}^{\infty} p(x) \log p(x) dx - \int_{-\infty}^{\infty} p(x) \log q(x) dx \\ &\approx -H_P - \frac{1}{N} \sum_{i=1}^N \log q(x_i), \end{aligned}$$

- (possibly) pull additional samples – sequential (active) learning



Neural Likelihood Estimation

- (possibly) pull additional samples – sequential (active) learning

Algorithm Sequential Neural Likelihood (SNL)

Input : observed data \mathbf{x}_o , estimator $q_\phi(\mathbf{x} | \boldsymbol{\theta})$,
number of rounds R , simulations per
round N

Output: approximate posterior $\hat{p}(\boldsymbol{\theta} | \mathbf{x}_o)$

set $\hat{p}_0(\boldsymbol{\theta} | \mathbf{x}_o) = p(\boldsymbol{\theta})$ and $\mathcal{D} = \{\}$

for $r = 1 : R$ **do**

for $n = 1 : N$ **do**

 sample $\boldsymbol{\theta}_n \sim \hat{p}_{r-1}(\boldsymbol{\theta} | \mathbf{x}_o)$ with MCMC

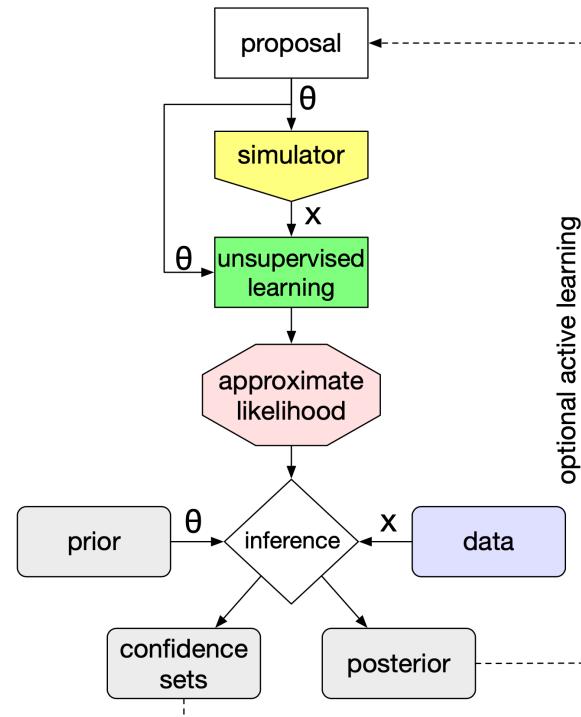
 simulate $\mathbf{x}_n \sim p(\mathbf{x} | \boldsymbol{\theta}_n)$

 add $(\boldsymbol{\theta}_n, \mathbf{x}_n)$ into \mathcal{D}

 (re-)train $q_\phi(\mathbf{x} | \boldsymbol{\theta})$ on \mathcal{D} and set

$\hat{p}_r(\boldsymbol{\theta} | \mathbf{x}_o) \propto q_\phi(\mathbf{x}_o | \boldsymbol{\theta}) p(\boldsymbol{\theta})$

return $\hat{p}_R(\boldsymbol{\theta} | \mathbf{x}_o)$



Neural Likelihood Estimation

- (possibly) pull additional samples – sequential (active) learning

Algorithm Sequential Neural Likelihood (SNL)

Input : observed data \mathbf{x}_o , estimator $q_\phi(\mathbf{x} | \boldsymbol{\theta})$,
number of rounds R , simulations per
round N

Output: approximate posterior $\hat{p}(\boldsymbol{\theta} | \mathbf{x}_o)$

set $\hat{p}_0(\boldsymbol{\theta} | \mathbf{x}_o) = p(\boldsymbol{\theta})$ and $\mathcal{D} = \{\}$

for $r = 1 : R$ **do**

for $n = 1 : N$ **do**

 sample $\boldsymbol{\theta}_n \sim \hat{p}_{r-1}(\boldsymbol{\theta} | \mathbf{x}_o)$ with MCMC

 simulate $\mathbf{x}_n \sim p(\mathbf{x} | \boldsymbol{\theta}_n)$

 add $(\boldsymbol{\theta}_n, \mathbf{x}_n)$ into \mathcal{D}

 (re-)train $q_\phi(\mathbf{x} | \boldsymbol{\theta})$ on \mathcal{D} and set

$\hat{p}_r(\boldsymbol{\theta} | \mathbf{x}_o) \propto q_\phi(\mathbf{x}_o | \boldsymbol{\theta}) p(\boldsymbol{\theta})$

return $\hat{p}_R(\boldsymbol{\theta} | \mathbf{x}_o)$

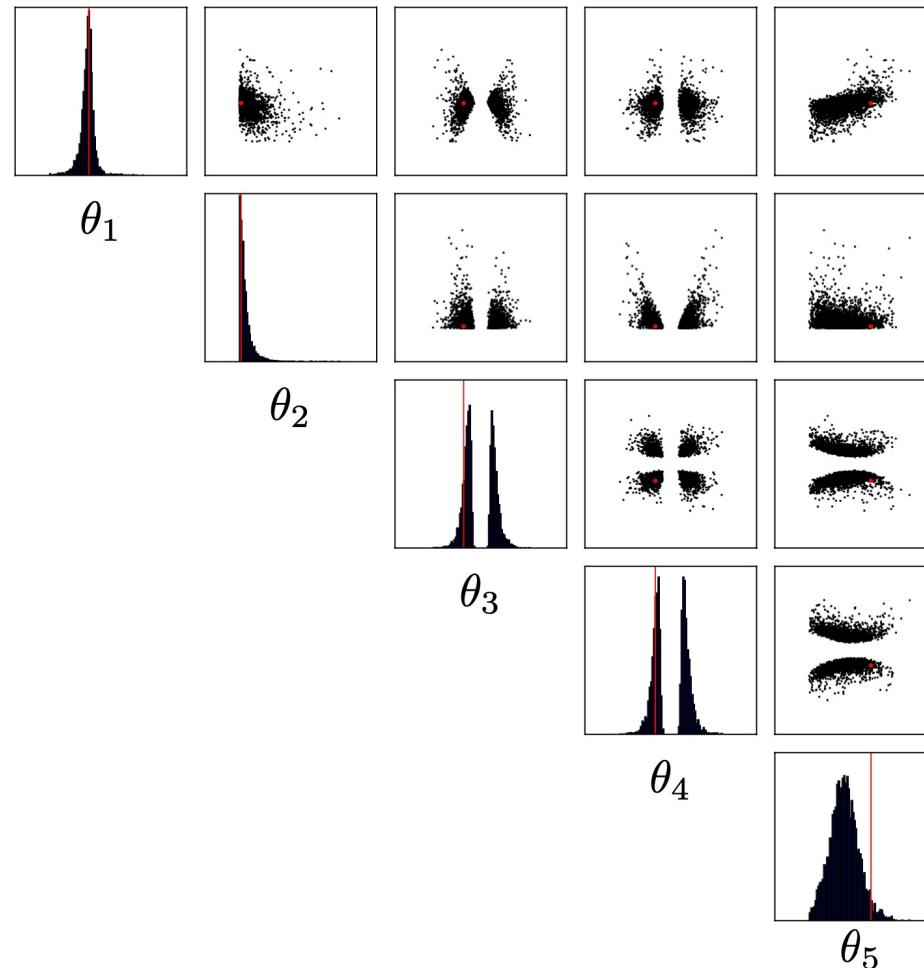
Pro:

- Provides better posterior as samples were chosen in the relevant parameter space

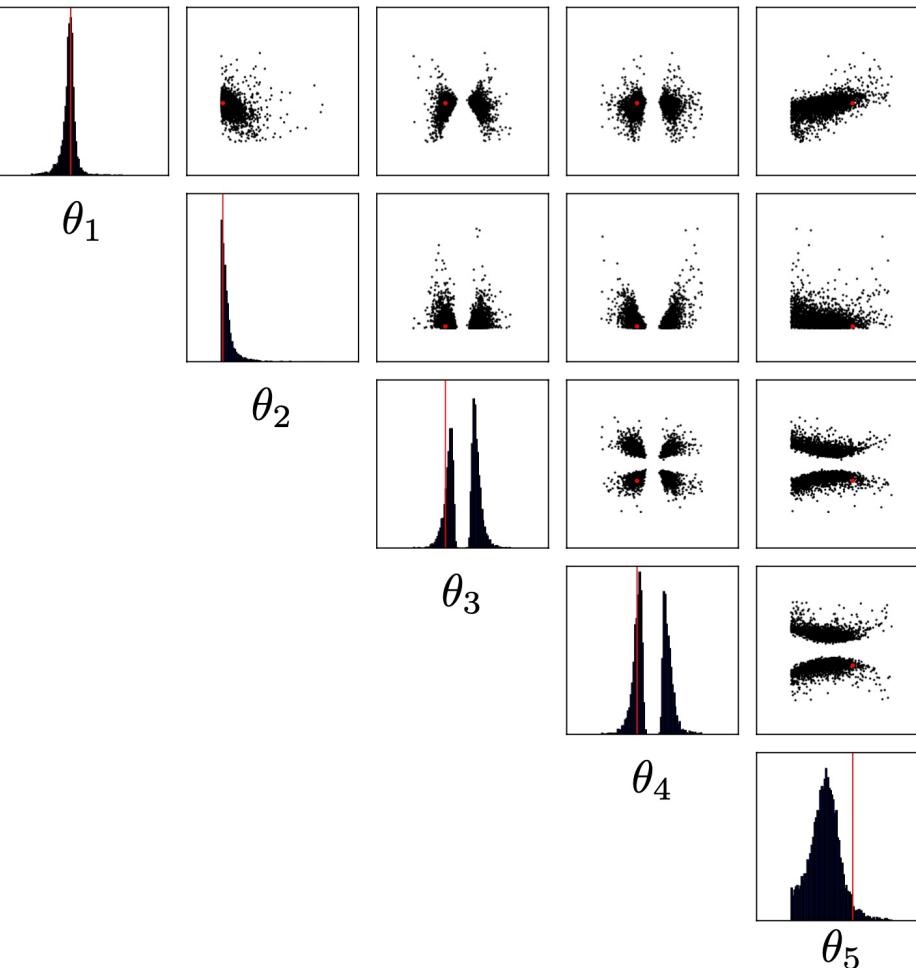
Con:

- It will not perform so well for samples other than \mathbf{x}_0

Neural Likelihood Estimation



(a) MCMC samples from exact posterior. True parameters indicated in red.



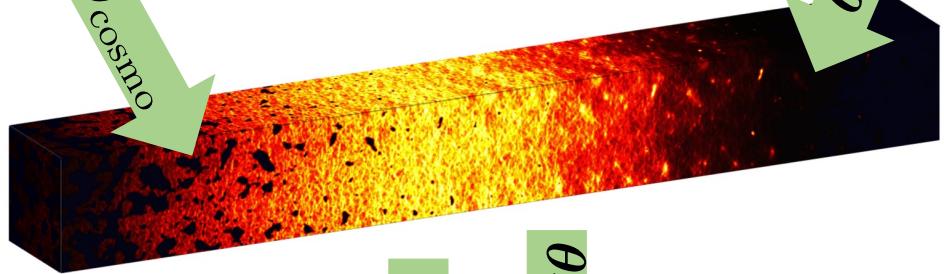
(b) MCMC samples from SNL posterior. True parameters indicated in red.

Example: SBI with NLE
for 21-cm PS

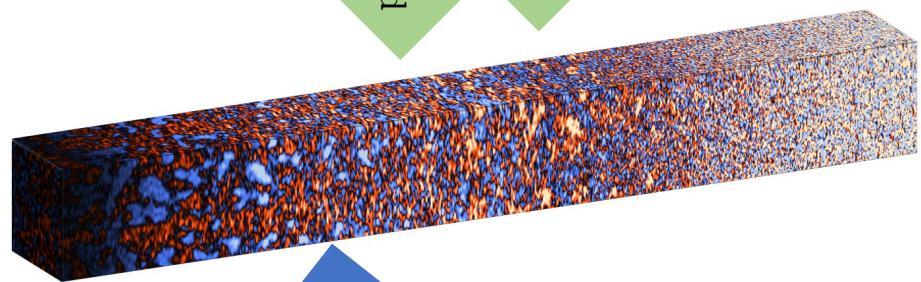
Astrophysics

Cosmology

Simulation

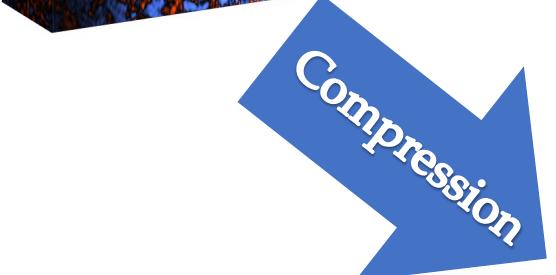


Foregrounds
simulator



Telescope
simulator

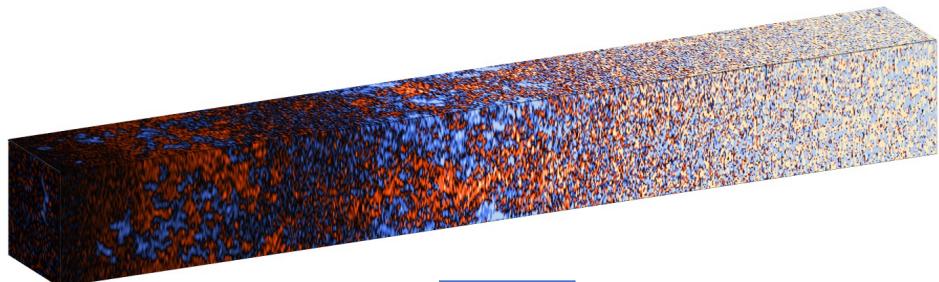
\mathcal{S}_{model}



Likelihood (analytic / LFI)

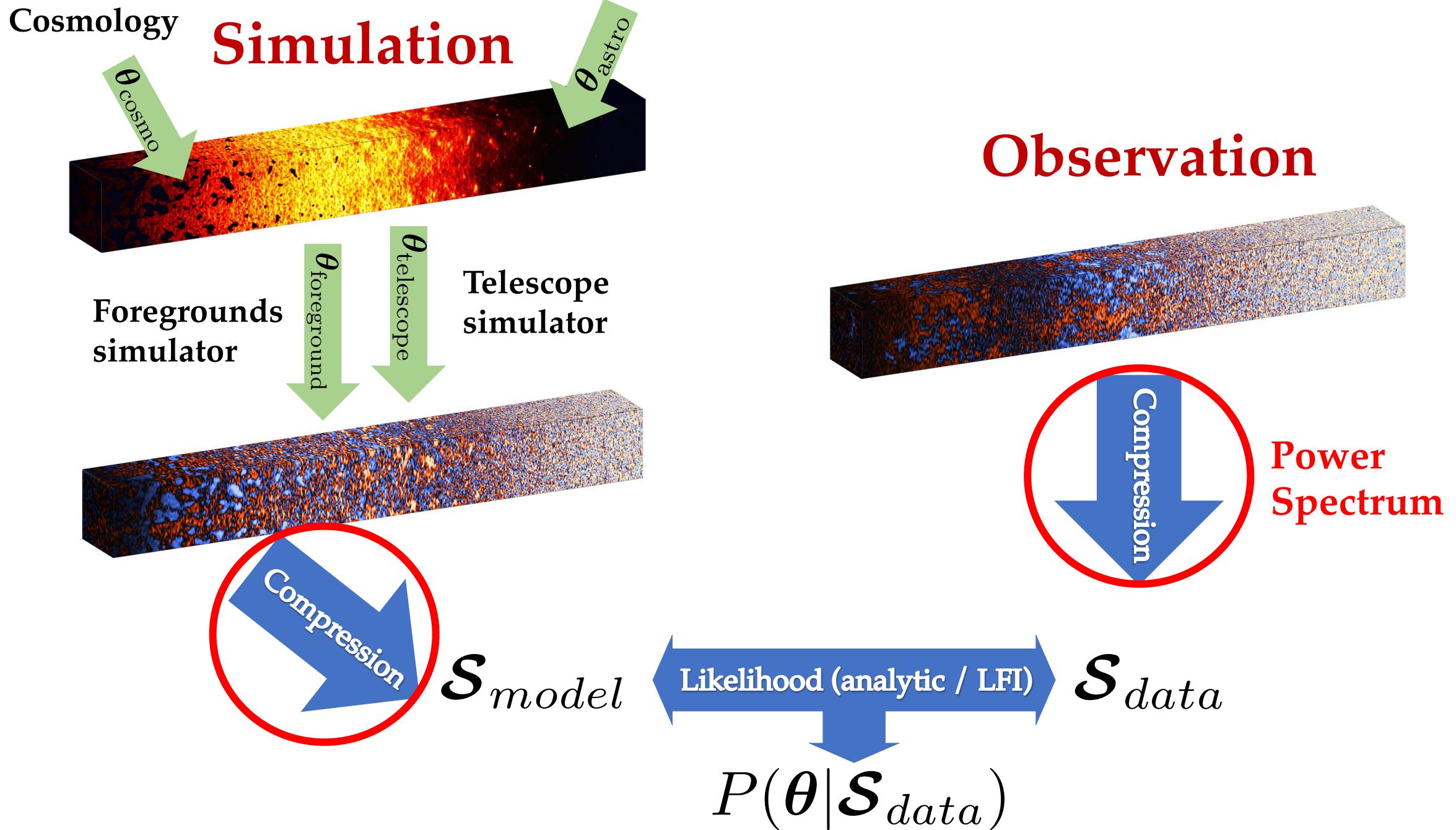
$$P(\theta | \mathcal{S}_{data})$$

Observation



\mathcal{S}_{data}

Astrophysics



Likelihood of the 21-cm PS

- Analytically not available
- Numerically non-tractable
- Eg. – power spectrum

$$\mathcal{L}(\boldsymbol{\sigma}, \boldsymbol{\mu}; \mathbf{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{\prod_i \sigma_i^2}} e^{-\frac{1}{2} \sum_i (x_i - \mu_i)^2 / \sigma_i^2}$$

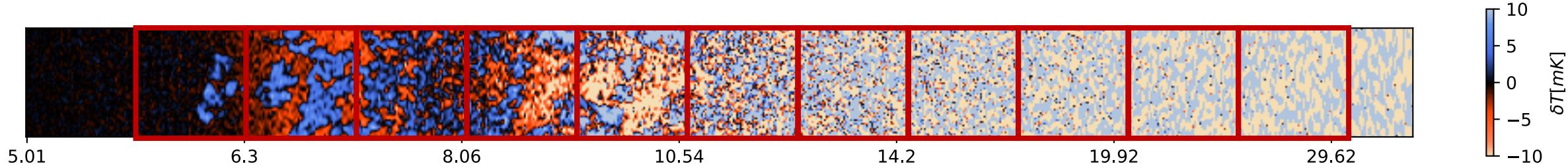
$$\mathcal{L}(\boldsymbol{\Sigma}, \boldsymbol{\mu}; \mathbf{x}) = \frac{1}{(2\pi)^{n/2} \sqrt{|\boldsymbol{\Sigma}|}} e^{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$$

??? higher order terms & non-constant $\sigma / \boldsymbol{\Sigma}$???

→ Neural Likelihood Estimation

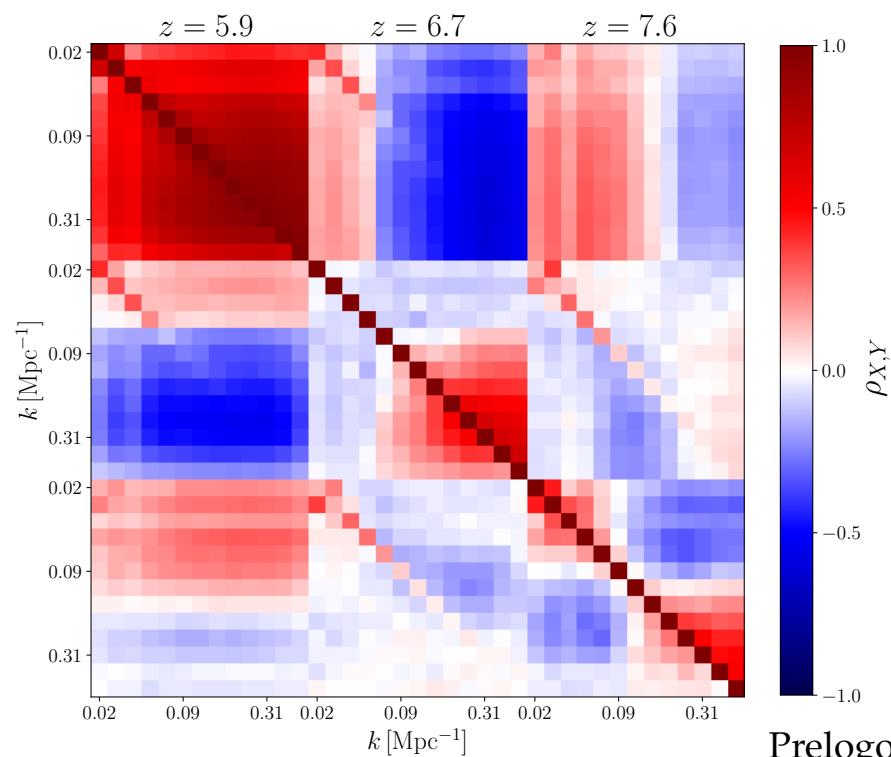
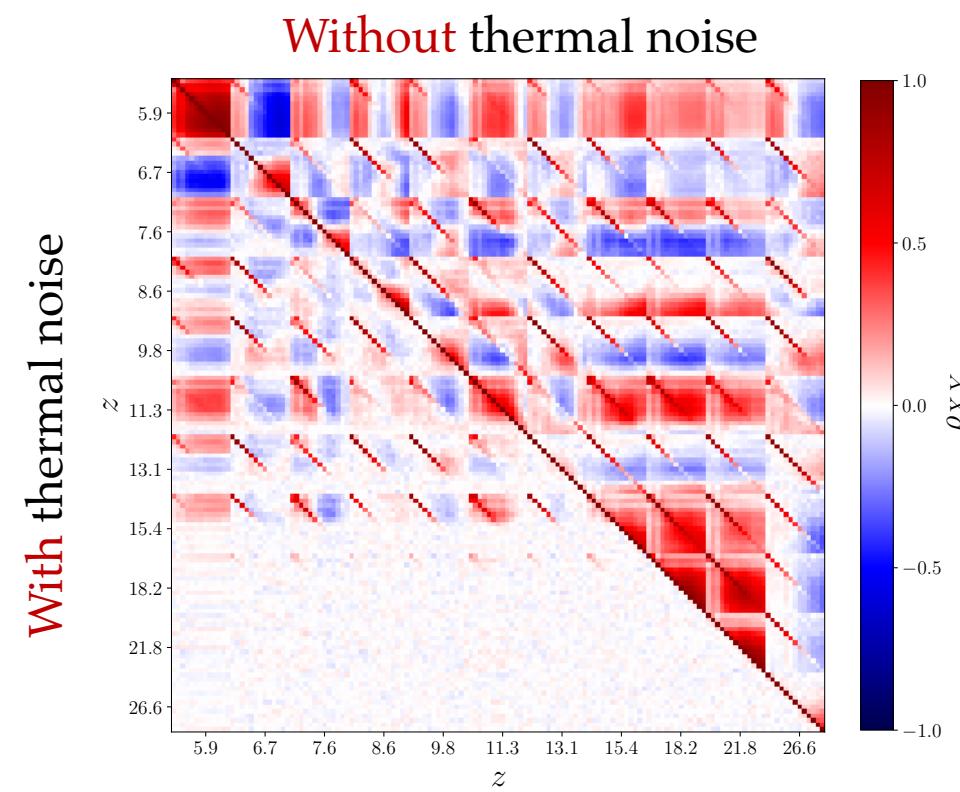
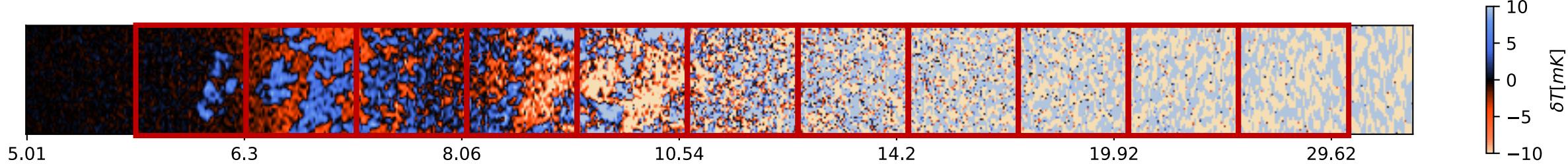
Likelihood of the 21-cm PS

- Is inclusion of the covariance worthwhile?



Likelihood of the 21-cm PS

- Is inclusion of the covariance worthwhile?



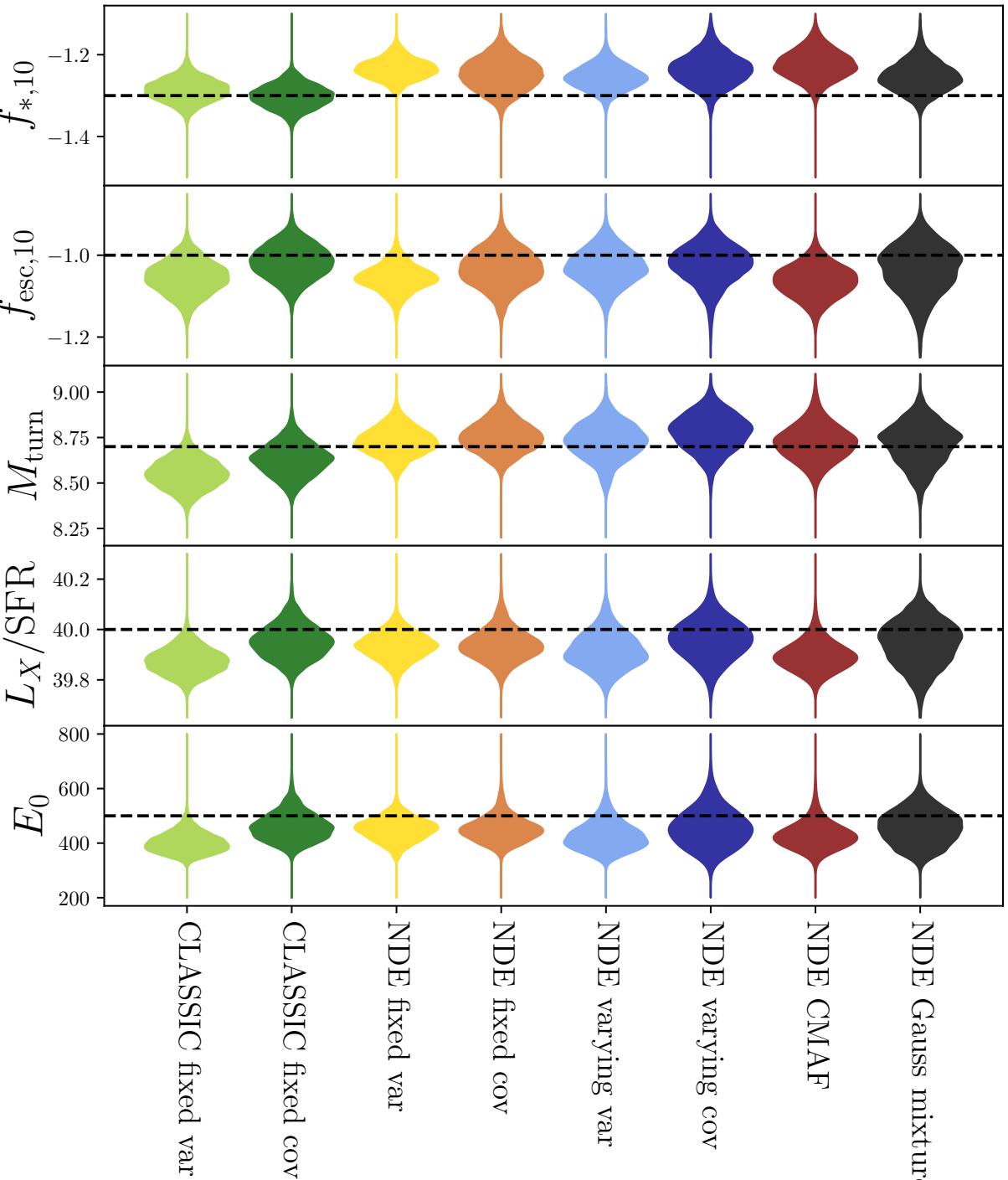
Likelihood of the 21-cm PS

- Classically, Σ is estimated only at the fiducial
 - Estimating it at every point in parameter space is too expensive
 - Even for a proper Gaussian, we rely on NLE
- Gaussian Mixture Density Network
- Conditional Masked Autoregressive Flows

Likelihood of the 21-cm PS

	CLASSIC		NDE				
	fixed var	fixed cov	fixed var	fixed cov	varying var	varying cov	CMAF
$f_{*,10}$	1.34	2.71	0.50	0.03	0.05	0.11	0.36
$f_{\text{esc},10}$	0.53	1.01	0.84	0.16	0.13	0.07	0.57
M_{turn}	2.37	0.49	0.32	0.38	0.03	0.19	0.05
L_X/SFR	1.19	0.50	0.59	0.26	0.16	0.06	0.48
E_0	1.38	0.15	0.26	0.18	0.17	0.02	0.21

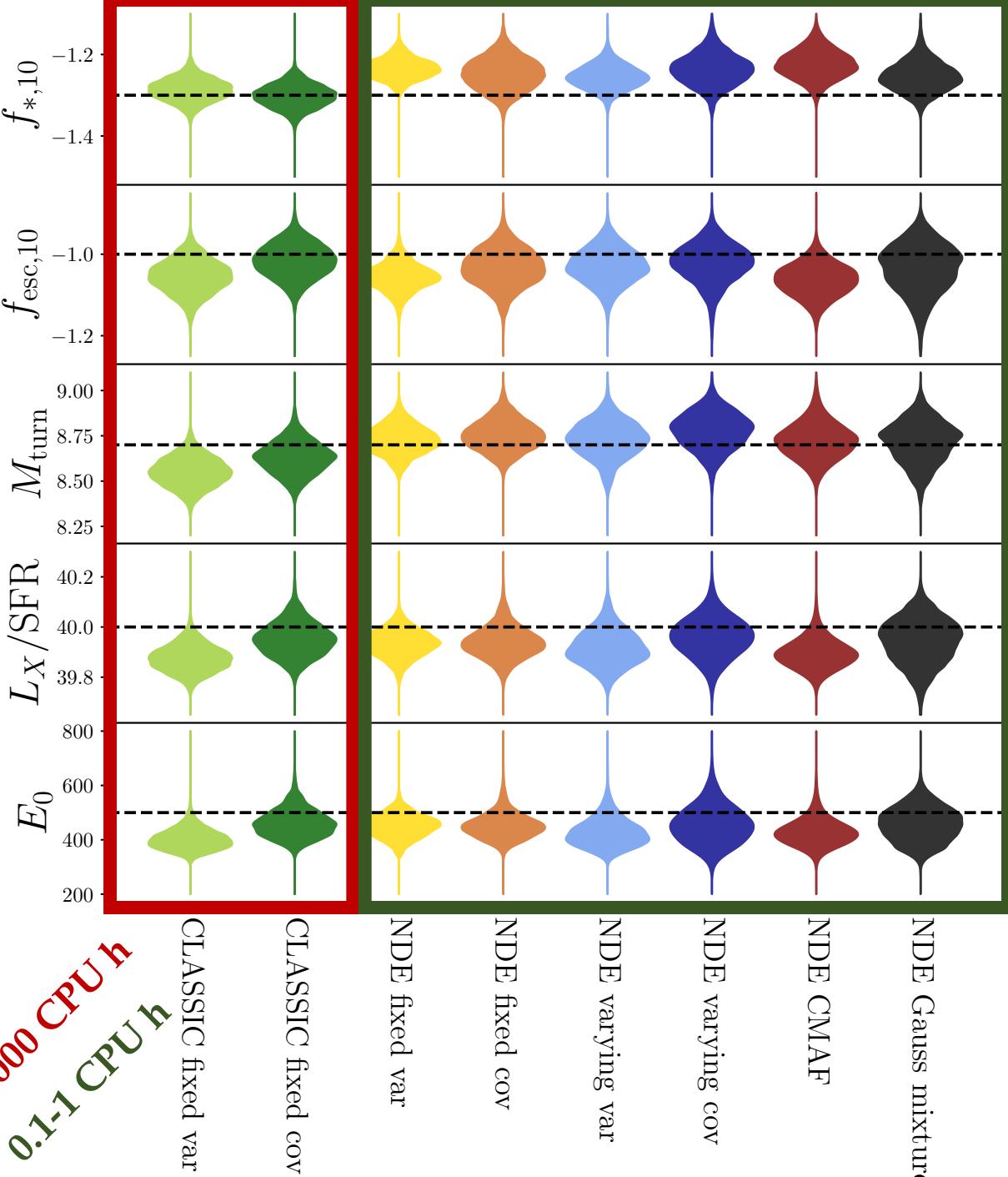
- Excluding correlations gives over-confident and biased posteriors
- NDE varying cov close to the NDE Gauss mixture
- NDE CMAF not performing very well
 - Possible issue is high expressivity, so it needs **sequential learning**



Likelihood of the 21-cm PS

	CLASSIC		NDE				
	fixed var	fixed cov	fixed var	fixed cov	varying var	varying cov	CMAF
$f_{*,10}$	1.34	2.71	0.50	0.03	0.05	0.11	0.36
$f_{\text{esc},10}$	0.53	1.01	0.84	0.16	0.13	0.07	0.57
M_{turn}	2.37	0.49	0.32	0.38	0.03	0.19	0.05
L_X/SFR	1.19	0.50	0.59	0.26	0.16	0.06	0.48
E_0	1.38	0.15	0.26	0.18	0.17	0.02	0.21

- Excluding correlations gives over-confident and biased posteriors
- NDE varying cov close to the NDE Gauss mixture
- NDE CMAF not performing very well
 - Possible issue is high expressivity, so it needs **sequential learning**



Likelihood of the 21-cm PS

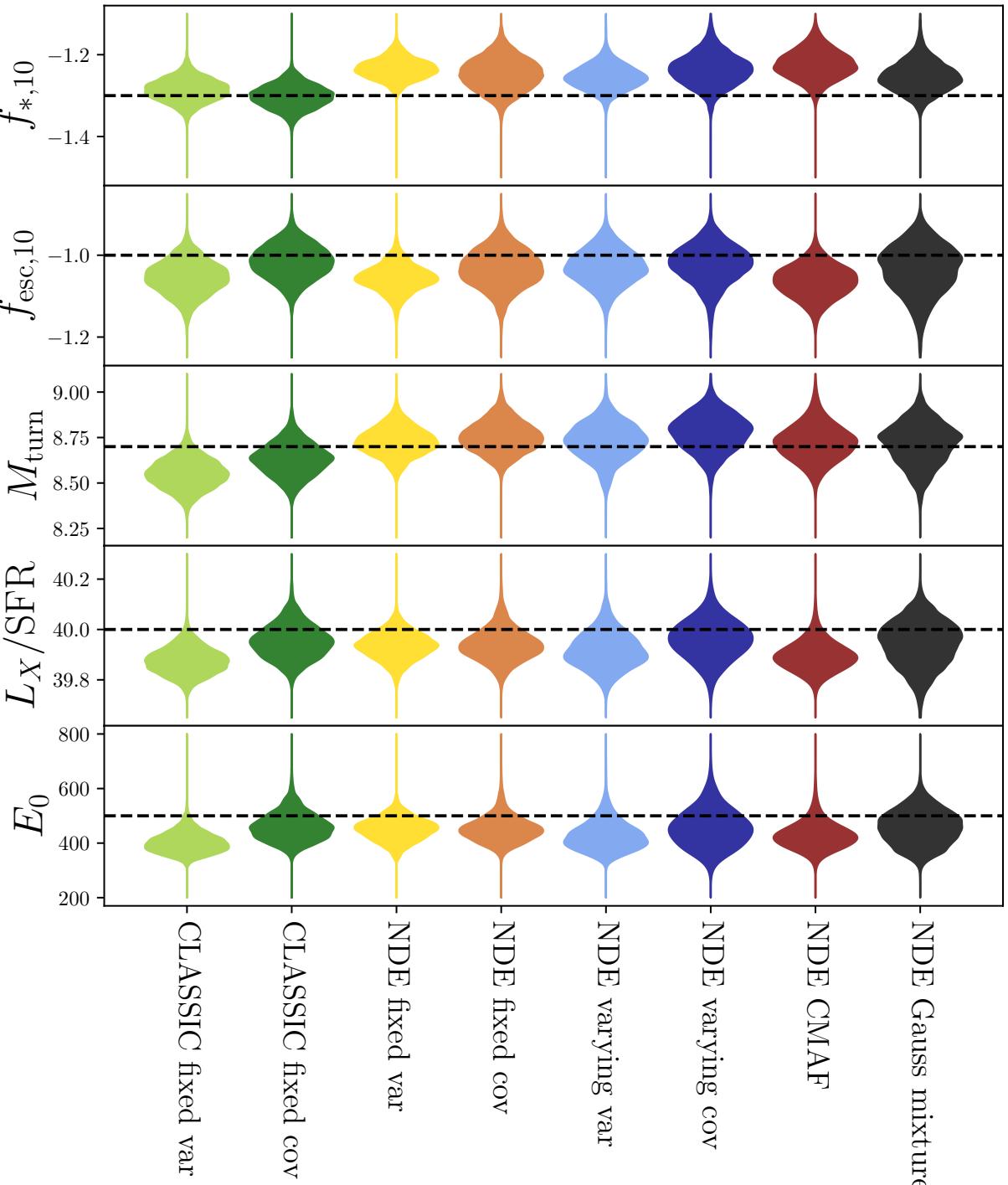
	CLASSIC		NDE				
	fixed var	fixed cov	fixed var	fixed cov	varying var	varying cov	CMAF
$f_{*,10}$	1.34	2.71	0.50	0.03	0.05	0.11	0.36
$f_{\text{esc},10}$	0.53	1.01	0.84	0.16	0.13	0.07	0.57
M_{turn}	2.37	0.49	0.32	0.38	0.03	0.19	0.05
L_X/SFR	1.19	0.50	0.59	0.26	0.16	0.06	0.48
E_0	1.38	0.15	0.26	0.18	0.17	0.02	0.21

BUT:

This is only qualitative description,
and only for the mock observation

- How does it perform for other points in the parameter space?
- Did the training converge?
- Can we quantify the best model?

→ Simulation-Based Calibration



Simulation Based Calibration (SBC)

- “prior” = “*data averaged posterior*” $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\boldsymbol{y}}) P(\tilde{\boldsymbol{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) d\tilde{\boldsymbol{y}} d\tilde{\boldsymbol{\theta}}$

Simulation Based Calibration (SBC)

- “prior” = “*data averaged posterior*” $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\mathbf{y}}) P(\tilde{\mathbf{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) d\tilde{\mathbf{y}} d\tilde{\boldsymbol{\theta}}$

1. Pull from prior

$$\tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$$

Simulation Based Calibration (SBC)

- “prior” = “*data averaged posterior*” $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\mathbf{y}}) P(\tilde{\mathbf{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) d\tilde{\mathbf{y}} d\tilde{\boldsymbol{\theta}}$

1. Pull from prior

$$\tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$$

2. Pull the data from the likelihood

$$\tilde{\mathbf{y}} \sim P(\mathbf{y}|\tilde{\boldsymbol{\theta}}) \Leftrightarrow \tilde{\mathbf{y}} = \text{simulator}(\tilde{\boldsymbol{\theta}})$$

Simulation Based Calibration (SBC)

- “prior” = “*data averaged posterior*” $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\mathbf{y}}) P(\tilde{\mathbf{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) d\tilde{\mathbf{y}} d\tilde{\boldsymbol{\theta}}$

1. Pull from prior $\tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$
2. Pull the data from the likelihood $\tilde{\mathbf{y}} \sim P(\mathbf{y}|\tilde{\boldsymbol{\theta}}) \Leftrightarrow \tilde{\mathbf{y}} = \text{simulator}(\tilde{\boldsymbol{\theta}})$
3. Calculate the posterior the sample $P(\boldsymbol{\theta}|\tilde{\mathbf{y}})$

Simulation Based Calibration (SBC)

- “prior” = “*data averaged posterior*” $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\mathbf{y}}) P(\tilde{\mathbf{y}}|\tilde{\boldsymbol{\theta}}) P(\tilde{\boldsymbol{\theta}}) d\tilde{\mathbf{y}} d\tilde{\boldsymbol{\theta}}$

1. Pull from prior

$$\tilde{\boldsymbol{\theta}} \sim P(\boldsymbol{\theta})$$

2. Pull the data from the likelihood

$$\tilde{\mathbf{y}} \sim P(\mathbf{y}|\tilde{\boldsymbol{\theta}}) \Leftrightarrow \tilde{\mathbf{y}} = \text{simulator}(\tilde{\boldsymbol{\theta}})$$

3. Calculate the posterior the sample

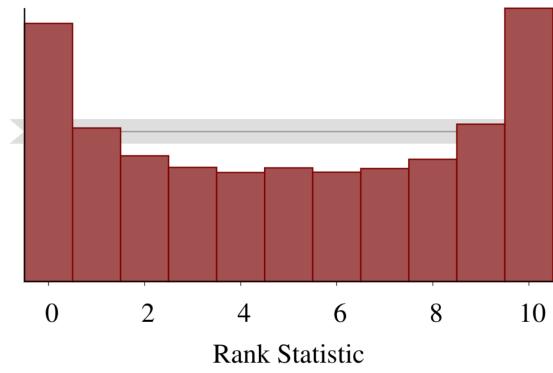
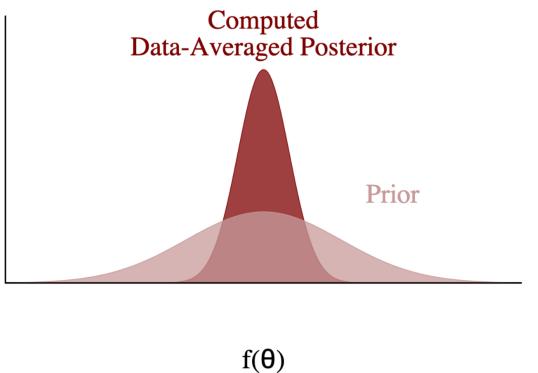
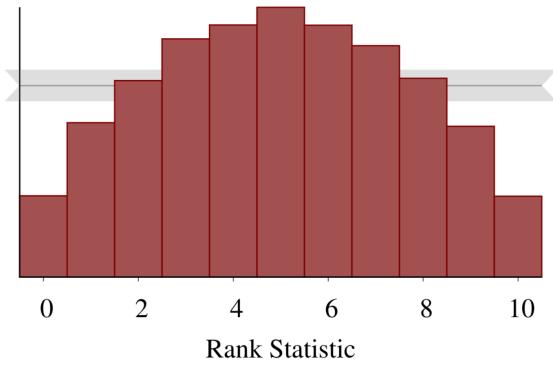
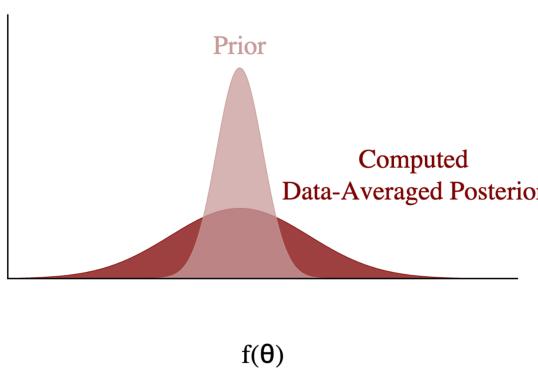
$$P(\boldsymbol{\theta}|\tilde{\mathbf{y}})$$

4. Repeat and average posteriors

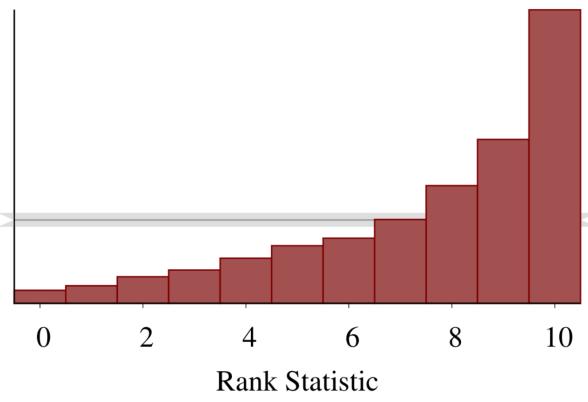
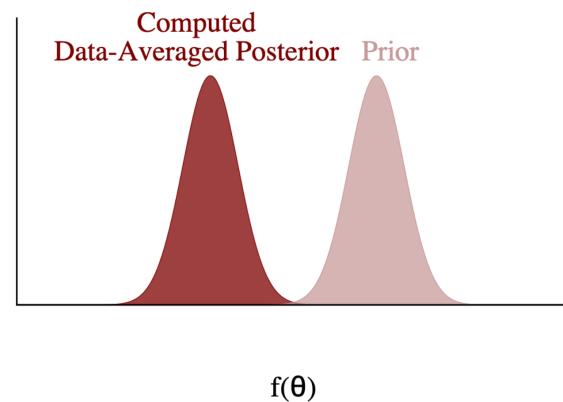
$$P(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{i=1}^N P_i(\boldsymbol{\theta}|\tilde{\mathbf{y}}_i)$$

Simulation Based Calibration (SBC)

- “prior” = “*data averaged posterior*” $P(\boldsymbol{\theta}) = \int P(\boldsymbol{\theta}|\tilde{\mathbf{y}}) P(\tilde{\mathbf{y}}|\boldsymbol{\theta})P(\boldsymbol{\theta}) d\tilde{\mathbf{y}} d\boldsymbol{\theta}$
- SBC – casting integral into 1D rank statistics distribution

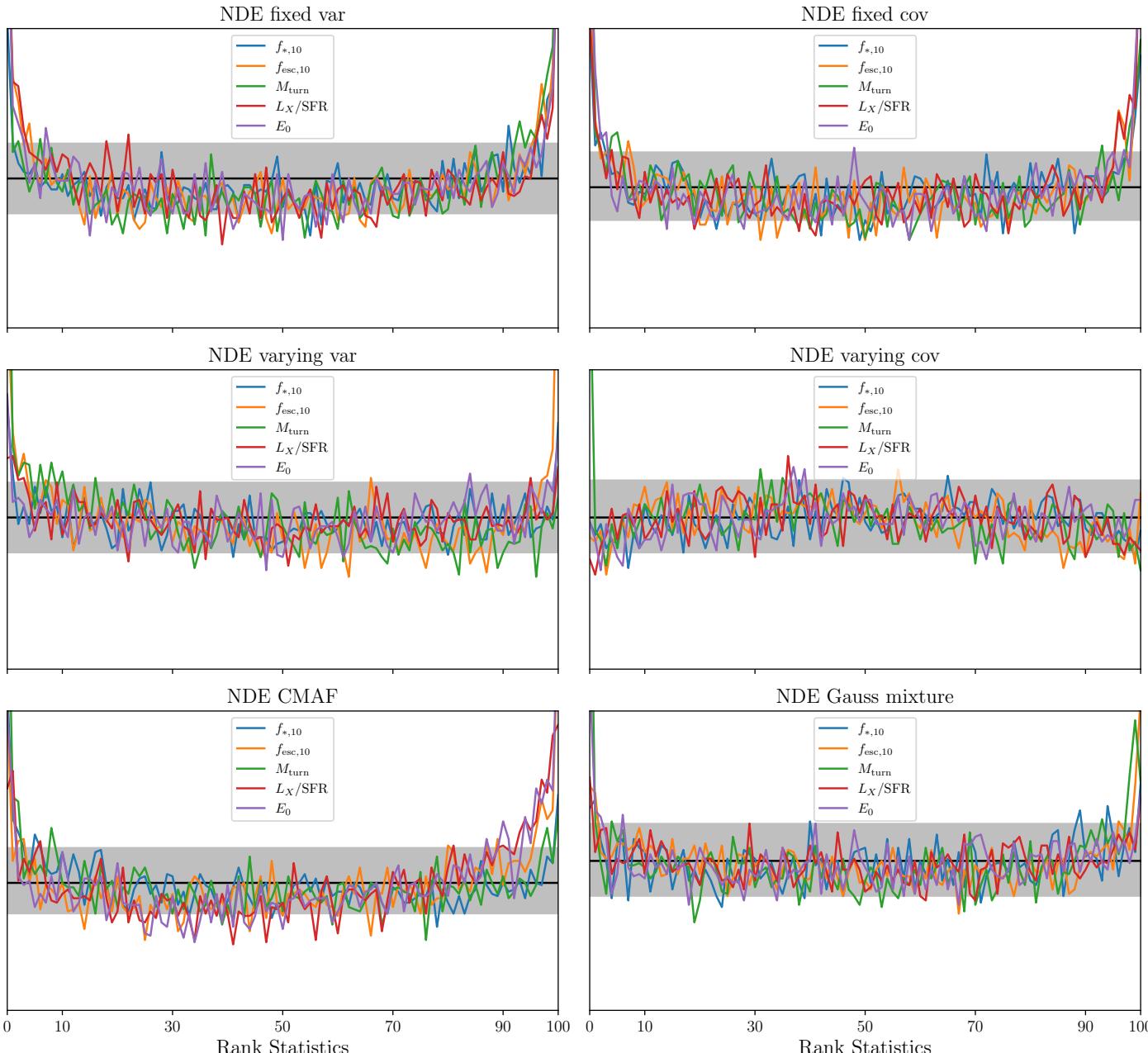


$$r \left(f(\boldsymbol{\theta}_1), \dots, f(\boldsymbol{\theta}_L), f(\tilde{\boldsymbol{\theta}}) \right) = \sum_{l=1}^L \mathbb{I} \left[f(\boldsymbol{\theta}_l) < f(\tilde{\boldsymbol{\theta}}) \right]$$



SBC for 21-cm PS

- *Expensive to compute!*
 - 10 000 posteriors
- However, once likelihood is trained, no new simulations are needed
 - “amortized inference”
- Procedure would be useful for classic inference too, but is impossible to compute
- NDE fixed var & cov – overconfident
- NDE varying var & cov – good
- NDE CMAF – overconfident
- NDE Gauss mixture – the best



Neural Ratio Estimation

and marginalization

NRE vs NLE/NPE

- No need to estimate a full distribution
- Ratio is a simple scalar
- Any neural network will do
- Doesn't suffer from the curse of dimensionality (in data)



Neural Ratio Estimation (NRE)

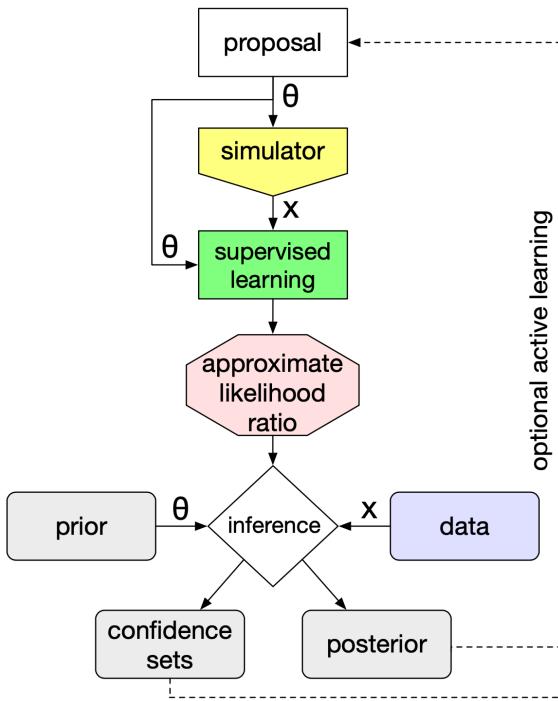
- The main focus:

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

- Introduce a binary random variable y , which classifies

$$\tilde{p}(x, \theta | y) = \begin{cases} p(x, \theta) & \text{if } y = 1 \\ p(x)p(\theta) & \text{if } y = 0 \end{cases}$$

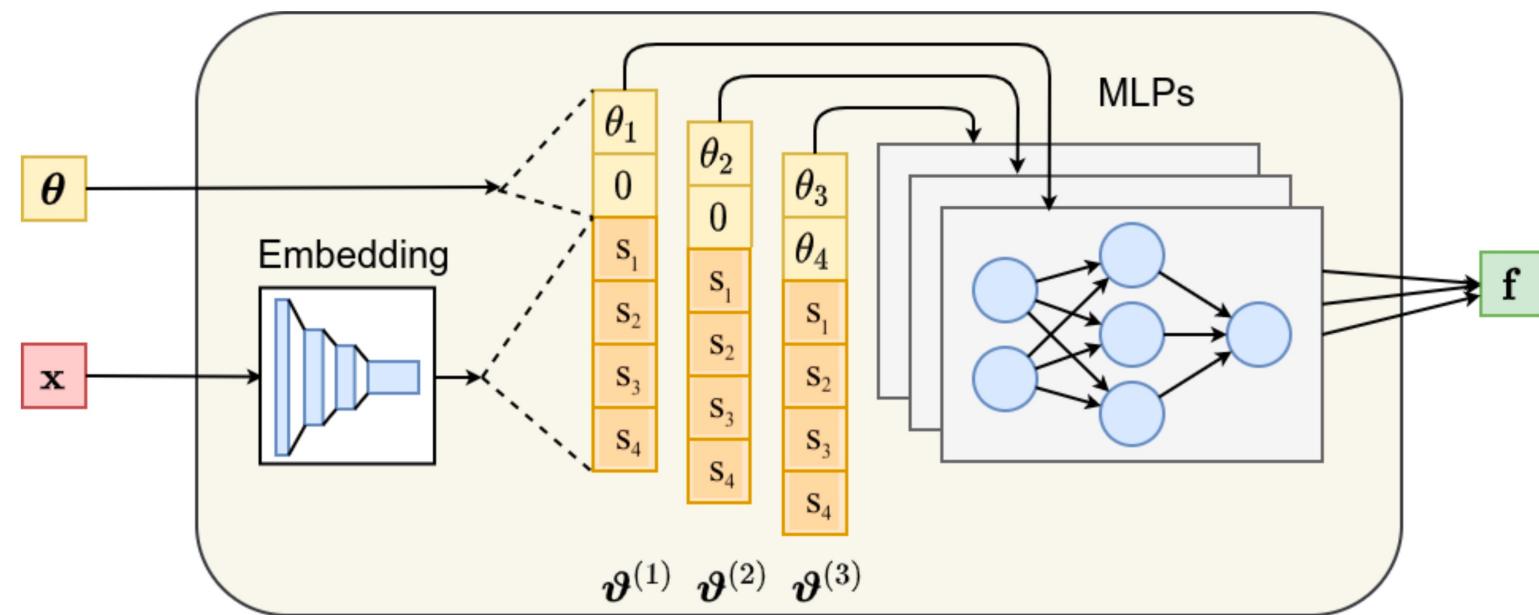
- Take a classifier $\tilde{p}(y = 1 | x, \theta)$
- Compute likelihood trick - *blackboard*
- Train it with all DL machinery developed for classifiers



Marginal NRE (MNRE)



- Nuisance parameters
 - Often we don't care about them
- Important parameters
 - Often we're interested only in 1D and 2D posteriors



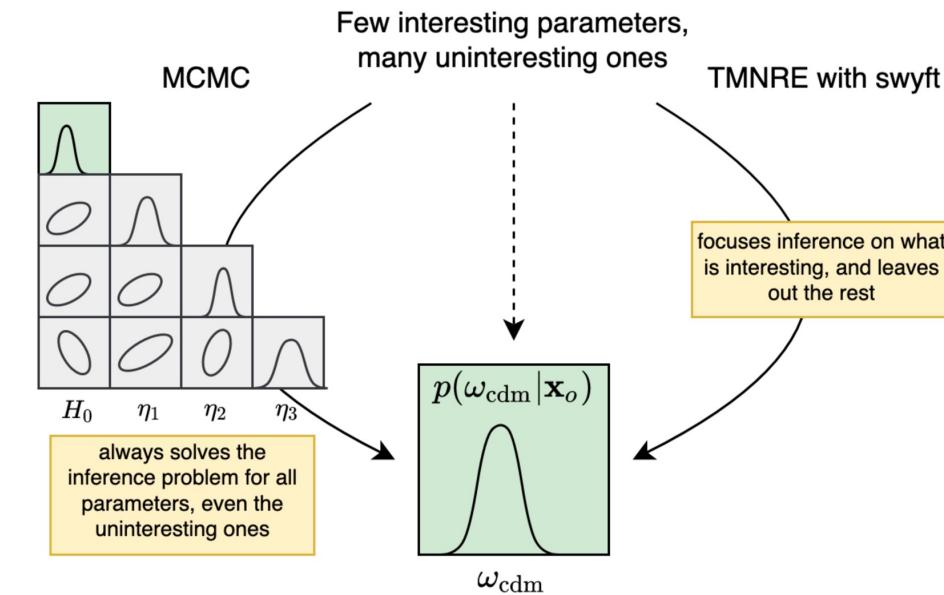
Marginal NRE (MNRE)

- For example, we're interested in only CDM energy density
- All other parameters are not important for us

$$p(\omega_{cdm} \mid \mathbf{x}) = \int d\boldsymbol{\eta} p(\omega_{cdm}, \eta_1, \dots, \eta_N \mid \mathbf{x})$$

- Running the full inference can be arbitrarily expensive, depending on the model complexity

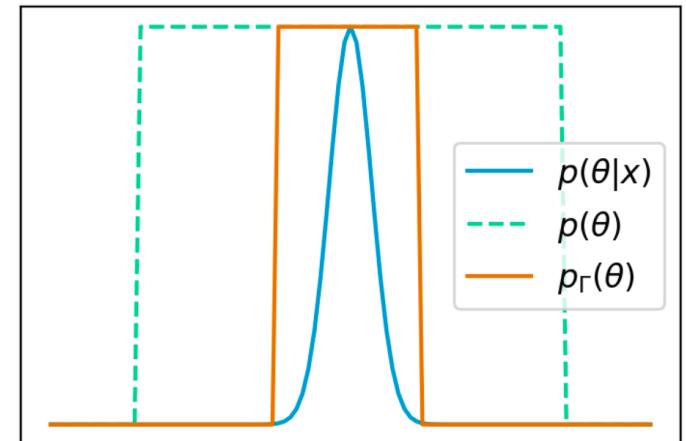
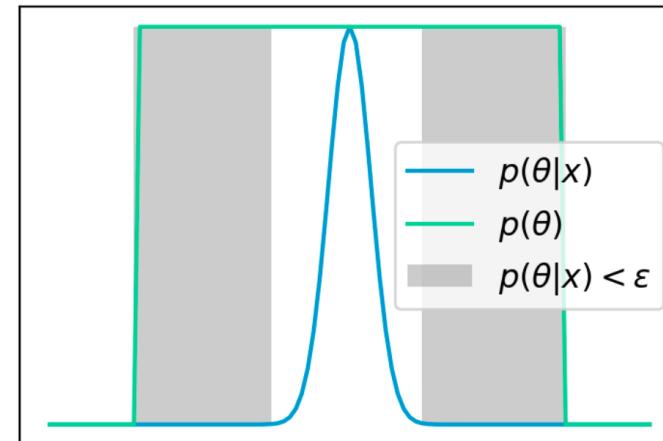
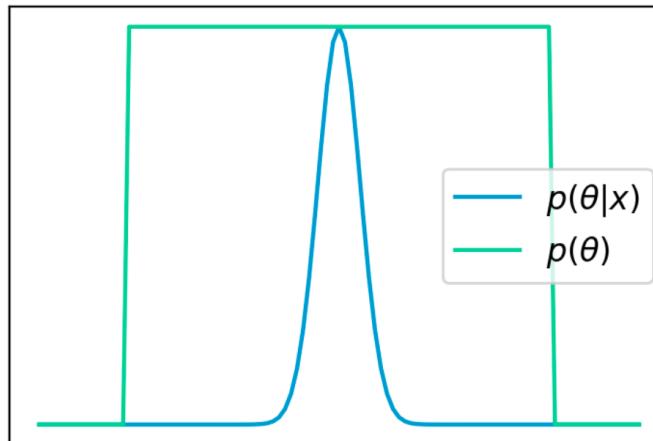
- Importantly, prior on η must be good!
- For other parameters important, but not crucial



Truncated Marginal NRE (TMNRE)

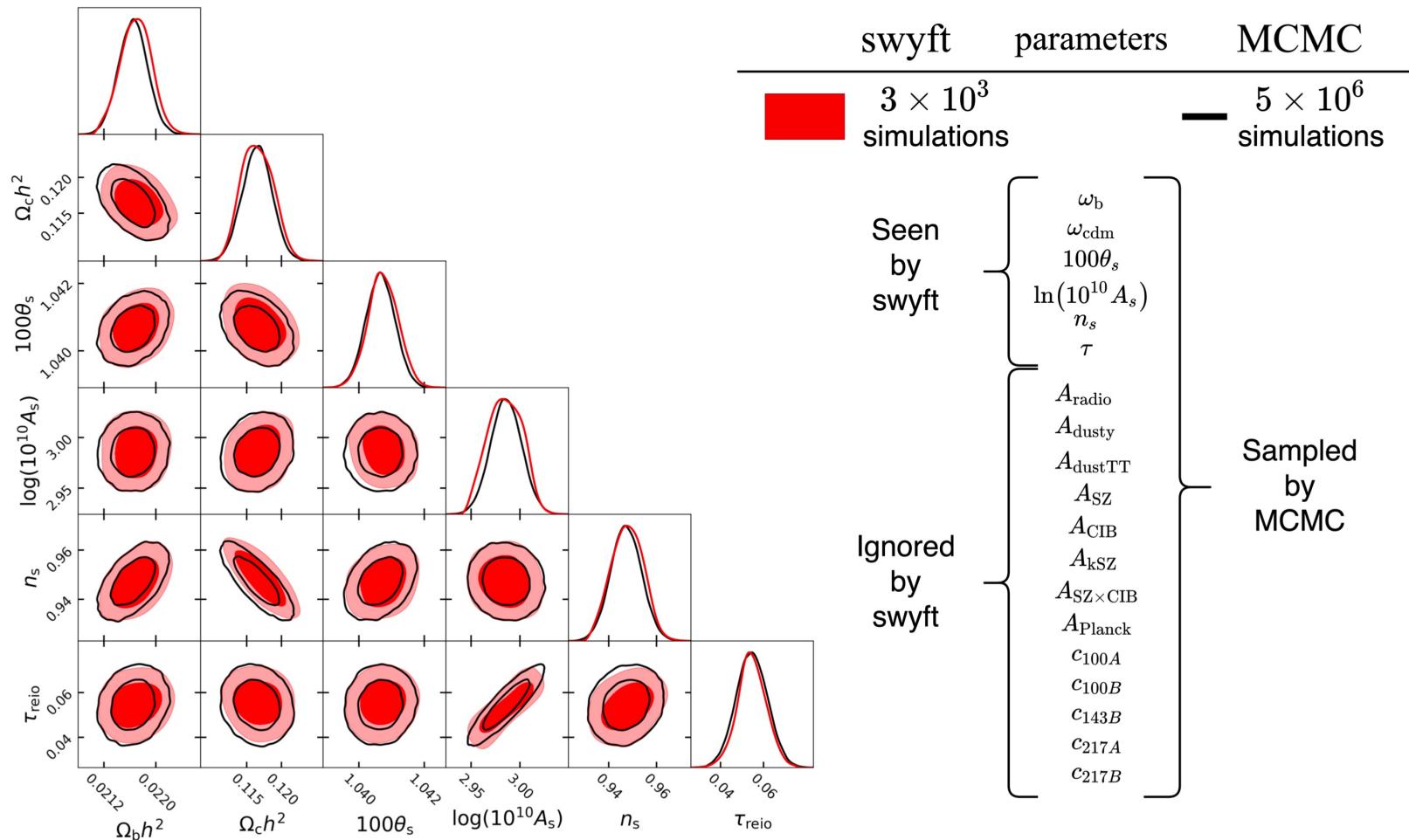


- Similar to the sequential (active) learning for NLE
- In several iterations focusing on the relevant prior volume
- Controlled by the posterior estimation around some x_0



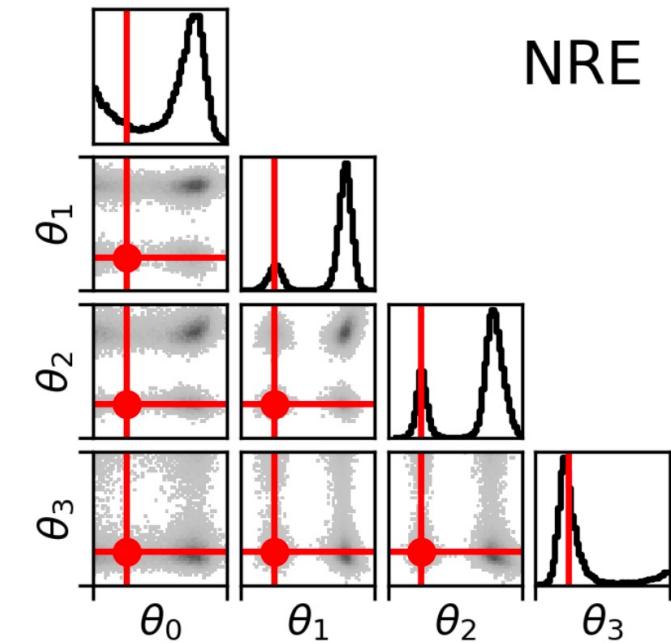
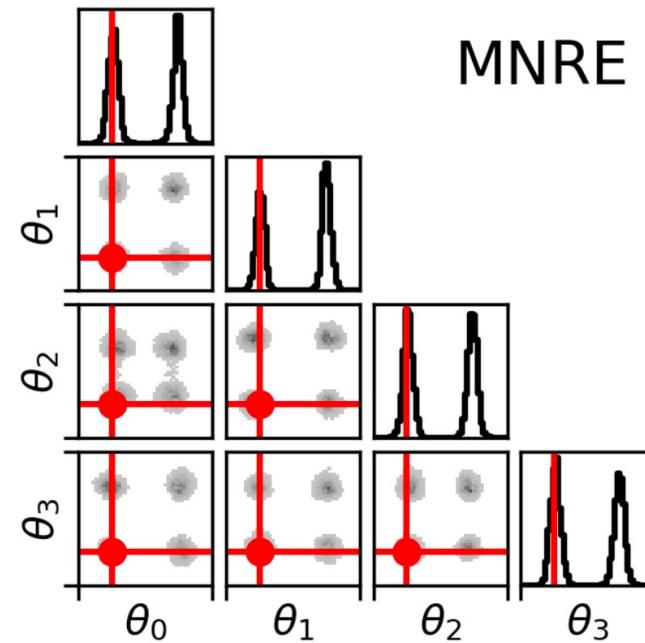
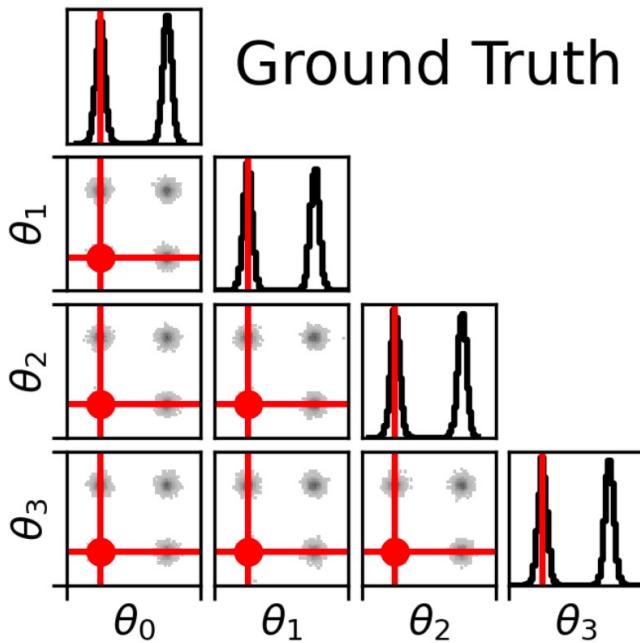
Truncated Marginal NRE (TMNRE)

- Example: MCMC vs TMNRE for CMB

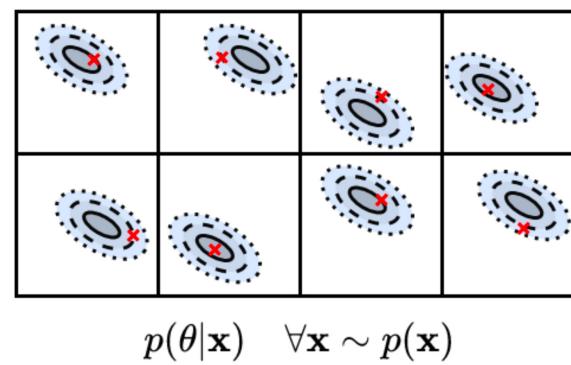


Truncated Marginal NRE (TMNRE)

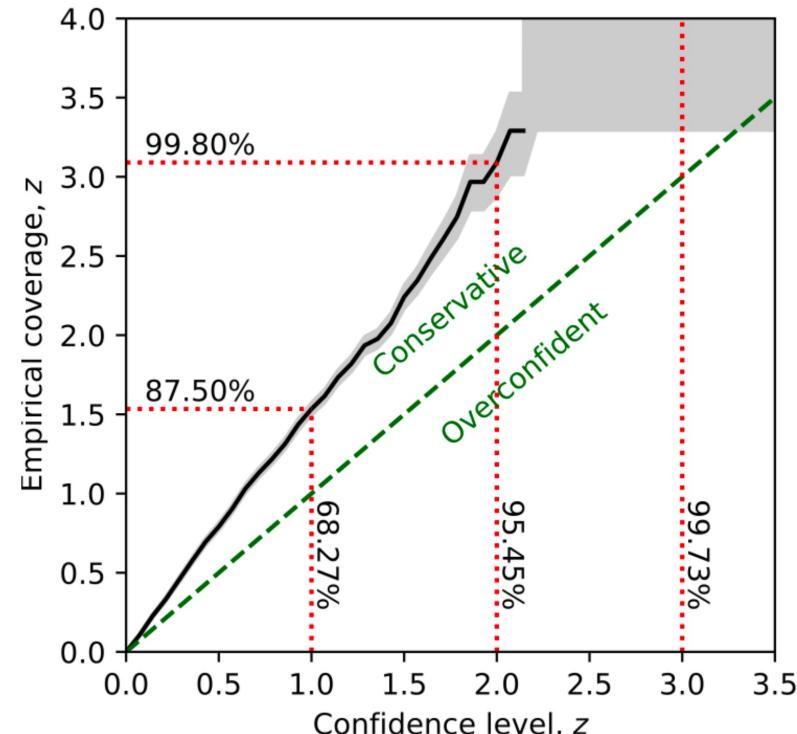
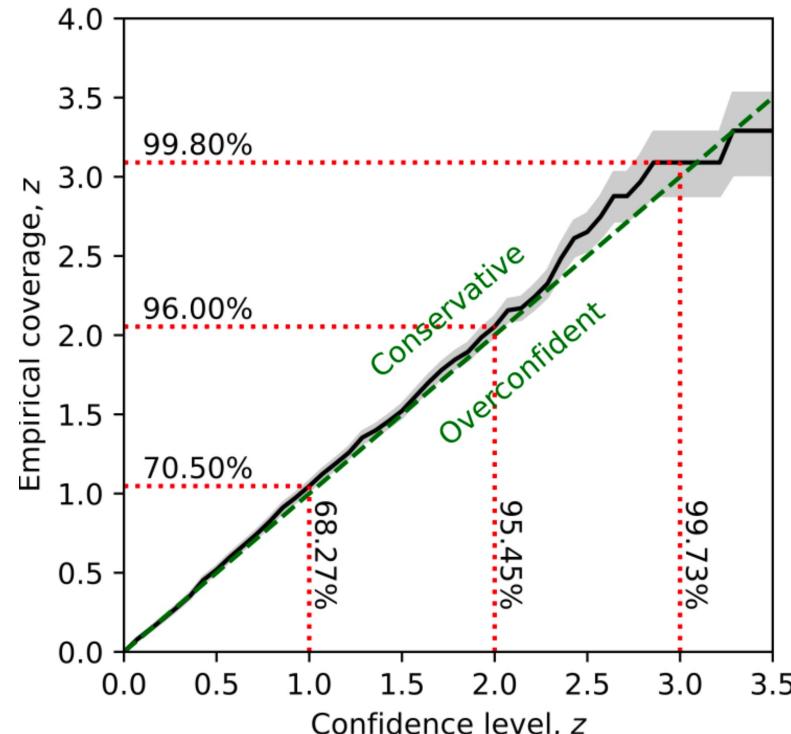
- Example: 10-dim eggbox



TMNRE coverage test



- SBC could be used to test posteriors in equivalent way as for NLE
- Another idea is to take all posteriors, and calculate confidence region for the true sample



TMNRE coverage test

- Example: CMB
 - It can be run for any 1D or 2D posterior



TMNRE vs. (Sequential) NLE

Similarities:

- *Amortization* – trained on full prior; once trained, any data can be passed
- *Sampling* – after training, both methods efficiently sample posterior (MCMC)
- *Compression* – both can have advantage of it, but trained differently
- *Marginalization* – both can target marginal posteriors (NLE with “nuisance-hardened” compression)
- *Active learning* – truncation vs. sequential posteriors / acquisition function

Differences:

- *Ratio vs. density estimation* – classifier is usually easier to train than a density estimator; only with a density estimator one can sample likelihood
- *Marginalization and active learning interplay* – truncation might work better for nuisance parameters than nuisance-hardened compression