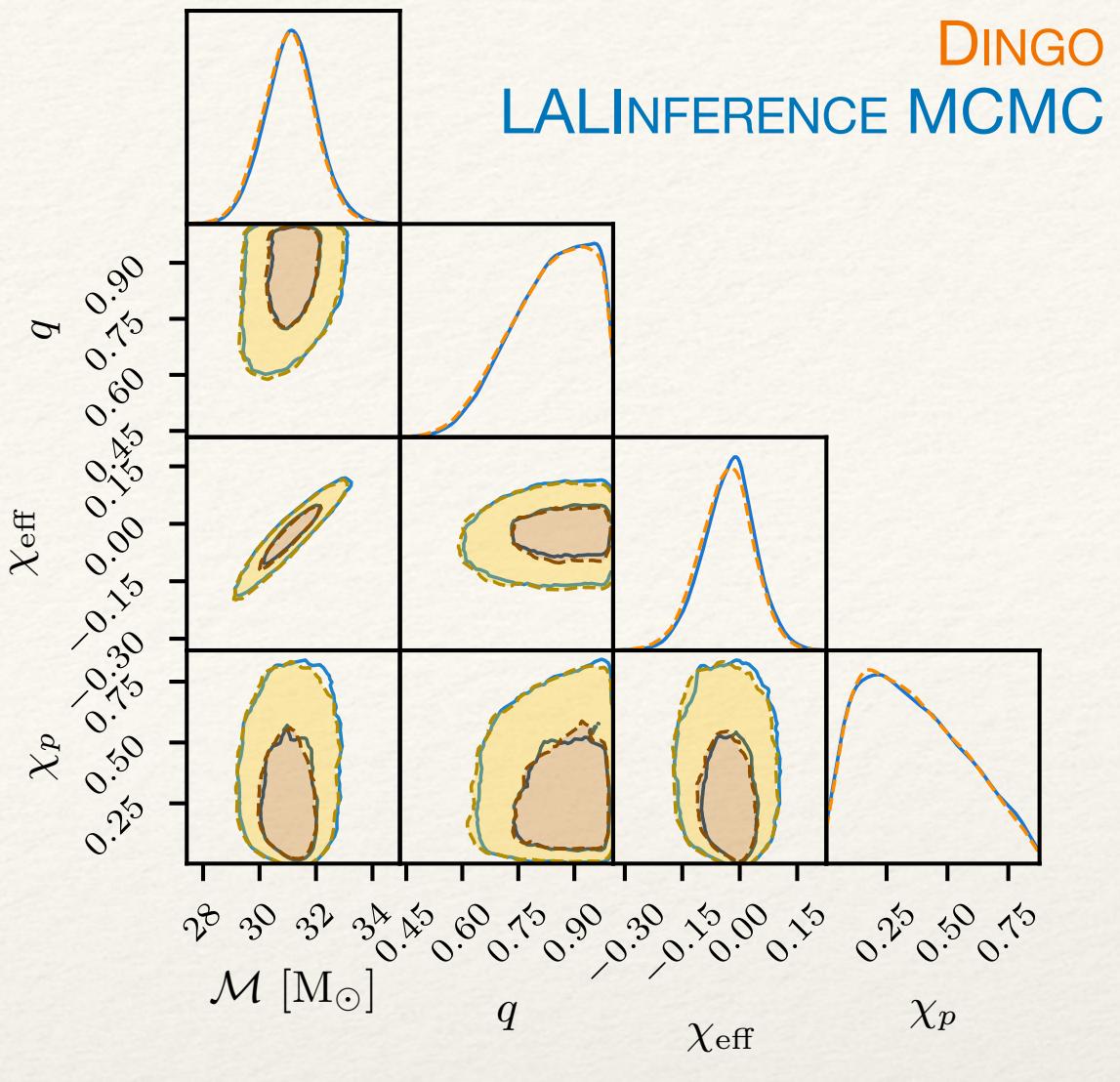


Scientific Machine Learning for Gravitational Wave Astronomy, ICERM
4 June 2025

Neural Posterior Estimation for Gravitational Waves

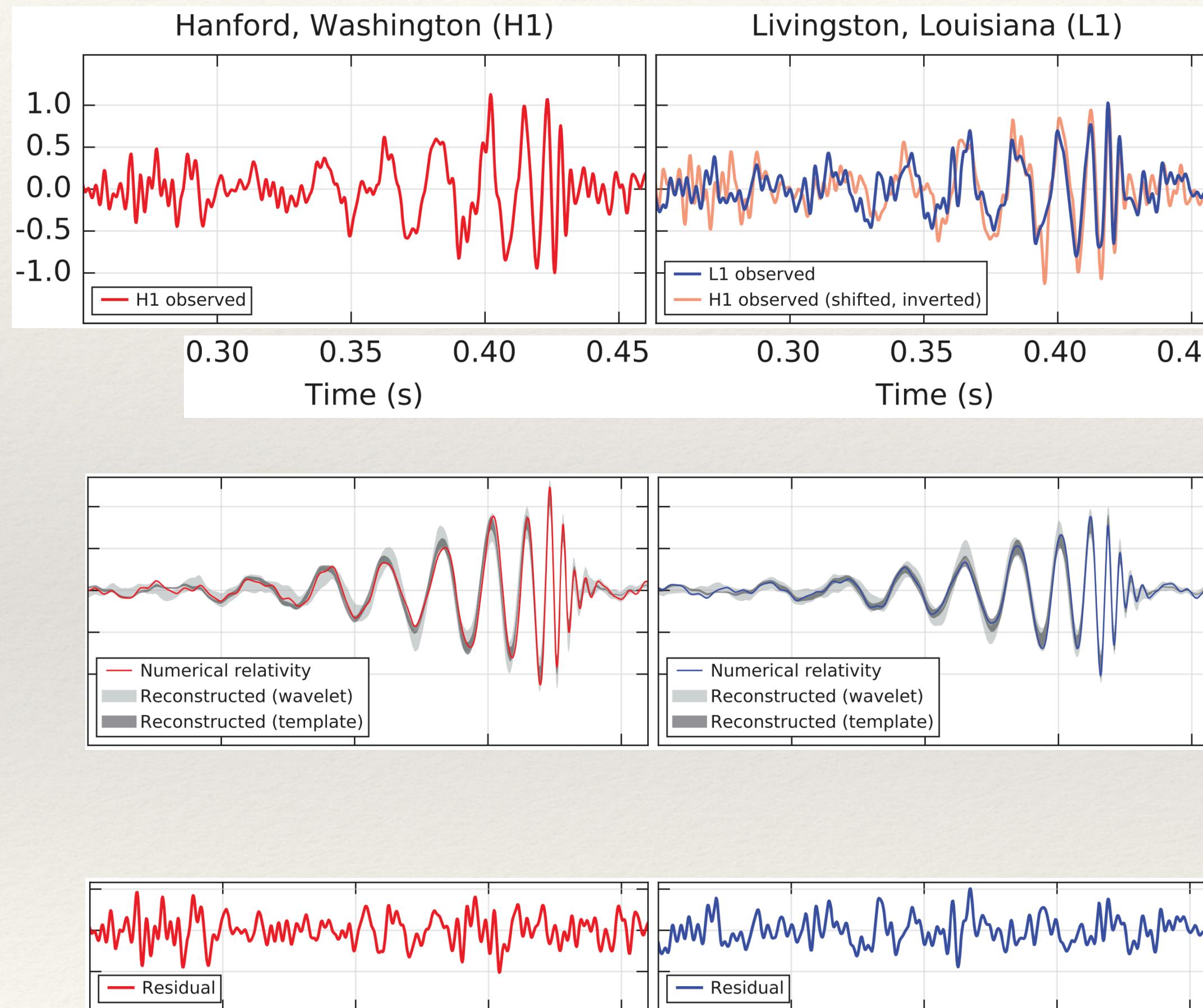
Stephen Green



Outline

- ❖ Simulation-based inference for gravitational waves
- ❖ Three case studies
 - ❖ Eccentricity
 - ❖ Binary neutron stars
 - ❖ Flow matching
- ❖ DINGO tutorial

Parameter estimation for gravitational waves



observed data

=

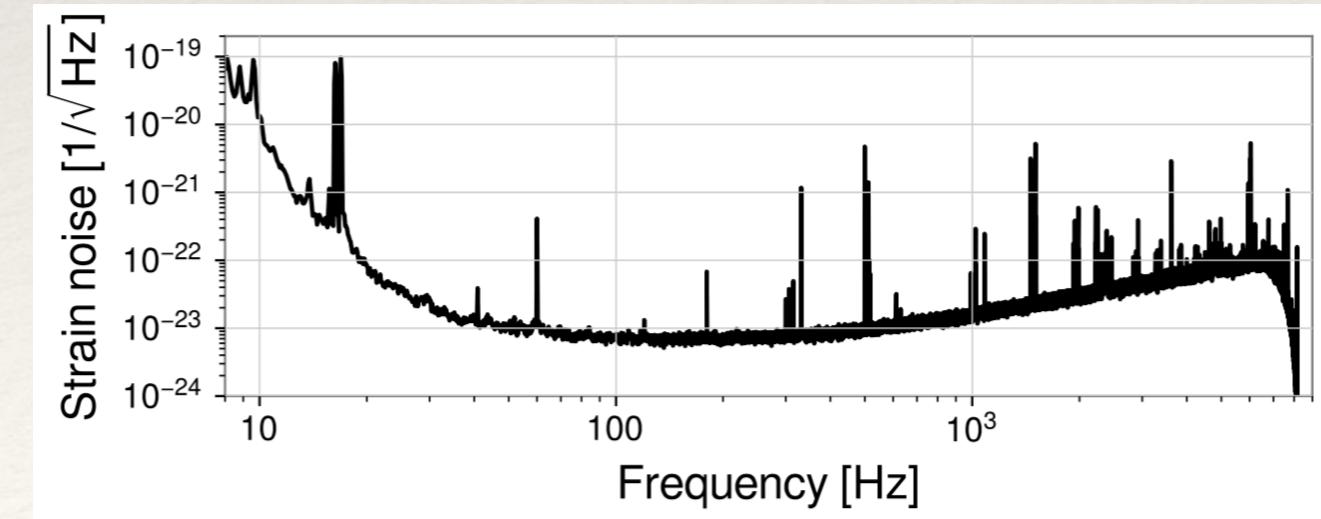
signal

+

noise

d

$h(\theta)$



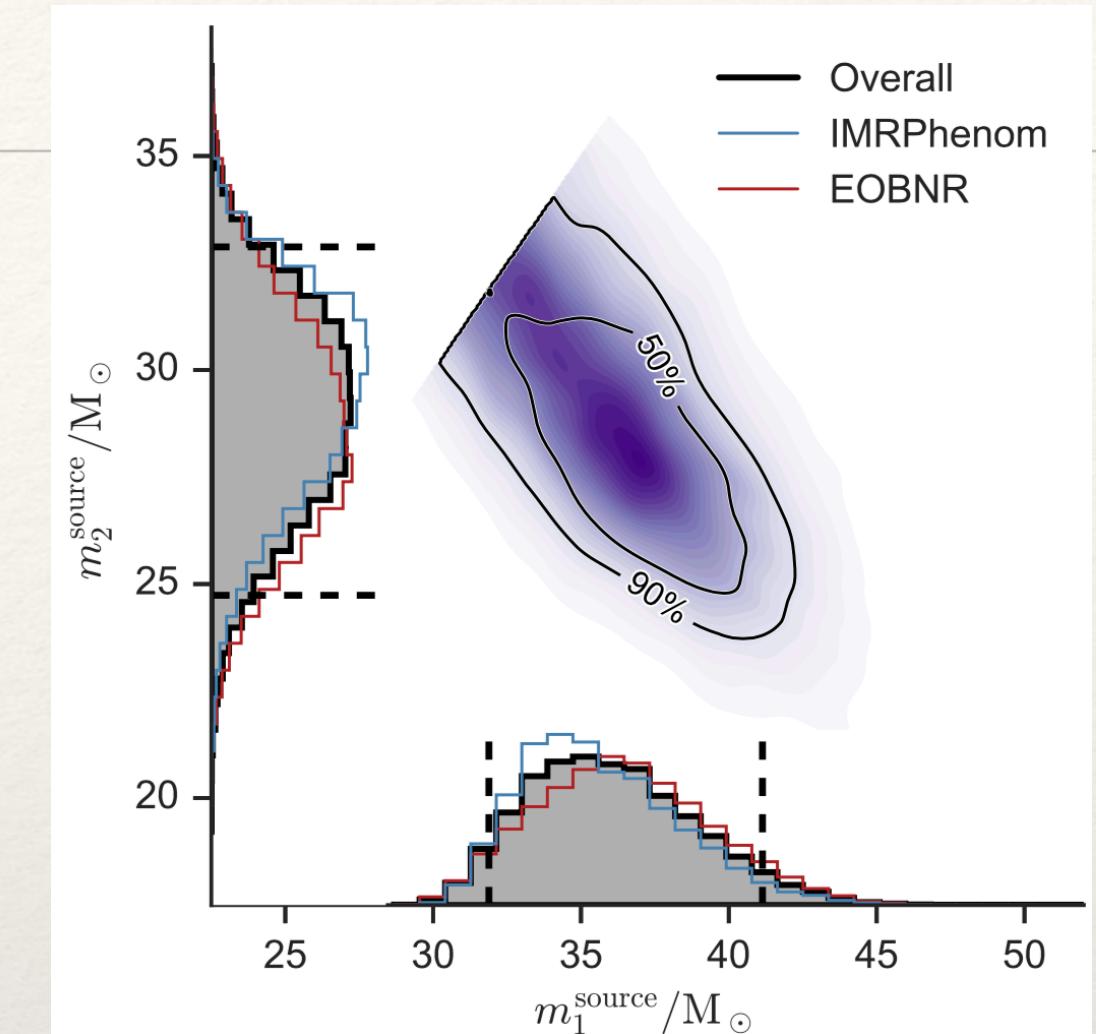
Bayesian inference

- ❖ Goal is to draw samples from the **posterior distribution**

$$\theta \sim p(\theta | d) = \frac{p(d | \theta)p(\theta)}{p(d)}$$

- ❖ **Gravitational waves**

- ❖ Likelihood assumes **stationary Gaussian** detector noise $p(d | \theta) = \mathcal{N}(h_I(\theta), S_{n,I})$
- ❖ Uninformative prior
- ❖ Typically use **stochastic sampler**, repeatedly evaluating right hand side
- ❖ **Can be expensive**, depending on cost of $h(\theta)$, and must be repeated for each event



Motivation

Speed

Handle the large number of events expected in the future

Enable fast alerts for electromagnetic observers

Accuracy

Move beyond approximations such as stationary Gaussian noise

Flexibility

Analyze data in most natural representation, e.g., time domain

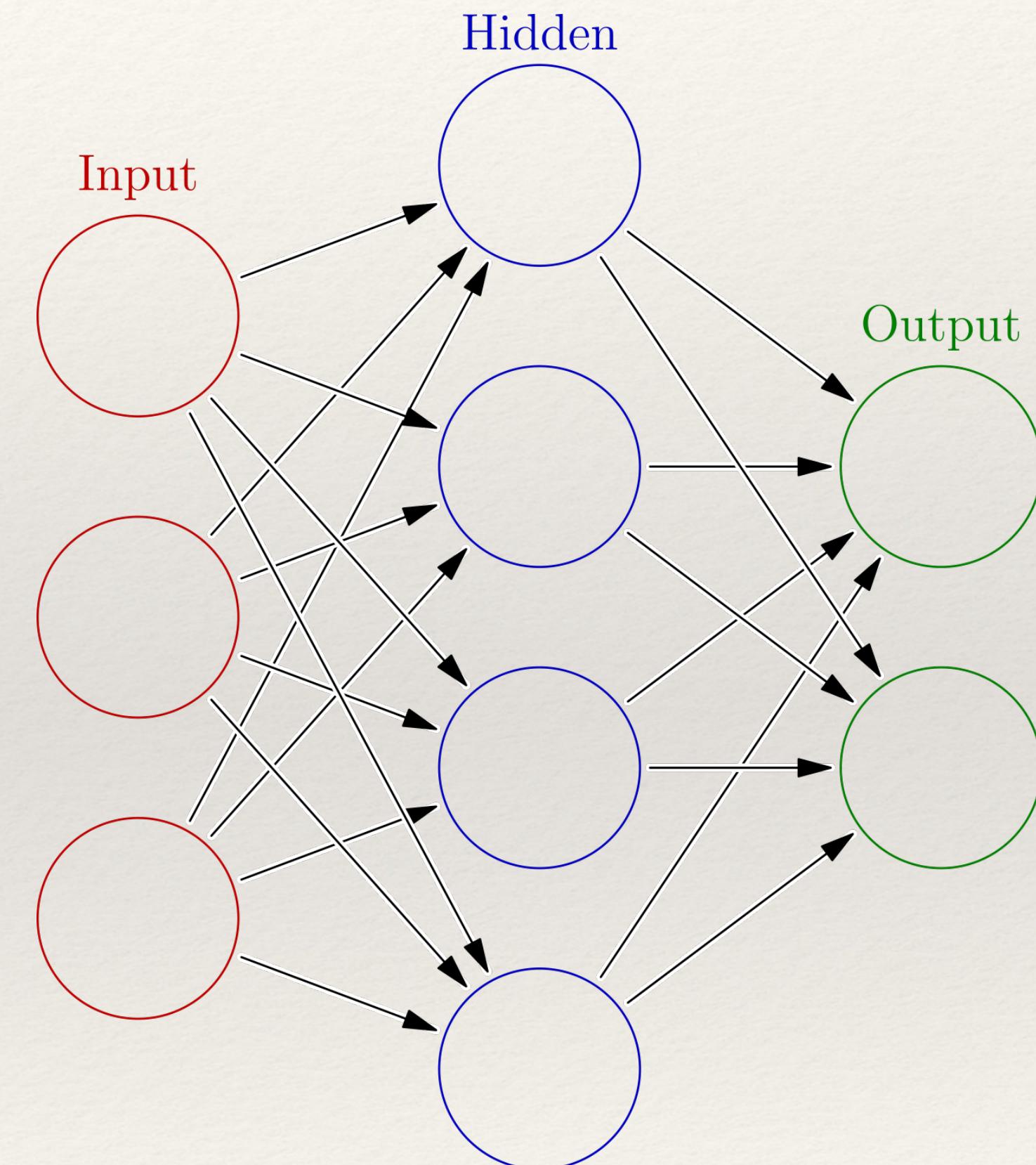
Incorporate other types of data, e.g., galaxy catalogs for populations

Marginalize unwanted parameters

Simulation-based inference

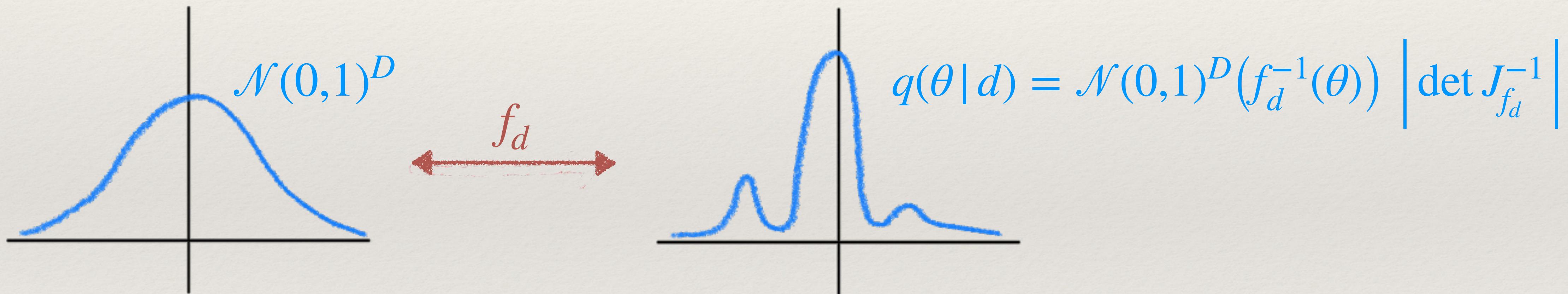
Two key facts

- ❖ Deep neural networks have tremendous capacity to **model complicated probability distributions**.
- ❖ Using **simulated data alone**, can train networks to learn Bayesian inference distributions (e.g., the posterior). No likelihood evaluations are required.



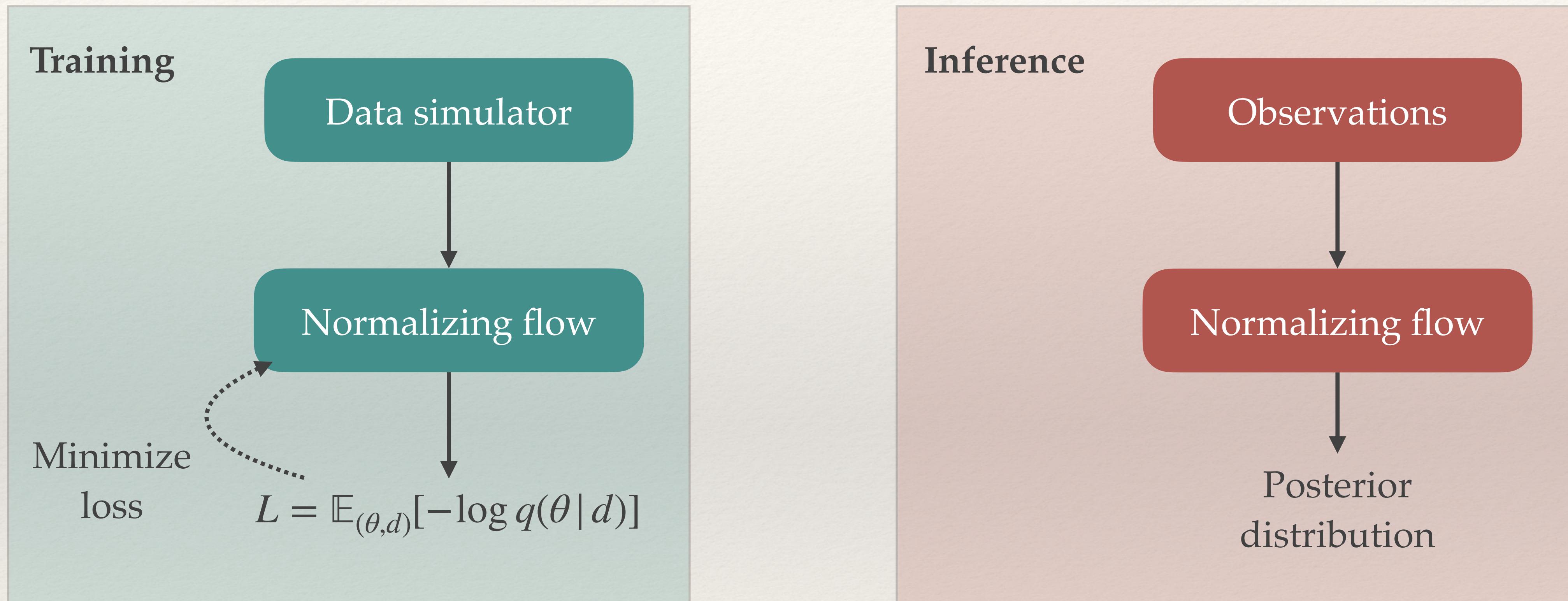
Distributions with neural networks

- For SBI, it is convenient to use a **normalizing flow**. Represents a complex distribution q using a **mapping** $f_d : u \mapsto \theta$ from simpler distribution:



- f_d is defined using neural networks. It must be invertible with simple Jacobian determinant.

Neural posterior estimation



More general than likelihood-based inference methods.

Inference is fast since it uses only forward neural network passes.

Loss function

- ❖ To train the network, specify a target loss function.
- ❖ Want $q_\phi(\theta | d) \rightarrow p(\theta | d)$
 - ❖ Take **Kullbeck-Liebler (KL) divergence** D_{KL} between these distributions

$$D_{\text{KL}}(p | q) = \int d\theta p(\theta | d) \log \frac{p(\theta | d)}{q_\phi(\theta | d)} \geq 0 \quad = 0 \text{ for identical distributions}$$

- ❖ This still depends on d so **marginalize** over it

$$\mathbb{E}_{p(d)} D_{\text{KL}}(p | q) = \int dd p(d) \int d\theta p(\theta | d) \log \frac{p(\theta | d)}{q_\phi(\theta | d)}$$

Loss function

$$\mathbb{E}_{p(d)} D_{\text{KL}}(p \mid q) = \int dd p(d) \int d\theta p(\theta \mid d) \log \frac{p(\theta \mid d)}{q_\phi(\theta \mid d)}$$

- ❖ To evaluate, re-order the integrals using **Bayes' theorem**

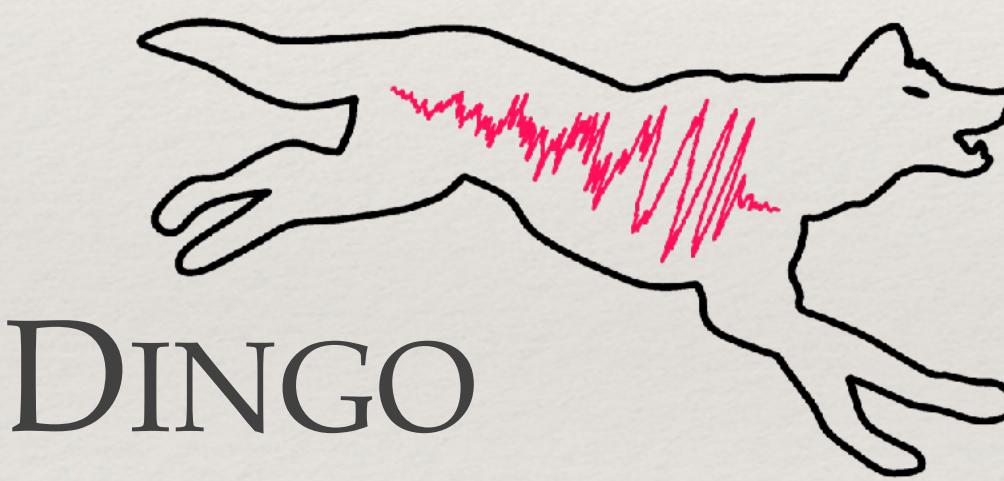
$$\mathbb{E}_{p(d)} D_{\text{KL}}(p \mid q) = \int d\theta p(\theta) \int dd p(d \mid \theta) \log \frac{p(\theta \mid d)}{q_\phi(\theta \mid d)}$$

- ❖ Finally

$$L[\phi] = \frac{1}{N} \sum_{\substack{\theta^{(i)} \sim p(\theta) \\ d^{(i)} \sim p(d \mid \theta^{(i)})}} -\log q_\phi(\theta^{(i)} \mid d^{(i)})$$

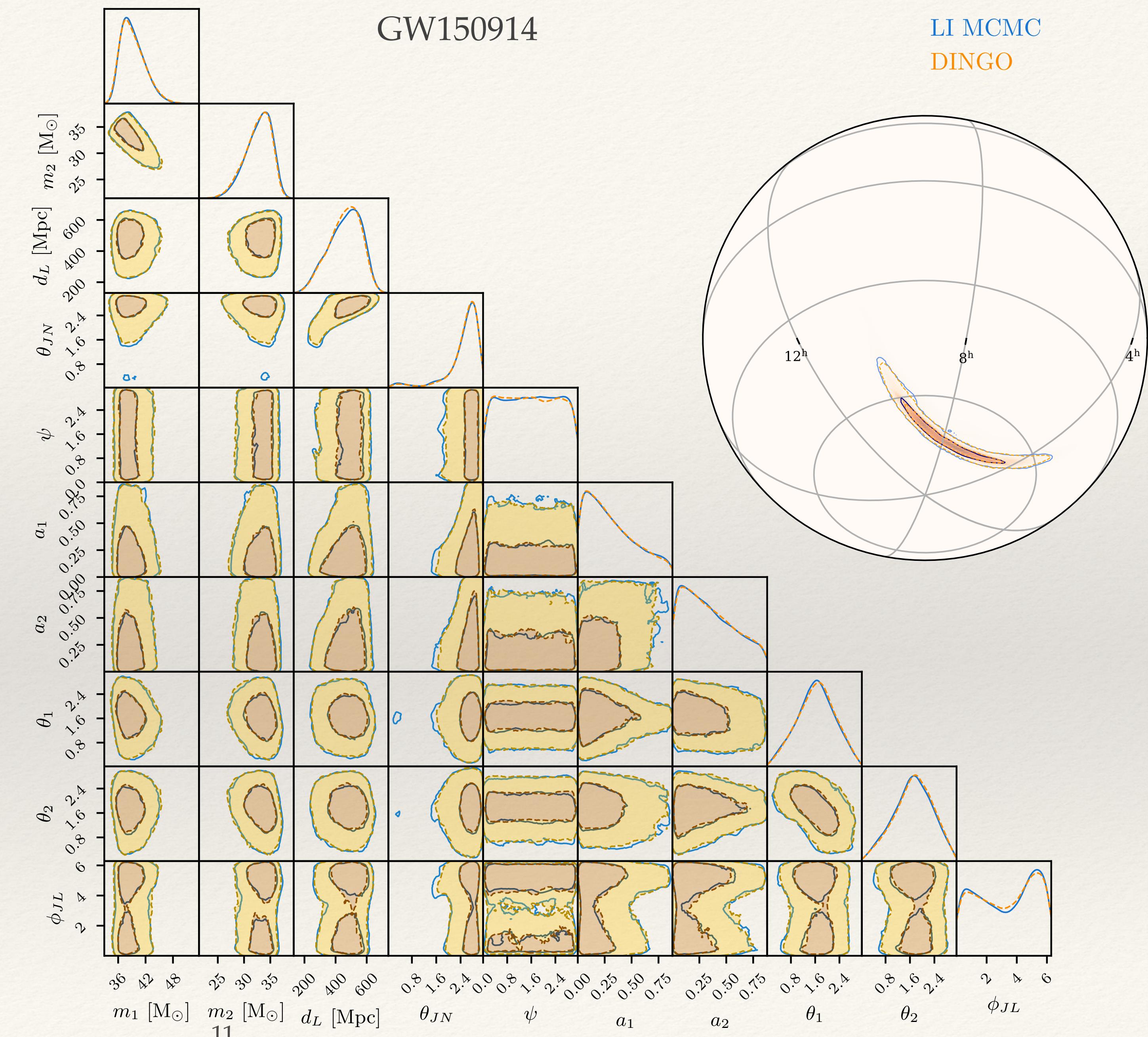
1. Sample $\theta^{(i)}$ from the prior
2. Simulate $d^{(i)} = \text{signal} + \text{noise}$

- ❖ Well-trained networks give extremely good agreement with standard samplers.
- ❖ Inference in seconds to minutes.

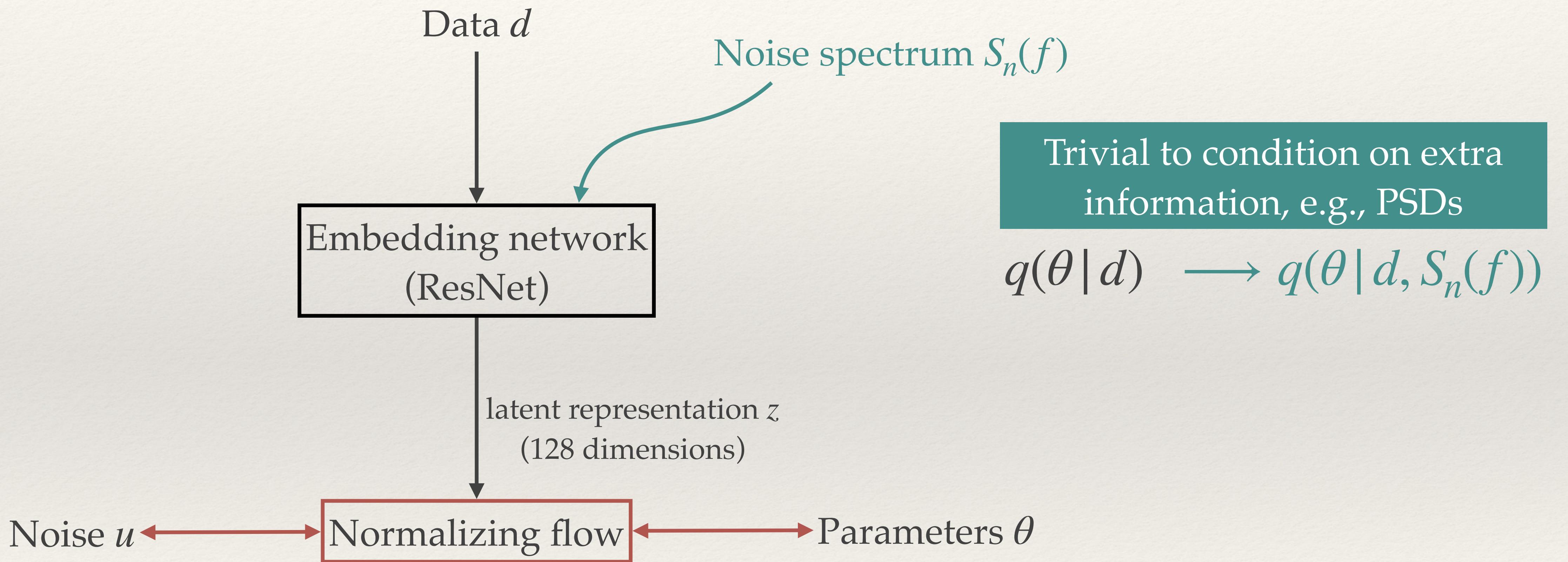


DINGO

<https://github.com/dingo-gw/dingo>



Architecture



Other SBI approaches

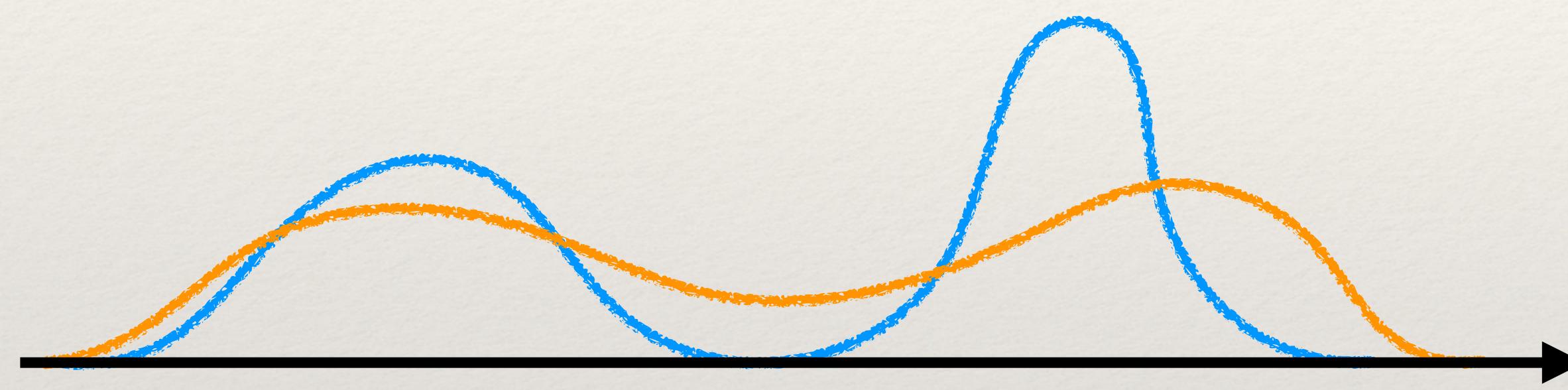
We can use neural density estimators in many ways: highly flexible!

- ❖ Neural likelihood estimation (NLE)
 - ❖ Learn likelihood $p(d | \theta)$
 - ❖ Requires MCMC to obtain posterior in the end
- ❖ Neural ratio estimation (NRE)
 - ❖ Train a classifier to distinguish samples from joint $p(\theta, d)$ and $p(\theta)p(d)$.
 - ❖ This gives ratio $p(d | \theta)/p(d)$. Use MCMC to obtain posterior samples

Trivial to marginalize parameters, add additional context, different data representations, etc.

Neural importance sampling

- By combining with classical likelihood-based techniques, we can correct SBI inaccuracies using **importance sampling**:



$$w_i \propto \frac{p(\theta_i)p(d | \theta_i)}{q(\theta_i | d)}$$

target (prior x likelihood)

proposal (NPE)

- The mean of the weights gives the **Bayesian evidence**:

$$p(d) = \frac{1}{n} \sum_{i=1}^n w_i$$

- Variance gives the **sampling efficiency**

$$\epsilon = \frac{\left(\sum_i w_i \right)^2}{n \sum_i w_i^2}$$

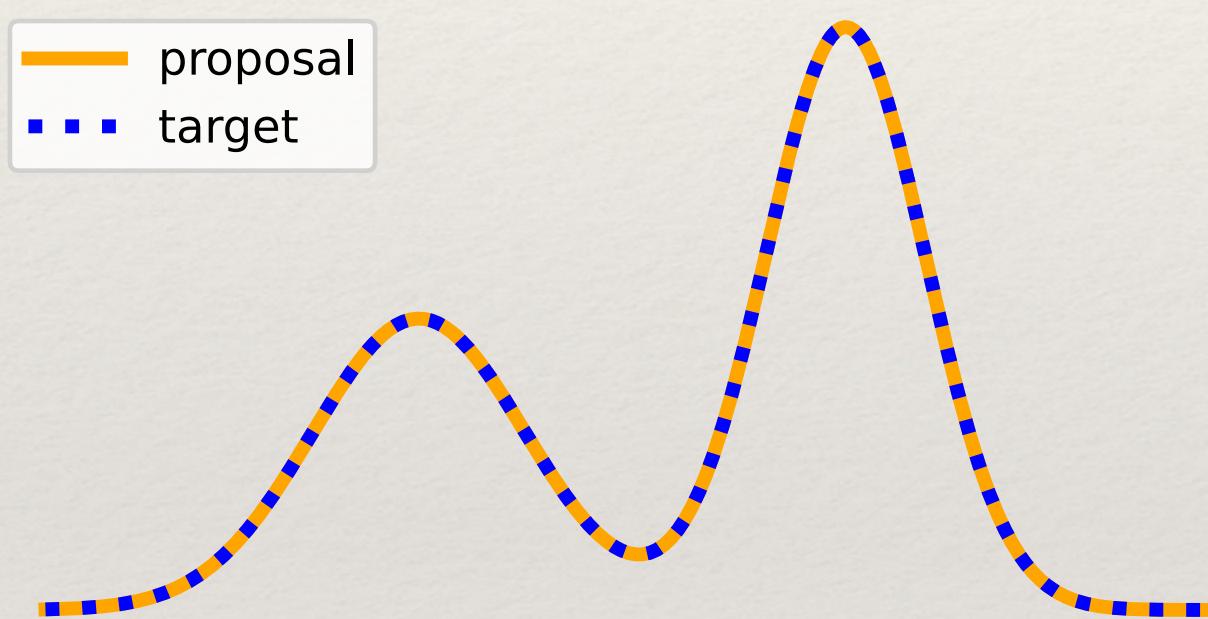
Neural importance sampling

- ❖ Performance depends on quality of the proposal:

- ❖ Perfect proposal

$$w_i = \text{constant} \implies \epsilon = 1$$

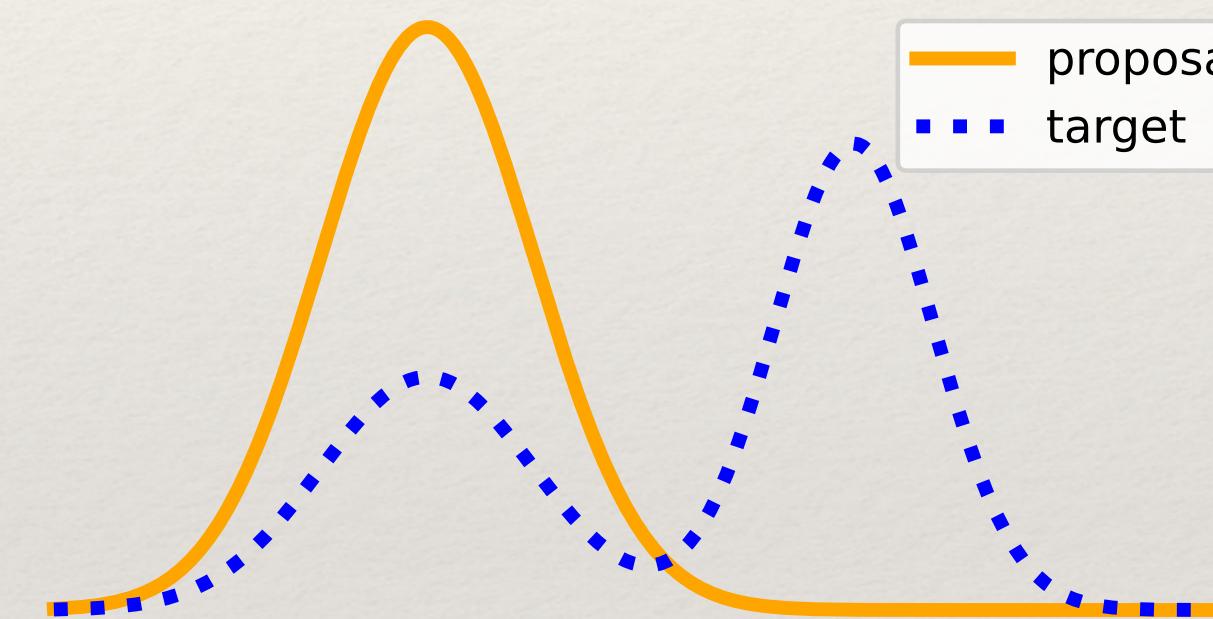
proposal
target



- ❖ Poor proposal

$$\text{large variance} \implies \epsilon \approx 0$$

proposal
target



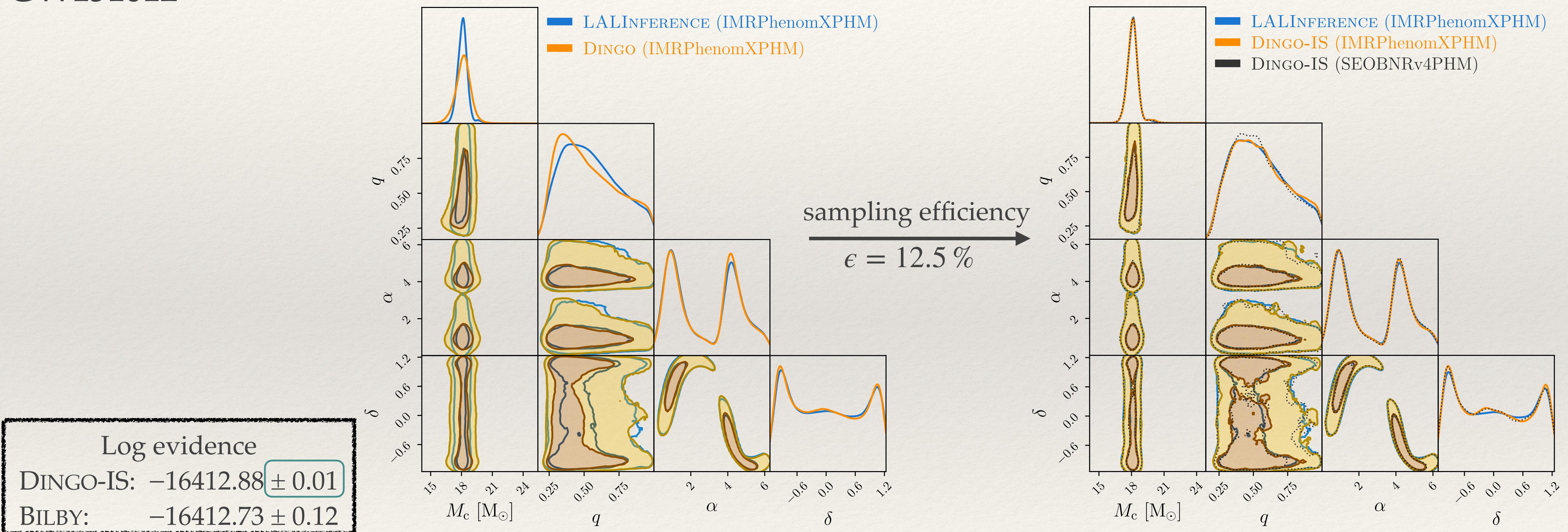
- ❖ NPE proposals have the **mass-covering property**, making them well-suited to IS

$$L[\phi] \sim D_{\text{KL}}(p \mid q) = \int d\theta p(\theta \mid d) \log \frac{p(\theta \mid d)}{q_\phi(\theta \mid d)}$$

Penalty in loss if q_ϕ misses a mode!

Neural importance sampling

GW151012



Validating results

- ❖ Good sampling efficiency for majority of events

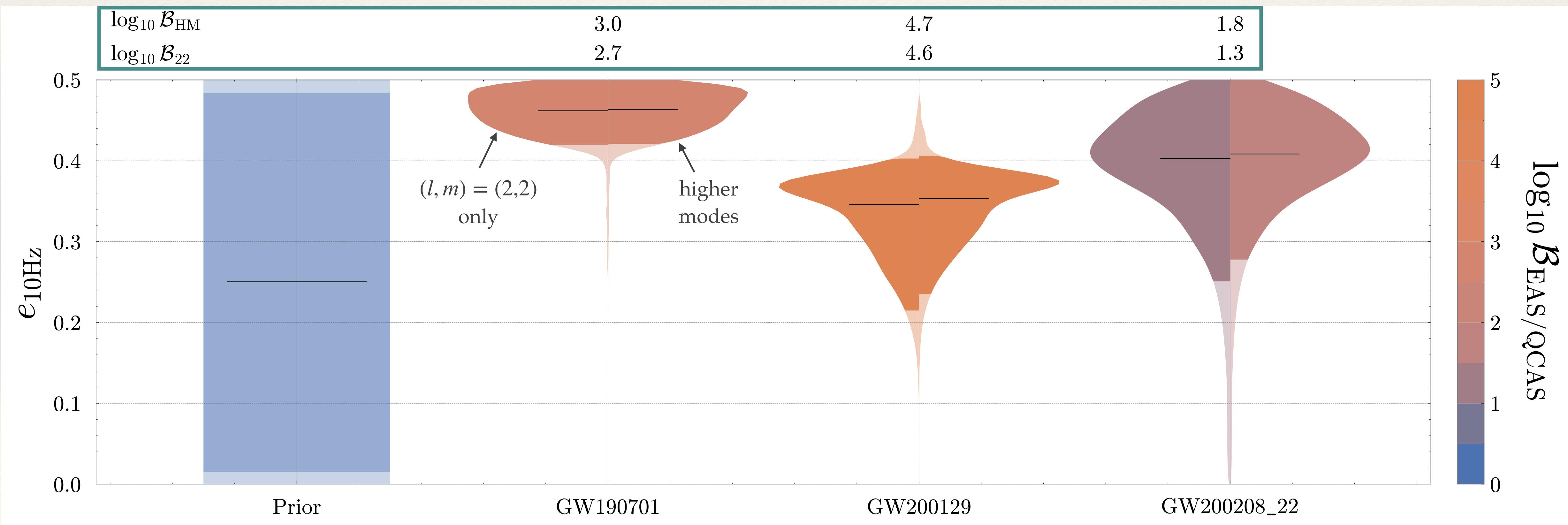
Event	$\log p(d)$	ϵ	Event	$\log p(d)$	ϵ	Event	$\log p(d)$	ϵ
GW190408_181802	-16178.332 ± 0.012	6.9%	GW190727_060333	-15992.017 ± 0.009	10.3%	GW191230_180458	-15913.798 ± 0.009	12.2%
GW190413_052954	-15571.413 ± 0.006	22.5%	GW190731_140936	-16376.777 ± 0.005	32.6%	GW200128_022011	-16305.128 ± 0.013	6.1%
GW190413_134308	-16399.331 ± 0.009	12.4%	GW190803_022701	-16132.409 ± 0.006	21.4%	‡GW200129_065458	-16226.851 ± 0.109	0.1%
GW190421_213856	-15983.248 ± 0.008	15.3%	GW190805_211137	-16073.261 ± 0.006	20.0%	GW200208_130117	-16136.381 ± 0.007	16.6%
GW190503_185404	-16582.865 ± 0.022	2.0%	GW190828_063405	-16137.220 ± 0.009	12.2%	GW200208_222617	-16775.200 ± 0.011	7.4%
GW190513_205428	-15946.462 ± 0.043	0.6%	GW190909_114149	-16061.634 ± 0.011	7.4%	GW200209_085452	-16383.847 ± 0.009	12.5%
GW190514_065416	-16556.466 ± 0.009	11.6%	GW190915_235702	-16083.960 ± 0.015	20.8%	GW200216_220804	-16215.703 ± 0.017	3.4%
GW190517_055101	-16271.048 ± 0.027	1.3%	GW190926_050336	-16015.813 ± 0.019	2.8%	GW200219_094415	-16133.457 ± 0.011	9.6%
GW190519_153544	-15991.171 ± 0.008	15.2%	GW190929_012149	-16146.666 ± 0.018	3.2%	GW200220_061928	-16303.782 ± 0.007	17.3%
GW190521_074359	-16008.876 ± 0.008	13.4%	GW191109_010717	-17925.064 ± 0.025	1.7%	GW200220_124850	-16136.000 ± 0.008	13.2%
GW190527_092055	-16119.012 ± 0.008	13.8%	GW191127_050227	-16759.328 ± 0.019	2.7%	GW200224_222234	-16138.613 ± 0.006	22.5%
GW190602_175927	-16036.993 ± 0.006	25.0%	‡GW191204_110529	-15984.455 ± 0.015	4.2%	‡GW200308_173609	-16173.938 ± 0.013	6.0%
GW190701_203306	-16521.381 ± 0.040	0.6%	GW191215_223052	-16001.286 ± 0.013	5.8%	GW200311_115853	-16117.505 ± 0.011	7.4%
GW190719_215514	-15850.492 ± 0.008	13.4%	GW191222_033537	-15871.521 ± 0.007	16.5%	‡GW200322_091133	-16313.568 ± 0.307	0.0%
		8.0%			25.8%			0.1%

Application: Eccentric binaries

- ❖ Use DINGO for a large study of eccentricity binaries using an expensive waveform model.
- ❖ “Effective one body” SEOBNRv4EHM model (Ramos-Buades+, 2022).
 - ❖ Two new parameters: **eccentricity** and **relativistic anomaly**
 - ❖ Aligned spin, but includes effect of higher-order multipoles
- ❖ Costs:
 - ❖ Nested sampling: $O(\text{week})$ per event w/ 320 cores.
 - ❖ DINGO: $O(\text{minute})$ for initial samples, and $O(\text{hour})$ for importance sampling.

Application: Eccentric binaries

- ❖ Three events with evidence for eccentricity ($\log_{10} \text{Bayes} \geq 1$ compared to aligned spin)

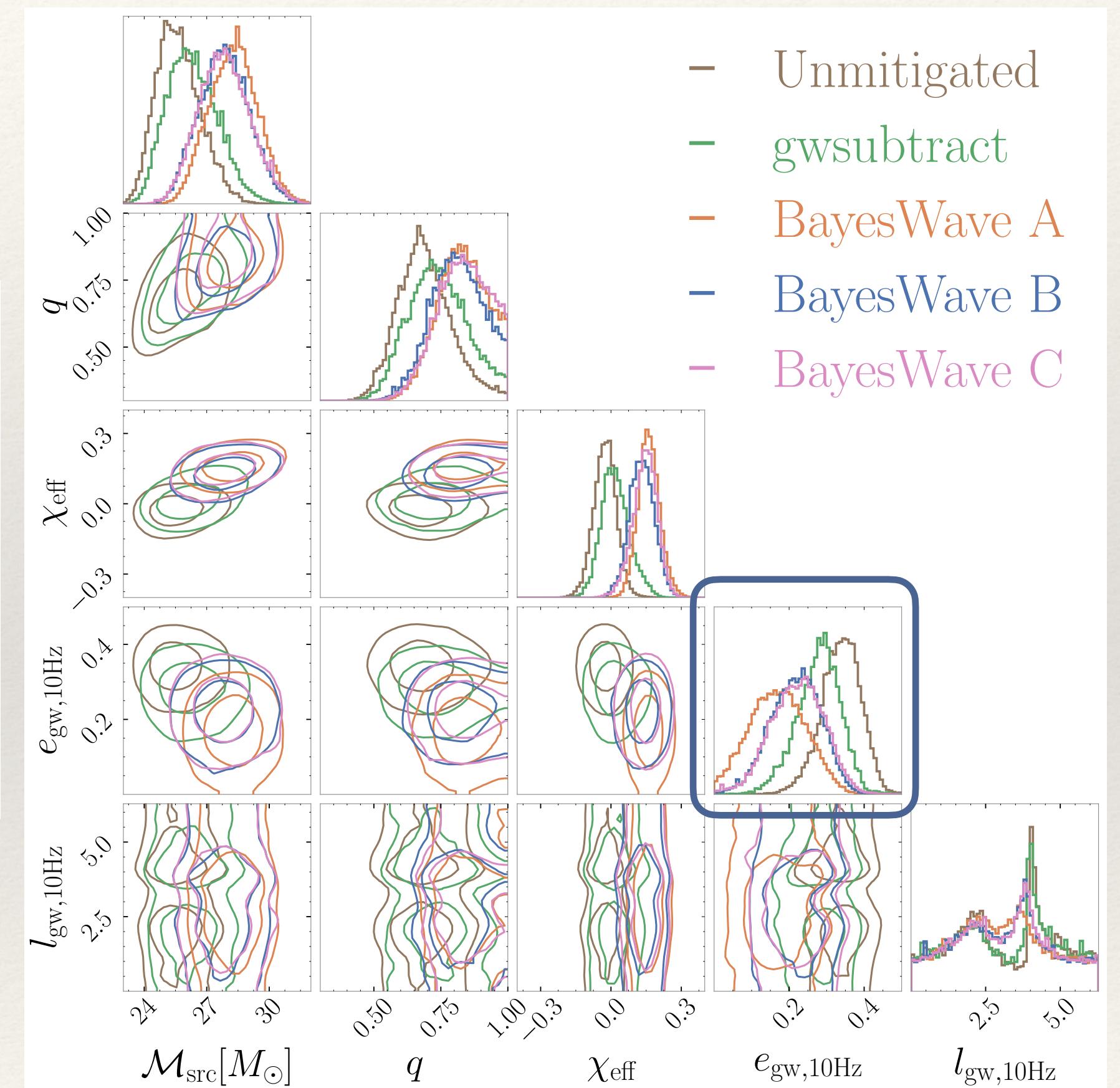


Application: Eccentric binaries

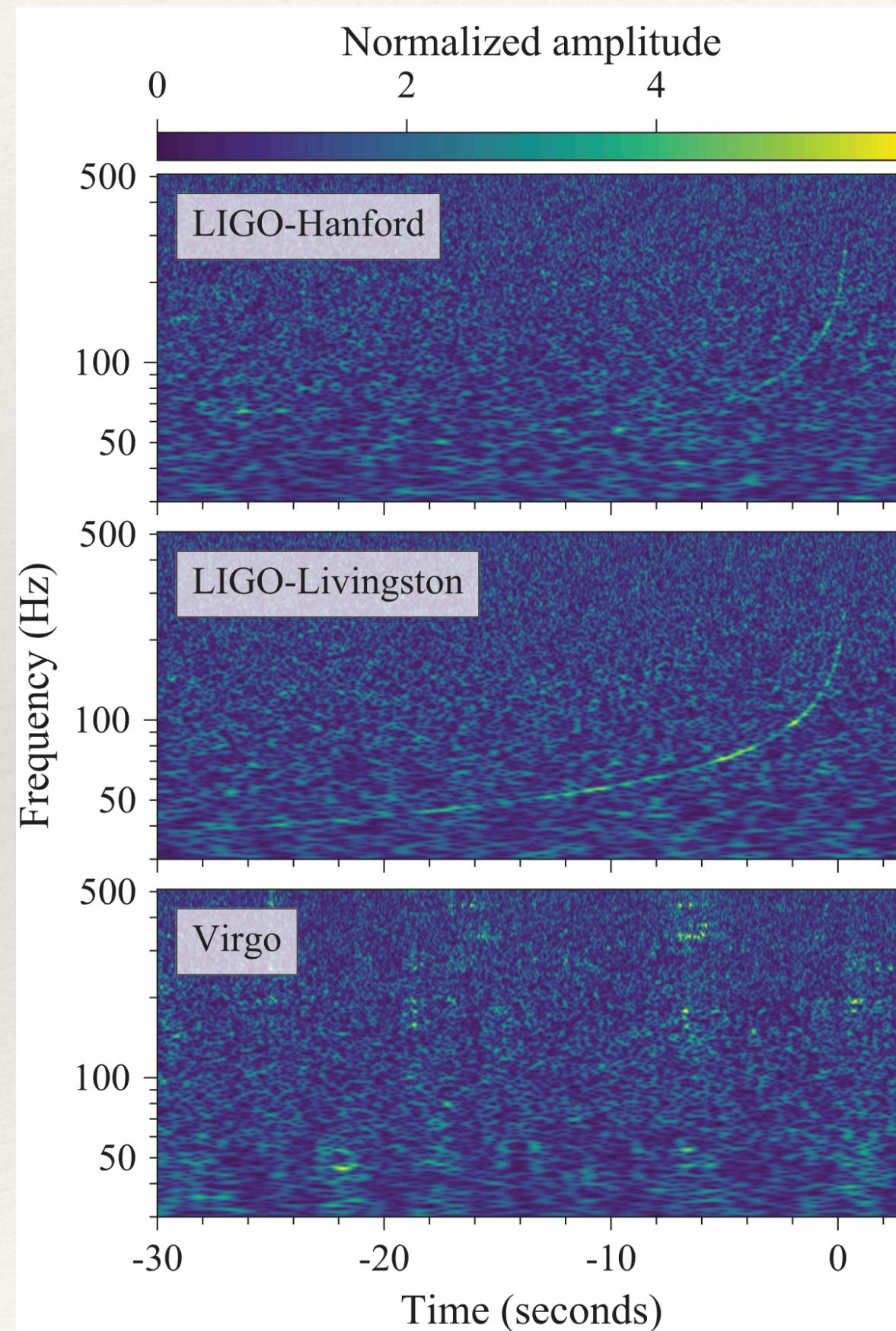
- ❖ Caveat for GW200129: Result highly dependent on glitch subtraction algorithm

Glitch Subtraction	e Prior	SEOBNRv4 $\log_{10} \mathcal{B}$	SEOBNRv4HM $\log_{10} \mathcal{B}$	SEOBNRv4PHM $\log_{10} \mathcal{B}$	NRSur7dq4 $\log_{10} \mathcal{B}$	$e_{10\text{Hz}}$	$e_{\text{gw}, 10\text{Hz}}$
GW200129							
gwsubtract	Uniform	4.57	4.75	4.92	4.0	$0.34^{+0.11}_{-0.06}$	$0.27^{+0.10}_{-0.12}$
BayesWave A	Uniform	1.7	1.84	2.20	1.53	$0.24^{+0.10}_{-0.10}$	$0.17^{+0.14}_{-0.13}$
BayesWave B	Uniform	2.92	3.08	3.43	2.35	$0.28^{+0.09}_{-0.11}$	$0.22^{+0.12}_{-0.13}$
BayesWave C	Uniform	2.85	2.93	2.63	1.43	$0.27^{+0.09}_{-0.10}$	$0.22^{+0.13}_{-0.14}$

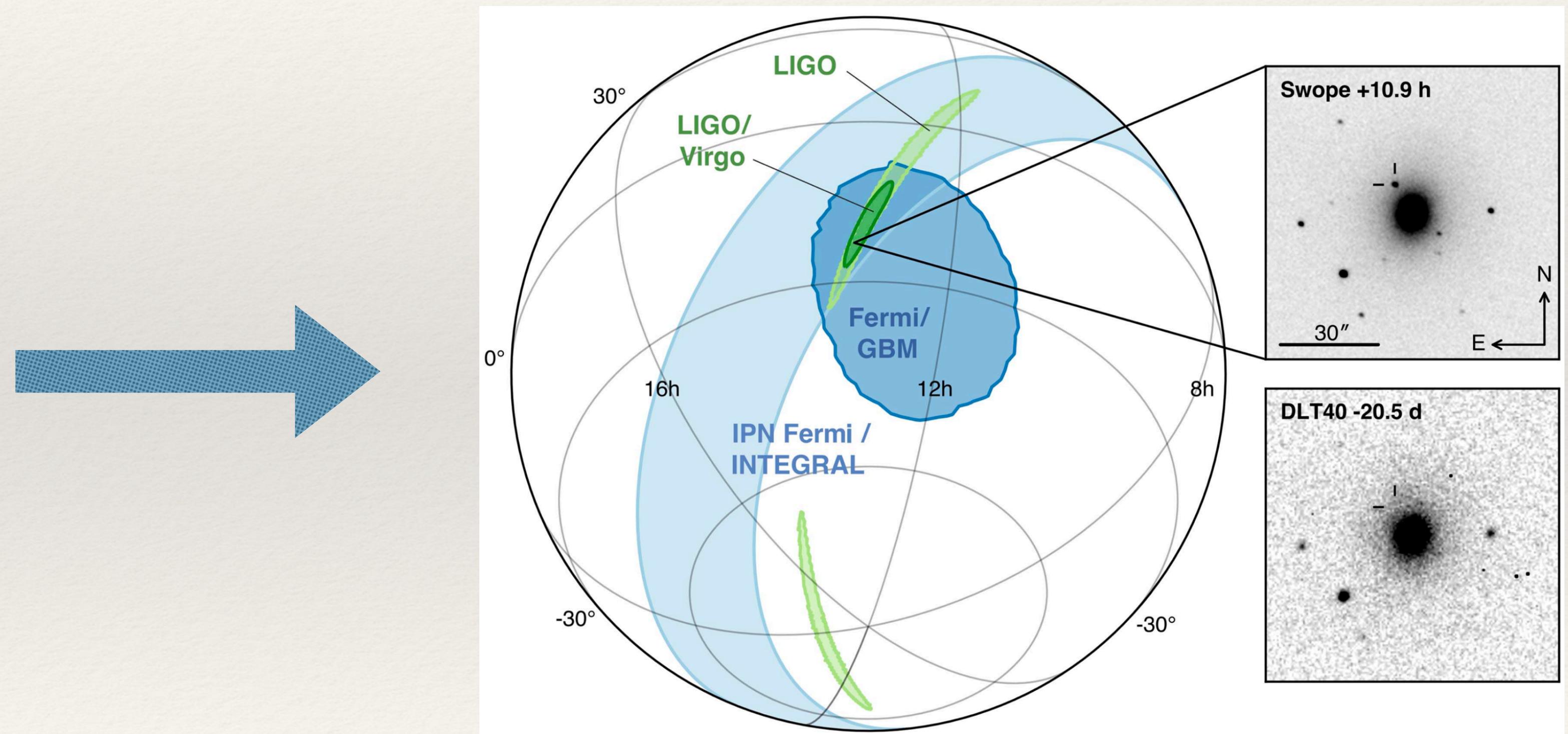
- ❖ However, all cases show eccentricity, even in comparison to precession!



Binary neutron stars



GW170817: Rapid sky localization enables multimessenger astrophysics

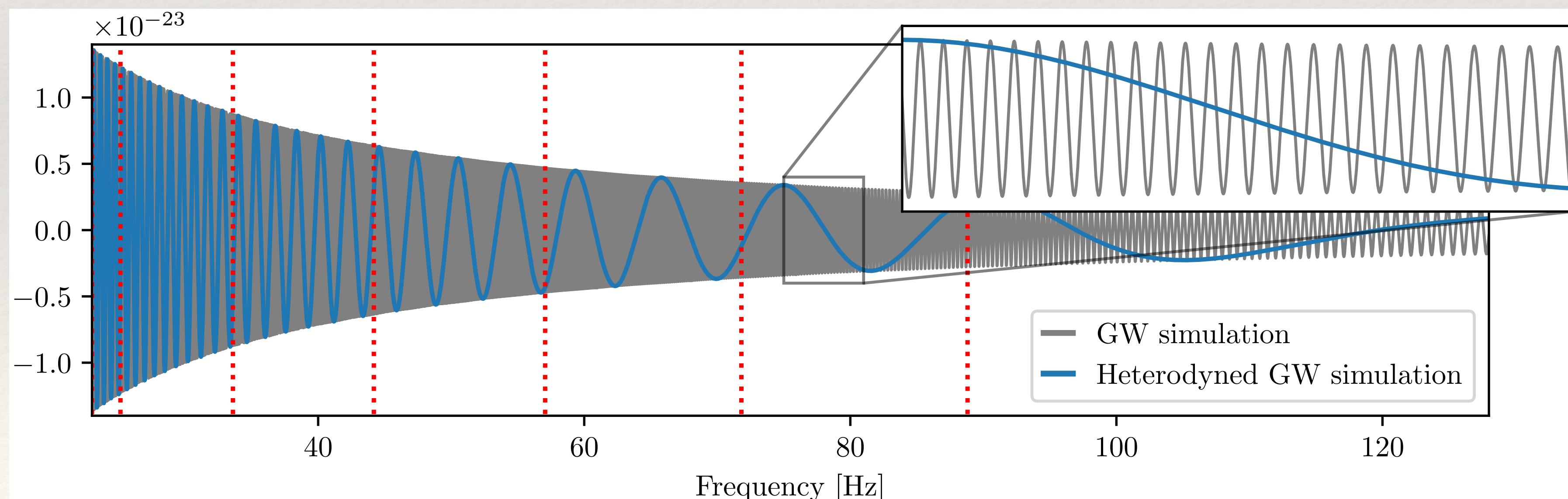


Binary neutron stars

❖ **Challenge:** BNS signals are longer and more complex than BBH

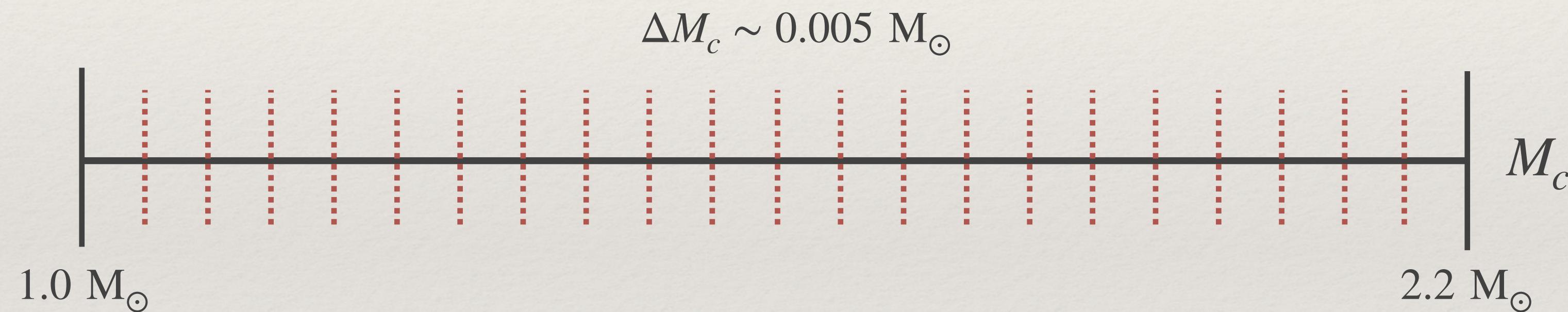
❖ **Solution:**

1. Heterodyning (Cornish 2010) — factor out overall phase $\propto (M_c f)^{-5/3}$
2. Multibanding (Vinciguerra+ 2017) — use reduced resolution at higher f



Binary neutron stars

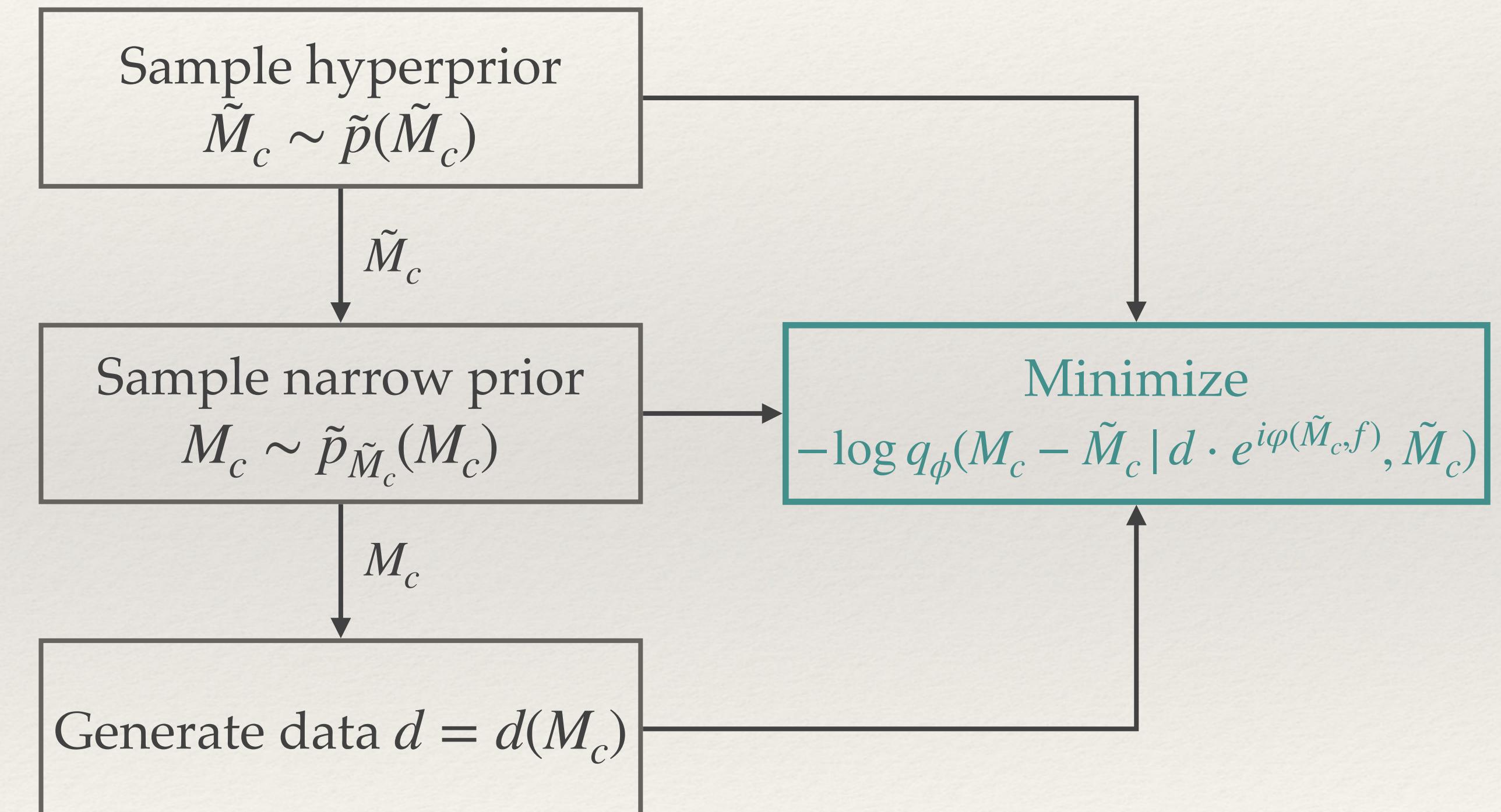
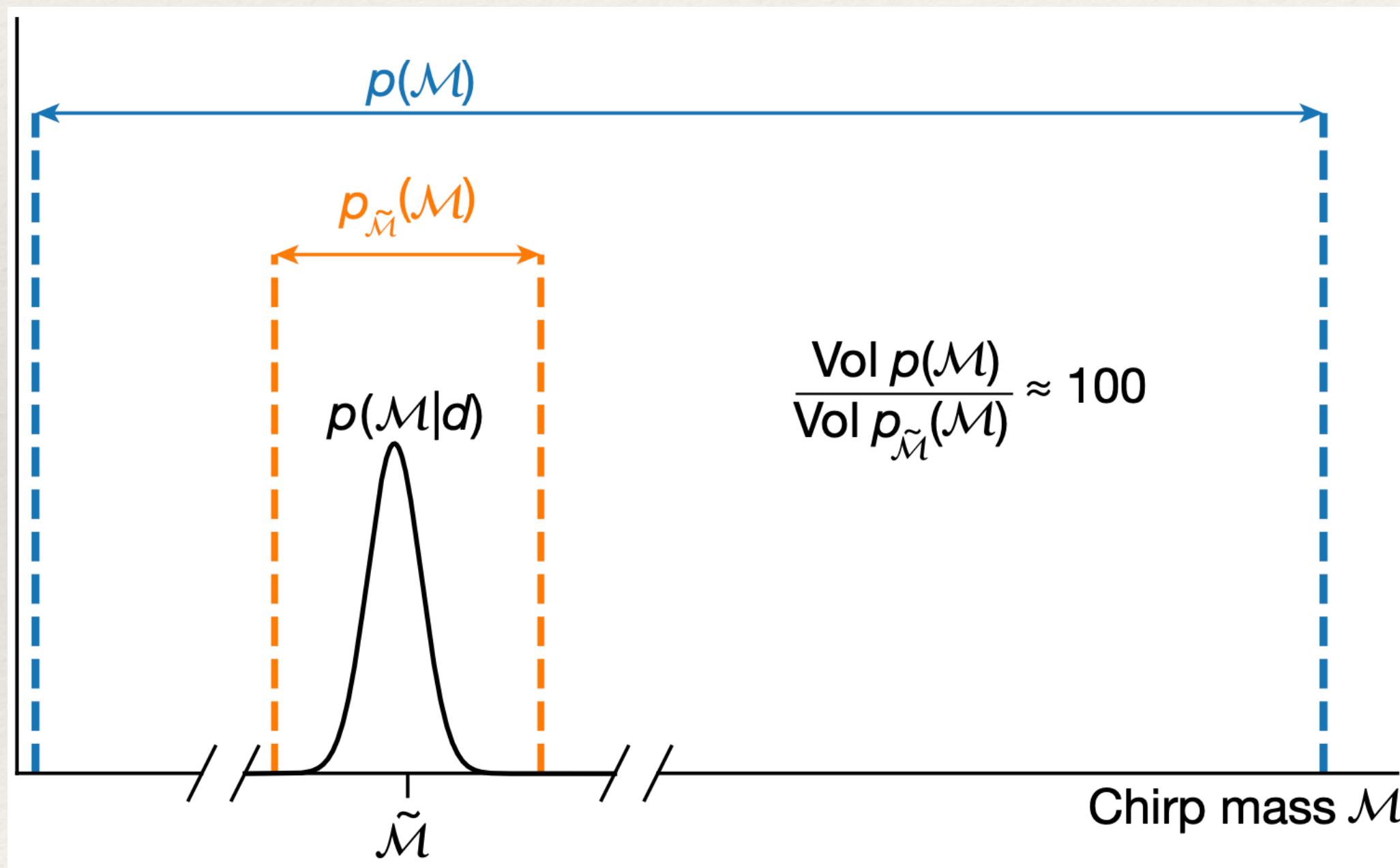
- ❖ But heterodyning depends on a **specific chirp mass** \tilde{M}_c . To achieve significant simplification this must be close to the true chirp mass.
- ❖ **One solution:** Divide up the prior, and train a network only over a narrow chirp mass range,



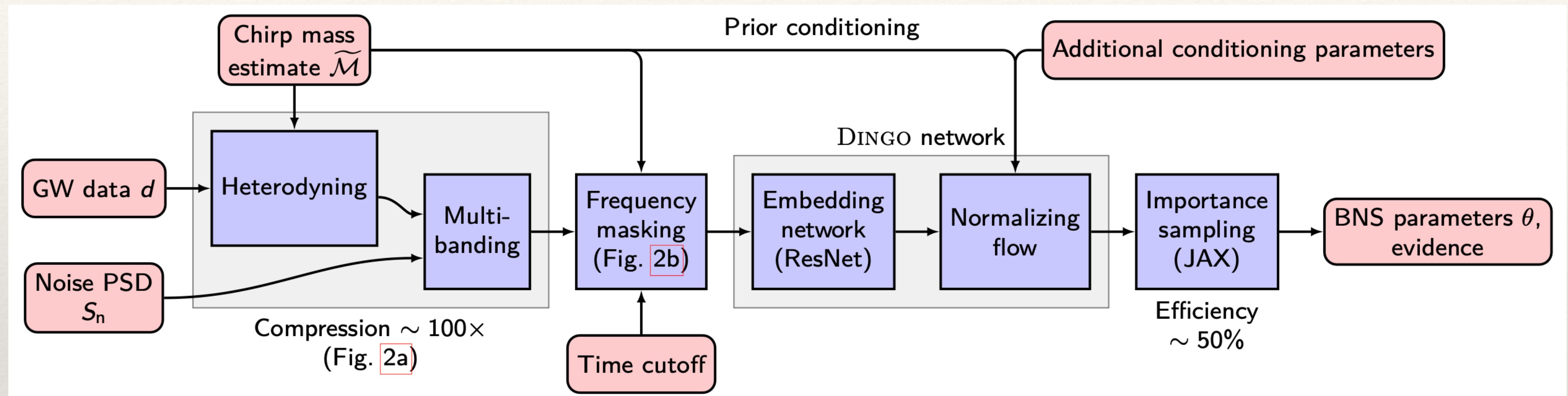
- ❖ **Impractical:** Would require too many networks.
- ❖ Instead use **prior conditioning**.

Prior conditioning

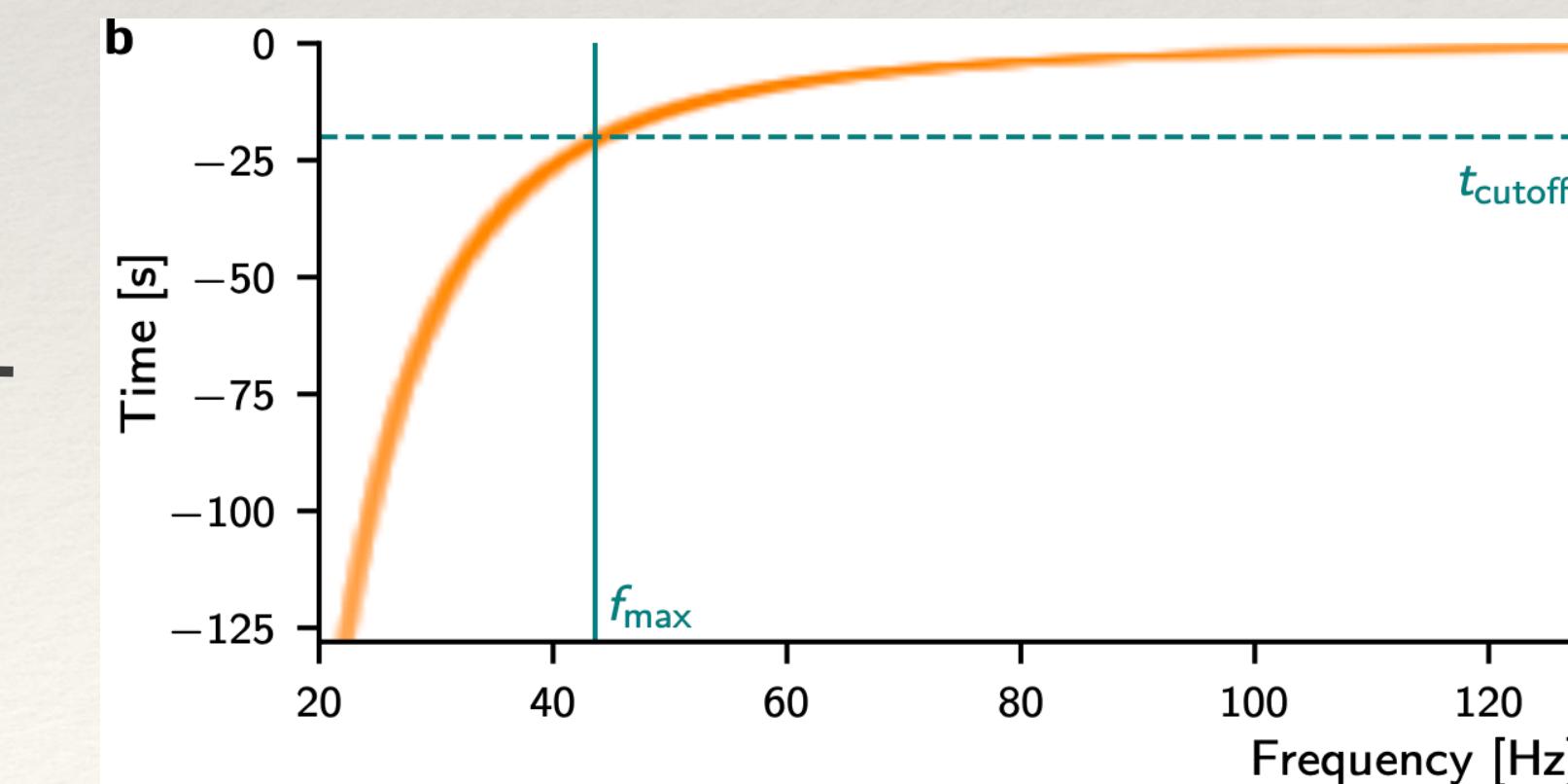
- ❖ Train network **conditioned** on narrow prior $p_{\tilde{M}_c}(M_c)$.



Binary neutron star inference



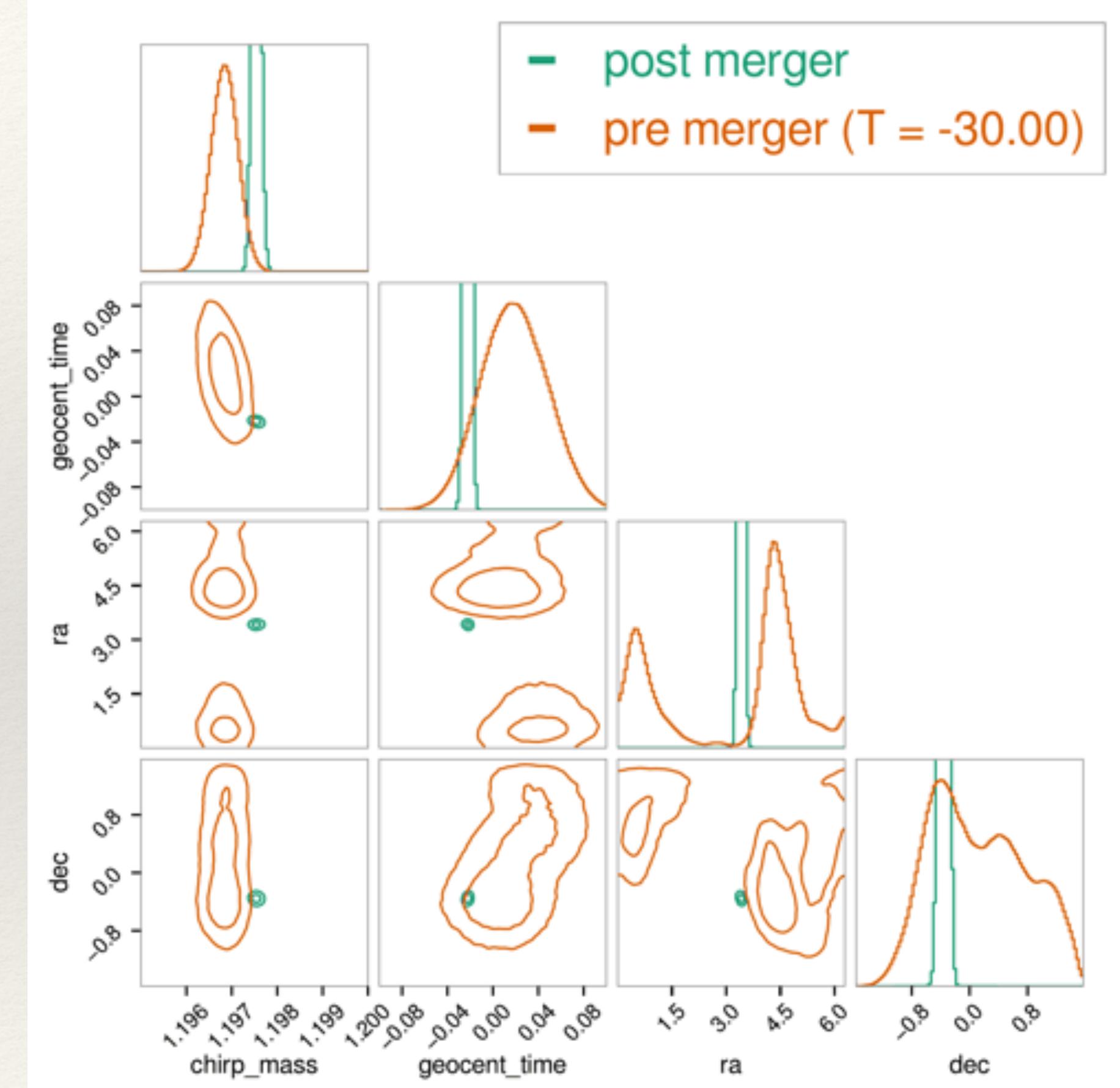
Enable pre-merger inference



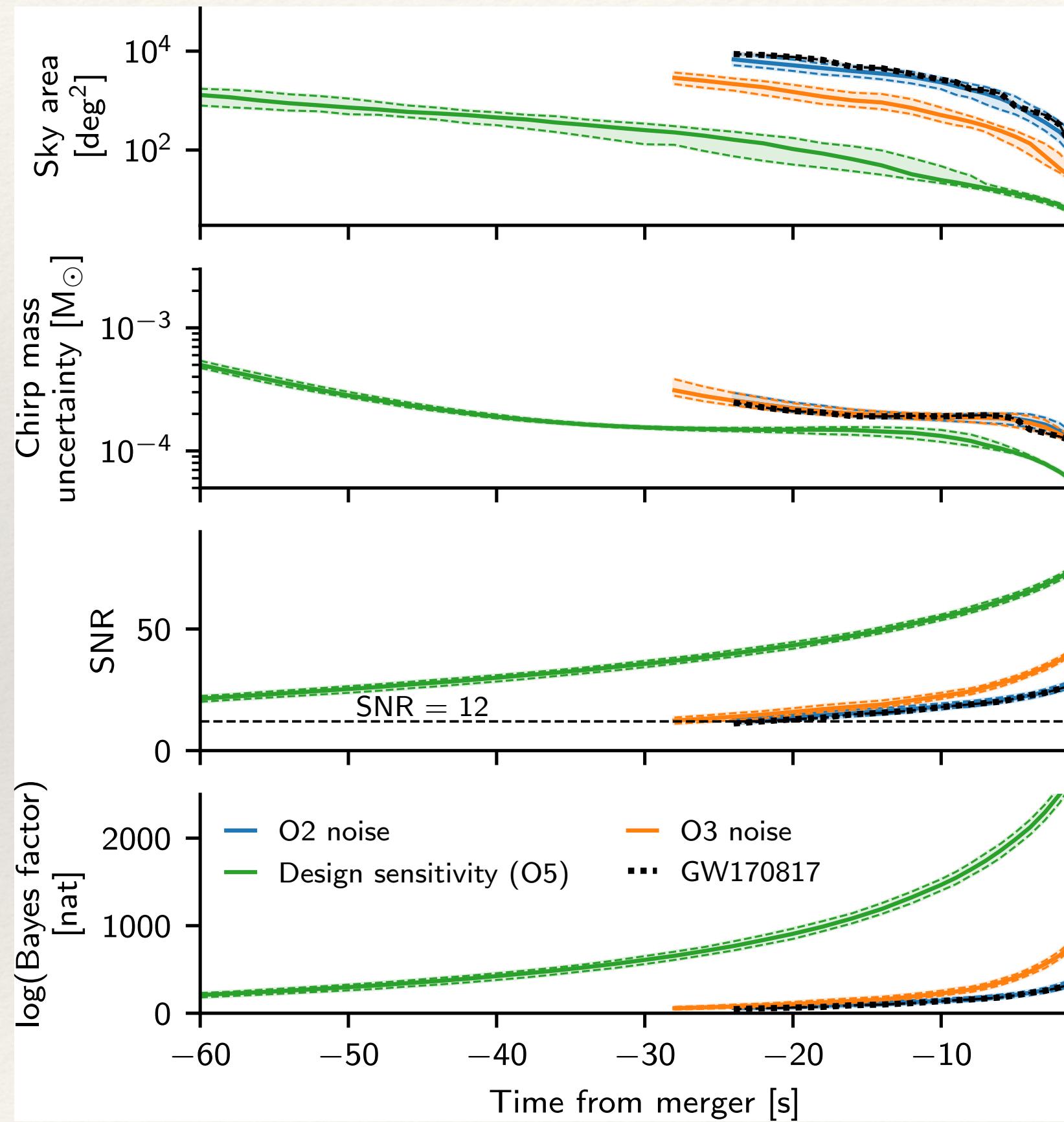
Binary neutron star inference

❖ Pre-merger inference:

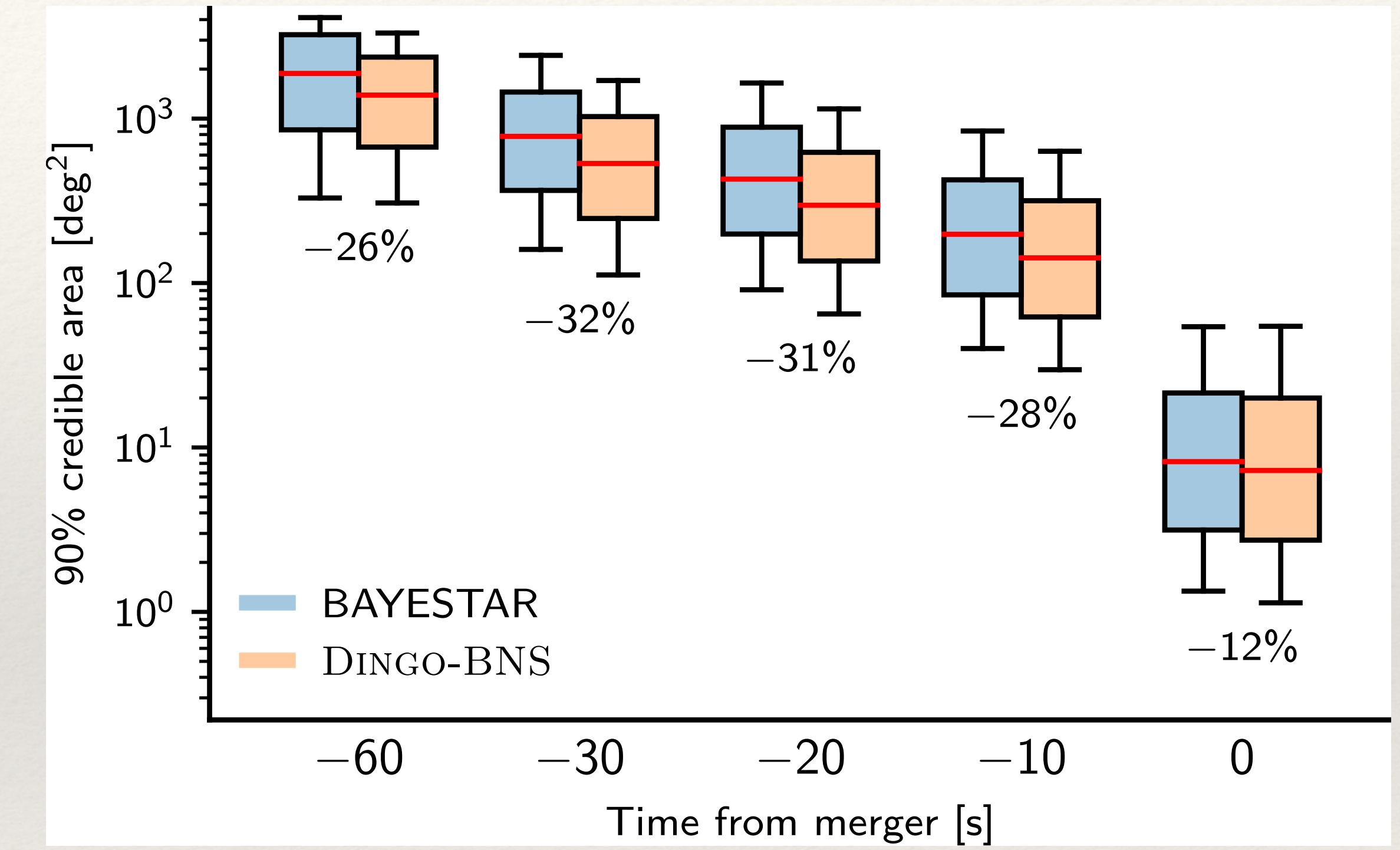
Full Bayesian analysis takes ~ 1 second!



Pre-merger BNS



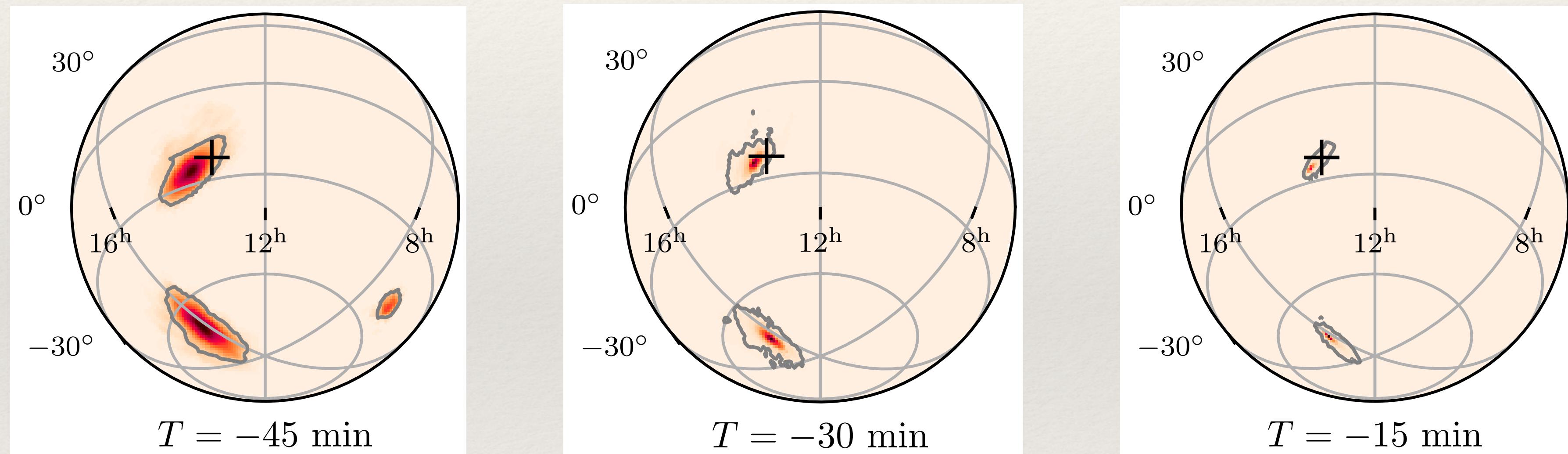
Evolution of pre-merger estimates



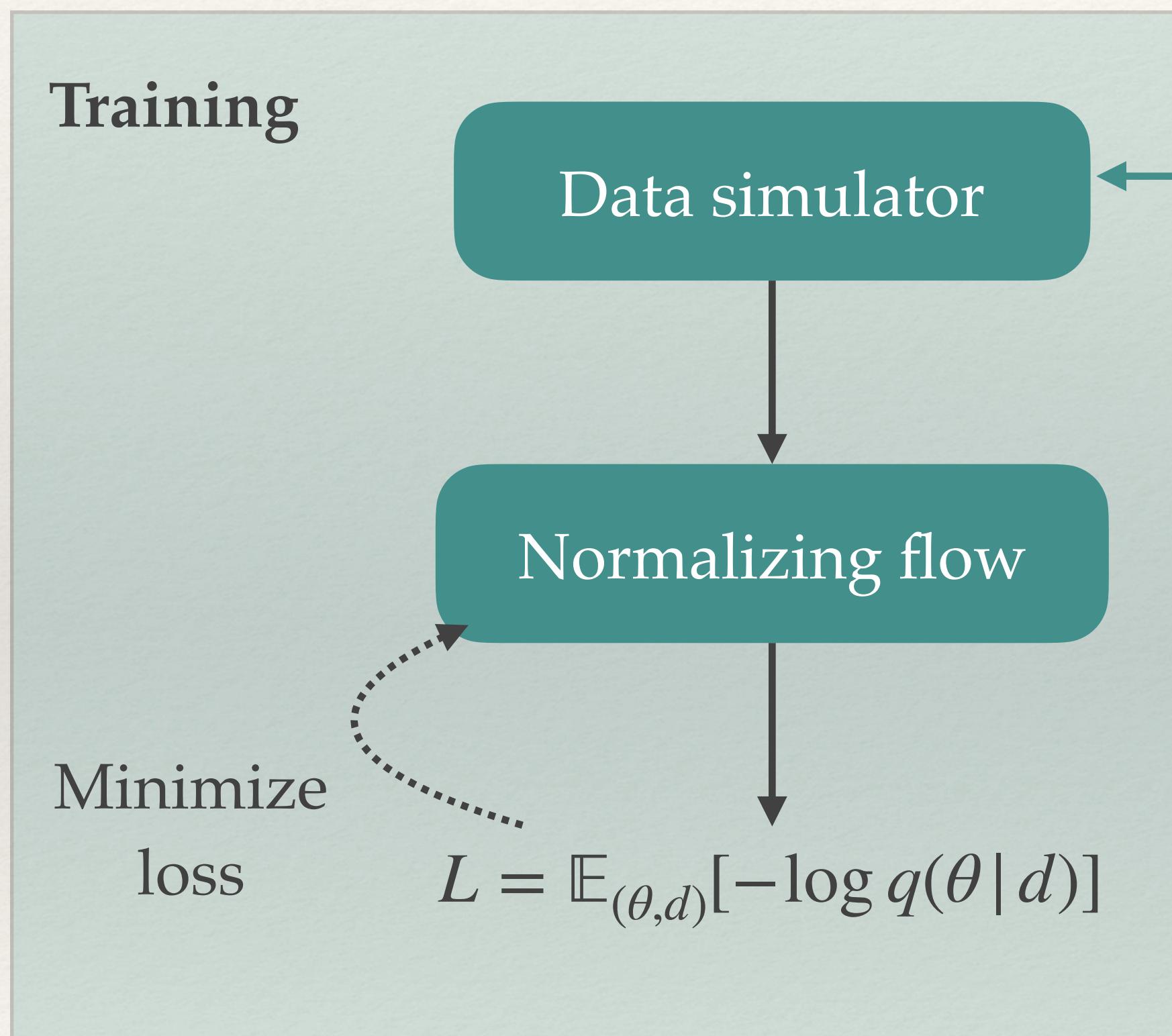
Improvements over BAYESTAR

Binary neutron stars

- ❖ With XG detectors, obtain sky position many minutes before merger



SBI is extremely general



Augment data simulator as desired, e.g.,

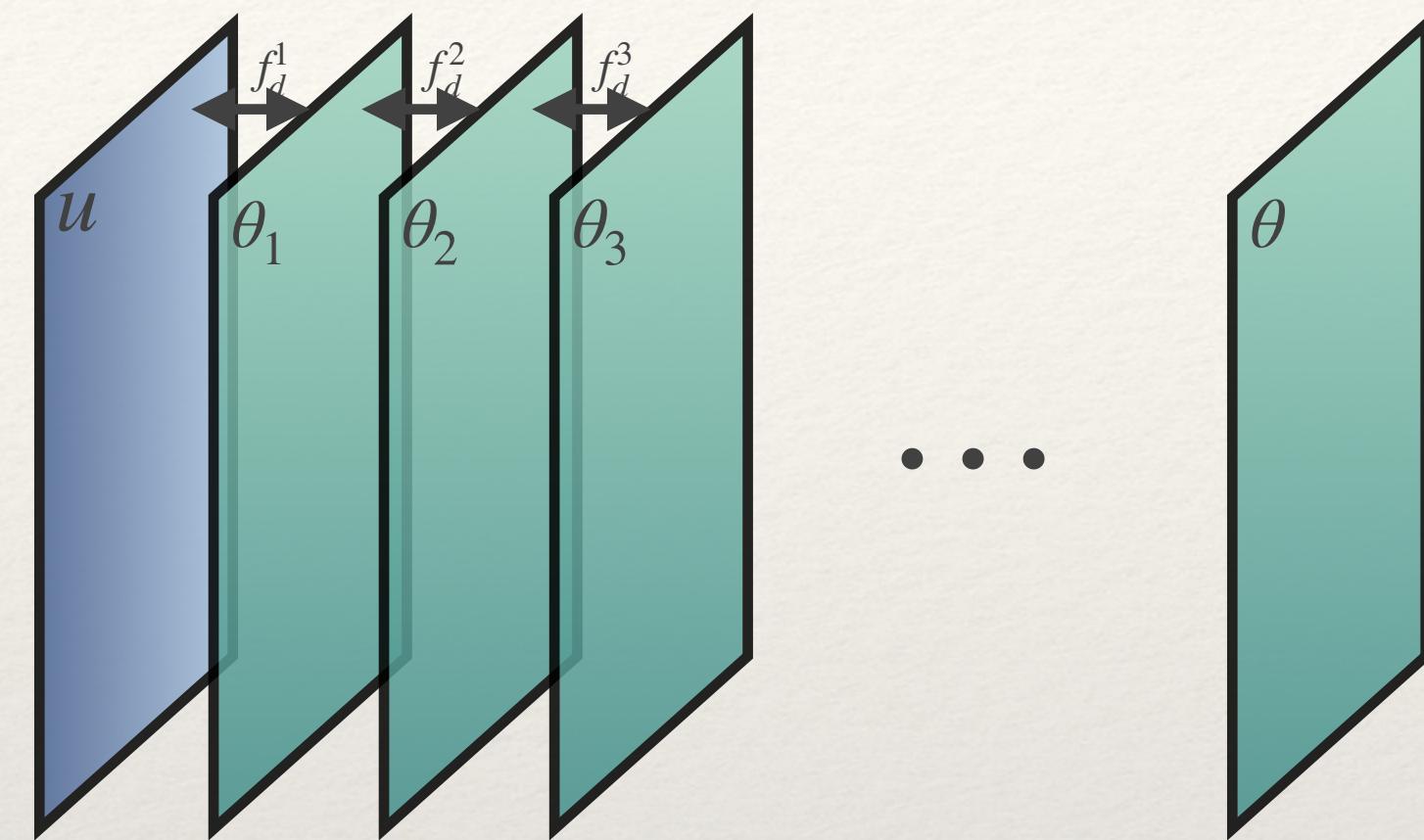
- ❖ Include varying noise spectrum to treat non-stationarity from event to event
- ❖ Inject signals into real noise
- ❖ Marginalize over unwanted parameters
- ❖ Use any data representation
- ❖ Include additional data channels

The network will learn how to incorporate this information.

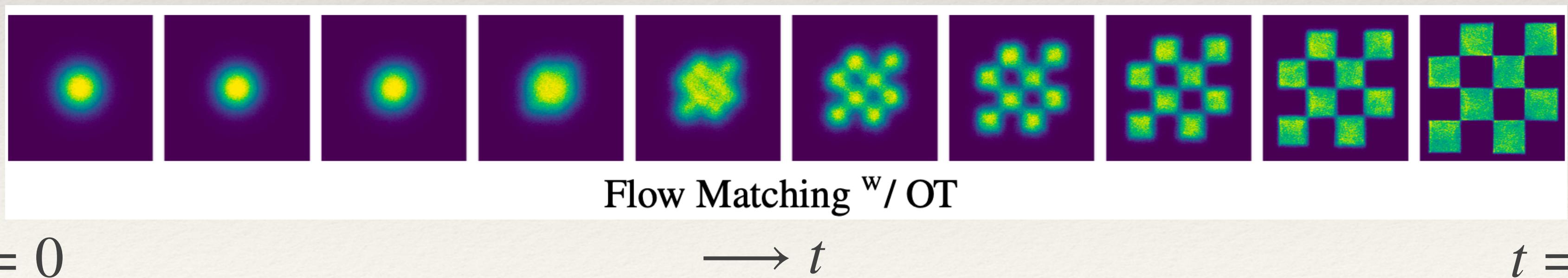
Challenge is making it work in practice!

Improved density estimators

- # ❖ Discrete normalizing flow



- ❖ → continuous normalizing flow / flow matching [Lipman+ (ICLR 2023)]

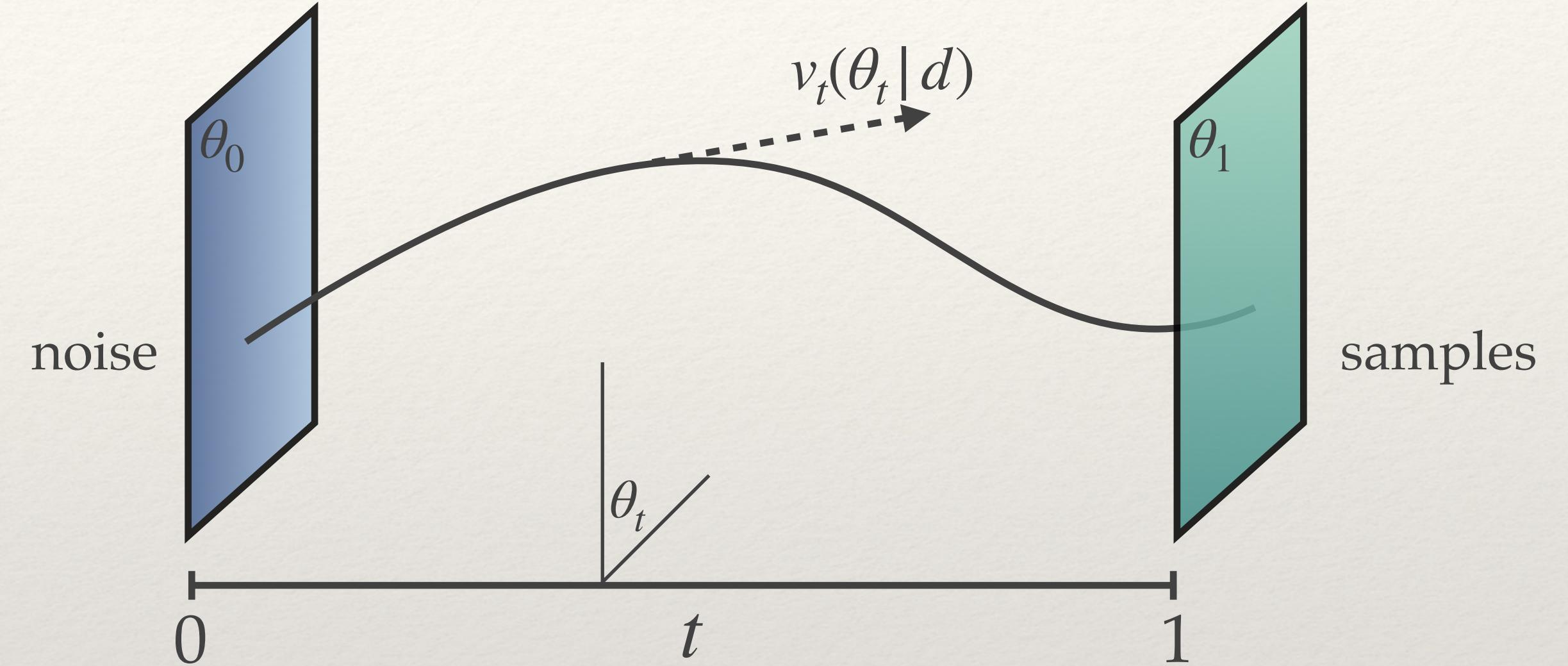


Continuous normalizing flows

- ❖ Sample trajectories defined by vector field

$$\frac{d\theta_t}{dt} = v_t(\theta_t | d)$$

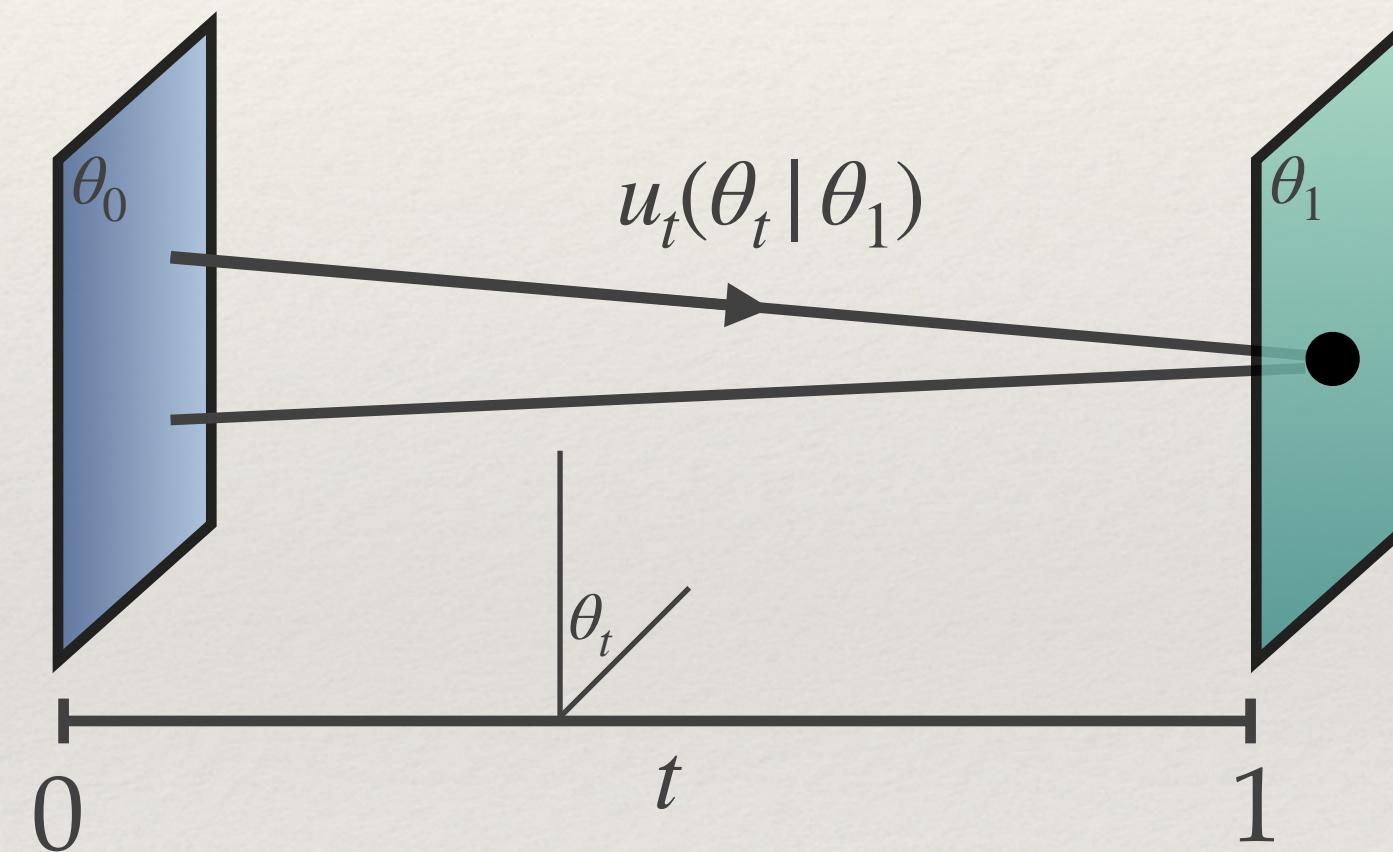
↑
neural network



- ❖ Density satisfies **transport equation** $\partial_t q_t = \nabla_{\theta_t} \cdot (q_t v_t)$ [cf. diffusion models - stochastic differential equations]
- ❖ However, expensive to sample and evaluate density, since many network evaluations required.
Makes training with loss $L = \mathbb{E}_{(\theta,d)}[-\log q(\theta | d)]$ impractical.

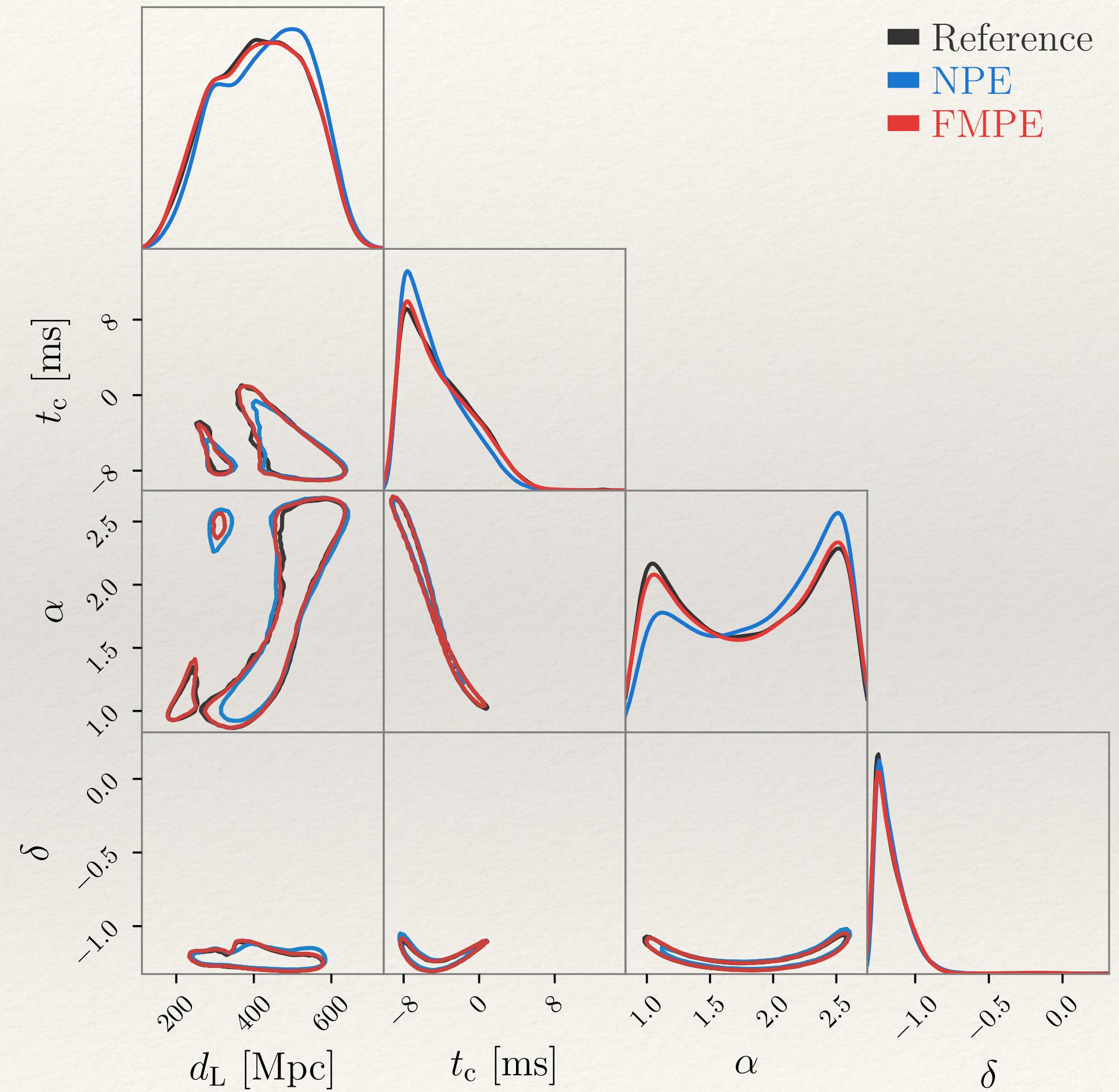
Flow matching

- ❖ Lipman+ (ICLR 2023): Directly regress vector field using **sample-conditional flow matching**
 - ❖ If we knew a target trajectory (u_t, p_t) that gives rise to an approximation to posterior $p(\theta_1 | d)$, **flow match** with $L_{\text{FM}} = \mathbb{E}_{(d, \theta_t, t)} \|v_t(\theta_t | d) - u_t(\theta_t | d)\|^2$
 - ❖ Instead, consider a single point θ_1 and a simple conditional trajectory $(u_t(\theta_t | \theta_1), p_t(\theta_t | \theta_1))$
- E.g., optimal transport
- ❖ Remarkably, **matching to the sample-conditional path** yields a vector field that gives the marginal path
$$L_{\text{SCFM}} = \mathbb{E}_{(d, \theta_1, t, \theta_t)} \|v_t(\theta_t | d) - u_t(\theta_t | \theta_1)\|^2$$



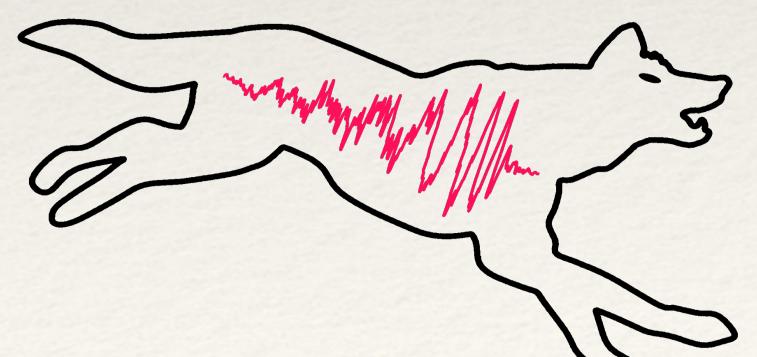
Flow matching posterior estimation (FMPE)

- ❖ We applied flow matching to posterior estimation.
- ❖ Proved that under reasonable assumptions, **mass covering property holds**.
- ❖ Outperforms NPE with discrete flows
 - ❖ **Faster training, better scaling to large networks**



Conclusions

- ❖ Simulation-based inference can deliver **fast and accurate** PE for gravitational waves.
 - ❖ Enabled new **eccentricity** studies, finding evidence in three events.
 - ❖ For **binary neutron stars**, prior conditioning enables O(second) pre- and post-merger inference.
 - ❖ New architectures promise to bring ever-improving performance.
- ❖ Going forwards, efforts focused on (1) extending to new sources and observatories, (2) training on realistic noise, and (3) building flexible networks.



Thank you!