GR-Amaldi – July 14th, 2025

Identifying and Mitigating Machine Learning Biases for the Gravitational Wave Detection Problem

arxiv: 2501.13846 (accepted for publication in PRD)





Narenraju Nagarajan Chris Messenger



<u>Problem Setup</u>: Detection of BBH GWs interpreted as a classification problem in the deep learning framework.

Generate BBH GW

signals given parameters randomised from the standard distributions

Get real detector noise for H1 and L1 from the third observing run

Sample length set based on lowest possible chirp mass





<u>Motivation</u>: Previous works were affected by ML biases. Finding the source of these biases and mitigating them might allow ML to move closer to the detection performance of matched-filtering pipelines.

M. B. Schäfer et al., Phys. Rev. D 107, 023021 (2023)

Koloniari, Alexandra E. et al. (2024) -- AresGW

Marx, Ethan et al. (2024) -- AFrame





- * Imagine your definition of a "perfect" search/classification analysis
 - * This is likely to be either a Bayesian of Frequentist or other?
 - * Remember, an ideal & realistic analysis <u>does not detect everything</u>
- * We don't have access to an ideal analysis for comparison
 - * So biases in detection e.g. as a function of chirp mass are measured relatively

Identifying and Mitigating Machine Learning Biases for the Gravitational **Wave Detection Problem**



Bias due to limited samples Imbalance in the sample distribution for a certain GW parameter might introduce biases.

Vokinger, K.N et al. Commun Med 1, 25 (2021) Tang, Kaihua et al. ArXiv abs/2009.12991 (2020)



$$\tau = t_c - 5(8\pi f_l)^{-8/3} \mathcal{M}_c^{-5/3}$$



Create image Clock with time set to 5:40

Image created



Is it true that chatGPT cannot create images of a clock with times other than 10:10?

No, that's not true — ChatGPT *can* create images of clocks showing any time, including times other than 10:10. However, some AI image generation systems (including DALL·E, which powers image generation here) sometimes default to 10:10 due to training data biases, as that's a common time shown in product photos of clocks. It's not a hard limit, and I can correct it on request — as I did when generating the 5:40 clock in the last image.

The time is 10:10 on that image. I asked for 5:40. Correct this.

Image created





Bias due to lack of variation Lack of augmentation along GW parameters that the model cannot interpolate trivially or lack of variation in noise characteristics.

Shah, Harshay et al. ArXiv abs/2006.07710 (2020) A. Mikołajczyk et al., IIPhDW, pp. 117-122 (2018)





Spectral bias

Lower frequency features are easier to learn than higher frequency features. Generalisation toward higher frequencies are typically brittle.

Rahaman, Nasim et al. ICML (2018). Cao, Y. et al., arXiv:1912.01198 (2019)



ML model was trained to learn sine waves at different frequencies but same signal duration. The colours represent how well the network learns samples of a certain frequency during training. 1 is perfect and 0 is non-existent learning. Rahaman et al. Proceedings of the 36th International Conference on Machine Learning, PMLR 97:5301-5310, 2019.



Other biases discussed in the paper [arxiv: 2501.13846 - accepted for publication in PRD]

- Bias due to class overlap 1.
- Bias due to class imbalance
- Bias due to train-test dataset overlap 3.
- Bias due to train-test mismatch 4.
- Bias due to disproportionate evaluation 5.
- Bias due to limited feature representation 6.
- Bias due to sample difficulty 7.
- Bias due to insufficient information 8.

and possibly more...





Sage Methodology

- * The primary features of SAGE
 - * Bias is very sensitive to the mass training distribution
 - Thoughtful network design addresses inductive bias

- Out-of-distribution PSDs aid classification
- Training on real data is key



Our aim was to use the MLGWSC-1 [Schäfer *et al.*, 2023] dataset for apples-to-apples comparison

1. Injection study with 1 month of O3a detector noise (H1, L1) 2. Injected signals generated using IMRPhenomXPHM 3. Mass priors are U[7, 50] solar masses with m1 > m24. Spins χ are distributed isotropically with magnitudes between |0.0, 0.99|5. Other GW parameters were randomised from the standard distributions









Comparing the histogram of found injections between Sage and PyCBC for the injection study at a False Alarm Rate (FAR) of 1 per month.

We find 443 signals more than matched-filtering and is equivalent to an increase of ~11.2% in the number of signals detected.

PyCBC results obtained from - M. B. Schäfer et al., Phys. Rev. D 107, 023021 (2023)







Comparing the histogram of found injections between Sage and AresGW for the injection study at a False Alarm Rate (FAR) of 1 per month.

We find 1397 signals more than AresGW on this dataset and is equivalent to an increase of ~48.29% in the number of signals detected.

AresGW results obtained from - P. Nousi et al., Phys. Rev. D 108, 024022 (2023)





$$V(\mathcal{F}) \approx \frac{V(d_{max})}{N_I} \sum_{i=1}^{N_{I,\mathcal{F}}} \left(\frac{\mathcal{M}_{c,i}}{\mathcal{M}_{c,max}}\right)^{5/2}$$
 Para as a s

ameter-dependent sensitive distance, $\propto V(\mathscr{F})^{\frac{1}{3}}$, function of false alarm rate per month.







- PSD estimation is not required for ML-based detection

- Mitigating biases might aid in glitch rejection

- OOD PSDs are good for classification

- Data / parameter efficiency via GW domain knowledge

(Empirical Contributions)





Identified sources of biases in ML for GW detection problem

Provided training strategies and mitigation tactics to address biases concurrently.

We detect ~11.2% more signals than the benchmark PyCBC result via the MLGWSC-1 injection study at an FAR of 1/ month in O3a noise.

Analysis of O3 is underway :)

arxiv: 2501.13846 - accepted for publication in PRD



Talk @ AIslands – May 13th, 2025

Supplemental Slides

arxiv: 2501.13846





Narenraju Nagarajan Chris Messenger Histograms of chirp mass of the found signals for the injection study at an FAR of 1/month



"Sage – Broad" has a broader distribution of noise PSDs than "Sage – Limited".

"Sage – D4 Metric" uses the template bank density, $U(\tau_0, \tau_3)$, to set the mass priors during training. "Sage – Annealed Train" transitions from U(m1, m2) to $U(\tau_0, \tau_3)$ during training.

"Sage – D3 Metric" uses $U(\tau_0, \tau_3)$ to set mass priors for training. Uses simulated coloured Gaussian noise for training and testing.





Ablation Study

All aspects of the methodology were kept the same except following changes















FAR for those triggers as a function of its corresponding glitch SNR for the H1 and L1 detectors. We conclude that Sage can reject glitches effectively in real detector noise.

Obtained all triggers by Sage related to glitch events in the O3a testing dataset noise (event times and durations obtained from GravitySpy). Scatter plot shows the ranking statistic and



21



Changing the Nyquist limit from solid red line to dashed red line as a form of time-series compression





Compression using multirate sampling

Reduce need for larger receptive field



2

Example of compression applied to a BBH gravitational-wave



Increased variance provided for noise class via augmentation

3







Augmentation Method 2



Augmentation Method 1



Finetuned training strategies to handle biases due to limited samples

5



Transitioning from U(m1, m2) to $U(\tau_0, \tau_3)$ during training









 \mathcal{M}_{c}

2D Histogram of $U(\tau_0, \tau_3)$ plotted in the mass1 vs mass2 space



Proper inductive bias used for two stage network *architecture*





Multi-scale Residual Block $C_{out}^{0} = p, k^{0} = q$ Multi-scale Residual Block $C_{out}^1 = p, k^1 = \lfloor q/2 \rfloor$ Max-Pooling S = k = 8Multi-scale Residual Block $C_{out}^2 = 2p, k^2 = \lfloor q/2 \rfloor$ Multi-scale Residual Block $C_{out}^3 = 2p, k^3 = \lfloor q/4 \rfloor$ Max-Pooling S = k = 4Multi-scale Residual Block $C_{out}^4 = 4p, k^4 = \lfloor q/4 \rfloor$ Multi-scale Residual Block

Input Time Series



6



Backend classifier powered by spatial and channel-wise attention









*Aframe is not optimised for the D4 mass distribution in MLGWSC-1. They operate on a broader mass prior [5, 100] Msun.









Generated a noise class of 250,000 noise realisations from the O3a noise PSDs. Generated a signal class dataset by injecting the same signal into these noise realisations.

Comparing the normalised histogram of network outputs for "Sage – Broad" and "Sage – Limited" suggests that training a model with a larger variance in noise PSDs allows it to be more confident about a given signal compared to a model that is not.









Actual value of GW parameters: chirp mass and time of coalescence, plotted against the network predicted value for a validation epoch in different Sage runs. The top row corresponds to the Sage -Annealed Train run and the bottom row to the Sage - Broad run.







Error bar obtained using 3 runs of Sage with different training seeds.





















