# ML4GW: An AI-based pipeline for Real-time Gravitational Wave Analysis

Erik Katsavounidis (MIT, LIGO Laboratory,A3D3 Institute)

on behalf of the ML4GW team: Will Benoit, Deep Chatterjee, Michael Coughlin, Malina Desai, Katya Govorkova, Alec Gunny, Phil Harris, Ethan Marx, Eric Moreno, Saleem Muhamed, Rafia Omer, Ryan Raikman

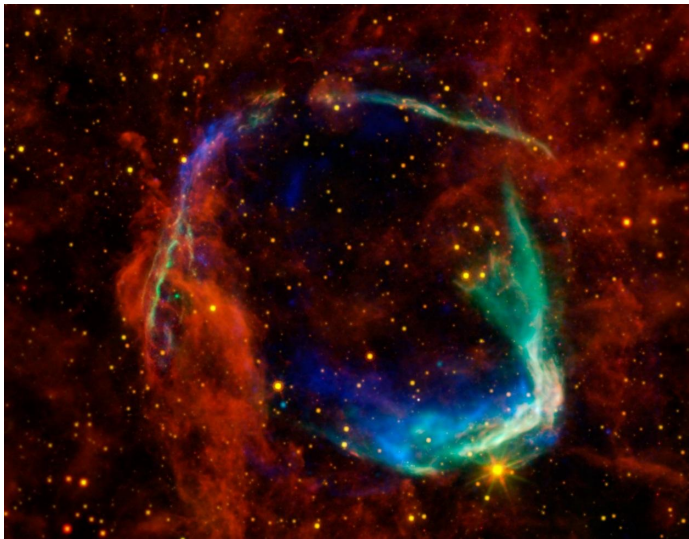at the University of Minnesota, MIT and the A3D3 Institute

17th Marcel Grossmann Meeting

July 7-12, 2024

# Observing our Universe



SN185

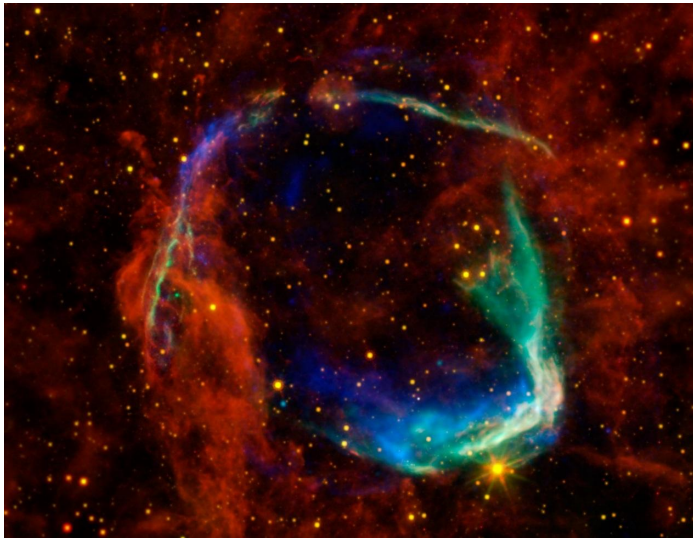*Image Credit: X-ray: NASA/CXC/SAO & ESA; Infared: NASA/JPL-Caltech/B. Williams (NCSU)*

SN1987A

Electromagnetic waves

# Observing our Universe



SN185

*Image Credit: X-ray: NASA/CXC/SAO & ESA; Infared: NASA/JPL-Caltech/B. Williams (NCSU)*



SN1987A    Koshiba, M. et al. 1988

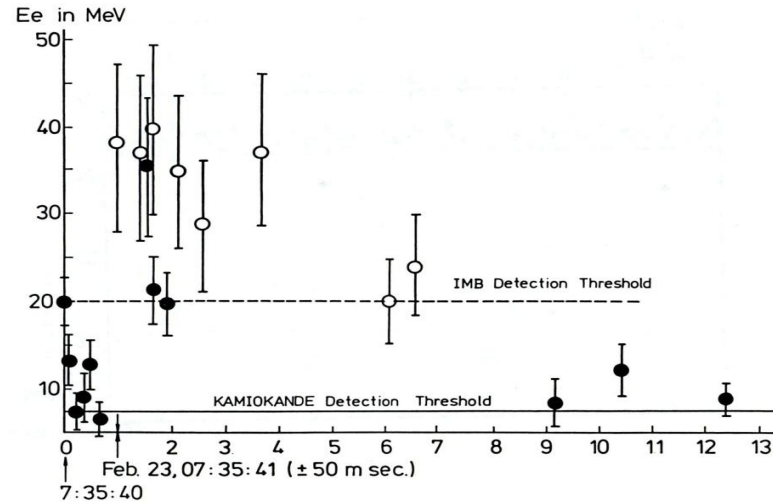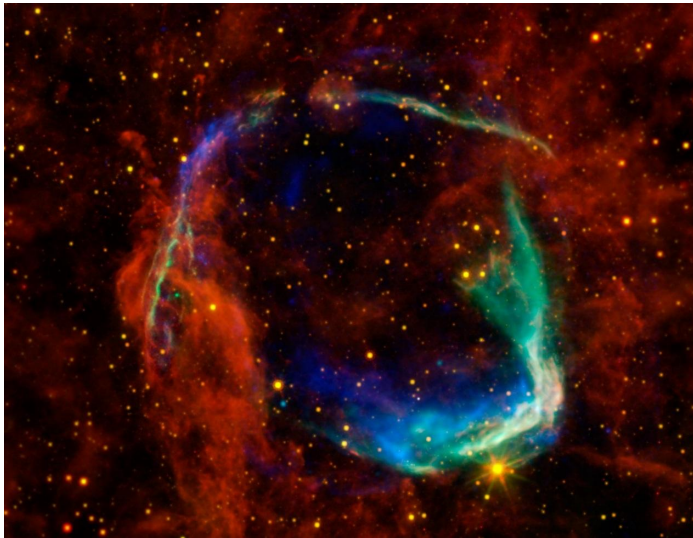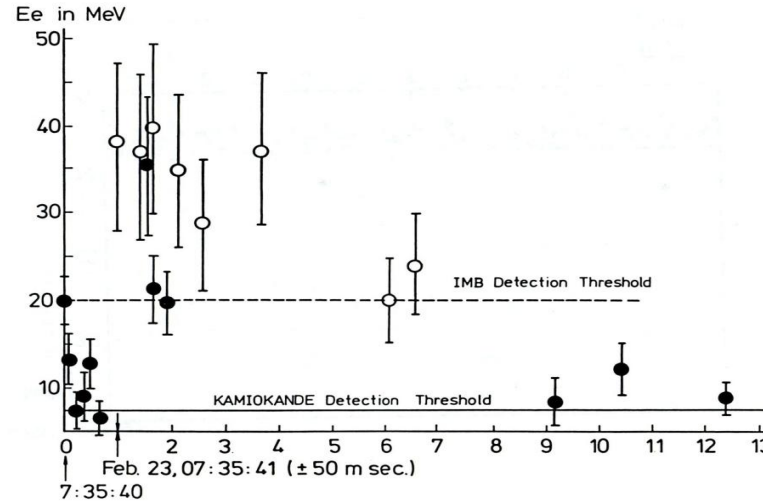Electromagnetic waves

Particles: neutrinos, cosmic rays

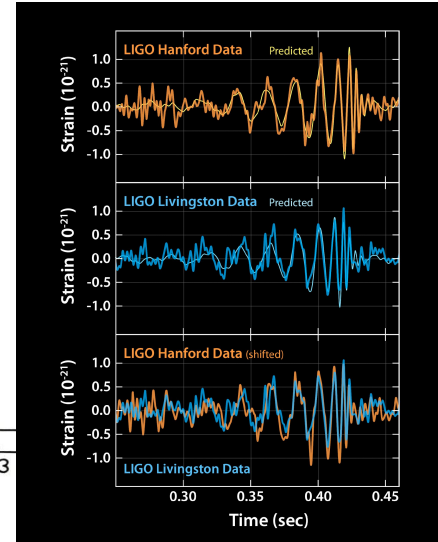3

# Observing our Universe



SN185

*Image Credit: X-ray: NASA/CXC/SAO & ESA; Infared: NASA/JPL-Caltech/B. Williams (NCSU)*



SN1987A

Koshiba, M. et al. 1988



GW150914
PRL 116, 061102, 2016

Electromagnetic waves

Particles: neutrinos, cosmic rays

Gravitational waves

4

# Gravitational wave detectors

Strain Noise
$$h = \Delta L \, / \, L$$

Strain 1/√Hz

Initial LIGO

Advanced LIGO

The lower noise gets, the better the ability to detect GWs

27.500 30.863 32.703 36.708 41.203 43.654 48.999 55.000 61.735 65.406 73.416 82.407 87.307 97.999 110.00 123.47 130.81 146.83 164.81 174.61 196.00 220.00 246.94 261.63 293.66 329.63 349.23 392.00 440.00 493.88 523.25 587.33 659.26 698.46 783.99 880.00 987.77 1046.5 1174.7 1318.5 1396.9 1568.0 1760.0 1975.5 2093.0 2349.3 2637.0 2793.0 3136.0 3520.0 3951.1 4186.0

A B C D E F G A B C D E F G A B C D E F G A B C D E F G A B C D E F G A B C D E F G A B C D E F G A B C

10Hz · 100Hz · 1kHz · 10kHz

test mass (mirror)

photodiode

LIGO Hanford Data (shifted)

Strain (10⁻²¹)

LIGO Livingston Data

Time (sec)

(after calibration)

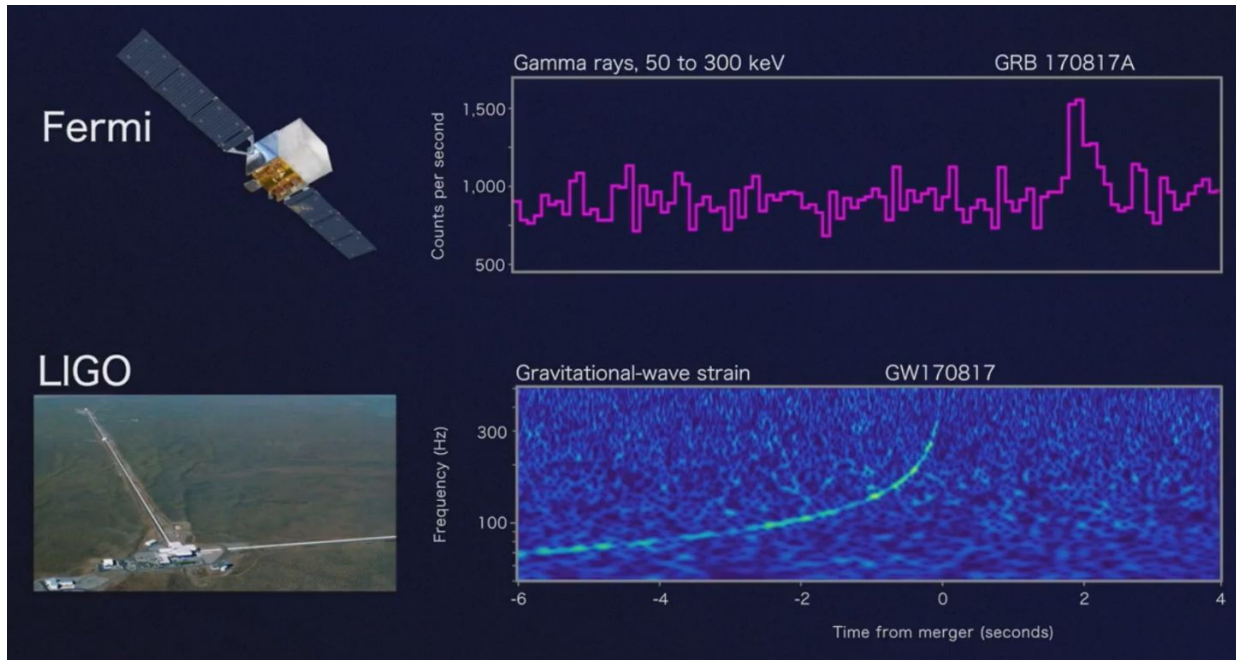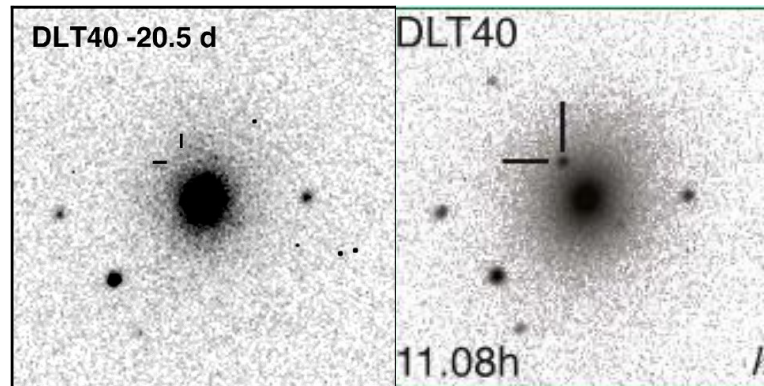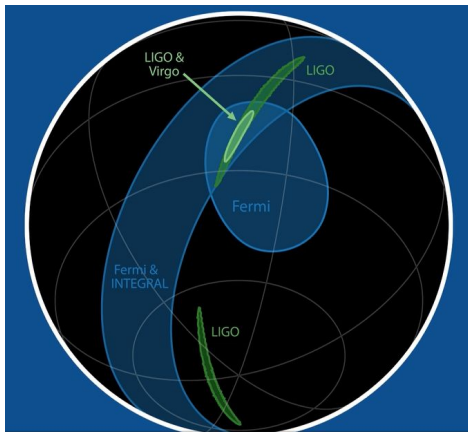# GW170817: The first Binary Neutron Star Merger



Image credit: NASA GSFC & Caltech/MIT/LIGO Lab

sGRB progenitors

Kilonova and the origins of heavy elements

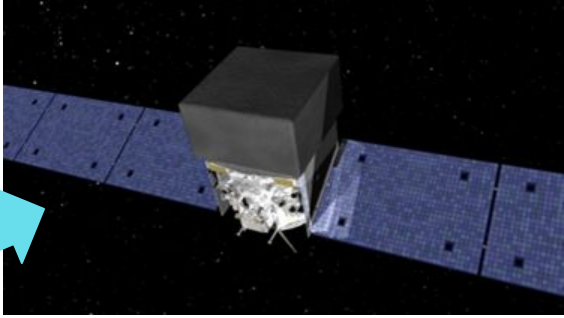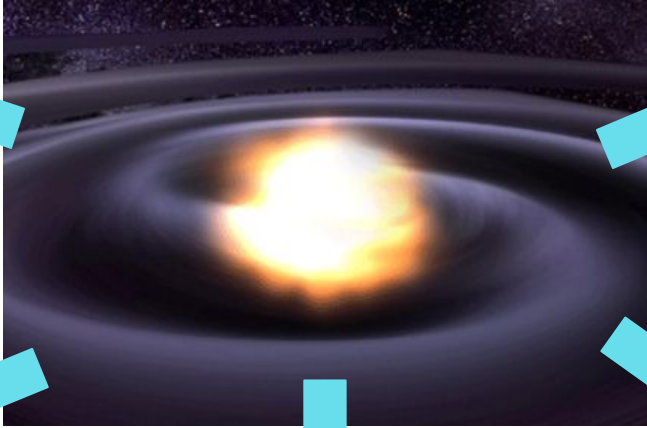'Standard siren' measurement of the Hubble constant

Speed of gravity

EM Partners with LIGO-Virgo, *Astrophys. J. Lett.* 848, L12 (2017)

# Multi-Messenger Astrophysics

Image credit: NASA Goddard Space Flight Center/ Dana Berry
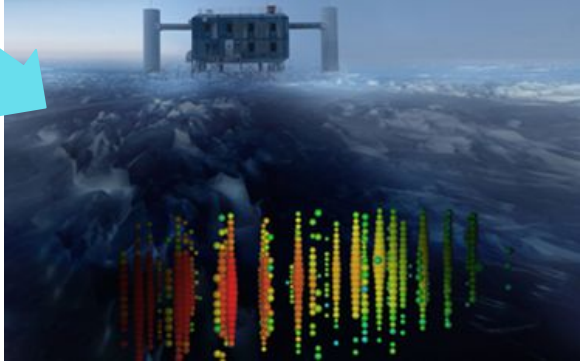
Gravitational waves

X-rays/Gamma-rays

Visible/infrared light

Radio waves

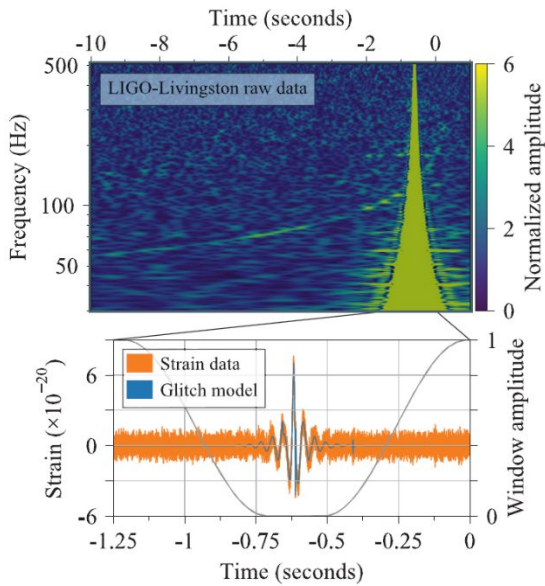Neutrinos

# The Challenge: the 3 deadly F's

## Fast:
need to identify GW transients as quickly as possible in order to have a chance to catch the earliest light
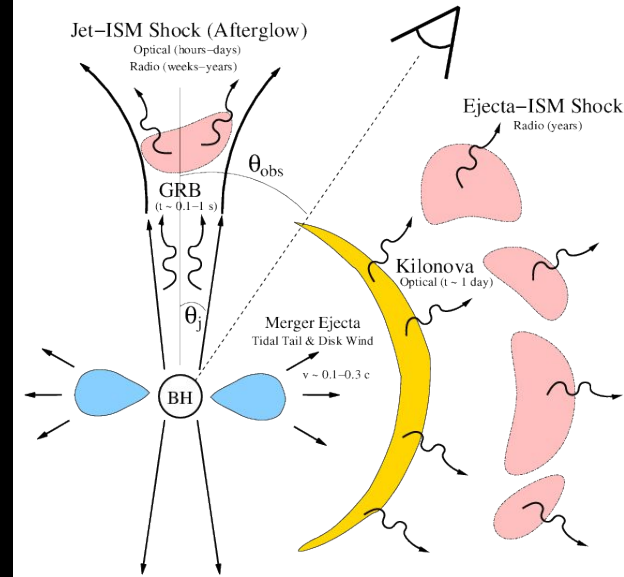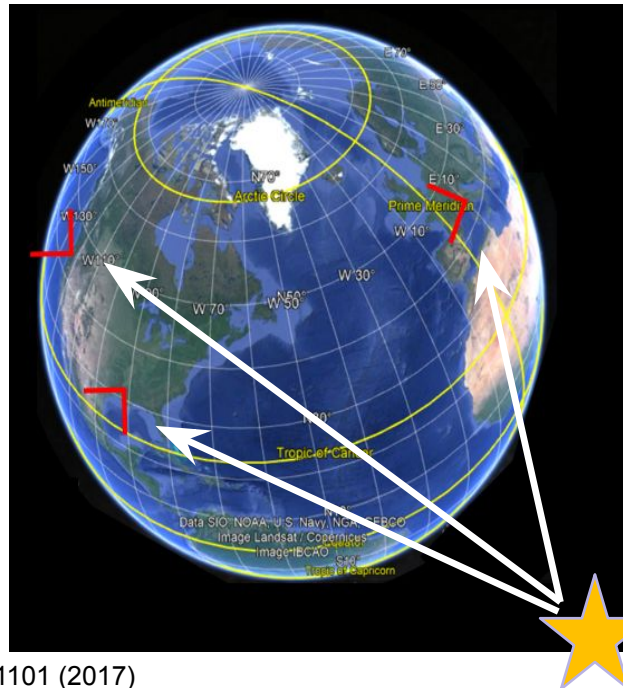
## Fuzzy:
gravitational-wave detectors are more like radio receivers than telescopes

## Faint:
for EM counterparts at the nominal BNS merger range of 200Mpc and BBH ranges out to Gpc



LIGO-Virgo Collaborations, *Phys. Rev. Lett.* 119, 161101 (2017)

Metzger and Berger, *Astrophys. J.* 746, 1 (2012)

8

# GW transient events



Image: https://www.zdnet.com/article/kafka-channels-the-big-data-firehose/

11 events from O1+O2

44 events in O3a, 55 total
1041 "subthreshold" events in O1,O2,O3a

35 events in O3b, 90 total
(catalogs are cumulative)
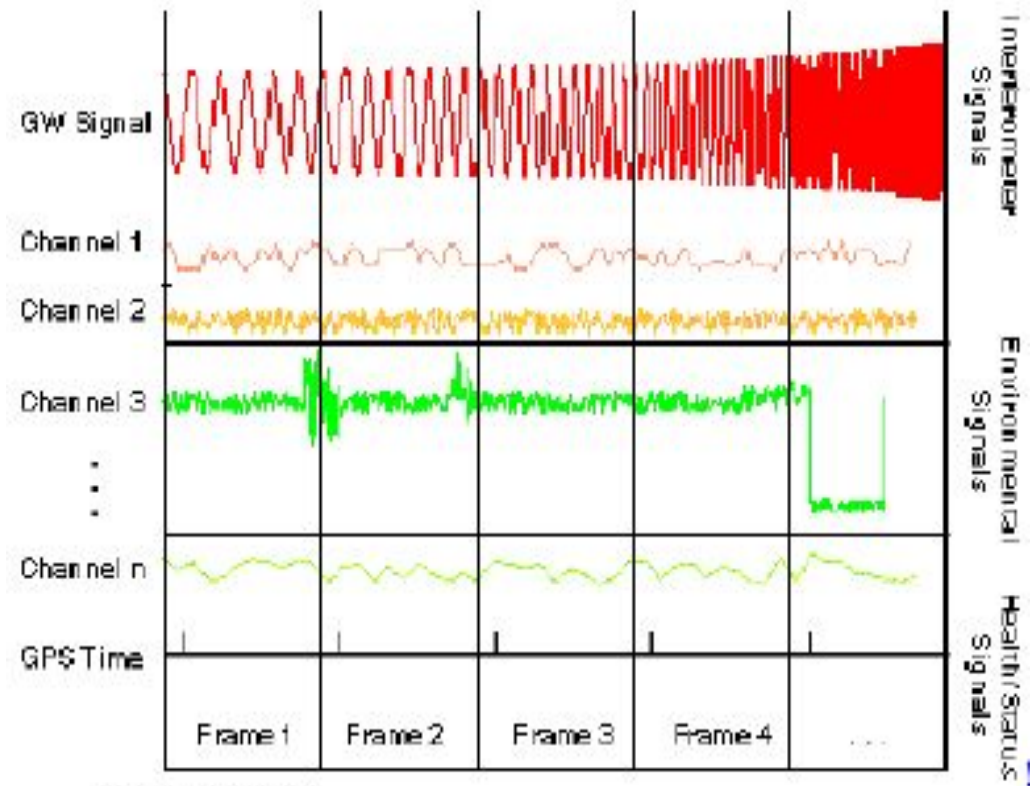
# Gravitational-wave detector data

Continuous **time series** (1Hz, 128Hz … 16kHz)

Gravitational Wave channel:
~20GB/day (per instrument)

Physical Environment
Monitors (seismometers,
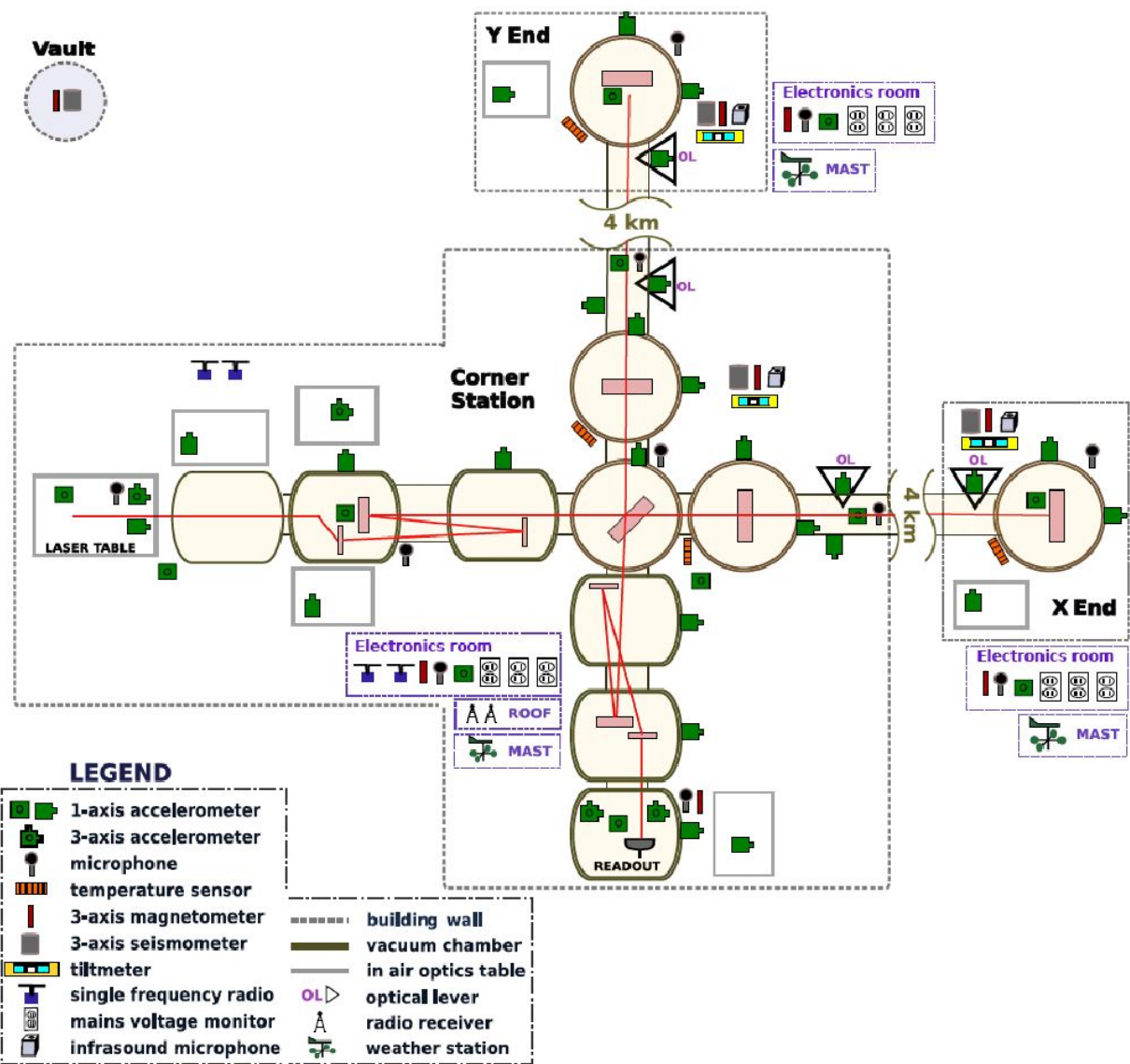accelerometers,
magnetometers, microphones
etc)

Internal Engineering Monitors
(sensing, housekeeping,
status etc)

Together with various
intermediate data products
>2TB/day (per instrument)

# Interefometric and environmental sensors



Each LIGO detector records over 200,000 auxiliary channels that monitor instrument (interferometric) behavior and environmental conditions.

Enables the study of correlations (couplings) of the gravitational wave channel with the environment (including global events, e.g. lightnings).

**LEGEND**

- 1-axis accelerometer
- 3-axis accelerometer
- microphone
- temperature sensor
- 3-axis magnetometer
- 3-axis seismometer
- tiltmeter
- single frequency radio
- mains voltage monitor
- infrasound microphone

- building wall
- vacuum chamber
- in air optics table
- OL ▷ optical lever
- radio receiver
- weather station

# Machine Learning for Gravitational-wave data

Lots of data

Rich, complex signal space

Rich, complex noise space

Low-latency/real-time requirements
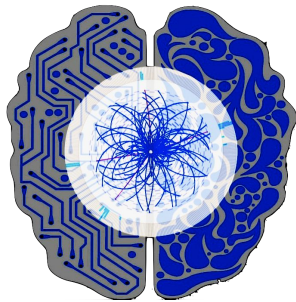
→ Neural Networks!

## Computing revolution:

Success of deep learning has led to sophisticated algorithms

Rise of heterogeneous computing has enabled deep learning

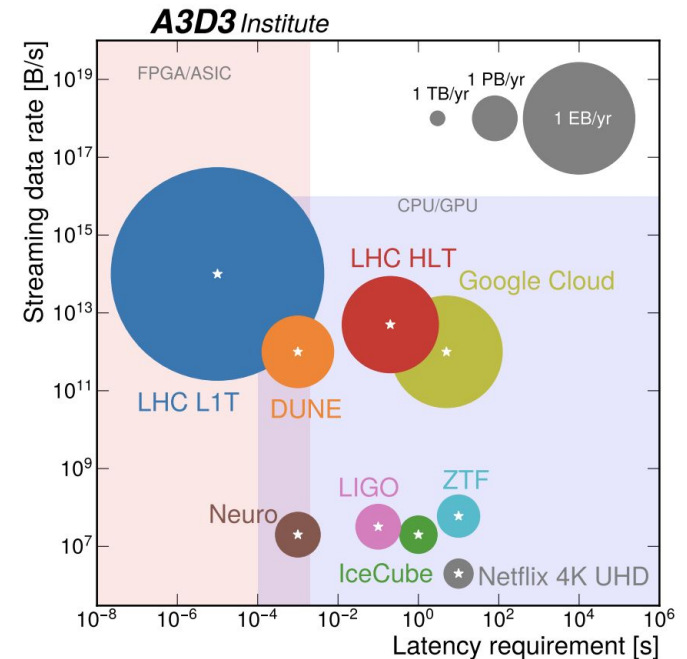Developing ML+GPU integration has enabled large throughput computing

Developing ML+FPGA/ASIC for low latency computing



Accelerated AI Algorithms for Data-Driven Discovery

FastML for Science        A3D3 Institute

# Requirements for ML deployment in GW searches

**Training**

Load time-series data from disk and efficiently move to GPU

Leverage simulations to create robust datasets

Implement signal processing operations on GPU

**Inference**

Offline - produce predictions on O(100+years) of background data

Online - produce transient detections on real-time data in O(1s) and estimate parameters in O(1s)

Stream time-series into NN

Heterogeneous computing backends/data-types

*Infrastructure design goals*

*Intuitive* - maps on to familiar, physically meaningful concepts
*Composable* - hierarchical layers of abstraction support new use cases seamlessly

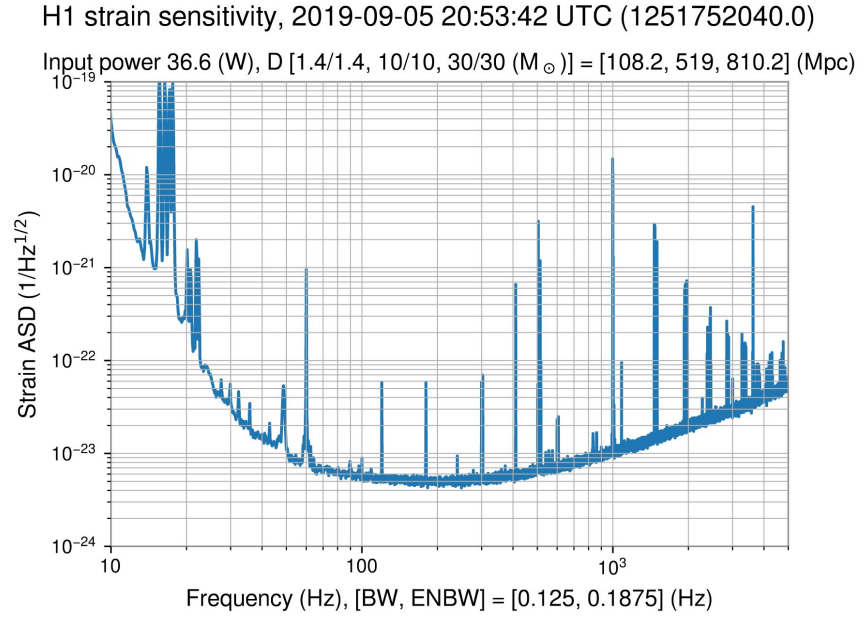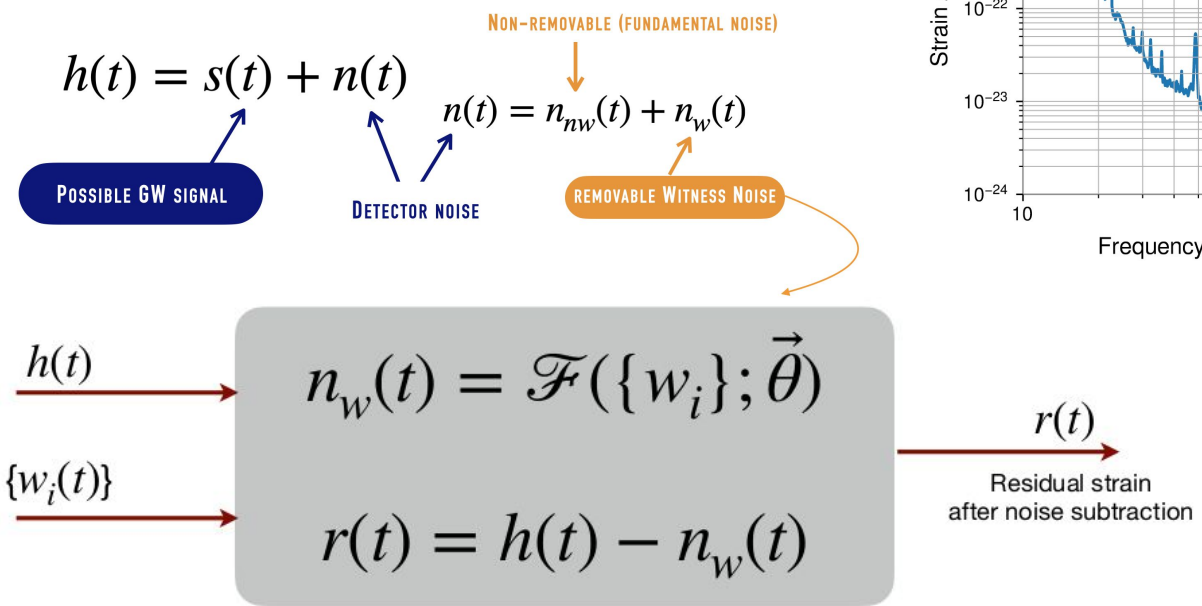*Integrated* - ecosystem of tools following same standards and nomenclature
*Efficient* - make the most out of parallel computing resources

# Gravitational-wave data analysis workflows

**Data quality:** identify and mitigate noise sources ("detector characterization")



**Noise subtraction** nonlinear regression

**Detection:** identity data instances that stand out statistically as deviating from noise



**Modelled transients** supervised
**Unmodelled transient** semi-supervised
**Vetoes** Glitch identification

**Inverse problem:** extract intrinsic and extrinsic signal/source parameters



**Parameter estimation** Normalizing flows

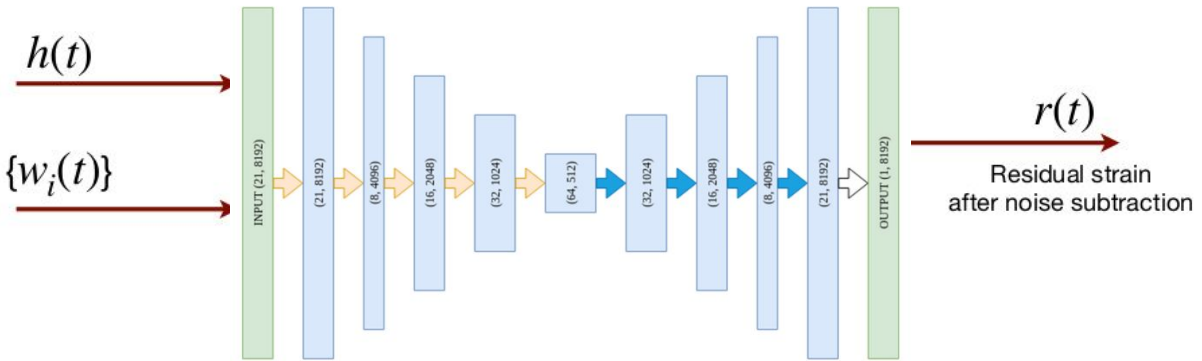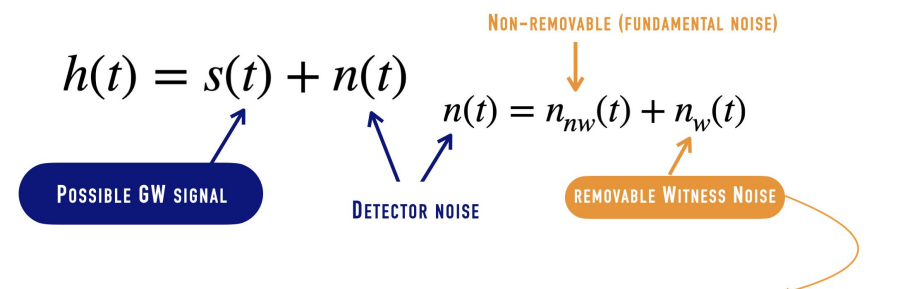ML4GW/HERMES: a new ecosystem for end-to-end ML-based GW searches

# DeepClean: a noise subtraction platform

Non-fundamental noise in interferometers can be subtracted, when such noise is "witnessed" by auxiliary channels:

$$h(t) = s(t) + n(t)$$

**NON-REMOVABLE (FUNDAMENTAL NOISE)**

$$n(t) = n_{nw}(t) + n_{w}(t)$$

**POSSIBLE GW SIGNAL**

**DETECTOR NOISE**

**REMOVABLE WITNESS NOISE**

H1 strain sensitivity, 2019-09-05 20:53:42 UTC (1251752040.0)

Input power 36.6 (W), D [1.4/1.4, 10/10, 30/30 (M$_\odot$)] = [108.2, 519, 810.2] (Mpc)



$h(t)$

$\{w_i(t)\}$

$$n_w(t) = \mathcal{F}(\{w_i\}; \vec{\theta})$$

$$r(t) = h(t) - n_w(t)$$

$r(t)$

Residual strain after noise subtraction

Ormiston et al. "Noise Reduction in Gravitational-Wave Data via Deep Learning."
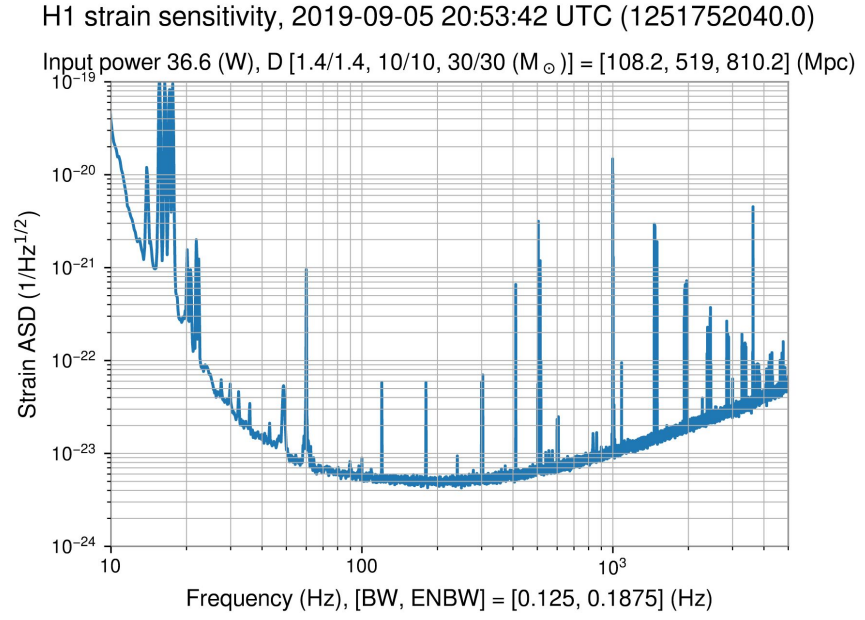
15

# DeepClean: a noise subtraction platform

Non-fundamental noise in interferometers can be subtracted, when such noise is "witnessed" by auxiliary channels:

$$h(t) = s(t) + n(t)$$

POSSIBLE GW SIGNAL

DETECTOR NOISE

$$n(t) = n_{nw}(t) + n_{w}(t)$$

NON-REMOVABLE (FUNDAMENTAL NOISE)

REMOVABLE WITNESS NOISE

H1 strain sensitivity, 2019-09-05 20:53:42 UTC (1251752040.0)

Input power 36.6 (W), D [1.4/1.4, 10/10, 30/30 (M$_\odot$)] = [108.2, 519, 810.2] (Mpc)

Frequency (Hz), [BW, ENBW] = [0.125, 0.1875] (Hz)

$h(t)$

$\{w_i(t)\}$

INPUT (21, 8192) | (21, 8192) | (8, 4096) | (16, 2048) | (32, 1024) | (64, 512) | (32, 1024) | (16, 2048) | (8, 4096) | (21, 8192) | OUTPUT (1, 8192)

$r(t)$

Residual strain after noise subtraction

Convolutional auto-encoder

Real-time implementation with ~1s latency

Provided to analyses downstream

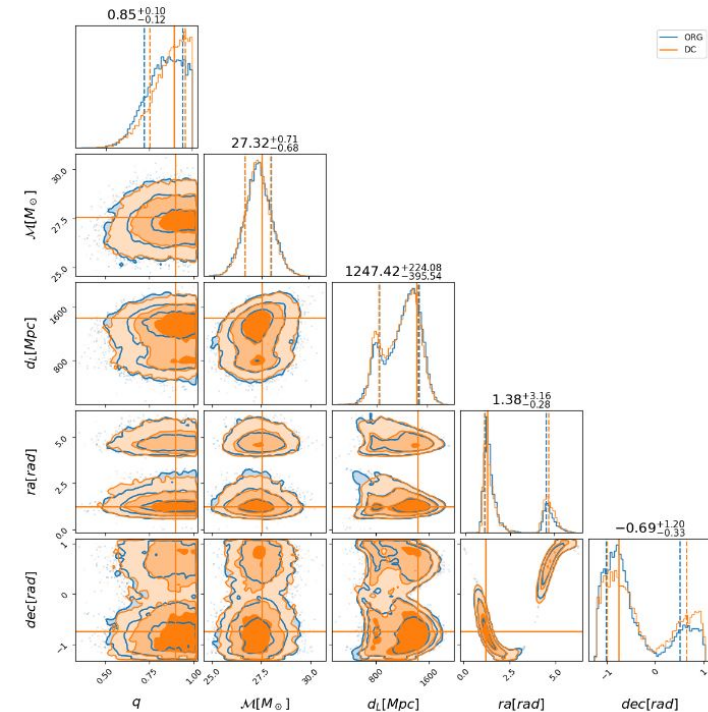Able to go beyond linear couplings/algorithms

Also implemented on an FPGA

Ormiston et al. "Noise Reduction in Gravitational-Wave Data via Deep Learning."

16

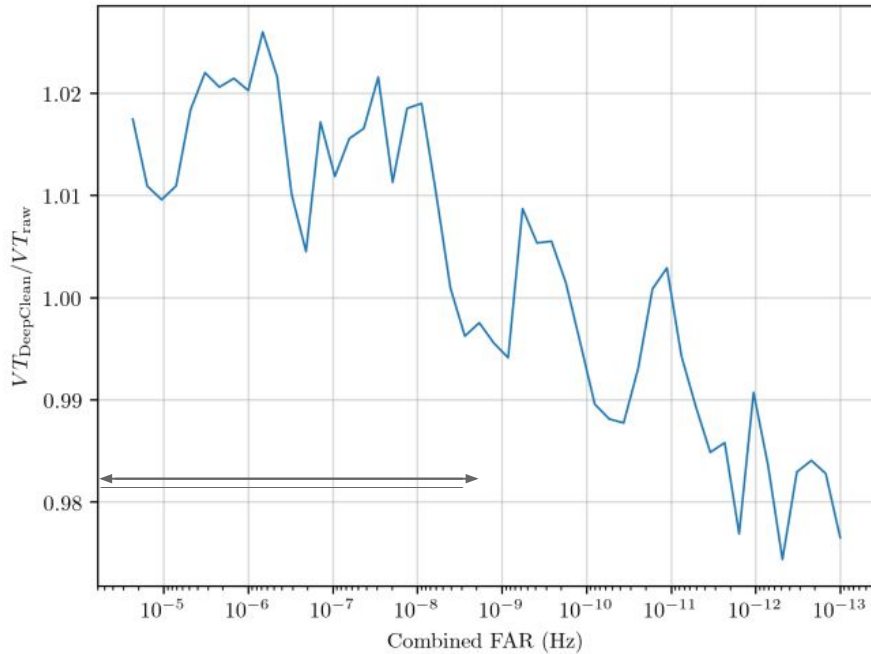# DeepClean performance in O3: Amplitude Spectral Densities and Parameter Estimation



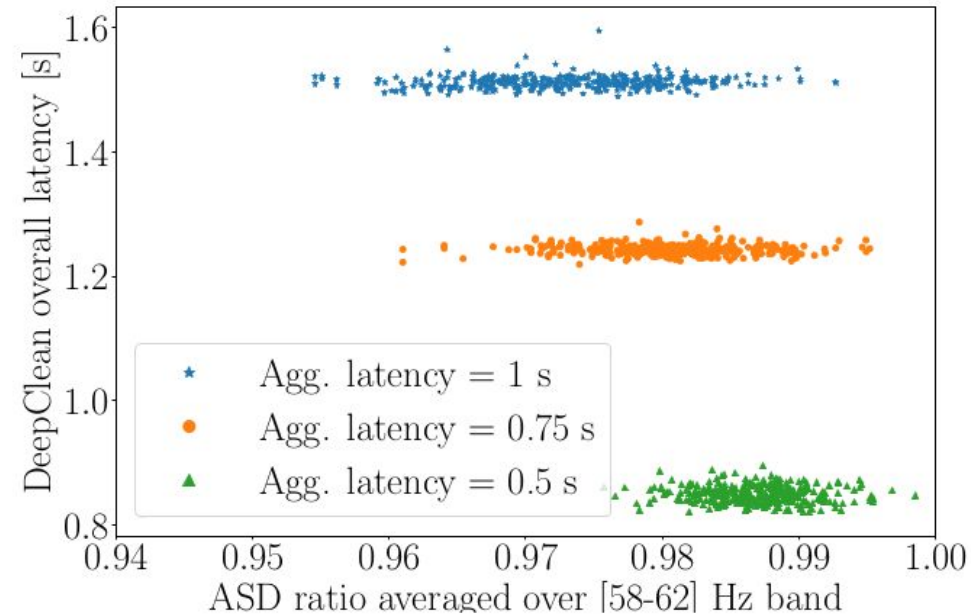Demonstrated non-linear subtraction on 60 Hz power lines and sidebands!

Multiple tests on BBH injections to demonstrate unbiased recovery of astrophysical signals

Ormiston et al. "Noise Reduction in Gravitational-Wave Data via Deep Learning."

Saleem et al (2023), **Demonstration of Machine Learning-assisted real-time noise regression in gravitational wave detectors**.

# DeepClean performance in O3: time-volume reach and latencies



Sensitive volume (V*T) fractional gain or loss (with/without DeepClean) as a function of the false alarm rate in a GstLAL search.

Trade-off between latency and quality: the ASD ratio improves with higher aggregation latency, at the cost of increased overall latency.
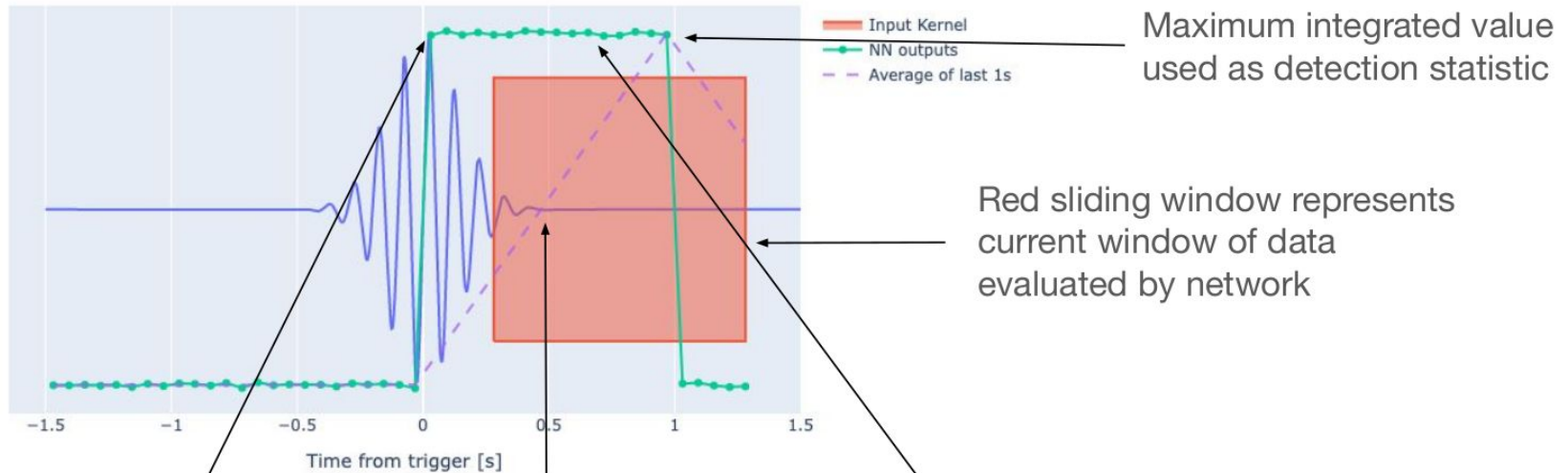
Saleem et al (2023), **Demonstration of Machine Learning-assisted real-time noise regression in gravitational wave detectors**.

# Aframe: detecting compact binary coalescences

Residual Network architecture trained to map 1.5 second windows of h(t) time-series data to scalar value that indicates likelihood of signal being present in the window.

Training on 10 days of coincident H1 and L1 strain from beginning of LIGO-Virgo-KAGRA's O3a run; 100,000 IMRPhenomPv2 waveforms from astrophysical prior used for training search space: 5 - 100 M_solar.

Extensive data augmentation to show the model as diverse a training set as possible.

Event identification:



Maximum integrated value used as detection statistic

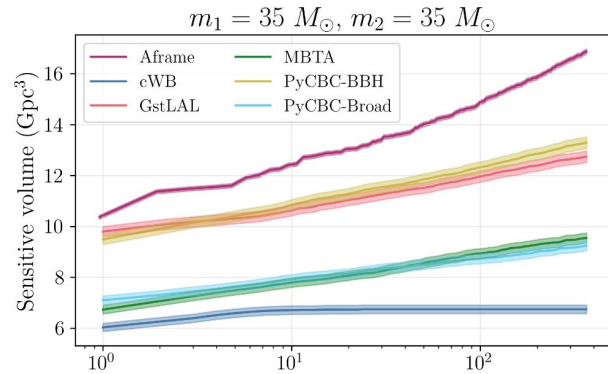Red sliding window represents current window of data evaluated by network

As coalescence time enters window, network begins to ring

Dotted purple line shows average of last 1 second of network outputs

Green dots show time-series of neural network outputs

19

# Aframe: pipeline sensitivity



Sensitive volume calculation is the same as the one used to measure performance of LVK pipelines in GWTC-3 catalog

Competitive performance on higher-mass catalog distributions

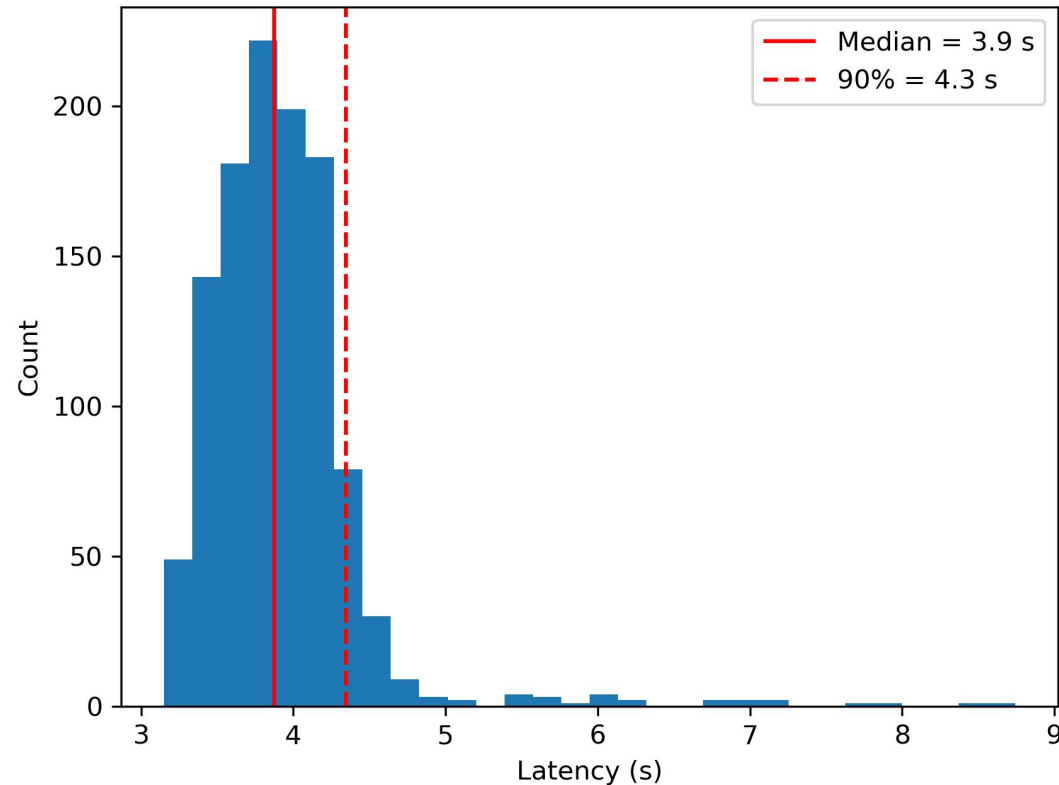Work remains to be done for lower masses – alternative architectures or smarter training techniques

https://zenodo.org/records/7890437

# Computational cost and throughput

*Training*: **~44 hours** on 1 GPU on a 16GB V100 GPU

*Evaluation*: Analysis of O3 took ~30 hours for 21 years of livetime on 8 GPUs, roughly **750** seconds of data processed per second per GPU

Total time scales roughly linearly with number of GPUs

Ran pipeline online over ~1 month of "replayed" O3 data emulating real-time environment: median (90%) latency is 8.4 (37.1) seconds faster than rest of CBC pipelines (as reported here: https://arxiv.org/abs/2308.04545)

# AMPLFI: from detections to astrophysics

Source parameter estimation:

$$d_i(t) = n_i(t) + F_{+,i} h_+(t; \theta) + F_{\times,i} h_\times(t; \theta)$$

Detector output    noise    **Antenna factor**   **Our goal (its parameters)**   **Antenna factor**   **Our goal (its parameters)**

Calculate posteriors via likelihood estimation:

$$p(\boldsymbol{\theta}|\mathbf{d}) = \frac{L(\mathbf{d}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{Z} = \frac{L(\mathbf{d}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\int L(\mathbf{d}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

Likelihood-Free (simulation-based) Inference: train neural network to estimate posteriors $p(\boldsymbol{\theta}|\mathbf{d})$
Model true posterior distribution with a normalizing flow:

$$p(\theta|d) \sim q_\phi(\theta|d) = p_u(T_d^{-1}(\theta)) |det J_{T_d^{-1}}(\theta)|$$

G Papamakarios et al, JMLR Vol. 22, Art. 57 2617-26 (2021)

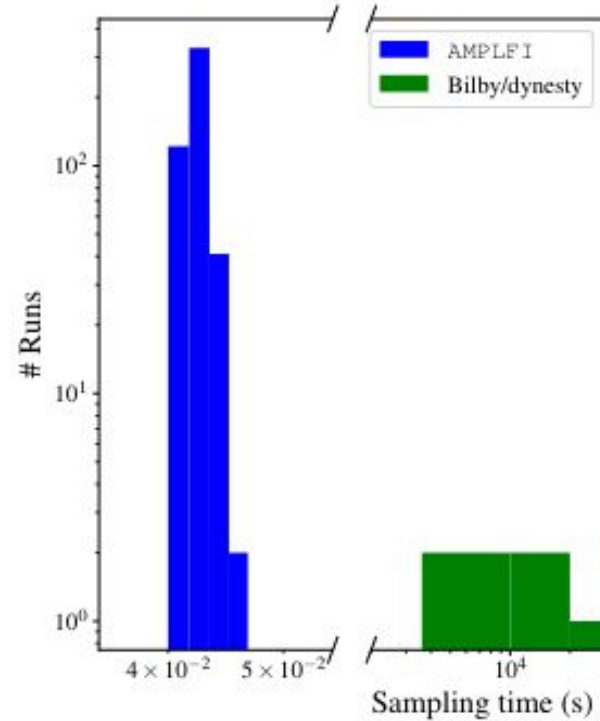The flow is parameterized via a neural network, and trained by minimizing the Kullback-Leibler divergence

$$L \approx -\frac{1}{N} \sum_{i=1}^{N} \log q_\phi(\theta^{(i)}|d)$$

Train neural network to estimate posteriors for any assumed signal morphology (e.g. sine-Gaussian, binary coalescences etc) embedded in <u>real instrument noise</u> and that's how you arrive at **AMPLFI: Accelerated Multi-messenger Parameter-estimation using LFI !**

Deep Chatterjee, Ethan Marx et al, submitted for publication (2024)

# AMPLFI performance



N=500, p-value=0.7674

Legend:
- $\mathcal{M}$ (0.840)
- $q$ (0.790)
- $d_L$ (0.263)
- $\phi_c$ (0.656)
- $\theta_{JN}$ (0.064)
- DEC (0.755)
- $\psi$ (0.609)
- RA (0.873)

P-P plot from inferences with AMPLFI over 500 Binary Black Hole systems

Sampling times for AMPLFI (1 GPU) vs. nested sampling runs (24 CPUs)

Deep Chatterjee, Ethan Marx et al, submitted for publication (2024)

# Putting all these together: "Who will bell the cat?"



Cartoon adopted from Alec Gunny

`ml4gw/HERMES`: https://github.com/ML4GW

# Bringing AI into GW and MMA data analyses



Scientist uses simulations to generate data, priors to regularize training

Models are distributed and versioned in centralized repositories

Dedicated inference applications host models, interacted with via simple client APIs

Dedicated tools make iteration/exploration frictionless

[ML4GW/HERMES](#): an ecosystem for ML applications in GW enabling fast deployment, fast inference, small computation footprint and optimized for computing heterogeneity

README.md

## ML4GW

Tools to make training and deploying neural networks in service of gravitational wave physics simple and accessible to all!

Includes a couple particular applications under active research.

⊘ View as: Public ▾

You are viewing the README and pinned repositories as a public user.

People

**Hardware-accelerated Inference for Real-Time Gravitational-Wave Astronomy**
Alec Gunny et al Nature Astronomy (2022)

# Summary and outlook

`ml4gw`: A new computing ecosystem for ML applications in GW data analyses

Emphasis on real-time processing for improving multi-messenger prospects of the GW observatories

> focus on latency

> minimal computational footprint (a

> couple of GPUs to keep up with

real-time)

Offline implementation

> portable, robust pipelines

> emphasis on throughput

> extensible

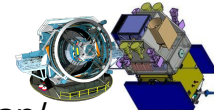End-to-end ML-based workflows have been implemented addressing:

> data cleaning
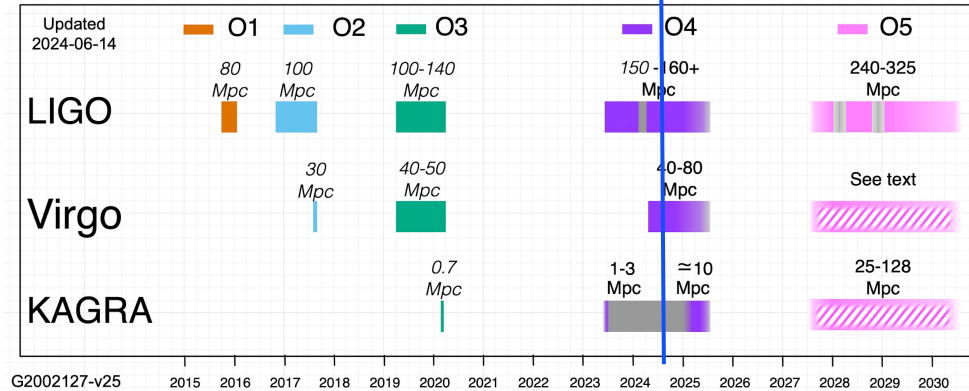
> transient event detections

> parameter estimation

Expecting to deploy in production for real-time use during LIGO-Virgo-KAGRA's current O4 observing run
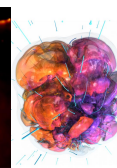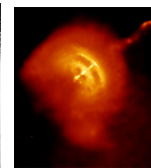
https://observing.docs.ligo.org/plan/ now

| Updated 2024-06-14 | O1 | O2 | O3 | | O4 | O5 |
|---|---|---|---|---|---|---|

LIGO — 80 Mpc, 100 Mpc, 100-140 Mpc, 150-160+ Mpc, 240-325 Mpc

Virgo — 30 Mpc, 40-50 Mpc, 40-80 Mpc, See text

KAGRA — 0.7 Mpc, 1-3 Mpc, ≃10 Mpc, 25-128 Mpc

G2002127-v25   2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030

**BBH populations**  **Known Unknowns**  **Unknown Unknowns**

**???**

26

# Summary and outlook

`ml4gw`: A new computing ecosystem for ML applications in GW data analyses

Emphasis on real-time processing for improving multi-messenger prospects of the GW observatories

  focus on latency

  minimal computational footprint (a

  couple of GPUs to keep up with

real-time)

Offline implementation

  portable, robust pipelines

  emphasis on throughput

  extensible

End-to-end ML-based workflows have been implemented addressing:

  data cleaning

  transient event detections

  parameter estimation

Expecting to deploy in production for real-time use during LIGO-Virgo-KAGRA's current O4 observing run

https://observing.docs.ligo.org/plan/ now



| | Updated 2024-06-14 | O1 | O2 | O3 | O4 | O5 |
|---|---|---|---|---|---|---|
| LIGO | | 80 Mpc | 100 Mpc | 100-140 Mpc | 150-160+ Mpc | 240-325 Mpc |
| Virgo | | | 30 Mpc | 40-50 Mpc | 40-80 Mpc | See text |
| KAGRA | | | | 0.7 Mpc | 1-3 Mpc  ≈10 Mpc | 25-128 Mpc |

G2002127-v25   2015 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025 2026 2027 2028 2029 2030

**BBH populations** **Known Unknowns** **Unknown Unknowns**

**???**

## Come join!
https://github.com/ML4GW

# EXTRA SLIDES

# Multi-Messenger Astrophysics



Image credit: Bill Saxton, NRAO

Gravitational waves

X-rays/Gamma-rays

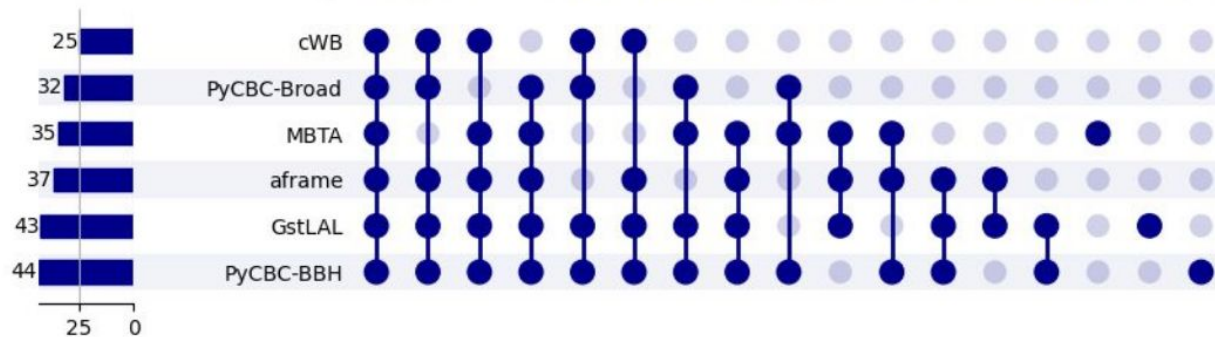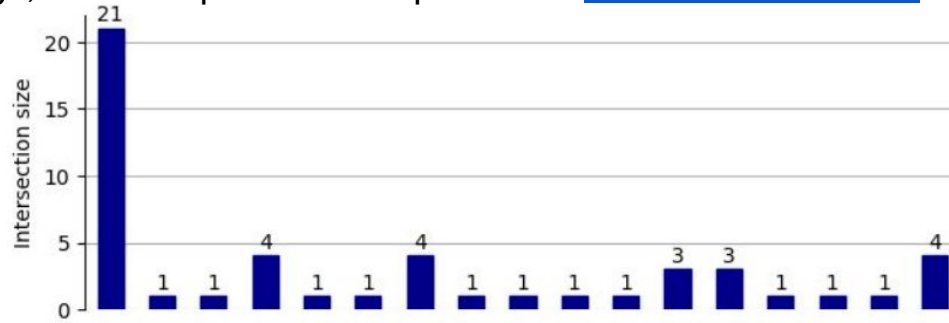Visible/infrared light

Radio waves

Neutrinos

# Comparison to LVK's O3 detections

Reanalyzing LVK's O3 observing run, 37/50 candidates detected by Aframe at false alarm threshold of 1 per 5 months; missed events have network matched filter SNR<13.1 or chirp mass < 10 M_solar. The latter is consistent with Aframe's sensitive volume measurements.

Also the next 10 most significant Aframe events have no overlaps with the LVK catalogs, but have partial overlap with the Olsen et al event list



| gpstime | FAR (1 / yr) |
|---|---|
| 1262635012.75 | 3.6 |
| 1246523564.75 | 4.0 |
| 1264333383.00 | 4.2 |
| 1238351045.00 | 4.4 |
| 1251010355.50 | 4.7 |
| 1264246793.25 | 5.9 |
| 1262163593.25 | 7.8 |
| 1249032684.75 | 11.0 |
| 1253452013.50 | 11.7 |
| 1259411705.25 | 12.0 |

E. Marx, W. Benoit et al, in preparation

# GWAK: an anomaly detection framework

Strong astrophysical motivation to look beyond modeled binary coalescences: supernova, neutron star glitches, magnetars, GRBs, FRBs, cosmic strings and cusps, unknown unknowns may emit GWs that we can not fully modeled currently and thus can not be searched with a matched-filter approach

We refer to them as anomalous and aim to develop a semi-supervised approach which would let us to discover such anomalous signals without explicit modeling
- use multiple autoencoders to create embedded space
- use real background and inject signals
- verify on anomalous signals that aren't included in training

GWAK is the Gravitational Wave Anomalous Knowledge, an algorithm using recurrent autoencoders inspired by similar approaches (QUasi Anomalous Knowledge, by Sang Eon Park et al. https://arxiv.org/abs/2011.03550) taken in performing anomaly detection in LHC data

Core idea: go beyond vanilla anomaly detection in 1-dimensional approach where the distance between the input and output is used as a metric for anomaly detection:
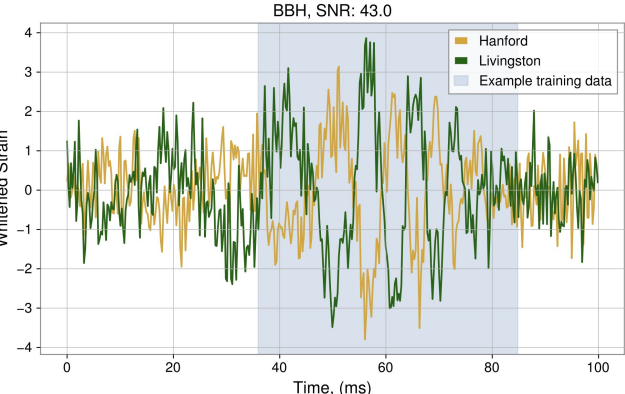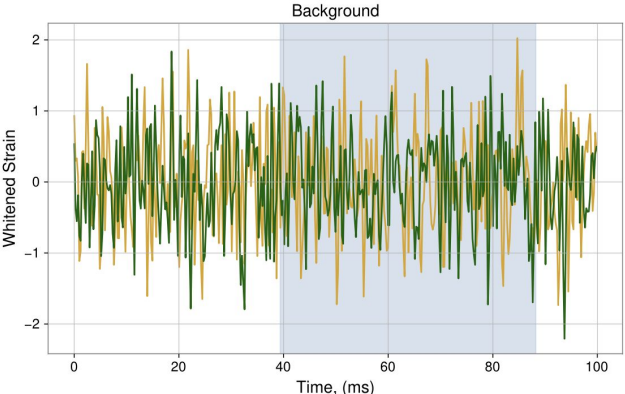


1-dim detection statistic

Introduce <u>multiples axes, for both signal and background</u> ⇒ allows to more efficiently select signal-like anomalies
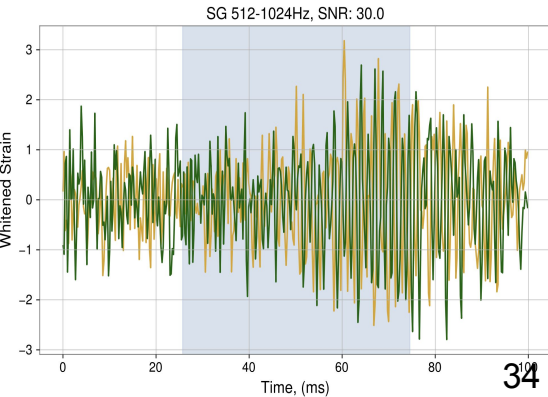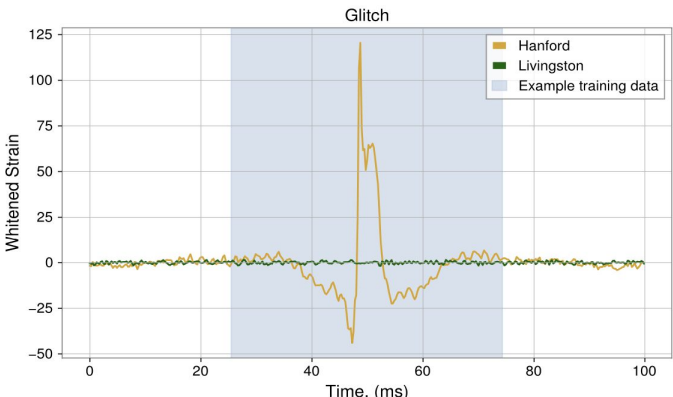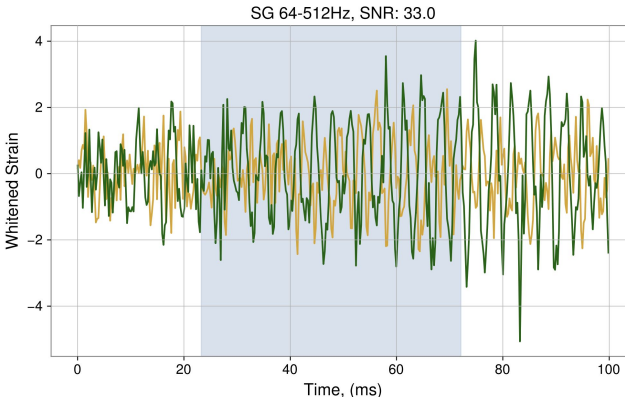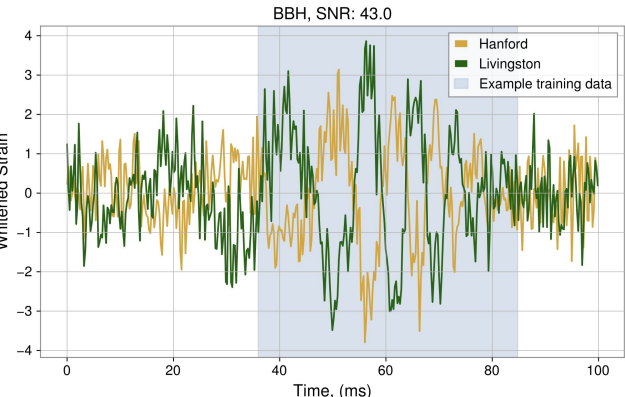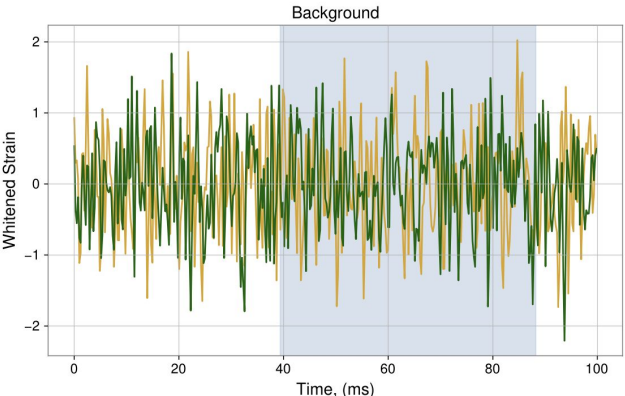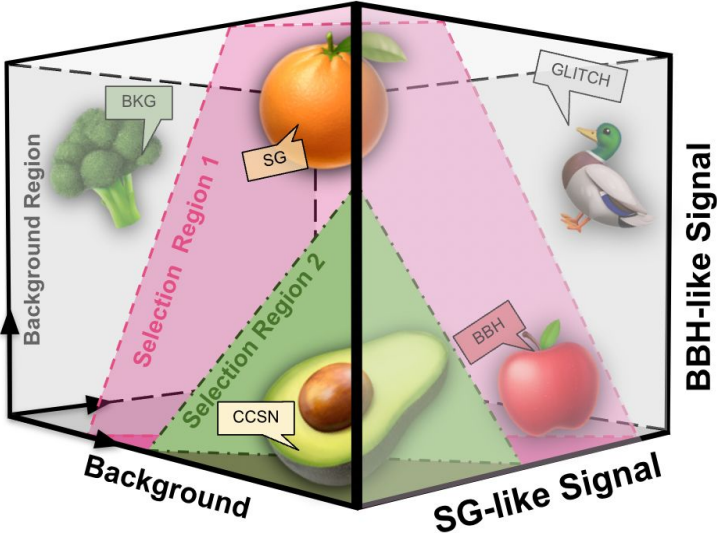
32

# GWAK: multi-dimensional approach



A 3-dim GWAK space example



Raikman et al (2023) , **GWAK: Gravitational-Wave Anomalous Knowledge with Recurrent Autoencoders**, https://arxiv.org/abs/2309.11537

# GWAK: multi-dimensional approach
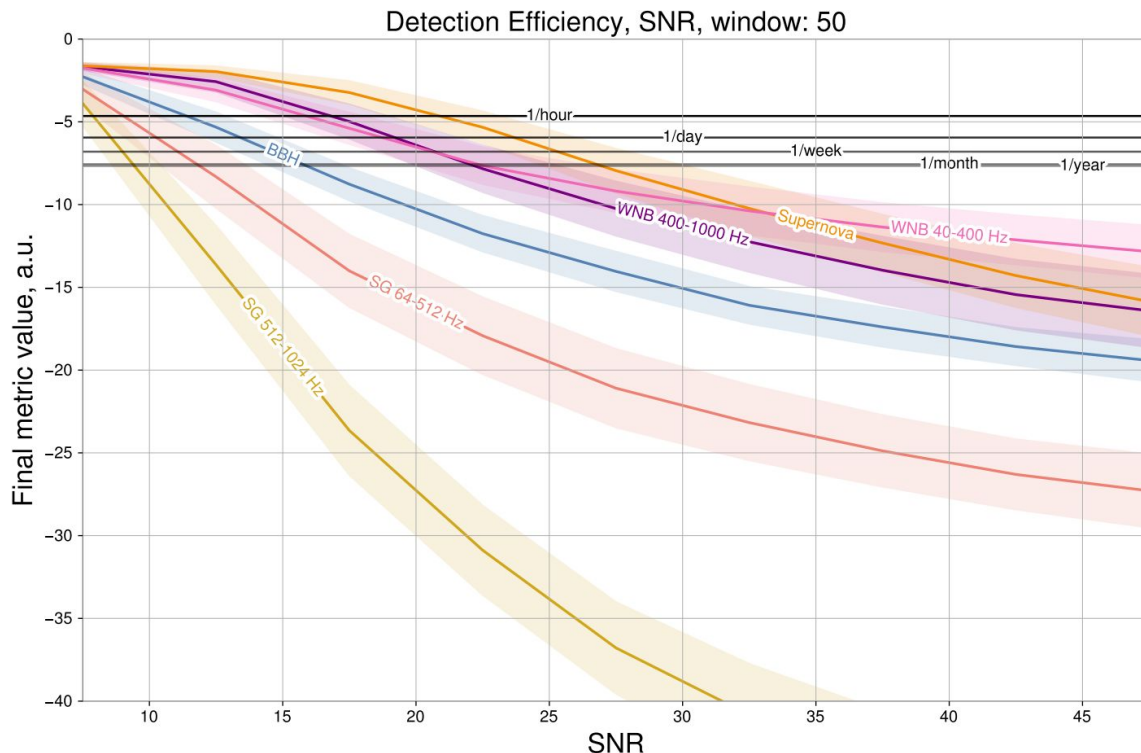
A 3-dim GWAK space example

# **GWAK** efficiency during O3

The final metric as a function of SNR for GWAK axes training signals, BBH (blue), SG 64-512 Hz (yellow), SG 512-1024 Hz (salmon)

and for potential (unseen) anomalies, WNB 40-400 Hz (pink), WNB 400-1000 Hz (purple), and Supernova (orange)

The black lines of varied width correspond to different FARs, from the FAR of 1 per hour to 1 per year

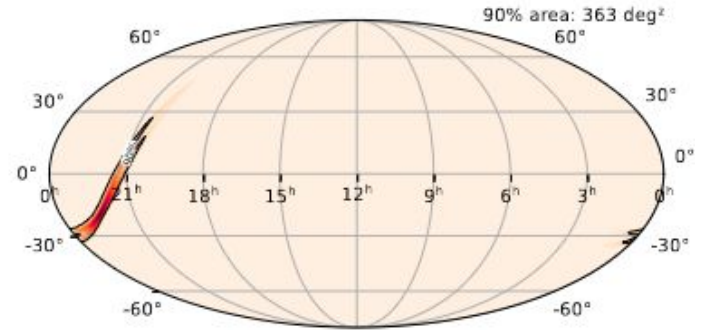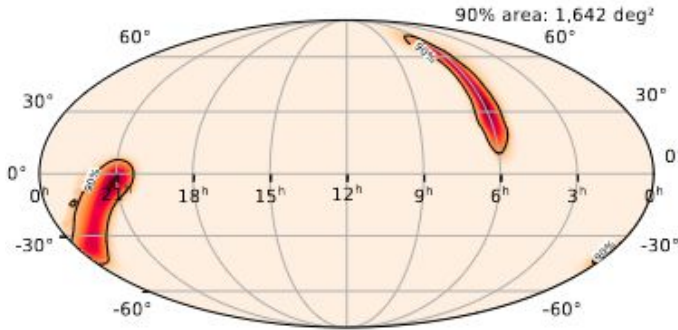For each of the lines, the events that are below that line would be detected.



Detection Efficiency, SNR, window: 50

Raikman et al (2023) , **GWAK: Gravitational-Wave Anomalous Knowledge with Recurrent Autoencoders**, https://arxiv.org/abs/2309.11537

# Low latency sky localization comparisons on O3 alerts

**AMPLFI (2-LIGO detector network)**

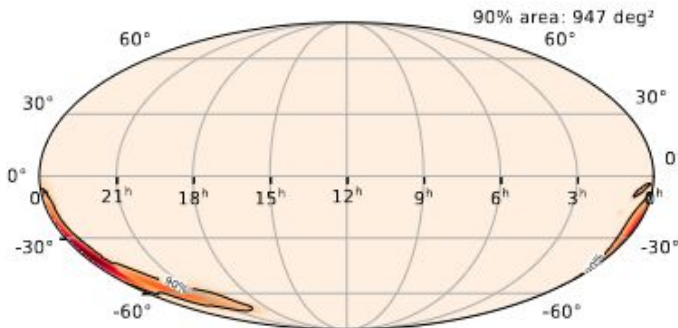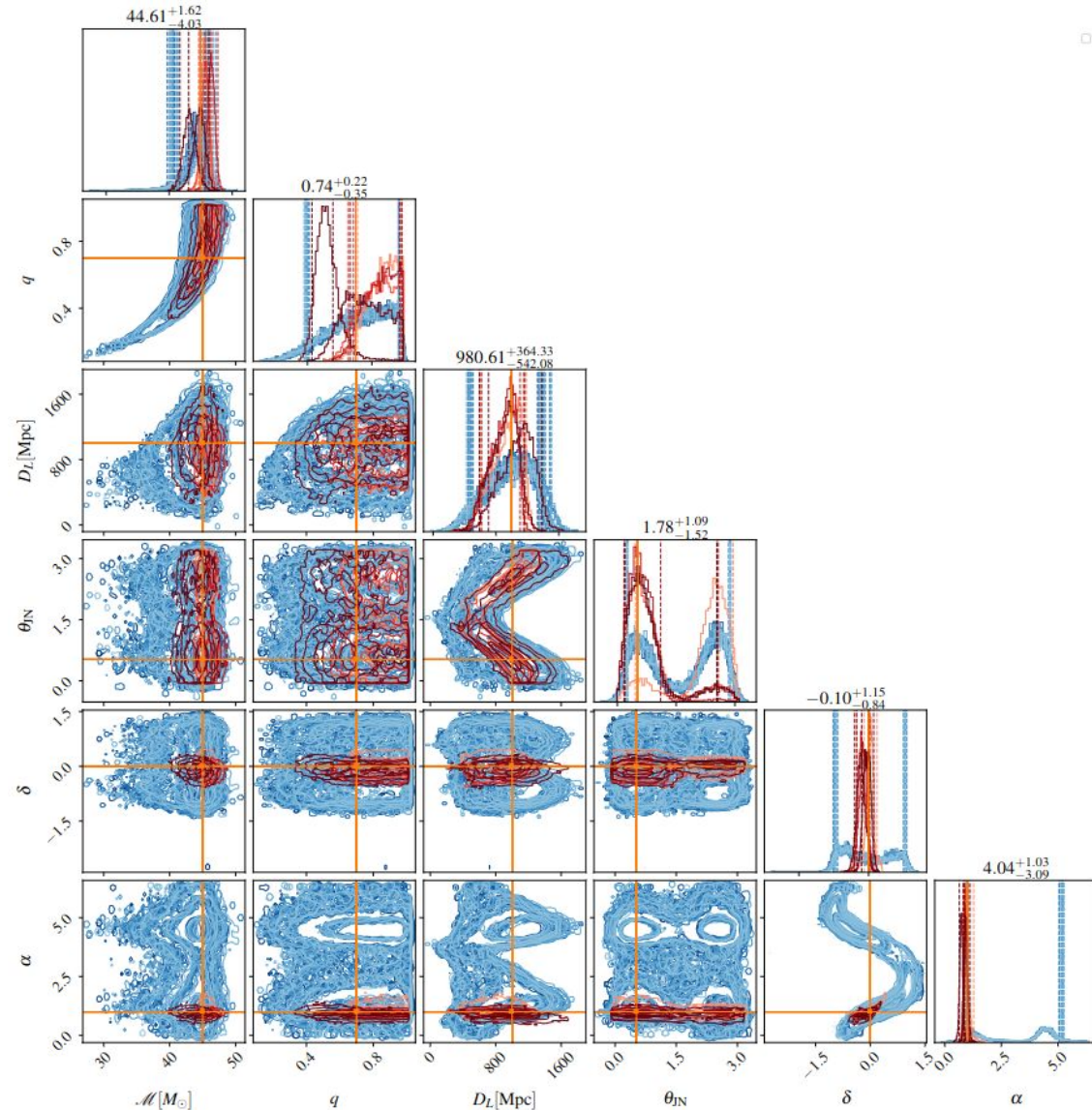**Bayestar (2-LIGO detector network)**

# Sample parameter estimation result

Injected signal: Mchirp = 45 solar masses, q=0.7, D_L = 1 Gpc and optimal SNR~20 added on 20 different background segments

Posterior samples: AMPLFI in blue, Bilby/Dynesty [Ashton et al. ApJS 241, 27 (2019)] in red

Parameter recovery is consistent with injections and stochastic samplers, although posterior widths are tighter with the later



Deep Chatterjee, Ethan Marx et al, submitted for publication (2024)

# The A3D3 Institute (www.a3d3.ai)

## Accelerated AI Algorithms for Data Driven Discovery

## Explore real-time AI in MMA, HEP and NeuroScience

New Types of Computing



LHC Physics

NeuroScience

Measured neural activity

Predicted future neural activity

Dynamical Model

Encoder → Latent Embedding → Generator