

# Scalable Bayesian inference for 3G: Leveraging hardware acceleration and normalizing flows

Thibeau Wouters



Utrecht  
University

Nikhef



# Contents

① Introduction

② Methods

③ Applications

④ Outlook and conclusion

# Parameter estimation in 3G

Parameter estimation is done with **Bayesian inference**:

$$\text{posterior} \propto \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

- Sample the posterior: MCMC or nested sampling
- $\mathcal{O}(10^6) - \mathcal{O}(10^8)$  likelihood evaluations per inference

# Parameter estimation in 3G

Parameter estimation is done with **Bayesian inference**:

$$\text{posterior} \propto \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

- Sample the posterior: MCMC or nested sampling
- $\mathcal{O}(10^6) - \mathcal{O}(10^8)$  likelihood evaluations per inference

What about **3G** detectors?

- ET will observe  $\mathcal{O}(10^5)$  events per year
- Signals will be longer ( $f_{\min} = 5$  Hz) and have higher SNRs

# Parameter estimation in 3G

Parameter estimation is done with **Bayesian inference**:

$$\text{posterior} \propto \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

- Sample the posterior: MCMC or nested sampling
- $\mathcal{O}(10^6) - \mathcal{O}(10^8)$  likelihood evaluations per inference

What about **3G** detectors?

- ET will observe  $\mathcal{O}(10^5)$  events per year
- Signals will be longer ( $f_{\min} = 5$  Hz) and have higher SNRs

**Premise:** Current software will not scale to 3G [1]

# Towards scalable inference

Ingredients to make inference scalable:

- Hardware acceleration: GPUs for parallel computations
- Normalizing flows (NFs): neural density estimators

# Towards scalable inference

Ingredients to make inference scalable:

- Hardware acceleration: GPUs for parallel computations
- Normalizing flows (NFs): neural density estimators

Both are used in two ways:

- Simulation-based inference [2–6]:
  - + Fast at inference ( $\sim$ seconds)
  - Needs pre-trained model
- See Lucia Papalini, Filippo Santoliquido,...

# Towards scalable inference

Ingredients to make inference scalable:

- Hardware acceleration: GPUs for parallel computations
- Normalizing flows (NFs): neural density estimators

Both are used in two ways:

- Simulation-based inference [2–6]:
  - + Fast at inference ( $\sim$ seconds)
  - Needs pre-trained model
    - See Lucia Papalini, Filippo Santoliquido,...
- Hybrid approach: faster likelihoods + NF proposals
  - + No pre-training: more flexibility
  - Re-implement software
    - Also see Luca Negri's poster: neural likelihood estimators

# Towards scalable inference

Ingredients to make inference scalable:

- Hardware acceleration: GPUs for parallel computations
- Normalizing flows (NFs): neural density estimators

Both are used in two ways:

- Simulation-based inference [2–6]:
  - + Fast at inference ( $\sim$ seconds)
  - Needs pre-trained model
  - See Lucia Papalini, Filippo Santoliquido,...
- **Hybrid approach: faster likelihoods + NF proposals → this talk**
  - + No pre-training: more flexibility
  - Re-implement software
  - Also see Luca Negri's poster: neural likelihood estimators

# Contents

① Introduction

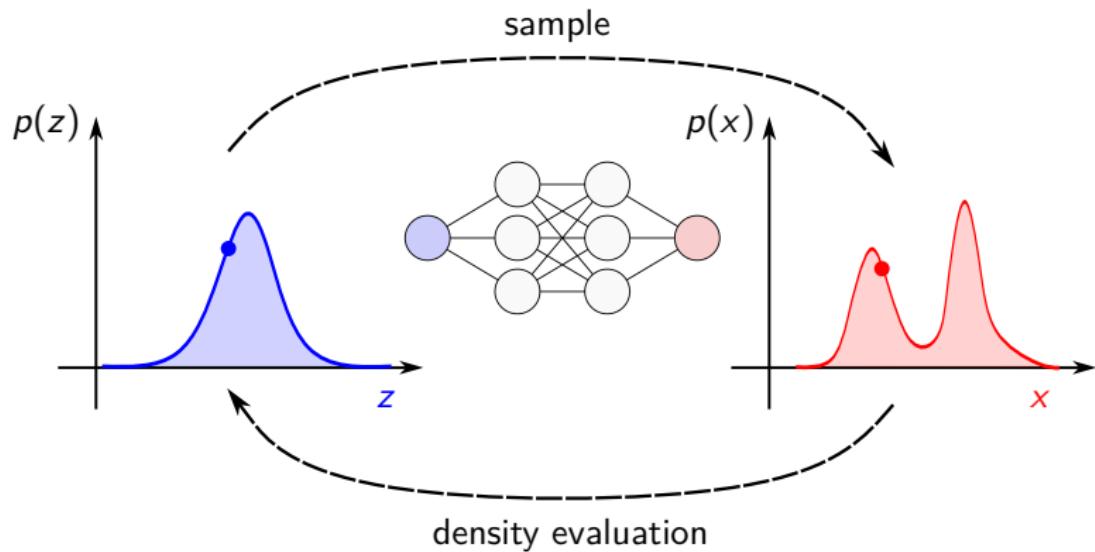
② Methods

③ Applications

④ Outlook and conclusion

# Normalizing flows (NFs)

- Trainable bijection between **latent** and **data** spaces
- Sample and evaluate complicated densities
- Used as **proposal distribution**, trained from MCMC chains



# JAX & FLOWMC

Acceleration with JAX features:

- Use GPU accelerators
- Automatic differentiation → gradient-based samplers
- Just-in-time (JIT) compilation



# JAX & FLOWMC

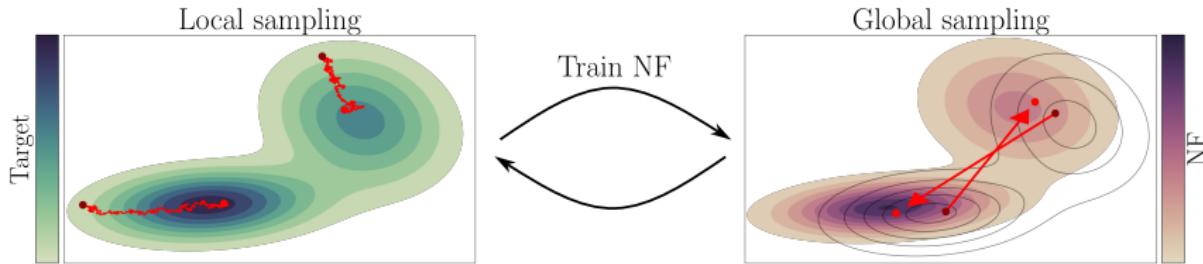
Acceleration with JAX features:

- Use GPU accelerators
- Automatic differentiation → gradient-based samplers
- Just-in-time (JIT) compilation



FLOWMC [7, 8]: MCMC + normalizing flows + JAX

- MCMC chains as training data: no pre-training
- Also see NESSAI  [9, 10]



# Contents

① Introduction

② Methods

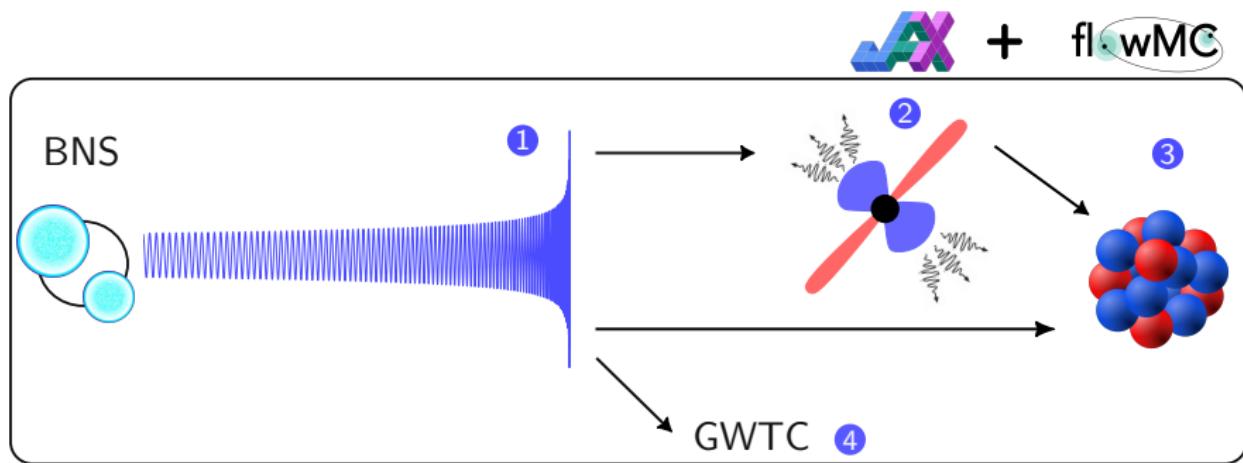
③ Applications

④ Outlook and conclusion

# Overview

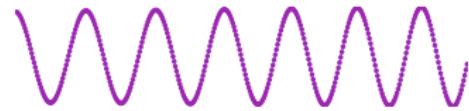
Analyzing a multi-messenger **binary neutron star** (BNS) signal:

- ① Gravitational waves
- ② Electromagnetic counterparts
- ③ Nuclear equation of state
- ④ Gravitational wave transient catalogue



# Gravitational waves

- Waveforms on GPU are  $\mathcal{O}(10^3)$  faster
- From LALSUITE to JAX: RIPPLE [11]
  - Also see SFTS



# Gravitational waves

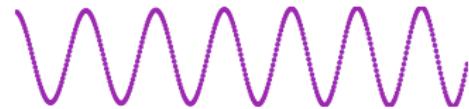
- Waveforms on GPU are  $\mathcal{O}(10^3)$  faster
  - From LALSUITE to JAX: RIPPLE  [11]
    - Also see SFTS 
- 



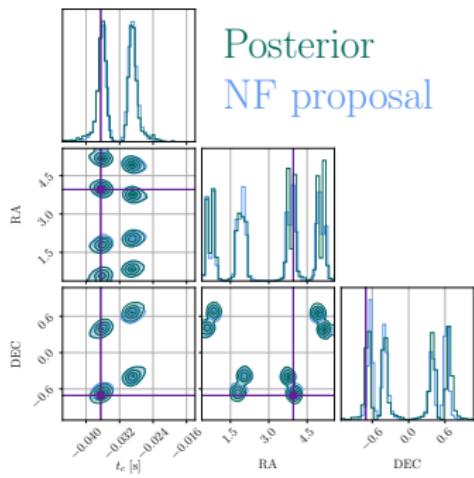
- Parameter estimation: JIM  [12, 13]
- BNS in LVK analyzed in  $\sim 15$  min

# Gravitational waves

- Waveforms on GPU are  $\mathcal{O}(10^3)$  faster
  - From LALSUITE to JAX: RIPPLE  [11]
    - Also see SFTS 
- 



- Parameter estimation: JIM  [12, 13]
- BNS in LVK analyzed in  $\sim 15$  min
- Ongoing work for ET – example:
  - BNS,  $f_{\min} = 20$  Hz, SNR = 21
  - ET- $\Delta$ , IMRPhenomD\_NRTidalv2
  - 30 mins on H100 GPU
- Evidence: HARMONIC [14]



# Overlapping signals

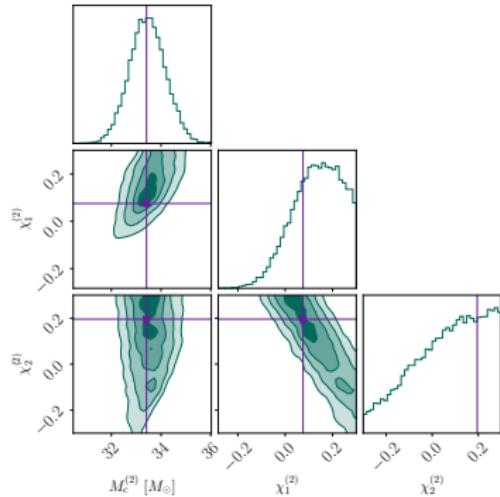
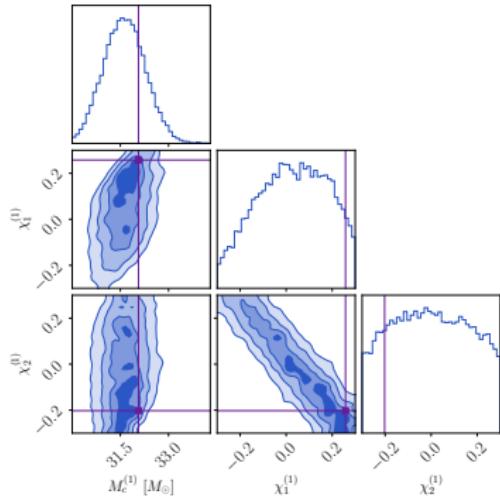
(Luca Negri, Justin Janquart, James Alvey, Uddipta Bhardwaj)

- Assess scaling of JIM: BBH+BBH in O5 with LVK
  - IMRPhenomD → 22 parameters (joint PE)
  - $M_c^{(1)} = 32M_\odot$ ,  $M_c^{(2)} = 33M_\odot$ ,  $\Delta t = 70$  ms
  - $\text{SNR}^{(1)} = 25.76$ ,  $\text{SNR}^{(2)} = 25.24$

# Overlapping signals

(Luca Negri, Justin Januart, James Alvey, Uddipta Bhardwaj)

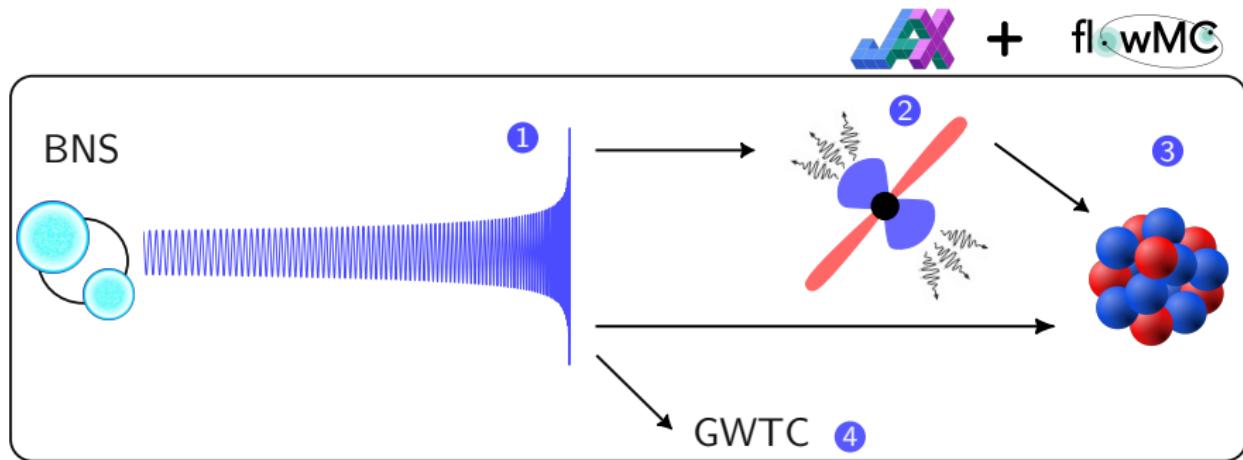
- Assess scaling of JIM: BBH+BBH in O5 with LVK
  - IMRPhenomD → 22 parameters (joint PE)
  - $M_c^{(1)} = 32M_\odot$ ,  $M_c^{(2)} = 33M_\odot$ ,  $\Delta t = 70$  ms
  - $\text{SNR}^{(1)} = 25.76$ ,  $\text{SNR}^{(2)} = 25.24$
  - **1h28m** on H100 (vs 23 days on 16 CPUs [15])



# Overview

Analyzing a multi-messenger **binary neutron star** (BNS) signal:

- ① Gravitational waves
- ② Electromagnetic counterparts
- ③ Nuclear equation of state
- ④ Gravitational wave transient catalogue



## Electromagnetic counterparts (Hauke Koehn, Tim Dietrich)

- BNS mergers lead to kilonovae, gamma-ray bursts
- Numerical models are too expensive (e.g. AFTERGLOWPY [16])

# Electromagnetic counterparts (Hauke Koehn, Tim Dietrich)

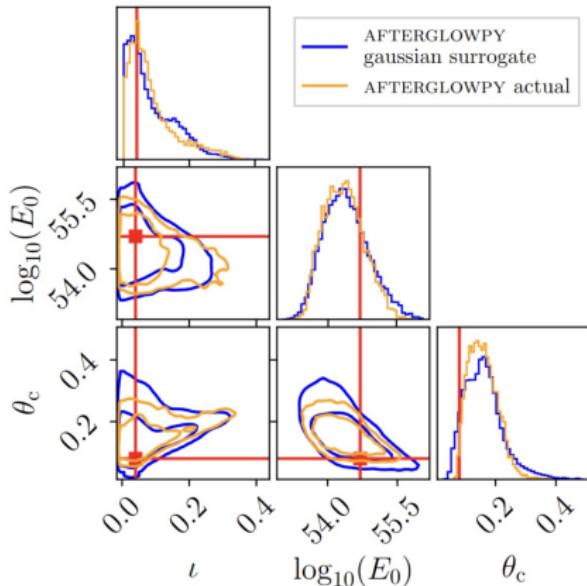
- BNS mergers lead to kilonovae, gamma-ray bursts
- Numerical models are too expensive (e.g. AFTERGLOWPY [16])
- Neural network surrogates for inference: FIESTA 

## FIESTA surrogates

- 1m36s
- 1 H100 GPU

## AFTERGLOWPY

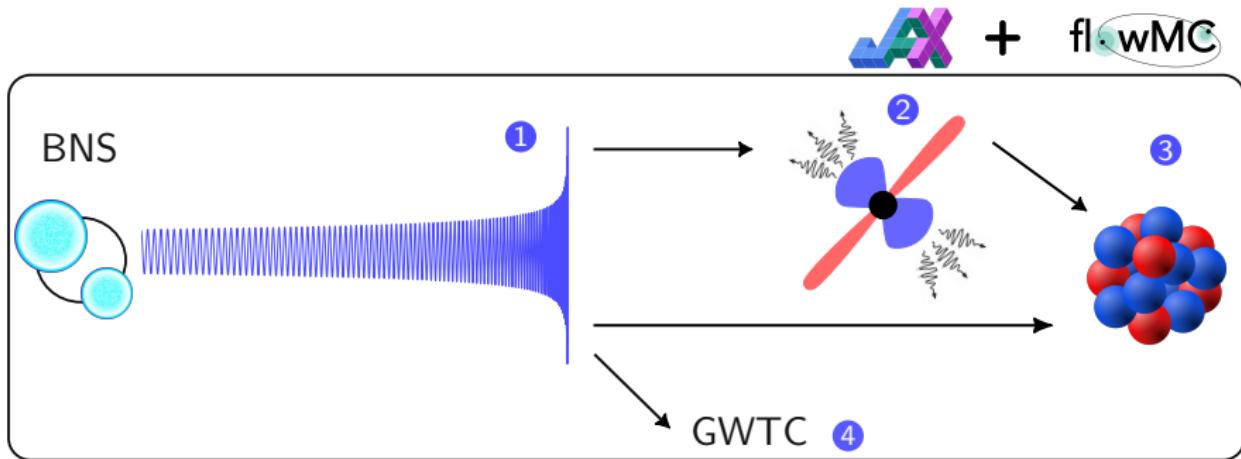
- 4 hours
- 30 CPUs



# Overview

Analyzing a multi-messenger **binary neutron star** (BNS) signal:

- ① Gravitational waves
- ② Electromagnetic counterparts
- ③ Nuclear equation of state
- ④ Gravitational wave transient catalogue

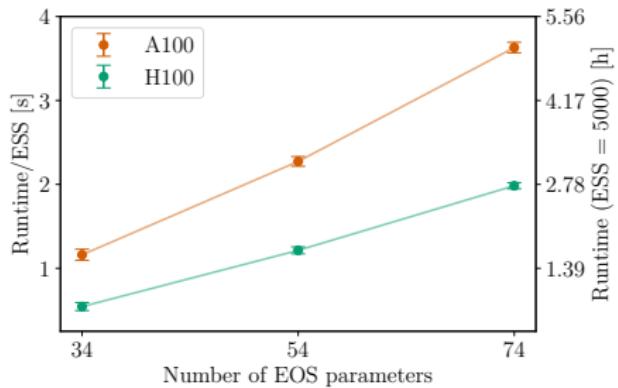


## Equation of state inference (Peter T.H. Pang)

- Goal: infer nuclear equation of state (EOS) of neutron stars [17]
- Computational bottleneck: solve TOV equations

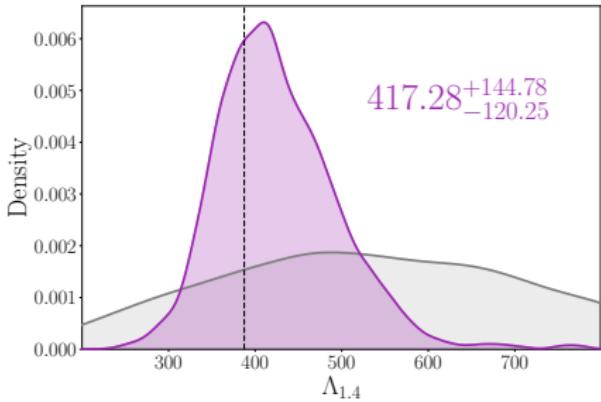
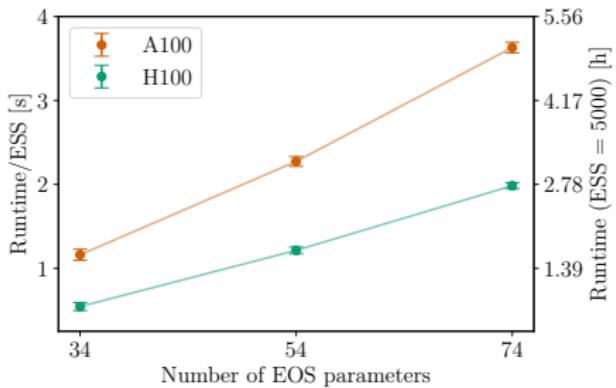
# Equation of state inference (Peter T.H. Pang)

- Goal: infer nuclear equation of state (EOS) of neutron stars [17]
- Computational bottleneck: solve TOV equations
- JESTER  [18]: JAX-based TOV solver
  - Full inference in  $\sim$ hours
  - No need for machine learning surrogates



# Equation of state inference (Peter T.H. Pang)

- Goal: infer nuclear equation of state (EOS) of neutron stars [17]
- Computational bottleneck: solve TOV equations
- JESTER  [18]: JAX-based TOV solver
  - Full inference in  $\sim$ hours
  - No need for machine learning surrogates
- End-to-end analysis: constrain EOS from 20 BNS in O5



# Contents

① Introduction

② Methods

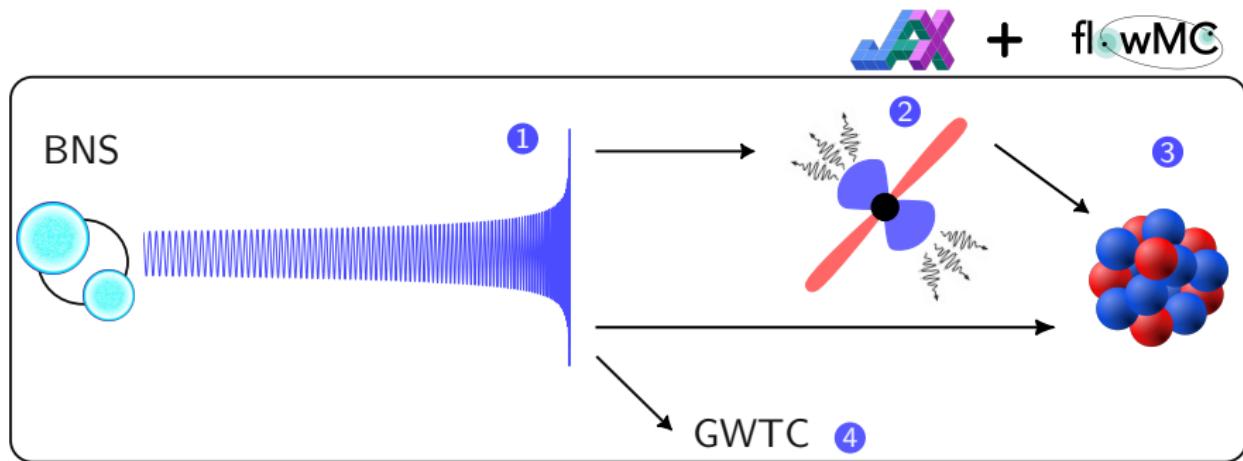
③ Applications

④ Outlook and conclusion

# Overview

Analyzing a multi-messenger **binary neutron star** (BNS) signal:

- ① Gravitational waves
- ② Electromagnetic counterparts
- ③ Nuclear equation of state
- ④ **Gravitational wave transient catalogue**



# Constructing GWTCs (Thomas Ng, Kaze Wong)

How can we use the trained normalizing flow proposal from FLOWMC?

Standard



FLOWMC



# Conclusion

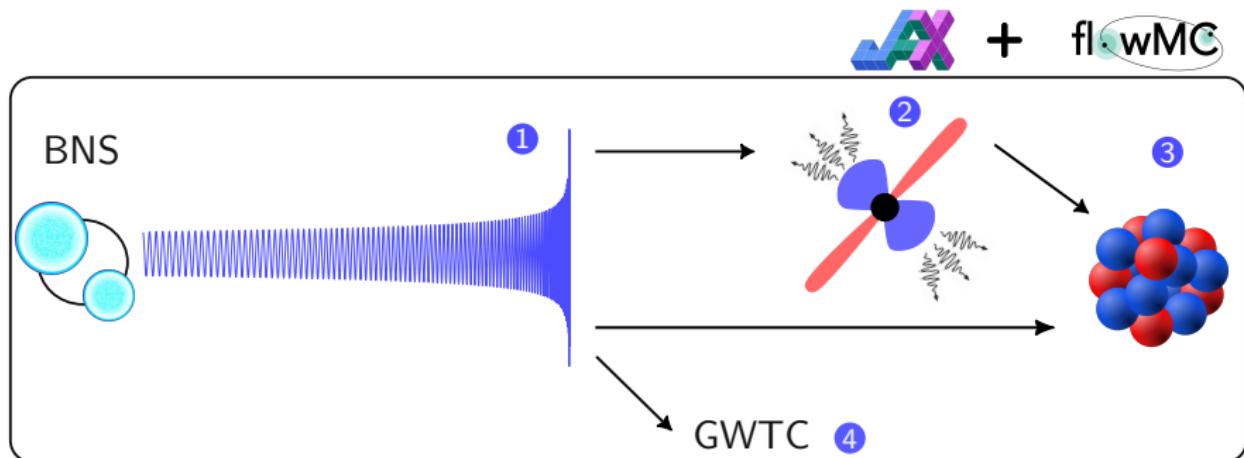
- Progress on scalable Bayesian inference software for 3G, with minimal amount of pre-training
- Hybrid acceleration: GPU + normalizing flow
  - JAX/GPU: likelihoods faster
  - FLOWMC: sampling converges faster
- Goal: joint multimessenger analyses in  $\sim$ hours (NMMA [19] in JAX)
- To do
  - GW injection studies for ET
  - More waveform models in JAX
  - Equation of state study with ET data

Let's talk!

# Thank you for your attention!

Software:

- FLOWMC  [7, 8]
- JIM  [12, 13]   
- FIESTA  
- JESTER  
- HARMONIC  [14, 20, 21]



# References I

- [1] Qian Hu and John Veitch. "Costs of Bayesian Parameter Estimation in Third-Generation Gravitational Wave Detectors: a Review of Acceleration Methods". In: (Dec. 2024). arXiv: [2412.02651 \[gr-qc\]](https://arxiv.org/abs/2412.02651).
- [2] Jurriaan Langendorff et al. "Normalizing Flows as an Avenue to Studying Overlapping Gravitational Wave Signals". In: *Phys. Rev. Lett.* 130.17 (2023), p. 171402. DOI: [10.1103/PhysRevLett.130.171402](https://doi.org/10.1103/PhysRevLett.130.171402). arXiv: [2211.15097 \[gr-qc\]](https://arxiv.org/abs/2211.15097).
- [3] Uddipta Bhardwaj et al. "Sequential simulation-based inference for gravitational wave signals". In: *Phys. Rev. D* 108.4 (2023), p. 042004. DOI: [10.1103/PhysRevD.108.042004](https://doi.org/10.1103/PhysRevD.108.042004). arXiv: [2304.02035 \[gr-qc\]](https://arxiv.org/abs/2304.02035).
- [4] Maximilian Dax et al. "Real-time inference for binary neutron star mergers using machine learning". In: *Nature* 639.8053 (2025), pp. 49–53. DOI: [10.1038/s41586-025-08593-z](https://doi.org/10.1038/s41586-025-08593-z). arXiv: [2407.09602 \[gr-qc\]](https://arxiv.org/abs/2407.09602).
- [5] Qian Hu et al. "Decoding Long-duration Gravitational Waves from Binary Neutron Stars with Machine Learning: Parameter Estimation and Equations of State". In: (Dec. 2024). arXiv: [2412.03454 \[gr-qc\]](https://arxiv.org/abs/2412.03454).
- [6] Filippo Santoliquido et al. "Fast and accurate parameter estimation of high-redshift sources with the Einstein Telescope". In: (Apr. 2025). arXiv: [2504.21087 \[astro-ph.HE\]](https://arxiv.org/abs/2504.21087).

## References II

- [7] Marylou Gabrié, Grant M. Rotskoff, and Eric Vanden-Eijnden. “Adaptive Monte Carlo augmented with normalizing flows”. In: *Proc. Nat. Acad. Sci.* 119.10 (2022), e2109420119. DOI: [10.1073/pnas.2109420119](https://doi.org/10.1073/pnas.2109420119). arXiv: [2105.12603 \[physics.data-an\]](https://arxiv.org/abs/2105.12603).
- [8] Kaze W. k. Wong, Marylou Gabrié, and Daniel Foreman-Mackey. “flowMC: Normalizing flow enhanced sampling package for probabilistic inference in JAX”. In: *J. Open Source Softw.* 8.83 (2023), p. 5021. DOI: [10.21105/joss.05021](https://doi.org/10.21105/joss.05021). arXiv: [2211.06397 \[astro-ph.IM\]](https://arxiv.org/abs/2211.06397).
- [9] Michael J. Williams, John Veitch, and Chris Messenger. “Nested sampling with normalizing flows for gravitational-wave inference”. In: *Phys. Rev. D* 103.10 (2021), p. 103006. DOI: [10.1103/PhysRevD.103.103006](https://doi.org/10.1103/PhysRevD.103.103006). arXiv: [2102.11056 \[gr-qc\]](https://arxiv.org/abs/2102.11056).
- [10] Michael J. Williams, John Veitch, and Chris Messenger. “Importance nested sampling with normalising flows”. In: *Mach. Learn. Sci. Tech.* 4.3 (2023), p. 035011. DOI: [10.1088/2632-2153/acd5aa](https://doi.org/10.1088/2632-2153/acd5aa). arXiv: [2302.08526 \[astro-ph.IM\]](https://arxiv.org/abs/2302.08526).
- [11] Thomas D. P. Edwards et al. “Differentiable and hardware-accelerated waveforms for gravitational wave data analysis”. In: *Phys. Rev. D* 110.6 (2024), p. 064028. DOI: [10.1103/PhysRevD.110.064028](https://doi.org/10.1103/PhysRevD.110.064028). arXiv: [2302.05329 \[astro-ph.IM\]](https://arxiv.org/abs/2302.05329).
- [12] Kaze W. K. Wong, Maximiliano Isi, and Thomas D. P. Edwards. “Fast Gravitational-wave Parameter Estimation without Compromises”. In: *Astrophys. J.* 958.2 (2023), p. 129. DOI: [10.3847/1538-4357/acf5cd](https://doi.org/10.3847/1538-4357/acf5cd). arXiv: [2302.05333 \[astro-ph.IM\]](https://arxiv.org/abs/2302.05333).

## References III

- [13] Thibeau Wouters et al. "Robust parameter estimation within minutes on gravitational wave signals from binary neutron star inspirals". In: *Phys. Rev. D* 110.8 (2024), p. 083033. DOI: [10.1103/PhysRevD.110.083033](https://doi.org/10.1103/PhysRevD.110.083033). arXiv: [2404.11397 \[astro-ph.IM\]](https://arxiv.org/abs/2404.11397).
- [14] Alicja Polanska et al. "Accelerated Bayesian parameter estimation and model selection for gravitational waves with normalizing flows". In: *38th conference on Neural Information Processing Systems*. Oct. 2024. arXiv: [2410.21076 \[astro-ph.IM\]](https://arxiv.org/abs/2410.21076).
- [15] Justin Janquart et al. "Analyses of overlapping gravitational wave signals using hierarchical subtraction and joint parameter estimation". In: *Mon. Not. Roy. Astron. Soc.* 523.2 (2023), pp. 1699–1710. DOI: [10.1093/mnras/stad1542](https://doi.org/10.1093/mnras/stad1542). arXiv: [2211.01304 \[gr-qc\]](https://arxiv.org/abs/2211.01304).
- [16] Geoffrey Ryan et al. "Gamma-Ray Burst Afterglows in the Multimessenger Era: Numerical Models and Closure Relations". In: *Astrophys. J.* 896.2 (2020), p. 166. DOI: [10.3847/1538-4357/ab93cf](https://doi.org/10.3847/1538-4357/ab93cf). arXiv: [1909.11691 \[astro-ph.HE\]](https://arxiv.org/abs/1909.11691).
- [17] Adrian Abac et al. "The Science of the Einstein Telescope". In: (Mar. 2025). arXiv: [2503.12263 \[gr-qc\]](https://arxiv.org/abs/2503.12263).
- [18] Thibeau Wouters et al. "Leveraging differentiable programming in the inverse problem of neutron stars". In: (Apr. 2025). arXiv: [2504.15893 \[astro-ph.HE\]](https://arxiv.org/abs/2504.15893).

## References IV

- [19] Peter T. H. Pang et al. "An updated nuclear-physics and multi-messenger astrophysics framework for binary neutron star mergers". In: *Nature Commun.* 14.1 (2023). Available at <https://github.com/nuclear-multimessenger-astronomy/nmma>, p. 8352. DOI: 10.1038/s41467-023-43932-6. arXiv: 2205.08513 [astro-ph.HE].
- [20] Jason D. McEwen et al. *Machine learning assisted Bayesian model comparison: learnt harmonic mean estimator*. 2023. arXiv: 2111.12720 [stat.ME]. URL: <https://arxiv.org/abs/2111.12720>.
- [21] Alicja Polanska et al. *Learned harmonic mean estimation of the marginal likelihood with normalizing flows*. 2024. arXiv: 2307.00048 [stat.ME]. URL: <https://arxiv.org/abs/2307.00048>.
- [22] Kurzgesagt. *Figures taken from “Neutron Stars - The Most Extreme Things that are not Black Holes”*. Accessed on May 14, 2025. 2019. URL: <https://www.youtube.com/watch?v=udFxKZRyQt4>.
- [23] Hergé. *Cover figure created with ChatGPT using this input figure from the comic Destination Moon*. Accessed on May 14, 2025. 2019. URL: <https://www.youtube.com/watch?v=udFxKZRyQt4>.

## Evidence calculation: HARMONIC I

Evidence  $Z$  can be computed from posterior samples with HARMONIC [20] with the **harmonic mean estimator**

$$\begin{aligned}\rho &\equiv \mathbb{E}_{P(\theta|d)} \left[ \frac{1}{L(\theta)} \right] \\ &= \int d\theta \frac{1}{L(\theta)} P(\theta|d) \\ &= \int d\theta \frac{1}{L(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{Z} = \frac{1}{Z}\end{aligned}$$

Therefore, estimate  $\rho$  with posterior samples:

$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{1}{L(\theta_i)}, \quad \theta_i \sim P(\theta|d)$$

## Evidence calculation: HARMONIC II

Can be interpreted as importance sampling

$$\rho = \int d\theta \frac{1}{Z} \frac{\pi(\theta)}{P(\theta|d)} P(\theta|d),$$

**but** with target = prior and sampling density = posterior. Therefore, importance sampling is inefficient – how to solve?

New proposal:

$$\begin{aligned}\rho &= \mathbb{E}_{P(\theta|d)} \left[ \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \right] \\ &= \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} P(\theta|d) \\ &= \int d\theta \frac{\varphi(\theta)}{\mathcal{L}(\theta)\pi(\theta)} \frac{\mathcal{L}(\theta)\pi(\theta)}{Z} = \frac{1}{Z}\end{aligned}$$

## Evidence calculation: HARMONIC III

Use the following estimator:

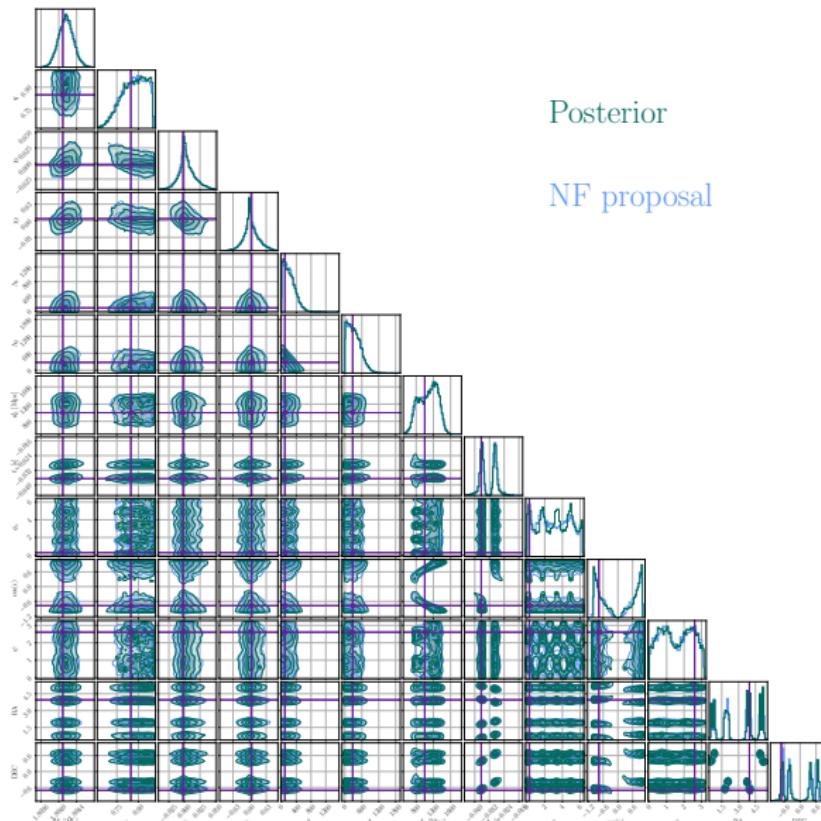
$$\hat{\rho} = \frac{1}{N} \sum_{i=1}^N \frac{\varphi(\theta_i)}{\mathcal{L}(\theta_i)\pi(\theta_i)}, \quad \theta_i \sim P(\theta|d)$$

Replace the target distribution  $\pi$  with  $\varphi$ : only requirement is that it is normalized

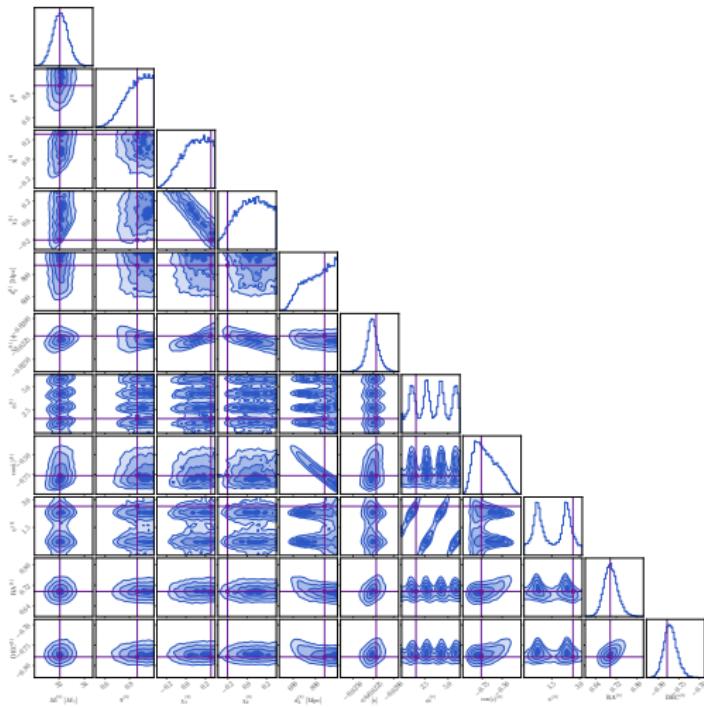
In practice, this can be achieved with a normalizing flow [21].

This has been verified to give accurate evidences (similar values as nested sampling) when GW posteriors are used [14].

# BNS in ET- $\Delta$ example: all parameters



# Overlapping signals: all parameters signal A



Overlapping signals: all parameters signal B

