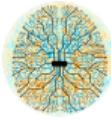


# Accelerated nested sampling with $\beta$ -flows

Metha Prathaban  
myp23@cam.ac.uk





## About Me

- ▶ 3rd year PhD student
- ▶ Work on Bayesian numerical method development in context of GWs

Current work is in collaboration with Will Handley and Harry Bevins.





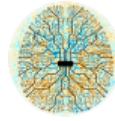
# Table of Contents

Bayesian Inference

Nested sampling

Accelerating NS

$\beta$ -flows



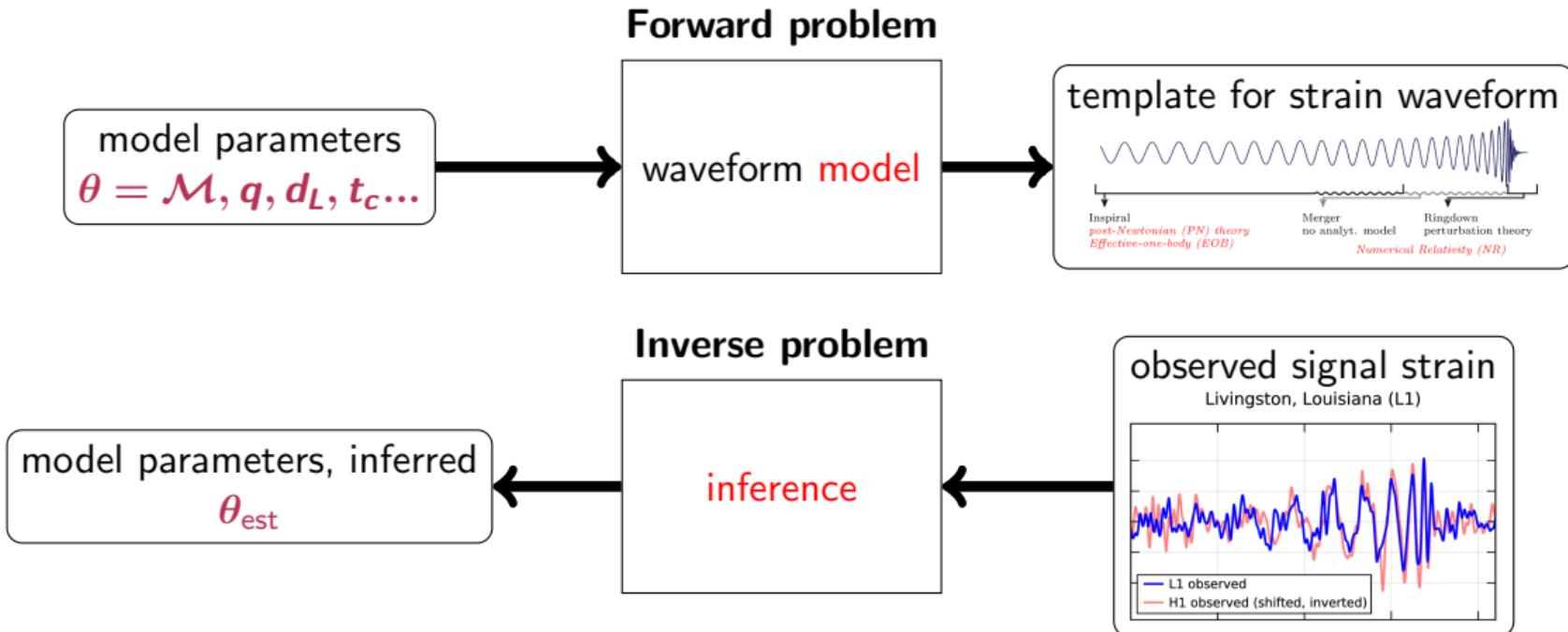
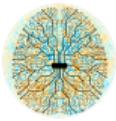
# Table of Contents

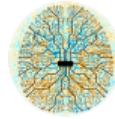
Bayesian Inference

Nested sampling

Accelerating NS

$\beta$ -flows





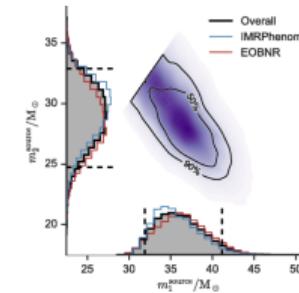
- ▶ One naive way to tackle this is to directly fit the waveform model to the data by minimizing the differences between the two.
- ▶ Produces a single set of best-fit parameters.

### Issues:

- ▶ **Degenerate solutions** - different combinations of parameters can produce near-identical signals (e.g. mass spin degeneracy).
- ▶ **Noisy data** - can lead to biased parameter estimates.

- ▶ Instead of single “best-fit” solution, obtain **posterior distribution** over possible parameters.

$$\{m_1, m_2\}^{\text{best-fit}} = \{36.3, 28.6\}$$



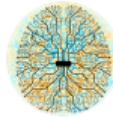
- ▶ Accounts for uncertainty in data and model.
- ▶ Even in the presence of noise or degeneracies, provides a complete probabilistic picture of parameter estimates.

Bayesian inference provides a robust framework for doing this!

Given some model  $\mathcal{M}$  and observed signal  $\mathcal{D}$ , Bayes' theorem enables us to relate the **posterior** probability of the set of parameters  $\theta$  which generated the signal to the **likelihood** of the  $\mathcal{D}$  given  $\theta$  and the **prior** probability of  $\theta$  given  $\mathcal{M}$ :

$$P(\theta|\mathcal{D}, \mathcal{M}) = \frac{P(\mathcal{D}|\theta, \mathcal{M})P(\theta|\mathcal{M})}{P(\mathcal{D}|\mathcal{M})} = \frac{\mathcal{L}(\mathcal{D}|\theta)\pi(\theta)}{\mathcal{Z}} \quad (1)$$

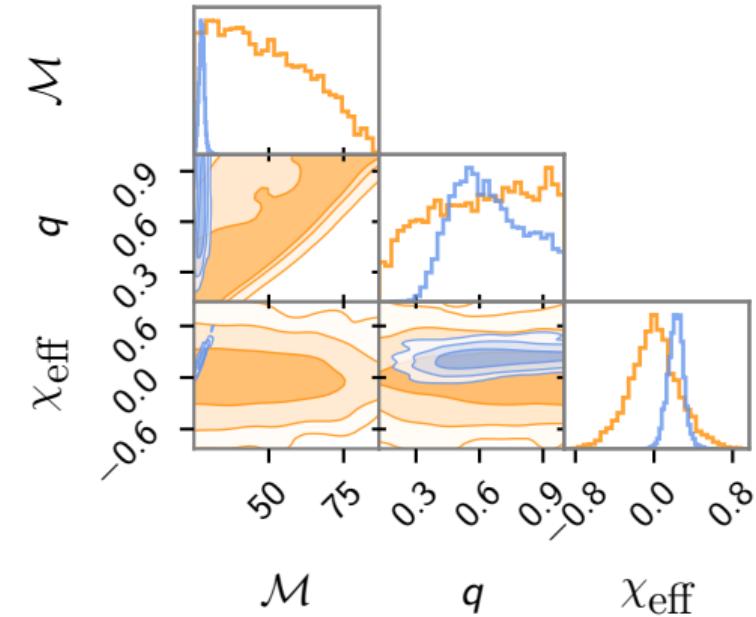
The **evidence**,  $\mathcal{Z}$ , plays a key role in model comparison.



- ▶ Define **prior**, **sample** (unnormalized) posterior ( $\mathcal{L}(D|\theta) \times \pi(\theta)$ ).

Challenges:

- ▶ High-dimensional parameter spaces  $\Rightarrow$  posterior occupies vanishingly small region of prior.
- ▶ Complex likelihoods with high costs



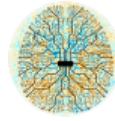


Goal is efficient exploration of parameter space, to do GW inference in feasible timescales.

**Posterior** samplers:

- ▶ Metropolis-Hastings
- ▶ Hamiltonian Monte-Carlo
- ▶ Ensemble samplers

None of these calculate the **evidence**,  $\mathcal{Z}$  - crucial for Bayesian model comparison (e.g. testing for precession vs. no precession)!



# Table of Contents

Bayesian Inference

Nested sampling

Accelerating NS

$\beta$ -flows



Nested sampling first and foremost calculates this **evidence**. The **evidence** is the integral of likelihood  $\times$  prior over the entire parameter space,

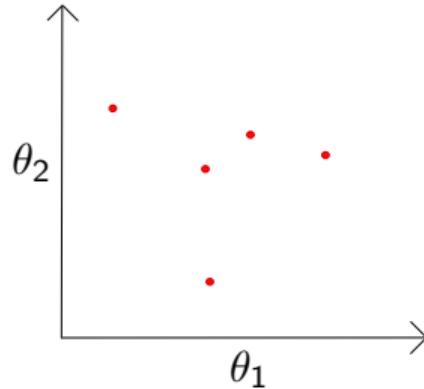
$$\mathcal{Z} = \int \mathcal{L}(\theta) \pi(\theta) d\theta, \quad (2)$$

which, in general, is a many dimensional integral.

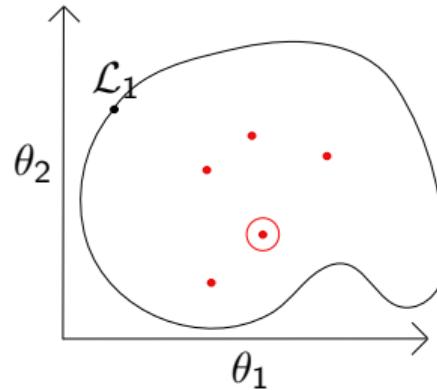
NS turn this into a 1D problem, performing this integral by summing over nested likelihood contours in the parameter space.

Nested sampling first and foremost calculates **evidence**,  $\mathcal{Z} = \int \mathcal{L}(\theta)\pi(\theta)d\theta$ .

- ▶ Prior is populated with set of ‘live points’.

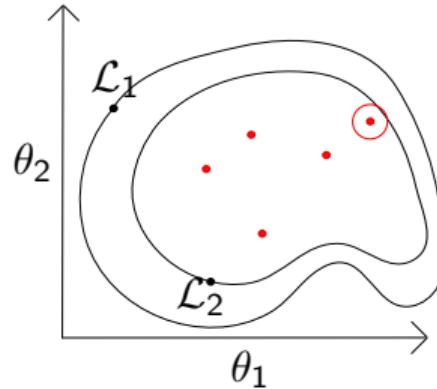


Nested sampling first and foremost calculates **evidence**,  $\mathcal{Z} = \int \mathcal{L}(\theta) \pi(\theta) d\theta$ .



- ▶ Prior is populated with set of ‘live points’.
- ▶ At each iteration  $i$ , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.

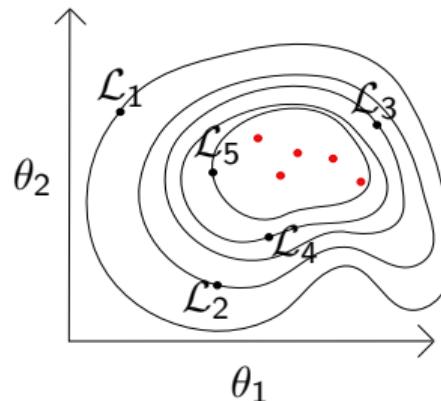
Nested sampling first and foremost calculates **evidence**,  $\mathcal{Z} = \int \mathcal{L}(\theta) \pi(\theta) d\theta$ .



- ▶ Prior is populated with set of ‘live points’.
- ▶ At each iteration  $i$ , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.

## Nested sampling (NS)

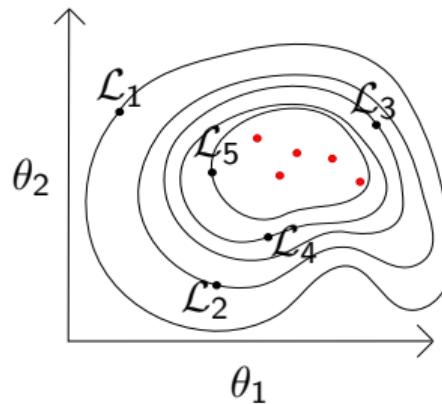
Nested sampling first and foremost calculates **evidence**,  $\mathcal{Z} = \int \mathcal{L}(\theta) \pi(\theta) d\theta$ .



- ▶ Prior is populated with set of ‘live points’.
- ▶ At each iteration  $i$ , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.
- ▶ Live points compress exponentially towards peak of likelihood.

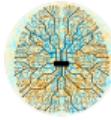
## Nested sampling (NS)

Nested sampling first and foremost calculates **evidence**,  $\mathcal{Z} = \int \mathcal{L}(\theta) \pi(\theta) d\theta$ .



- ▶ Prior is populated with set of ‘live points’.
- ▶ At each iteration  $i$ , point is lowest likelihood is deleted and new live point is drawn, which must have a likelihood higher than that of the deleted point.
- ▶ Live points compress exponentially towards peak of likelihood.
- ▶ **Evidence** is calculated as weighted sum over deleted ('dead') points.

- ▶ This is not trivial!
- ▶ Different samplers implement this step in different ways, but most samplers can broadly be split into 2 categories:
  1. Region samplers (e.g. MULTINEST, NESSAI, ULTRANEST)
    - ▶ Construct a bounding region for the deleted point's likelihood contour and sample within this.
    - ▶ Generates invalid points that have to be discarded.
  2. Path samplers (e.g. POLYCHORD, CPNEST, DYNESTY)
    - ▶ Run a Markov chain from the deleted point.
    - ▶ Generates correlated (but valid) points that have to be discarded.

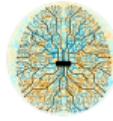


- ▶ The Kullback Liebler divergence between the **prior** and **posterior** is defined as:

$$\mathcal{D}_{\text{KL}} = \left\langle \log \frac{\mathcal{P}}{\pi} \right\rangle_{\mathcal{P}} \approx \log \frac{V_{\mathcal{P}}}{V_{\pi}} \quad (3)$$

- ▶ “Amount of compression from **prior** to **posterior**”





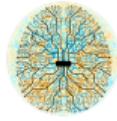
# Table of Contents

Bayesian Inference

Nested sampling

Accelerating NS

$\beta$ -flows



likelihood evaluation time

$$T \propto T_{\mathcal{L}} \times n_{\mathcal{L}}$$

(4)

number of likelihood evaluations



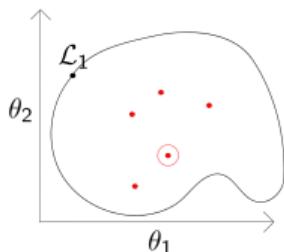
likelihood evaluation time

number of iterations

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times n_{\text{iter}}$$

(4)

drawing new live point





likelihood evaluation time

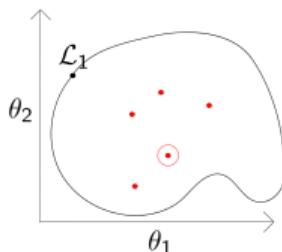
 $T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times \mathcal{D}_{\text{KL}} \times n_{\text{live}}$ 

resolution

(4)

compression between prior and posterior ( $\approx \log \frac{V_{\pi}}{V_{\mathcal{P}}}$ )

drawing new live point





likelihood evaluation time

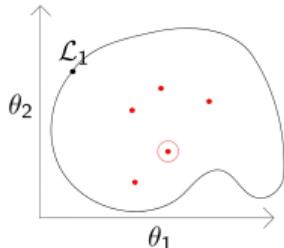
$$T \propto T_{\mathcal{L}}$$

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times \mathcal{D}_{\text{KL}} \times n_{\text{live}} \quad (4)$$

resolution

focus of this talk!

drawing new live point





likelihood evaluation time

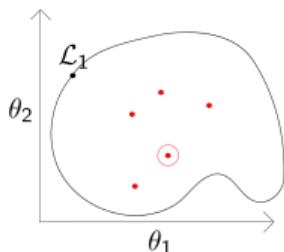
$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times$$

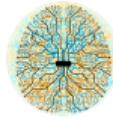
$$\mathcal{D}_{\text{KL}} \times n_{\text{live}}$$

(5)

compression between prior and posterior ( $\approx \log \frac{V_{\pi}}{V_{\mathcal{P}}}$ )

drawing new live point





likelihood evaluation time

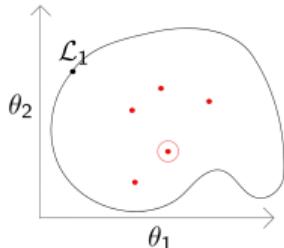
resolution (baked in)

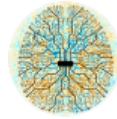
$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times \mathcal{D}_{\text{KL}} \times n_{\text{live}}$$

(5)

compression between prior and posterior ( $\approx \log \frac{V_{\pi}}{V_{\mathcal{P}}}$ )

drawing new live point





faster waveform models

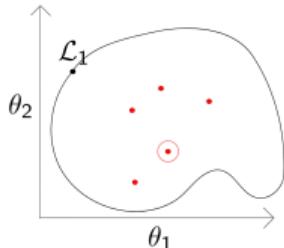
resolution (baked in)

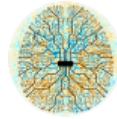
$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times \mathcal{D}_{\text{KL}} \times n_{\text{live}}$$

(5)

compression between prior and posterior ( $\approx \log \frac{V_{\pi}}{V_{\mathcal{P}}}$ )

drawing new live point





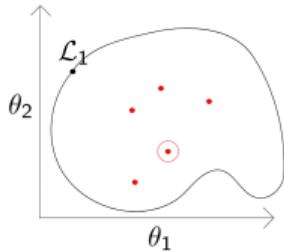
$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times \mathcal{D}_{\text{KL}} \times n_{\text{live}} \quad (5)$$

faster waveform models

better samplers

resolution (baked in)

compression between prior and posterior ( $\approx \log \frac{V_{\pi}}{V_{\mathcal{P}}}$ )





## Time of convergence of NS

$$T \propto T_{\mathcal{L}} \times f_{\text{sampler}} \times D_{\text{KL}} \times n_{\text{live}} \quad (6)$$

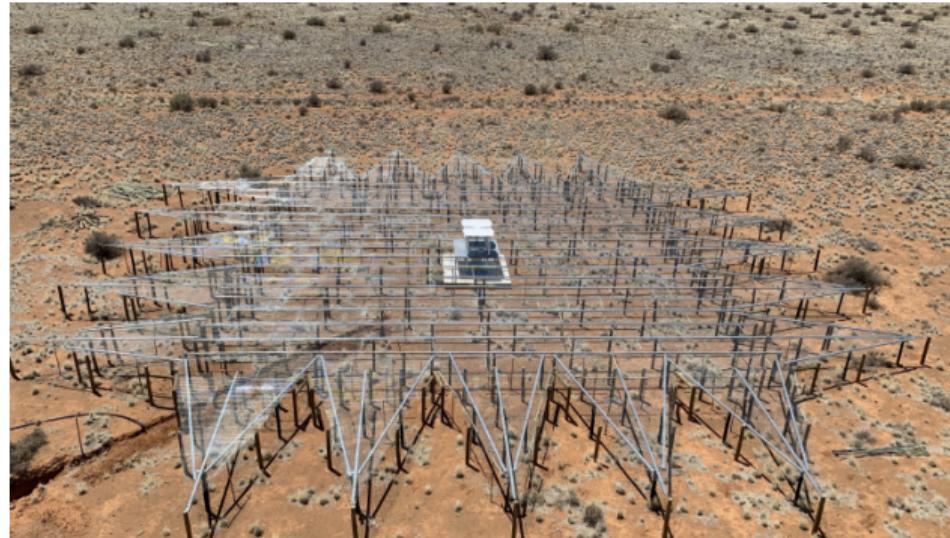
## Uncertainty in $\log \mathcal{Z}$

$$\sigma \propto \sqrt{D_{\text{KL}} / n_{\text{live}}} \quad (7)$$

Precision-normalized runtime has quadratic dependence on KL divergence. 2212.01760



One way to do this (REACH):



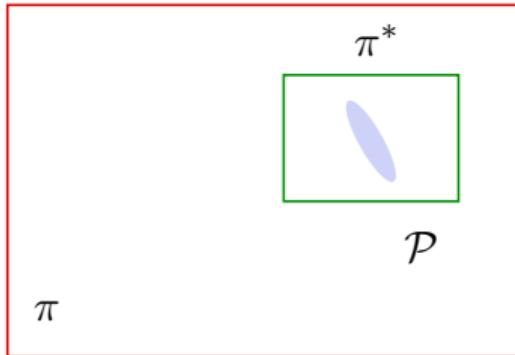


One way to do this (REACH):



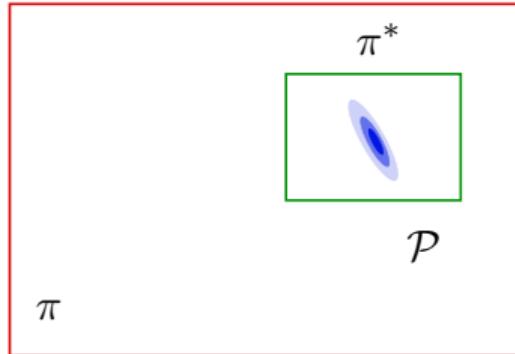
- ▶ Perform low resolution (low live points) run first to roughly identify where **posterior** lies.

One way to do this (REACH):

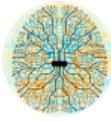


- ▶ Perform low resolution (low live points) run first to roughly identify where **posterior** lies.
- ▶ Then set off second, high resolution, run with **narrower** box **prior** (much quicker).

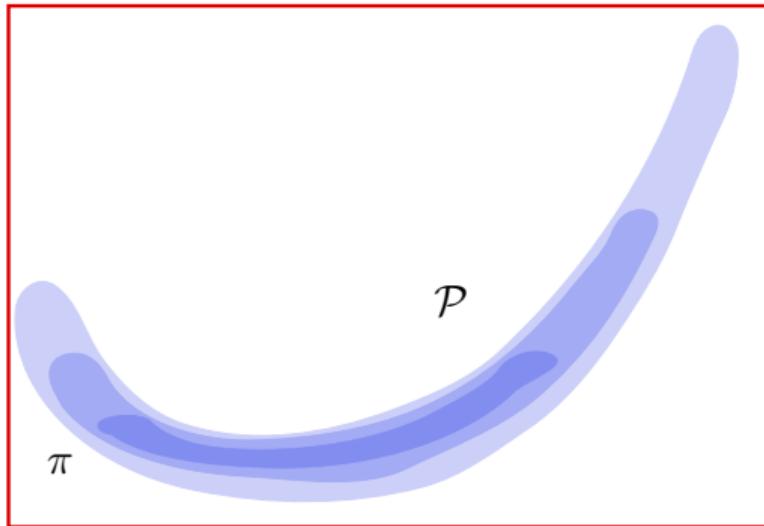
One way to do this (REACH):



- ▶ Perform low resolution (low live points) run first to roughly identify where **posterior** lies.
- ▶ Then set off second, high resolution, run with **narrower** box **prior** (much quicker).
- ▶ **Evidence** has **changed** (since different prior), but easy to correct (multiply new evidence by  $\frac{V_{\pi^*}}{V_\pi}$ )



# When does this break down?

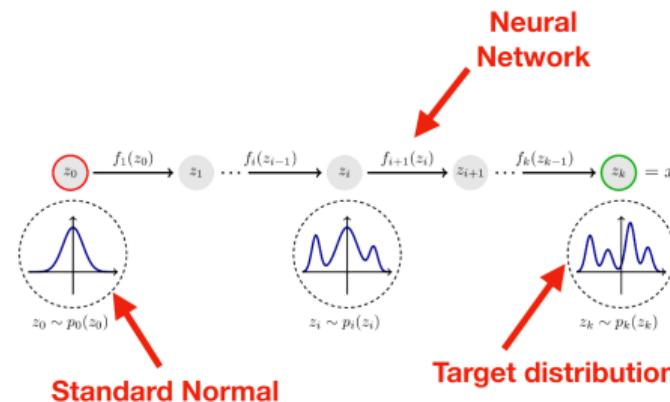


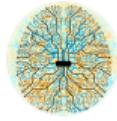
- ▶ Banana distributions, multi-modality etc.
- ▶ Precludes its use in most realistic GW cases...



- ▶ Can iterate on this by using **normalizing flows** (NF) to learn the rough **posterior**.

- ▶ Can iterate on this by using **normalizing flows** (NF) to learn the rough **posterior**.
- ▶ NFs perform density estimation, by learning a series of invertible mappings from the standard normal distribution to the target (posterior).





- ▶ Use **normalizing flows** (NF) to learn the rough **posterior**, and use this as our updated prior,  $\pi^*$ .
- ▶ In this case, can't do our trick of correcting the second **evidence** by volume ratio,  $\frac{V_{\pi^*}}{V_\pi}$ !
- ▶ Must rely on another technique to get around this!

- ▶ Use **normalizing flows** (NF) to learn the rough **posterior**, and use this as our updated prior,  $\pi^*$ .
- ▶ In this case, can't do our trick of correcting the second **evidence** by volume ratio,  $\frac{V_{\pi^*}}{V_\pi}$ !
- ▶ Must rely on another technique to get around this!

**Posterior repartitioning** (PR) can help us with this! (see e.g. 2212.01760)

Bayesian Analysis (0000)

00, Number 0, pp. 1

## Bayesian posterior repartitioning for nested sampling

Xi Chen<sup>\*†</sup>, Farhan Feroz<sup>‡</sup> and Michael Hobson<sup>†</sup>

Improving the efficiency and robustness of nested sampling using posterior repartitioning

Xi Chen · Michael Hobson · Saptarshi Das · Paul Gelderblom



Article

SuperNest: accelerated nested sampling applied to astrophysics and cosmology<sup>†</sup>

Aleksandr Petrosyan<sup>1,2,3†</sup> & Will Handley<sup>1,2,4†</sup>



- Evidence and posterior only depend on product of  $\mathcal{L}$  and  $\pi$ :

$$\mathcal{Z} = \int \mathcal{L}(\theta) \pi(\theta) d\theta \quad (8)$$

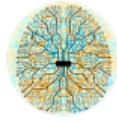
$$\mathcal{P}(\theta) = \frac{\mathcal{L}(\theta) \pi(\theta)}{\mathcal{Z}} \quad (9)$$

We are free to redefine the likelihood and prior however we like - as long as the product is the same! arXiv:1908.04655

$$\tilde{\mathcal{Z}} = \int \tilde{\mathcal{L}}(\theta) \tilde{\pi}(\theta) d\theta = \int \mathcal{L}(\theta) \pi(\theta) d\theta = \mathcal{Z} \quad (10)$$

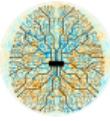


- ▶ Many sampling algorithms do not distinguish between  $\mathcal{L}$  and  $\pi$  at the algorithmic level.
- ▶ e.g. Metropolis-Hastings acceptance ratio only depends on the **joint distribution**,  $\mathcal{L}(\theta)\pi(\theta)$ .
- ▶ Nested sampling does distinguish between prior and likelihood at the algorithmic level, by '**sampling from the prior  $\pi$ , subject to the hard likelihood constraint,  $\mathcal{L}$** '.
- ▶  $\mathcal{Z}$  and  $\mathcal{P}$  will not change if we repartition  $\mathcal{L}$  and  $\pi$ , **but  $\mathcal{D}_{KL}$  will**.



PR-NS w/ NFs

$$\pi(\theta) \longrightarrow \text{NF}(\theta)$$



$$\pi(\theta) \longrightarrow \text{NF}(\theta)$$

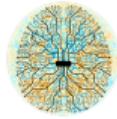
$$\mathcal{L}(\theta) \longrightarrow \frac{\mathcal{L}(\theta)\pi(\theta)}{\text{NF}(\theta)}$$



$$\pi(\theta) \longrightarrow \text{NF}(\theta)$$

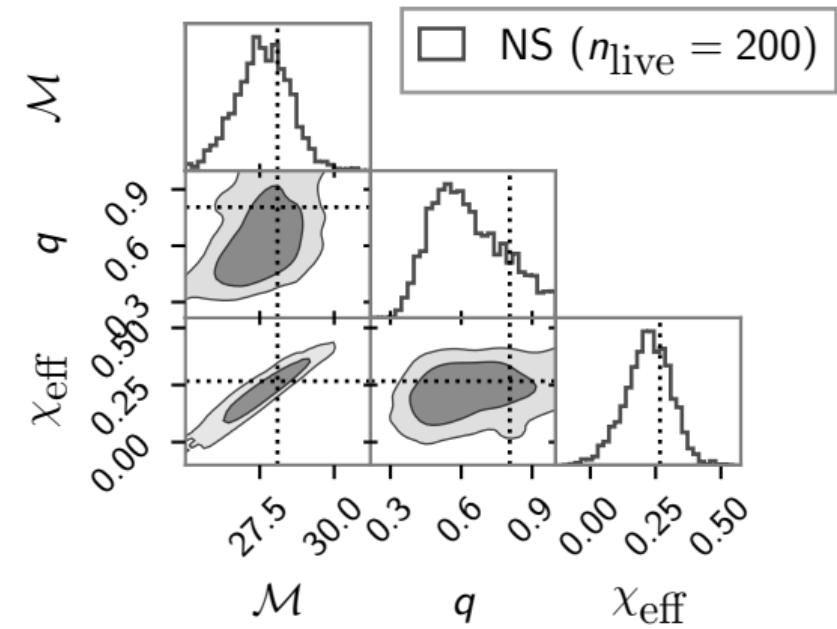
$$\mathcal{L}(\theta) \longrightarrow \frac{\mathcal{L}(\theta)\pi(\theta)}{\text{NF}(\theta)}$$

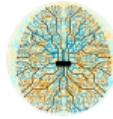
$$\mathcal{D}_{\text{KL}} \approx \log \frac{V_{\text{NF}}}{V_{\mathcal{P}}}$$



## Demo on simulated example

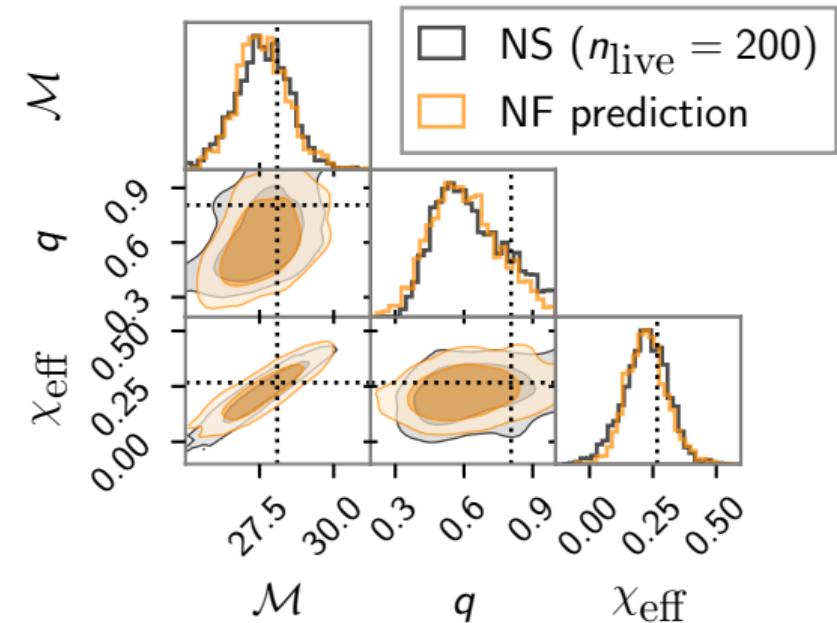
- ▶ Perform low resolution run on simulated data.





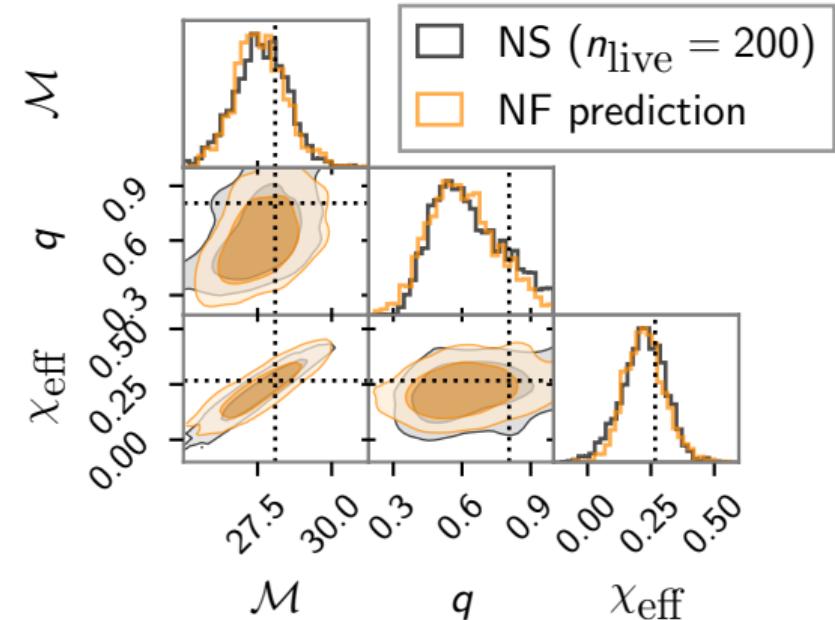
## Demo on simulated example

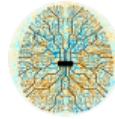
- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.



## Demo on simulated example

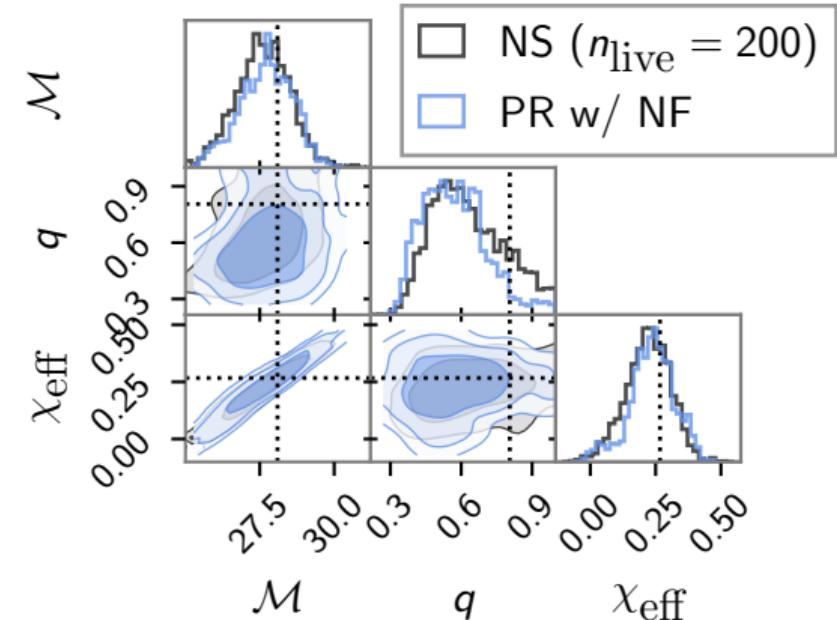
- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.
- ▶ Use this as ‘repartitioned prior’ for new high resolution run (using PR to also update likelihood accordingly to same evidences and posteriors out).





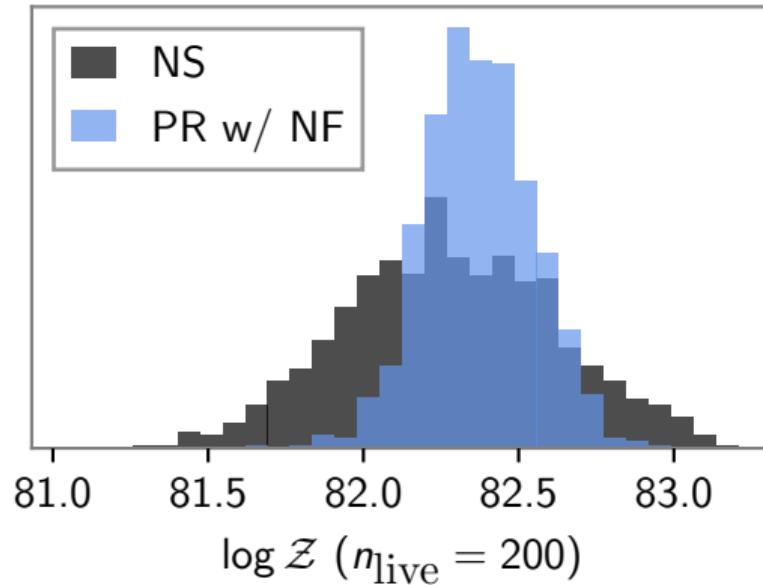
## Demo on simulated example

- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.
- ▶ Use this as ‘repartitioned prior’ for new high resolution run (using PR to also update likelihood accordingly to same evidences and posteriors out).

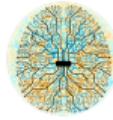


## Demo on simulated example

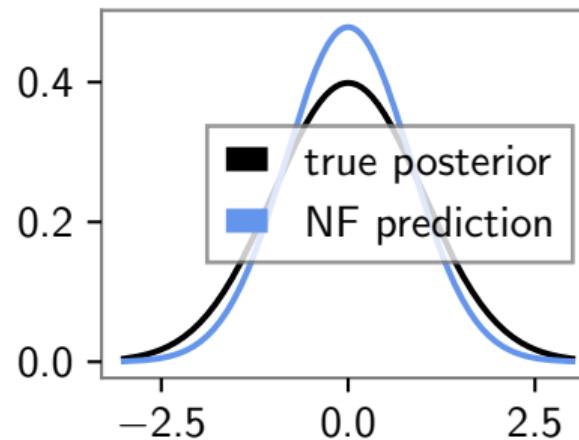
- ▶ Perform low resolution run on simulated data.
- ▶ Train NF on the weighted samples.
- ▶ Use this as ‘repartitioned prior’ for new high resolution run (using PR to also update likelihood accordingly to same evidences and posteriors out).

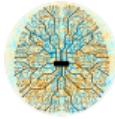


Same answer as doing a full resolution pass of NS, but **7x faster** (precision-normalized).

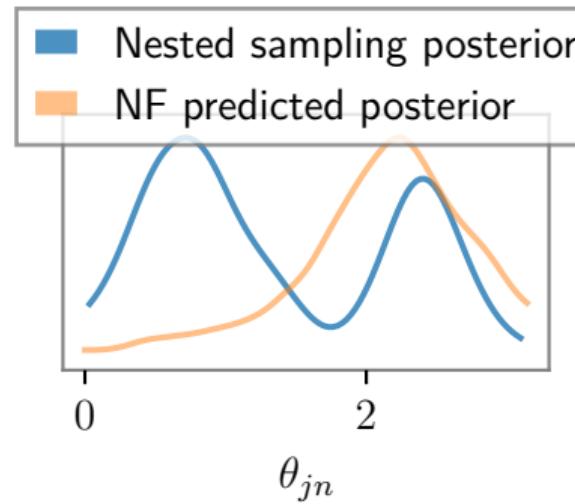


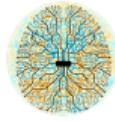
- ▶ If NF learns something **narrower** than true posterior, we can make the problem harder!



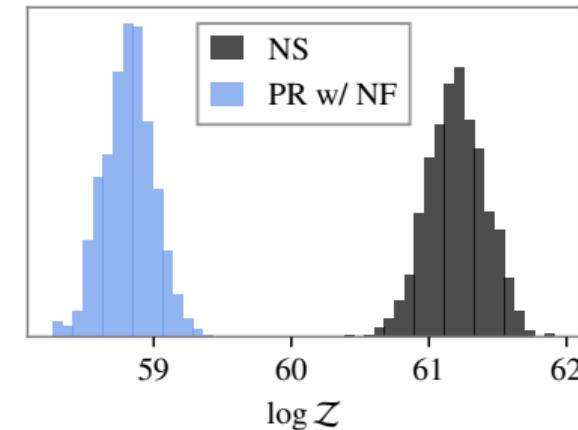
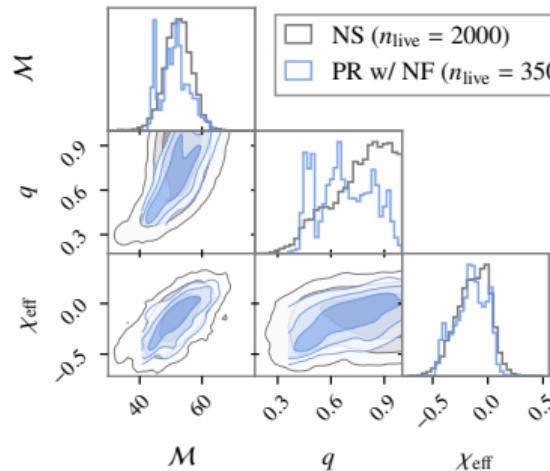


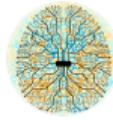
- ▶ Still have an issue with multi-modality (if NF only learns one mode, the others are cut off at the prior level in the high resolution run).





When the NF has been unable to properly learn the multi-modality, we can get biased posteriors and evidences (GW191222):

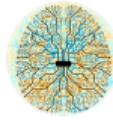




## Improving the method

In order to improve the robustness of the method:

- ▶ Repartitioned prior should ideally be able to **widen itself adaptively at runtime** to mitigate missed modes and badly learned posteriors.



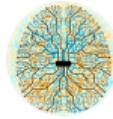
# Table of Contents

Bayesian Inference

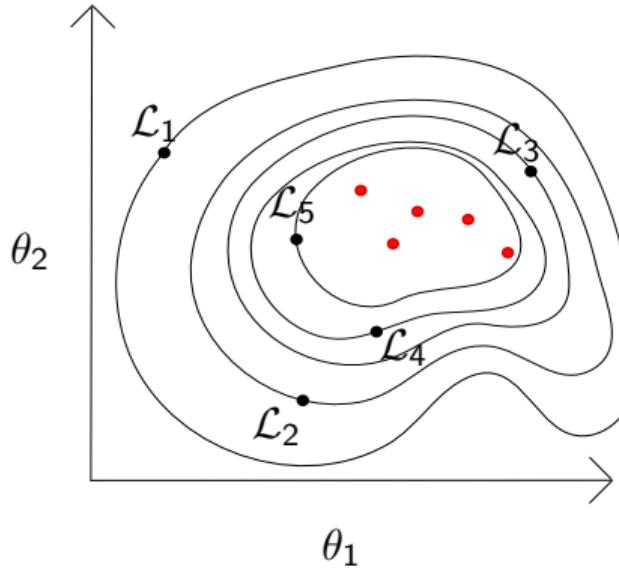
Nested sampling

Accelerating NS

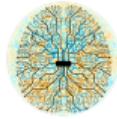
$\beta$ -flows



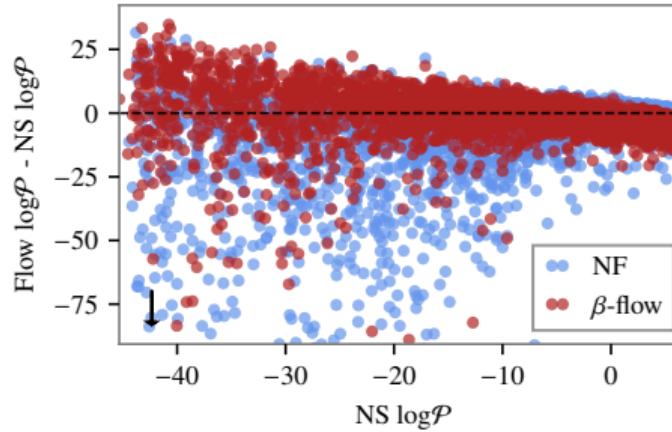
## $\beta$ -flows vs typical NFs



- ▶ Nested sampling sees tip to tail of the posterior in a systematic way, and NS has deep tails.
- ▶ NS can be used to train a specialized form of **conditional NFs** that can better learn these deep tail events.
- ▶  $\beta$ -flows are new and only used in this work so far, though broadly applicable.

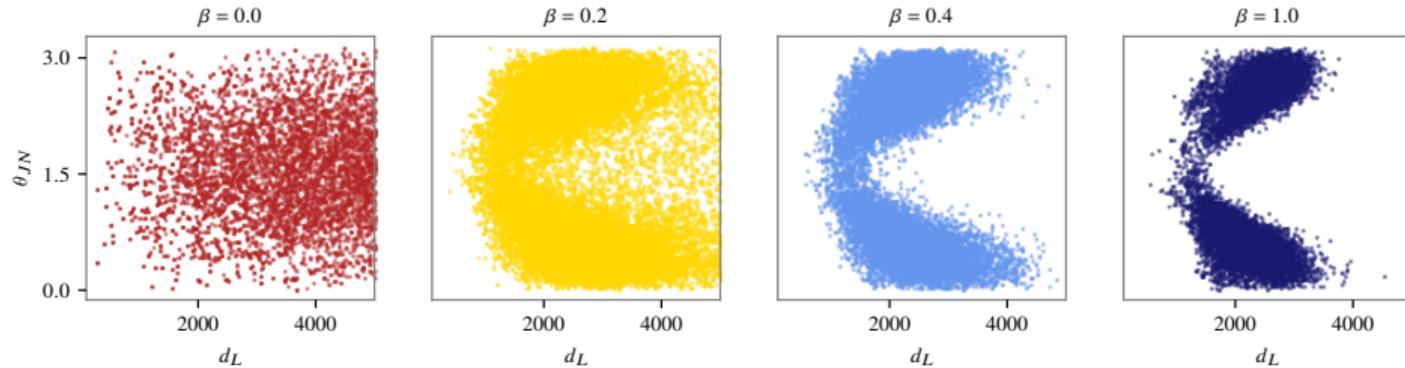


## Better at deep tail probabilities



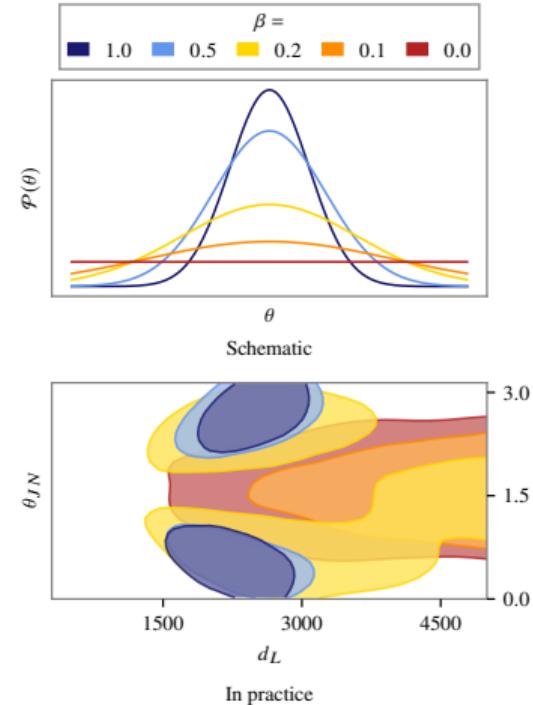
- ▶ For simulated example shown before, the  $\beta$ -flow is able to better predict the NS posterior probabilities.
- ▶  $\beta$ -flows exhibit less scatter in the tails (low posterior probabilities) than the NFs.

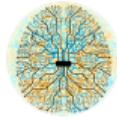
- ▶ NS compresses step by step from prior to posterior.
- ▶ We can label these stages by a parameter  $\beta$  (akin to inverse temperature  $\beta$  in e.g. materials science).
- ▶ Sliding scale from  $\beta = 0$  as the prior and  $\beta = 1$  as the posterior.



- ▶ We **sample over**  $\beta$  at runtime
- ▶ We can set the repartitioned prior to be anywhere between the posterior ( $\beta = 1$ ) and the original prior ( $\beta = 0$ ).
- ▶ Can **widen themselves adaptively** at runtime

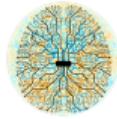
$$p(\beta) \propto \mathcal{L}^\beta \pi, \quad \beta \in [0, 1] \quad (11)$$



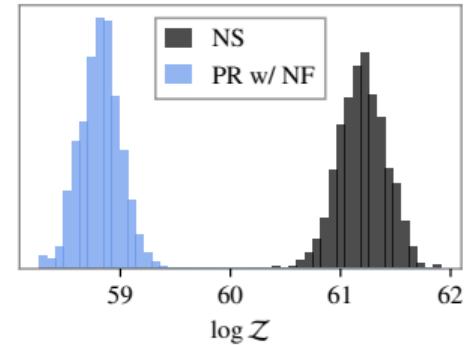
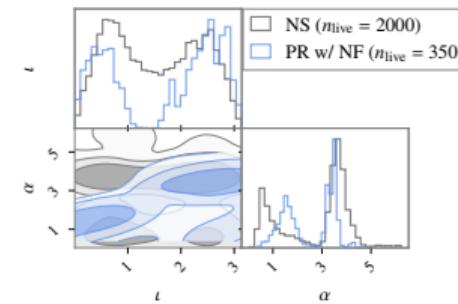
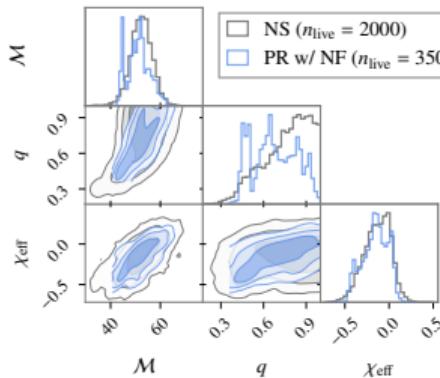


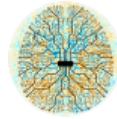
## GW191222 (again)

- ▶ Using  $\beta$ -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over  $\beta$  now fixes the problem.

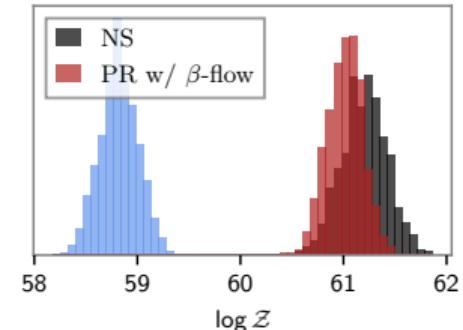
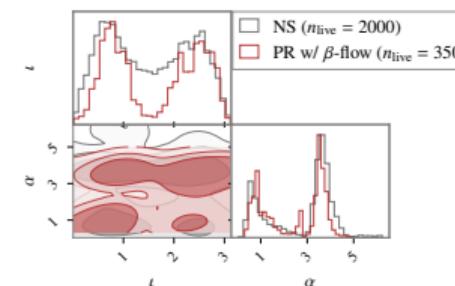
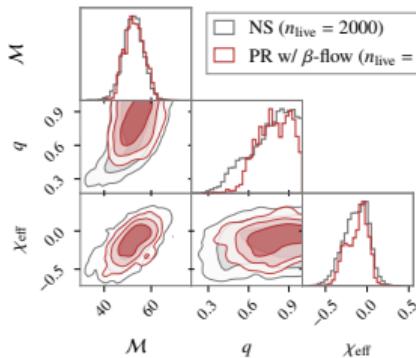


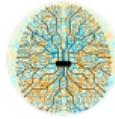
- ▶ Using  $\beta$ -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over  $\beta$  now fixes the problem.



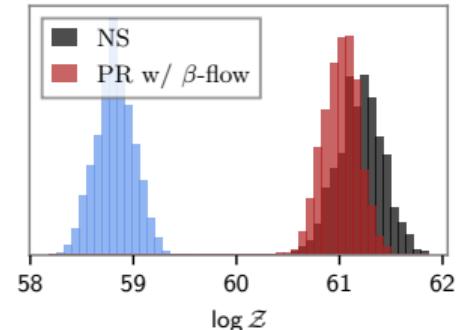
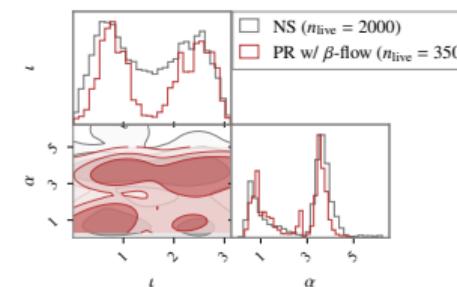
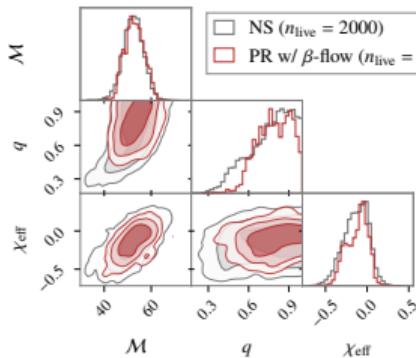


- ▶ Using  $\beta$ -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over  $\beta$  now fixes the problem.

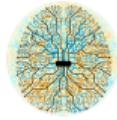




- ▶ Using  $\beta$ -flows to set the repartitioned prior for the high resolution run, instead of a typical NF, and sampling over  $\beta$  now fixes the problem.



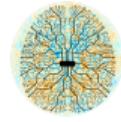
Only 2x (precision-normalized) as fast as normal NS for this real example, but robust.



## Summary

- ▶ Several ways to reduce NS runtime, including reducing amount of compression from prior to posterior.
  - ▶ Can perform low resolution run to identify rough posterior, learn distribution with a flow, and perform high resolution run with updated prior.
  - ▶ Use posterior repartitioning to get correct **evidences** out, despite changed prior.
  - ▶ Can achieve order of magnitude speedups on realistic GW examples.

- ▶ Several ways to reduce NS runtime, including reducing amount of compression from prior to posterior.
  - ▶ Can perform low resolution run to identify rough posterior, learn distribution with a flow, and perform high resolution run with updated prior.
  - ▶ Use posterior repartitioning to get correct **evidences** out, despite changed prior.
  - ▶ Can achieve order of magnitude speedups on realistic GW examples.
- ▶ Introduced  $\beta$ -flows:
  - ▶ Conditional normalizing flow, trained with whole NS run
  - ▶ Better at deep tail events
  - ▶ First application is in this paper, but their use is much broader!



---

Thank you for listening!