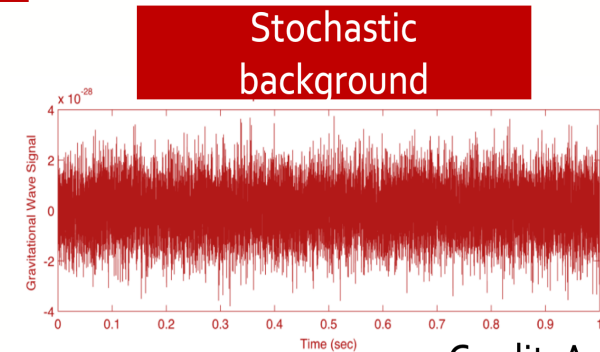
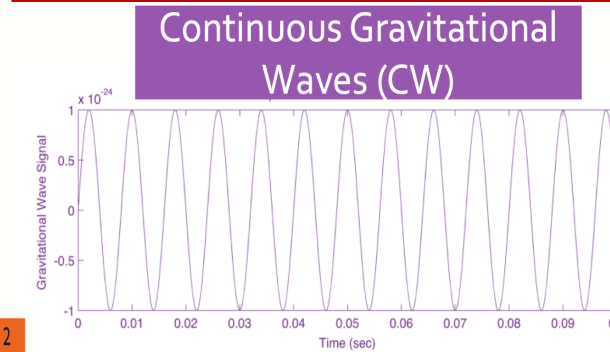
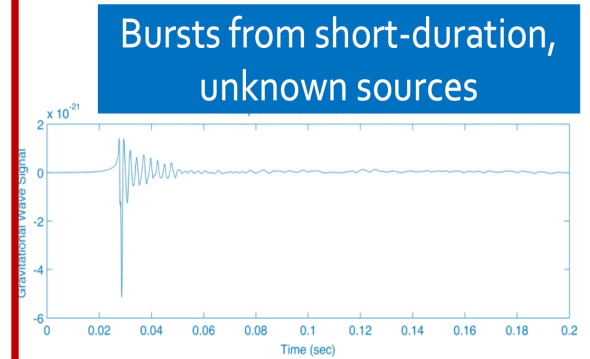
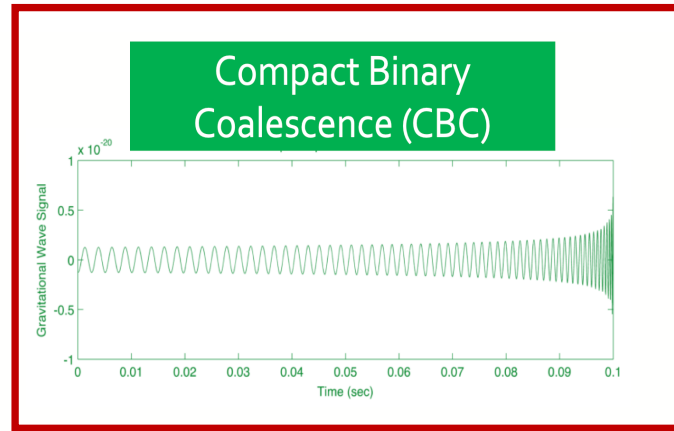


A brief Introduction to Data Analysis procedures in LVK

Ester Ruiz Morales

UPM & IFT-Madrid Virgo Group

- Data analysis is a huge topic!
- I present a very limited scope, reduced to:
 - CBC data
 - Mainly in off-line data processing for the GWTC.
- For Low Latency information, see for instance Cardiff 22 LVK meeting slides:



2

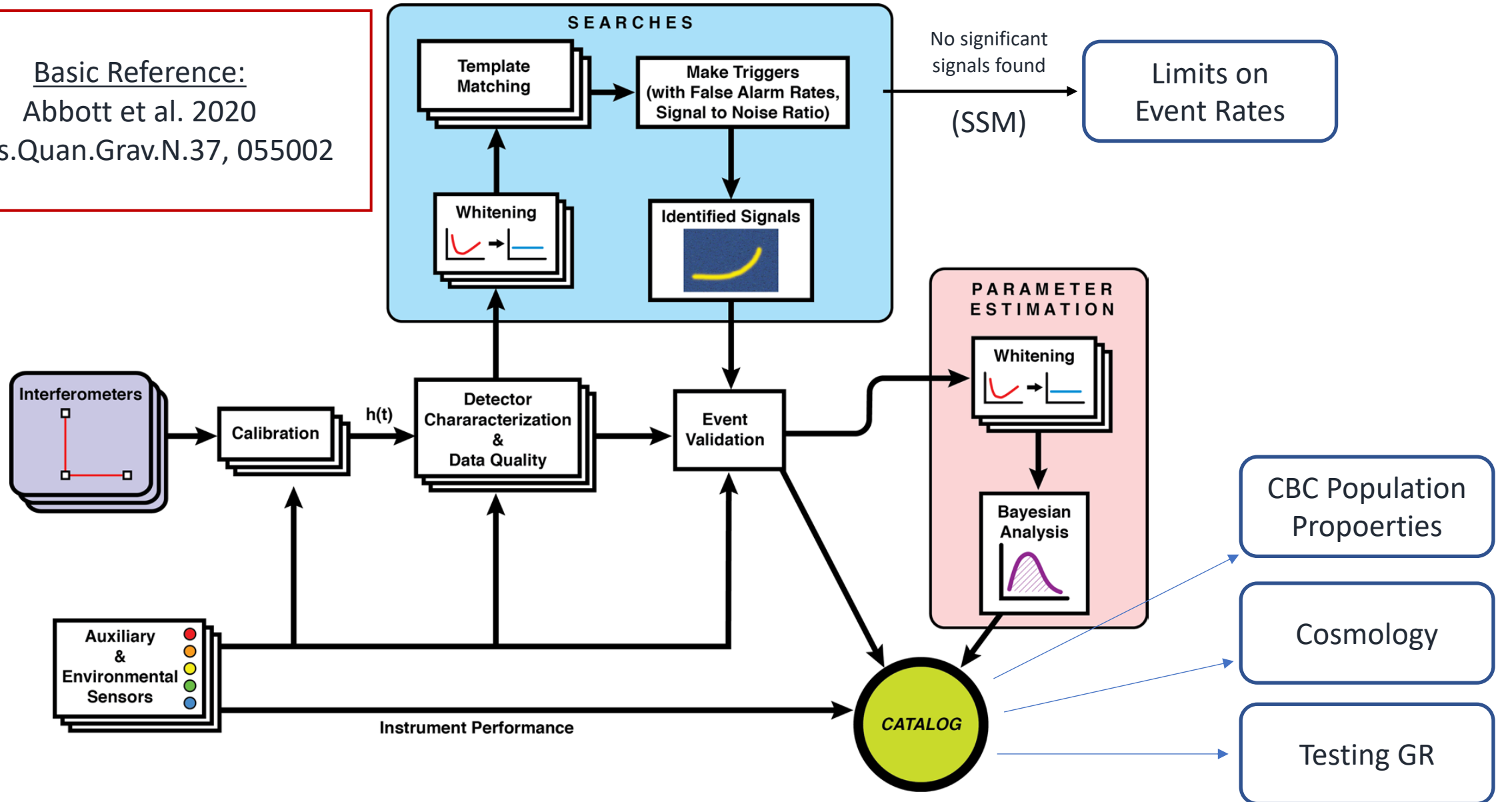
Credit: Amber Stuver

https://dcc.ligo.org/DocDB/0184/G2201664/003/LowLatency_Plenary.pdf

- Describe the basic methods/procedures, main properties & limitations.
- Many of them are automatized, but I won't describe that either.

Main steps in data processing

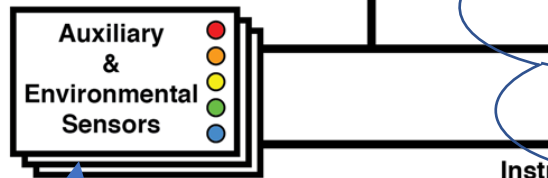
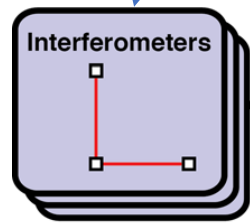
Basic Reference:
Abbott et al. 2020
Class.Quan.Grav.N.37, 055002



Data processing – Step 1: What is really measured?

raw data:
time-varying intensity of the
laser light measured at the
interferometer output

data used for analysis:
gravitational-wave
strain amplitude $d(t)$

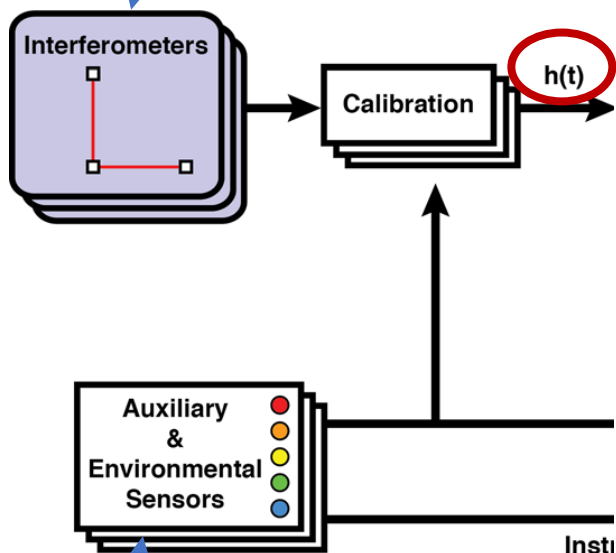


complex process ($\Delta L \sim 1 \times 10^{-19}$ m)
converts the changes of arm length difference into
the time series of the strain of a passing gravitational wave $h(t)$,
sampled at ~ 16 kHz,
Calibration valid in (10Hz – 5kHz) frequency range

Record the time series of the detector
& environment state

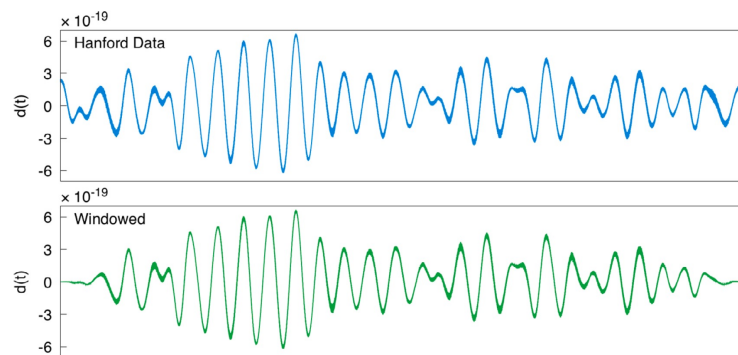
Data processing – Step 1: What is really measured?

raw data:
time-varying intensity of the
laser light measured at the
interferometer output



Record the time series of the detector
& environment state

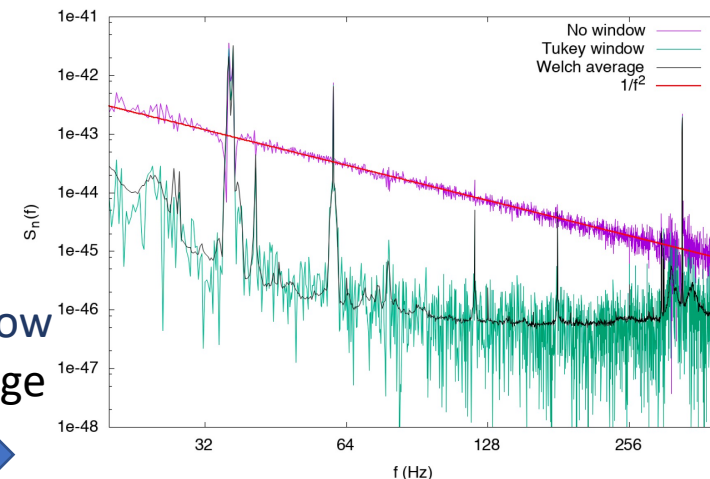
data used for analysis:
gravitational-wave
strain amplitude $d(t)$



FFT
Window
Average

PSD of the noise is not known,
must be estimated from data

Most of the strain amplitude is simply
NOISE
ONLY Aprox. stationary
with a Gaussian component + glitches



PSD: Power Spectral Density

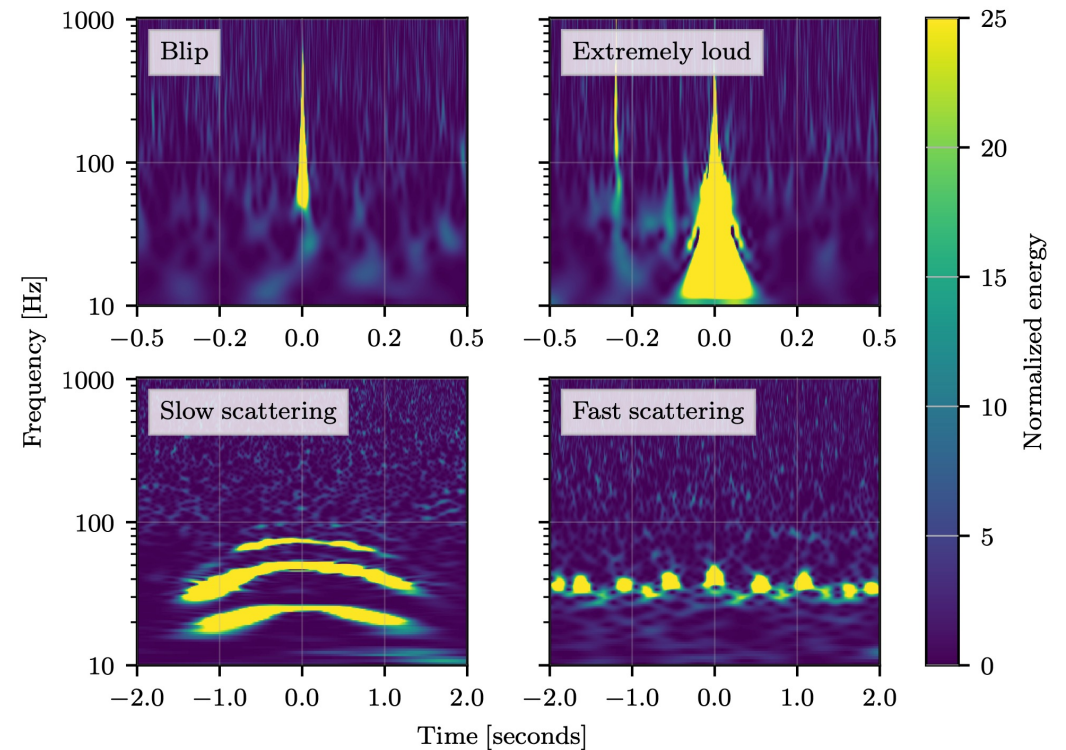
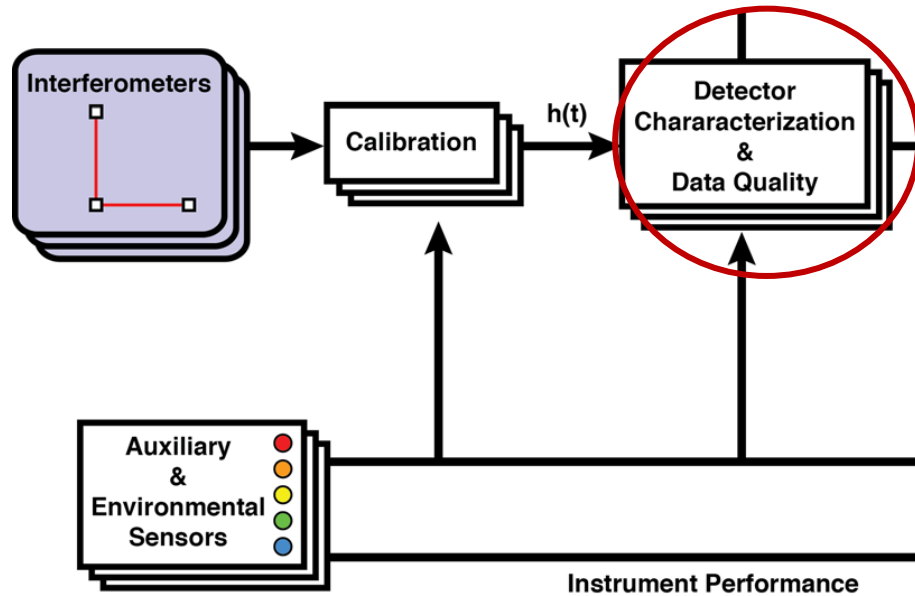
ESSENTIAL TO CHARACTERIZE NOISE WELL
PSD needed for the matched filtered search
and PE

Data processing – Step 2: Check data quality

DetChar & DQ functions:

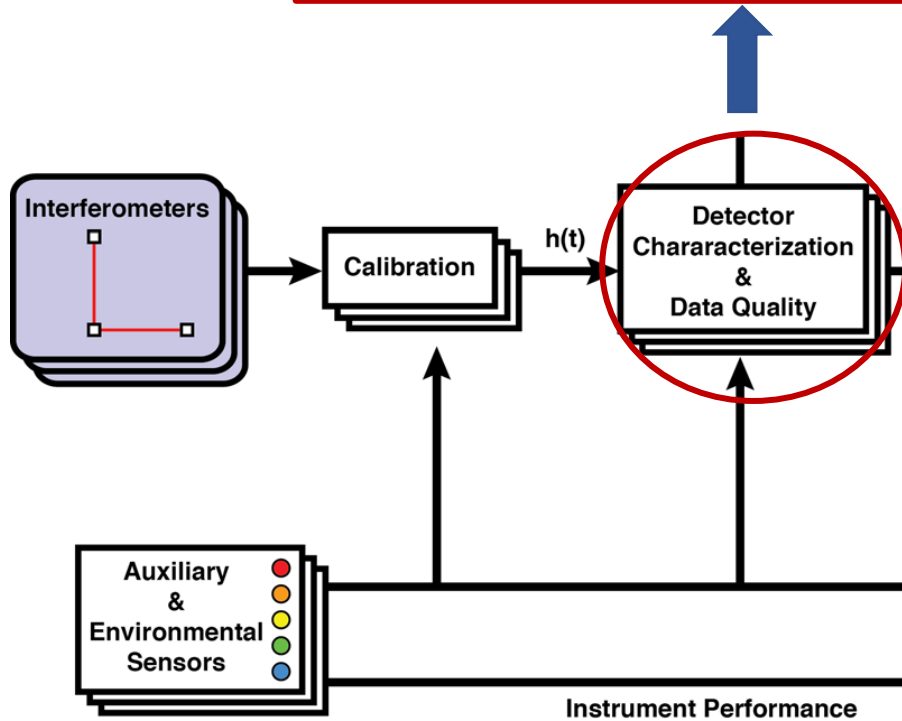
1.- Noise characterization and mitigation

- Find sources of noise and work in on-site mitigation
- Characterization of transient noise in the detectors



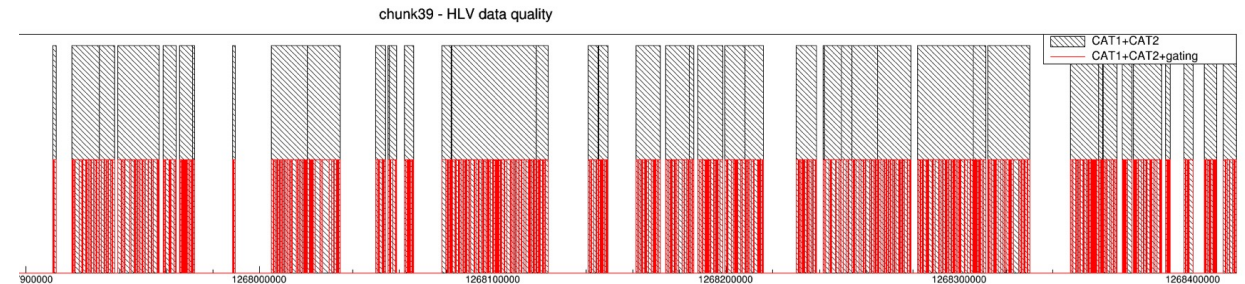
Data processing – Step 2: Check data quality

Output: Cat2 + gatedData
ready for searches &
public access



DetChar functions:

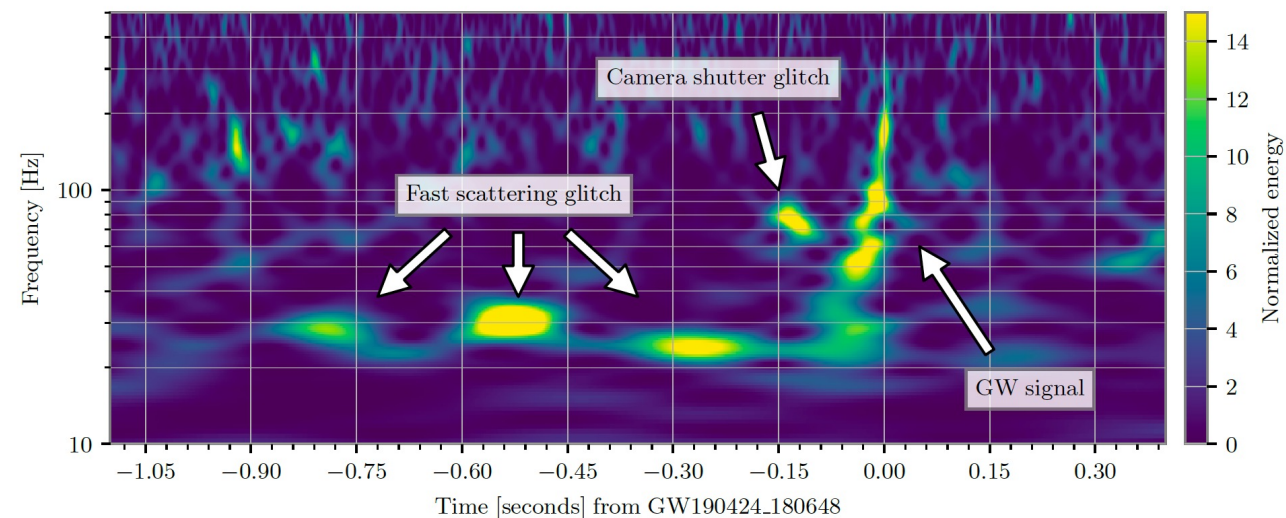
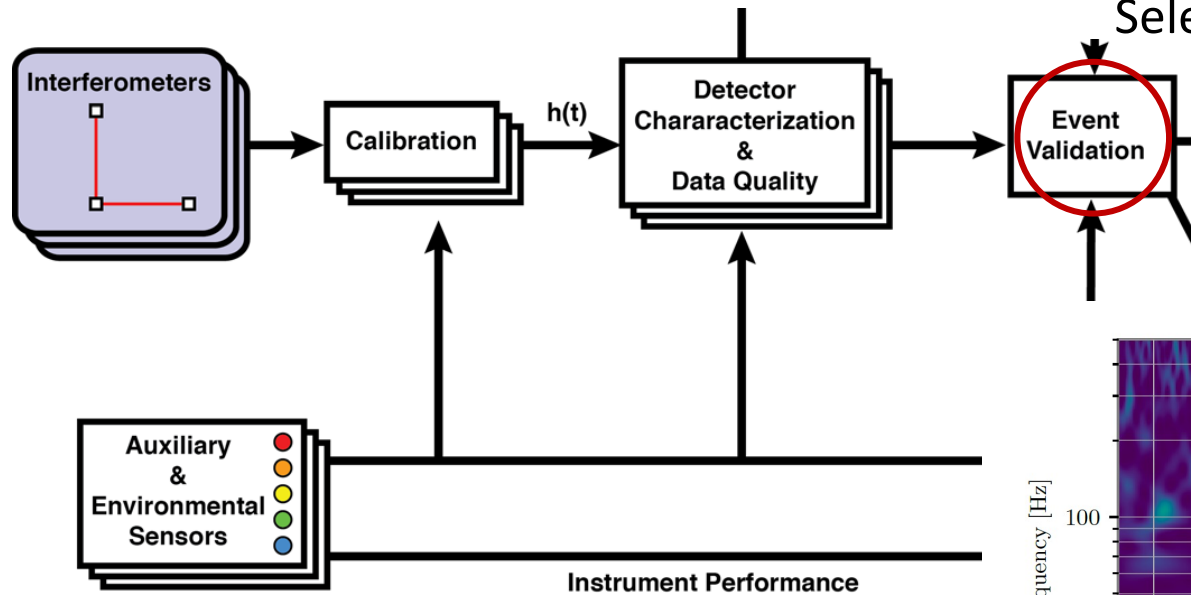
- 1.- Noise characterization and mitigation
 - Find sources of noise and work in on-site mitigation
 - Characterization of transient noise in the detectors
- 2.- Data Quality vetoes on strain time series
 - Select t when detectors are working without problems



Data processing – Step 2: Check data quality

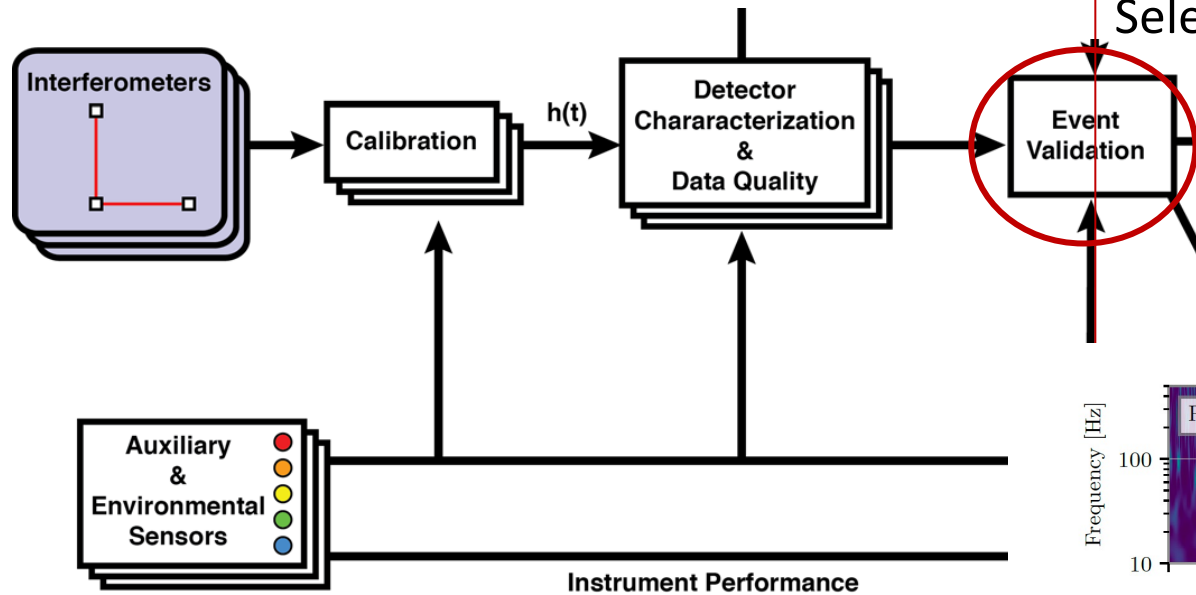
DetChar functions:

- 1.- Noise characterization and mitigation
 - Find sources of noise and work in on-site mitigation
 - Characterization of transient noise in the detectors
- 2.- Data Quality vetoes on strain time series
Select t when detectors are working without problems
- 3.- Event validation for candidates found by the search



Detector Characterization and Mitigation of Noise in Ground-Based Gravitational-Wave Interferometers
Published in: *Galaxies* 10 (2022) 1, 12

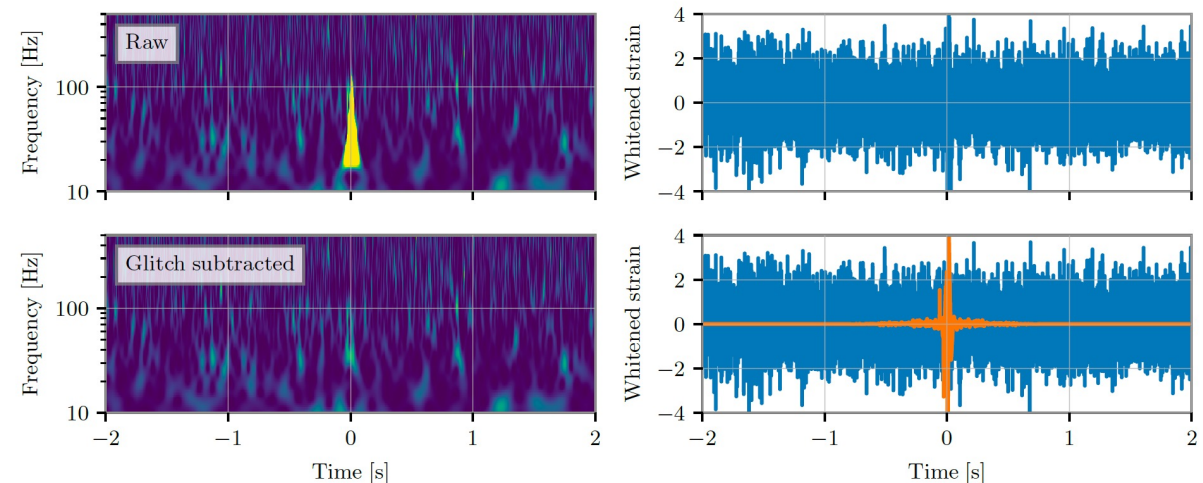
Data processing – Step 2: Check data quality



DetChar functions:

- 1.- Noise characterization and mitigation
 - Find sources of noise and work in on-site mitigation
 - Characterization of transient noise in the detectors
- 2.- Data Quality vetoes on strain time series
 - Select t when detectors are working without problems
- 3.- Event validation for candidates found by the search
- 4.- Noise subtraction: deglitching, gating...
20% of events, **“Delicate” processes...**

Noise subtraction software: BayesWave
New glitch + signal modelling methods:
<https://arxiv.org/abs/2205.13580>



Data processing – Step 3: Matched filtered searches

Signal identification: Match filtering of $d(t)$

1.- A common bank of templates for filtering is designed for each category: BBH, NSBH & BNS.

2.- For each template $h(t)$ w. param μ

$$\mathbf{h}(\theta) = A\mathbf{p}(t, \mu) \cos \phi + A\mathbf{q}(t, \mu) \sin \phi$$

calculate the SNR time series

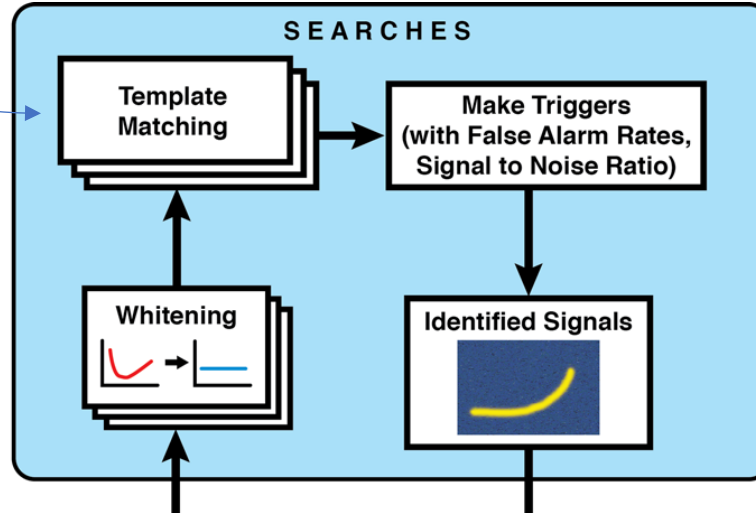
$$\rho(t, \mu) \equiv \sqrt{(\mathbf{d} | \mathbf{p}(t, \mu))^2 + (\mathbf{d} | \mathbf{q}(t, \mu))^2}$$

$$(\mathbf{a} | \mathbf{b}) = 2 \int_0^\infty \frac{\tilde{a}(f)\tilde{b}^*(f) + \tilde{a}^*(f)\tilde{b}(f)}{S_n(f)} df.$$

SNR quantifies the likelihood that the observed data contains a GW signal.

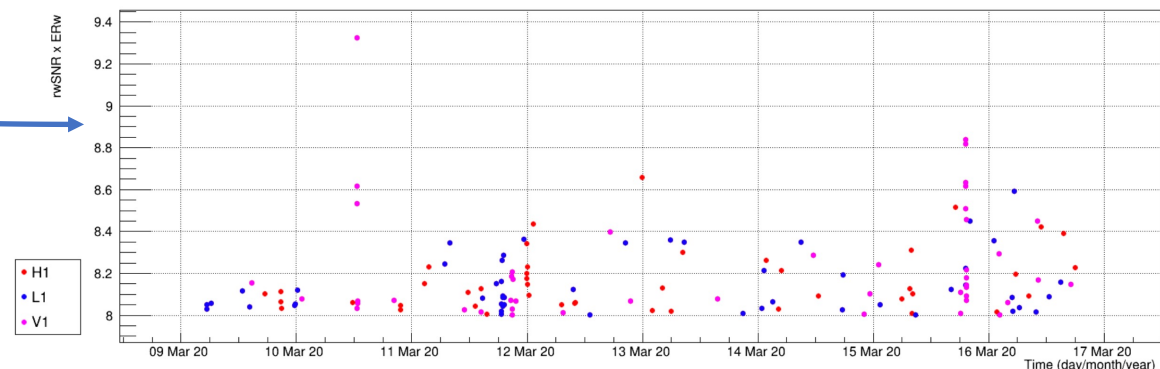
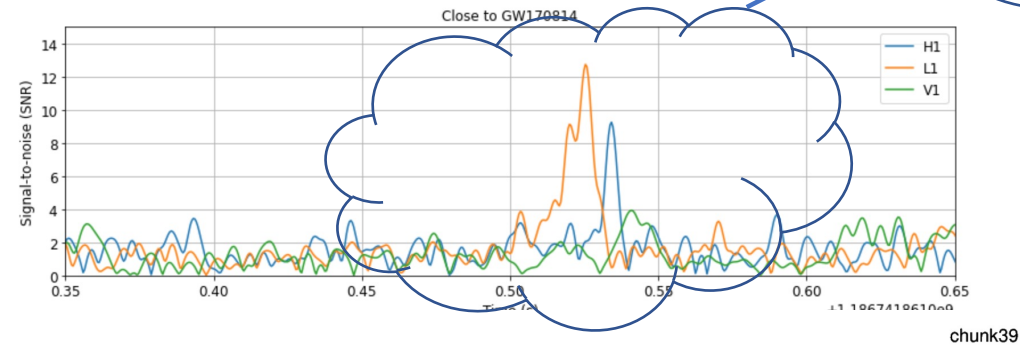
3.- As result of filtering, we get a collection of triggers:

4.- Cluster triggers in time in each detector, choose the most representative, look for coincident triggers in more than one detector to select signal candidates.



GraceDB

Candidate identified!



Data processing – Step 3': Event Ranking

SEARCHES

We get a HUGE amount of triggers, most of them caused by NOISE !!

Need to solve 2 problems:

1.- How do we rank the triggers to select the most “signal like” ones?

- SNR is optimal ranking statistics in Gaussian Noise: the higher the SNR the higher the likelihood that data contains a signal
- But non-gaussian glitches produce HIGH SNR!
- A reweighted-SNR (using χ^2 methods) is used as RE
- Each pipeline uses his own RE.

2.- How do we assign an statistical significance to each candidate event in terms of its ranking statistic?



Data processing – Step 3: Event Ranking

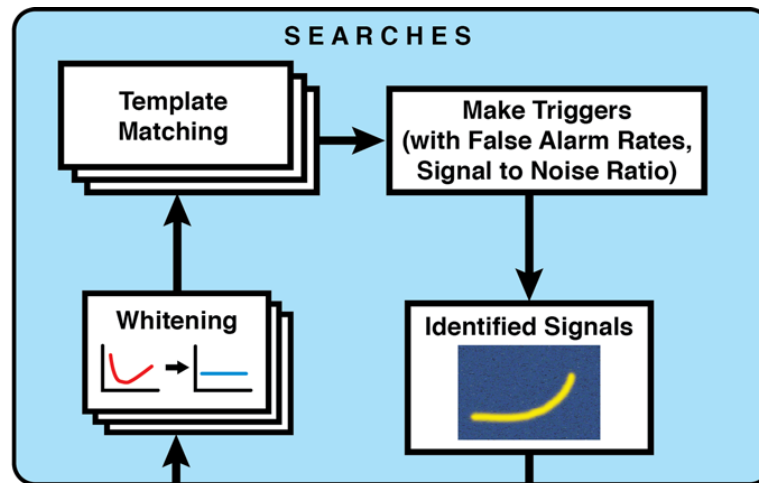
1.- Assign a FAR value to the event = the rate of background triggers with ranking statistics value equal to or greater than the RE of the event.

2.- The background distribution of the ranking statistic is estimated in a data driven way, by running the search over time-shifted detector data, so that coincidences become not physical.

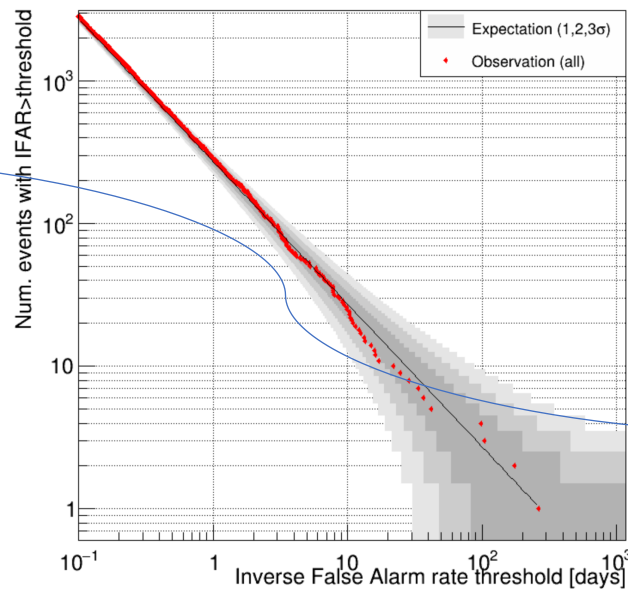
3.- Plot the cumulative FAR distribution for background and data, outstanding events would clearly appear, as in GWCT1 plots

4.- The FARxTobs = estimate of the probability of there being at least one noise trigger with a FAR this low or lower in the observed time.

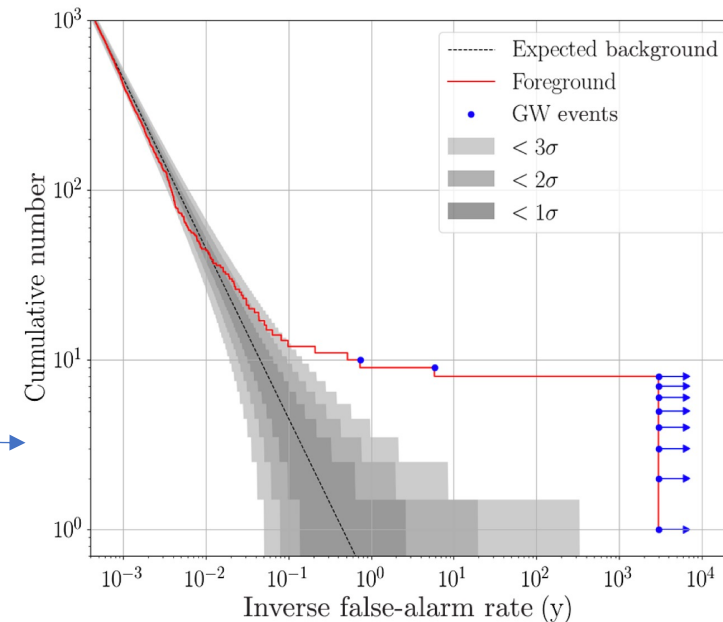
Typically: FAR < 2/year for O3a/b runs



O3 - 273.78 effective days of coincidences



MBTA in SSMO3b



PyCBC in GWTC1

Data processing – Step 3: Event Ranking

Problem:

Since FAR depends exponentially on the RE values, FAR values can differ by orders of magnitude among pipelines!

Name	Inst.	cWB			GstLAL			MBTA			PyCBC-broad			PyCBC-BBH		
		FAR (yr ⁻¹)	SNR	<i>p</i> _{astro}	FAR (yr ⁻¹)	SNR	<i>p</i> _{astro}	FAR (yr ⁻¹)	SNR	<i>p</i> _{astro}	FAR (yr ⁻¹)	SNR	<i>p</i> _{astro}	FAR (yr ⁻¹)	SNR	<i>p</i> _{astro}
GW191103_012549	HL	–	–	–	–	–	–	27	9.0	0.13	4.8	9.3	0.77	0.46	9.3	0.94
GW191105_143521	HLV	–	–	–	24	10.0	0.07	0.14	10.7	> 0.99	0.012	9.8	> 0.99	0.036	9.8	> 0.99
GW191109_010717	HL	< 0.0011	15.6	> 0.99	0.0010	15.8	> 0.99	1.8 × 10 ⁻⁴	15.2	> 0.99	0.096	13.2	> 0.99	0.047	14.4	> 0.99
GW191113_071753	HLV	–	–	–	–	–	–	26	9.2	0.68	1.1 × 10 ⁴	8.3	< 0.01	1.2 × 10 ³	8.5	< 0.01
GW191126_115259	HL	–	–	–	80	8.7	0.02	59	8.5	0.30	22	8.5	0.39	3.2	8.5	0.70
GW191127_050227	HLV	–	–	–	0.25	10.3	0.49	1.2	9.8	0.73	20	9.5	0.47	4.1	8.7	0.74
GW191129_134029	HL	–	–	–	< 1.0 × 10 ⁻⁵	13.3	> 0.99	0.013	12.7	> 0.99	< 2.6 × 10 ⁻⁵	12.9	> 0.99	< 2.4 × 10 ⁻⁵	12.9	> 0.99
GW191204_110529	HL	–	–	–	21	9.0	0.07	1.3 × 10 ⁴	8.1	< 0.01	980	8.9	< 0.01	3.3	8.9	0.74
GW191204_171526	HL	< 8.7 × 10 ⁻⁴	17.1	> 0.99	< 1.0 × 10 ⁻⁵	15.6	> 0.99	< 1.0 × 10 ⁻⁵	17.1	> 0.99	< 1.4 × 10 ⁻⁵	16.9	> 0.99	< 1.2 × 10 ⁻⁵	16.9	> 0.99
GW191215_223052	HLV	0.12	9.8	0.95	< 1.0 × 10 ⁻⁵	10.9	> 0.99	0.22	10.8	> 0.99	0.0016	10.3	> 0.99	0.28	10.2	> 0.99
GW191216_213338	HV	–	–	–	< 1.0 × 10 ⁻⁵	18.6	> 0.99	9.3 × 10 ⁻⁴	17.9	> 0.99	0.0019	18.3	> 0.99	7.6 × 10 ⁻⁴	18.3	> 0.99
GW191219_163120	HLV	–	–	–	–	–	–	–	–	–	4.0	8.9	0.82	–	–	–

From GWTC3

Data processing – Step 3: Event Ranking

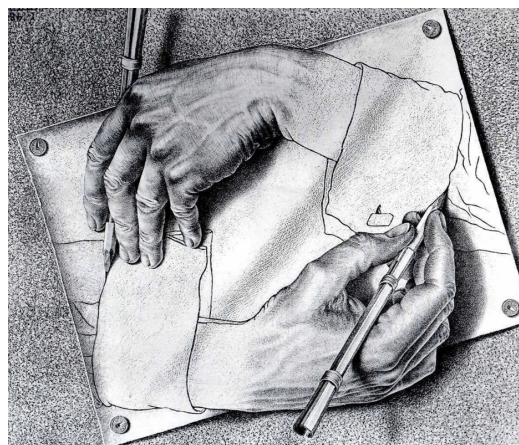
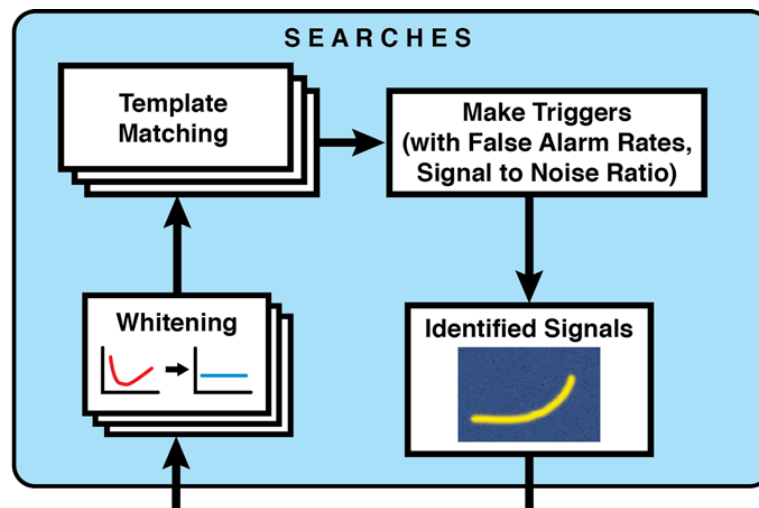
1.- Assign a FAR value to the event = the rate of background triggers with ranking statistics value equal to or greater than the RE of the event.

2.- The background distribution of the ranking statistic is estimated in a data driven way, by running the search over time-shifted detector data, so that coincidences become not physical.

3.- Plot the cumulative FAR distribution for background and data, outstanding events would clearly appear, as in GWCT1 plots

4.- The FAR_{xTobs} = estimate of the probability of there being at least one noise trigger with a FAR this low or lower in the observed time.

Typically: $FAR < 2/\text{year}$ for O3a/b runs



¿Does Pastro have this risk?

• Present situation:

FAR eliminated as statistical criteria from the catalogs after GWTC1.

For an evaluation of the FAR from Gaussian Noise in GD Detectors, see our recent paper: [arXiv:2209.05475](https://arxiv.org/abs/2209.05475)

• Since GWTC2.1 Pastro introduced to quantify the “probability of astrophysical origin”.

Characteristics:

- Event goes to catalog if $Pastro > 0.5$.
- Each pipeline computes its own.
- Based on prior knowledge of the population properties and rates!

Main steps in data processing

See Jose Francisco's talk
for the rest of topics!

Thank You!

