

AI/ML Computing for Gravitational Waves

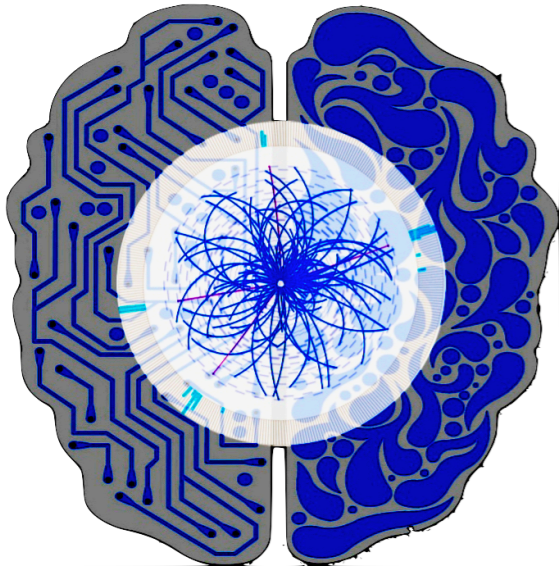


Phil Harris(MIT)
A3D3 deputy director

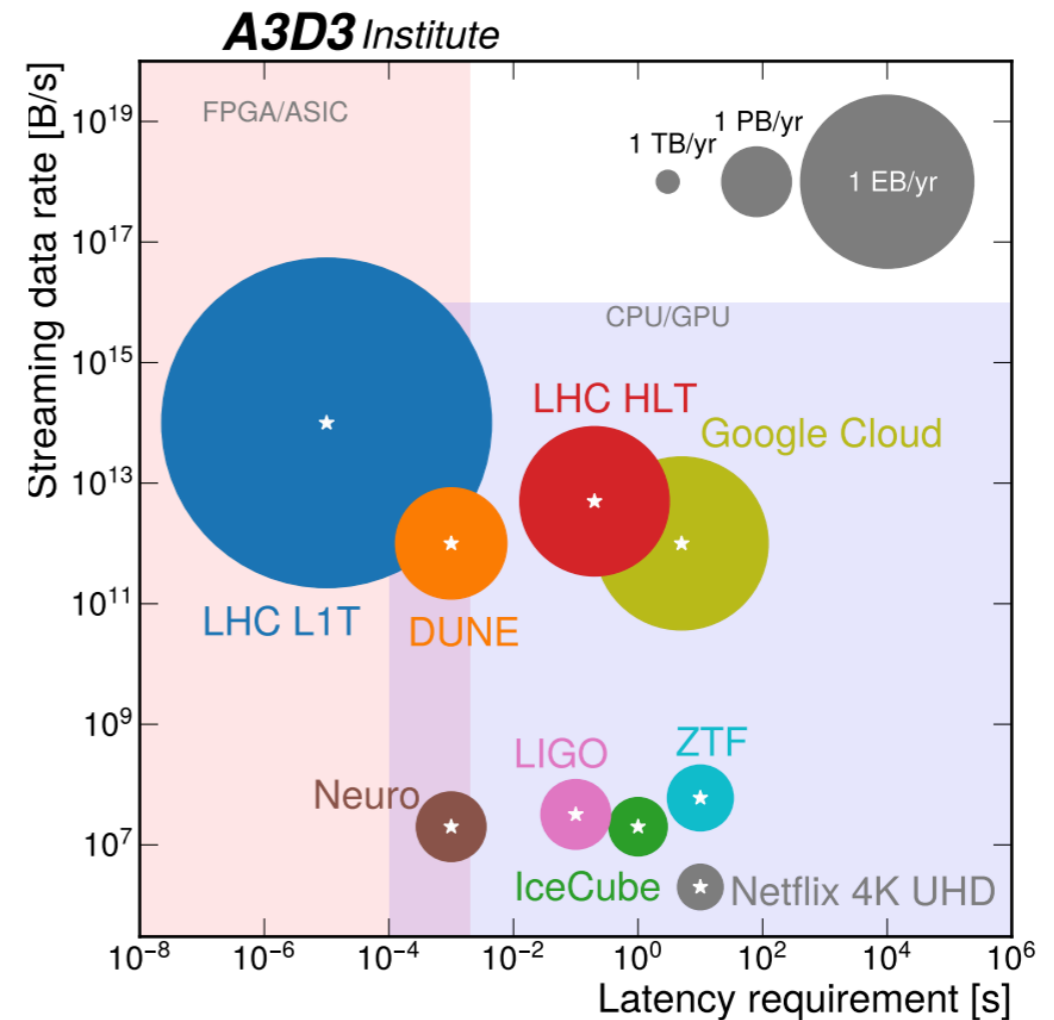
Alec Gunny¹, Ethan Marx¹, Will Benoit², Deep Chatterjee¹, Michael Coughlin², Katya Govorkova¹, Erik Katsavounidis¹, Eric Moreno¹, Rafia Omer², Ryan Raikman¹, Muhammed Saleem²

1 - Massachusetts Institute of Technology
2 - University of Minnesota

Understanding this Problem



FastML For Science

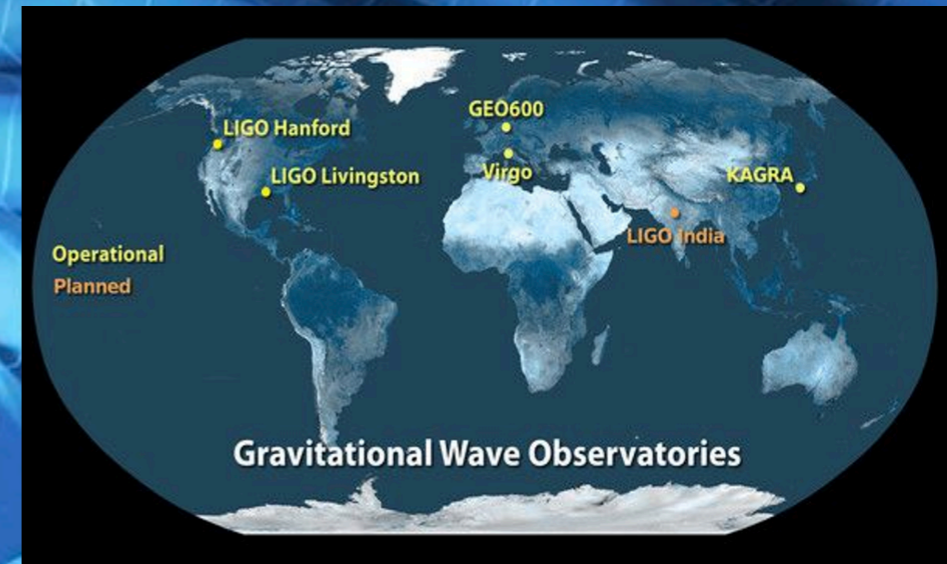
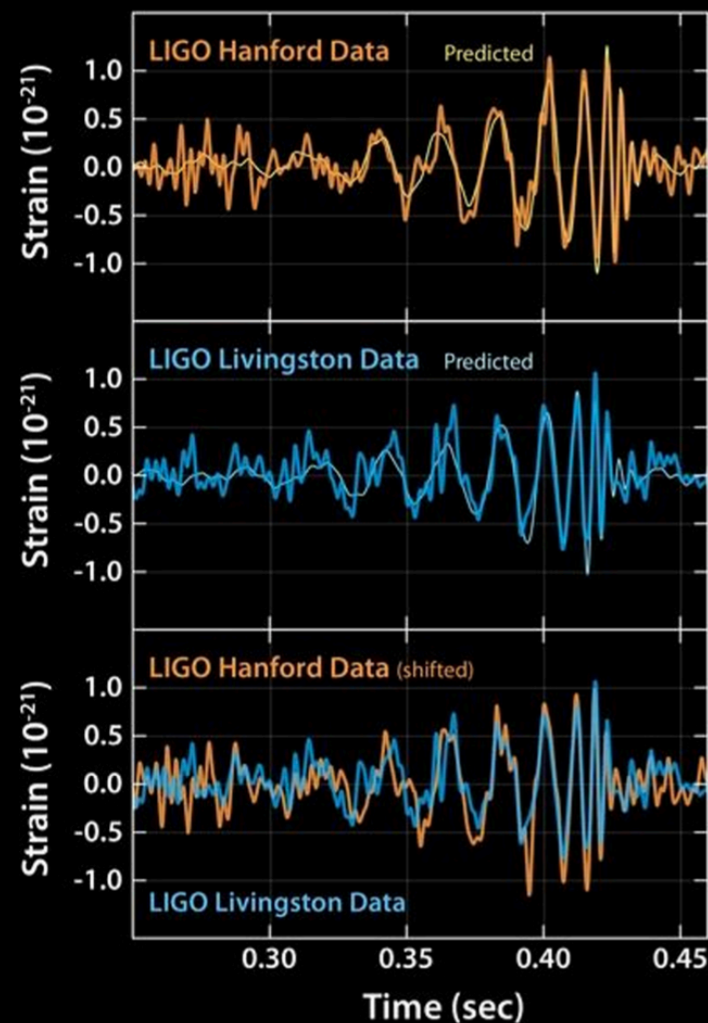


- FastML/A3D3 created to address real-time AI for science
 - Developing ML + GPU integration for large throughput computing
 - Developing ML+ FPGA/ASIC for low latency computing
- Science benchmarks are competitive with the rest of AI world

Gravitational Waves

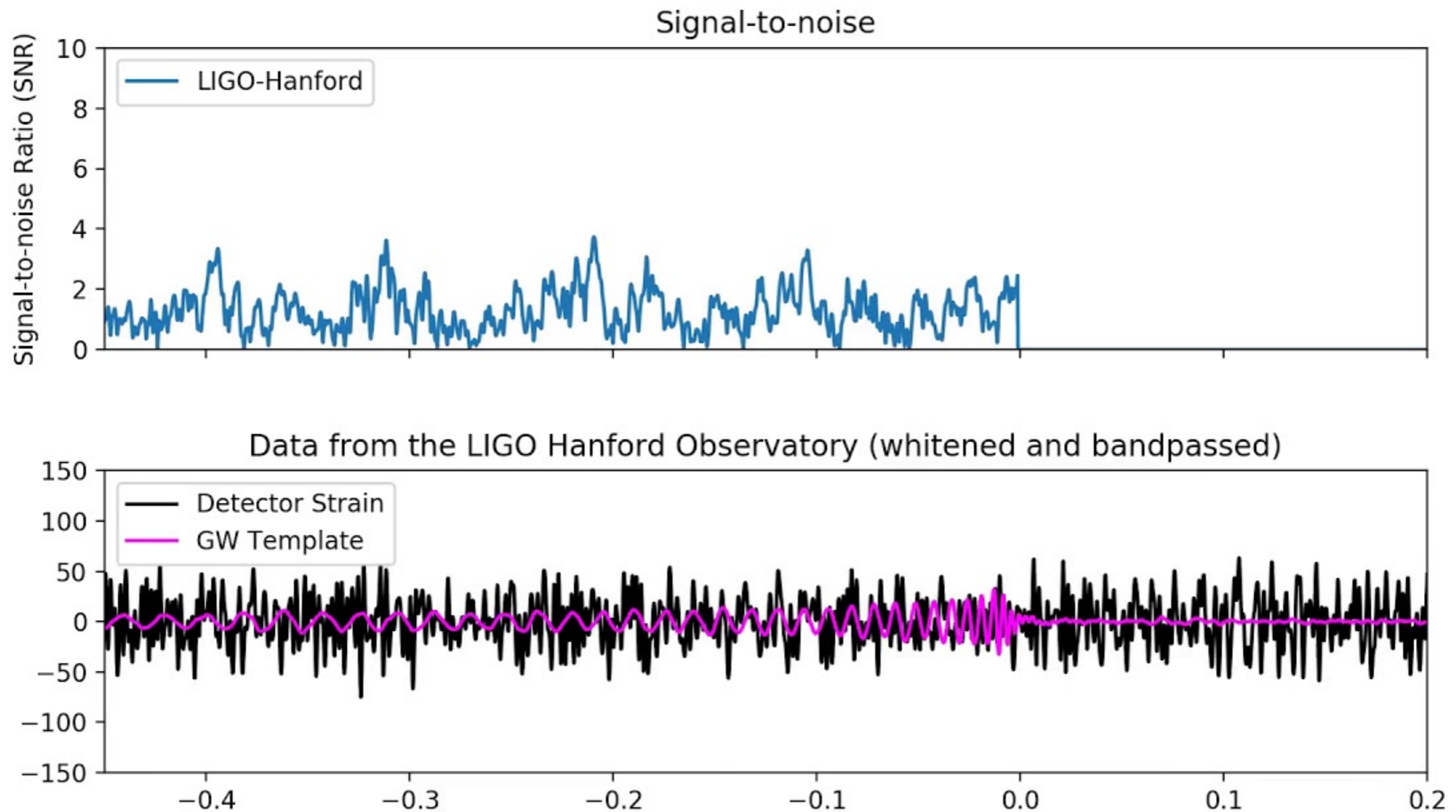
Ripples in Spacetime

International Gravitational Wave Network to characterize Gravitational Wave



Grand Challenge:
Can we identify GWs fast for downstream telescopes?

Typical LIGO Signal

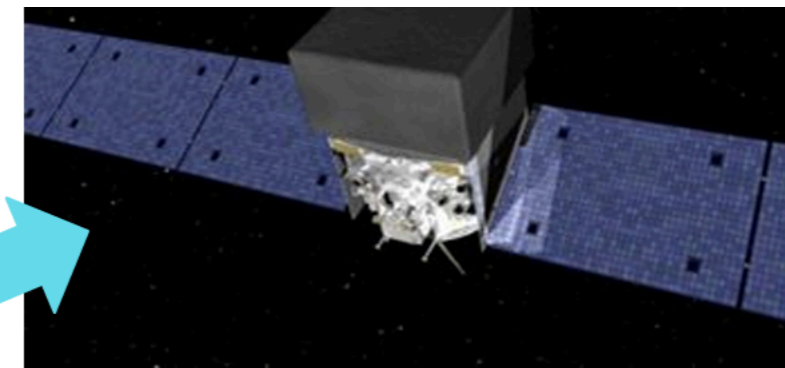
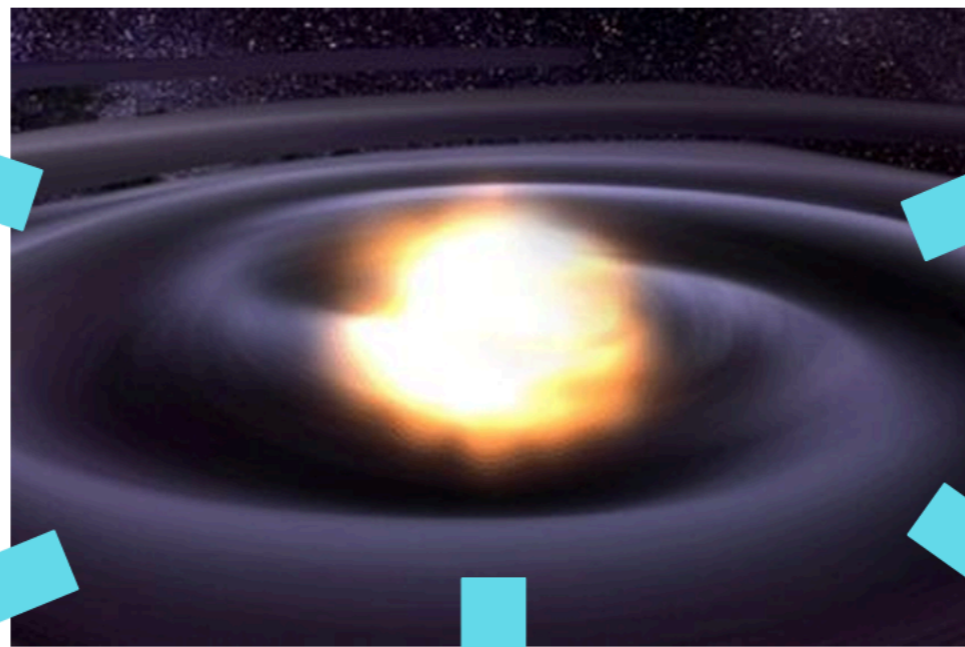


To ensure we aren't seeing a glitch we use at multiple detectors
2 LIGO detectors in US + Virgo detector in Europe + Kagra in Japan

Multi-Messenger Astronomy



Gravitational waves



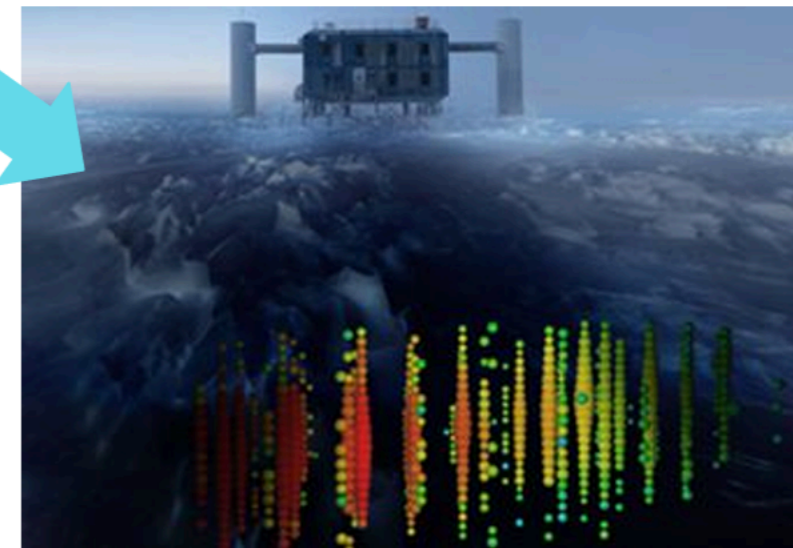
X-rays/Gamma-rays



Visible/infrared light



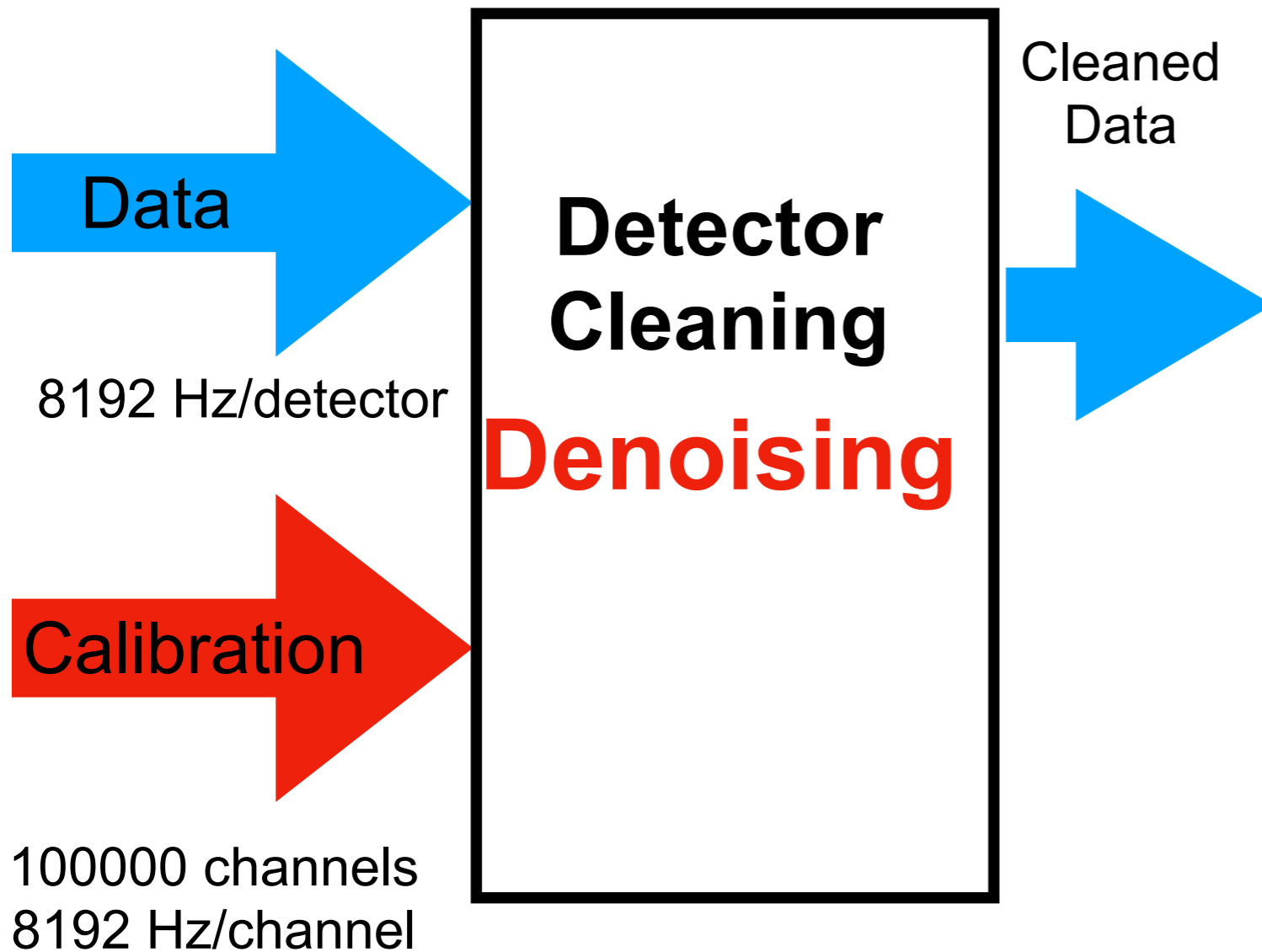
Radio waves



Neutrinos

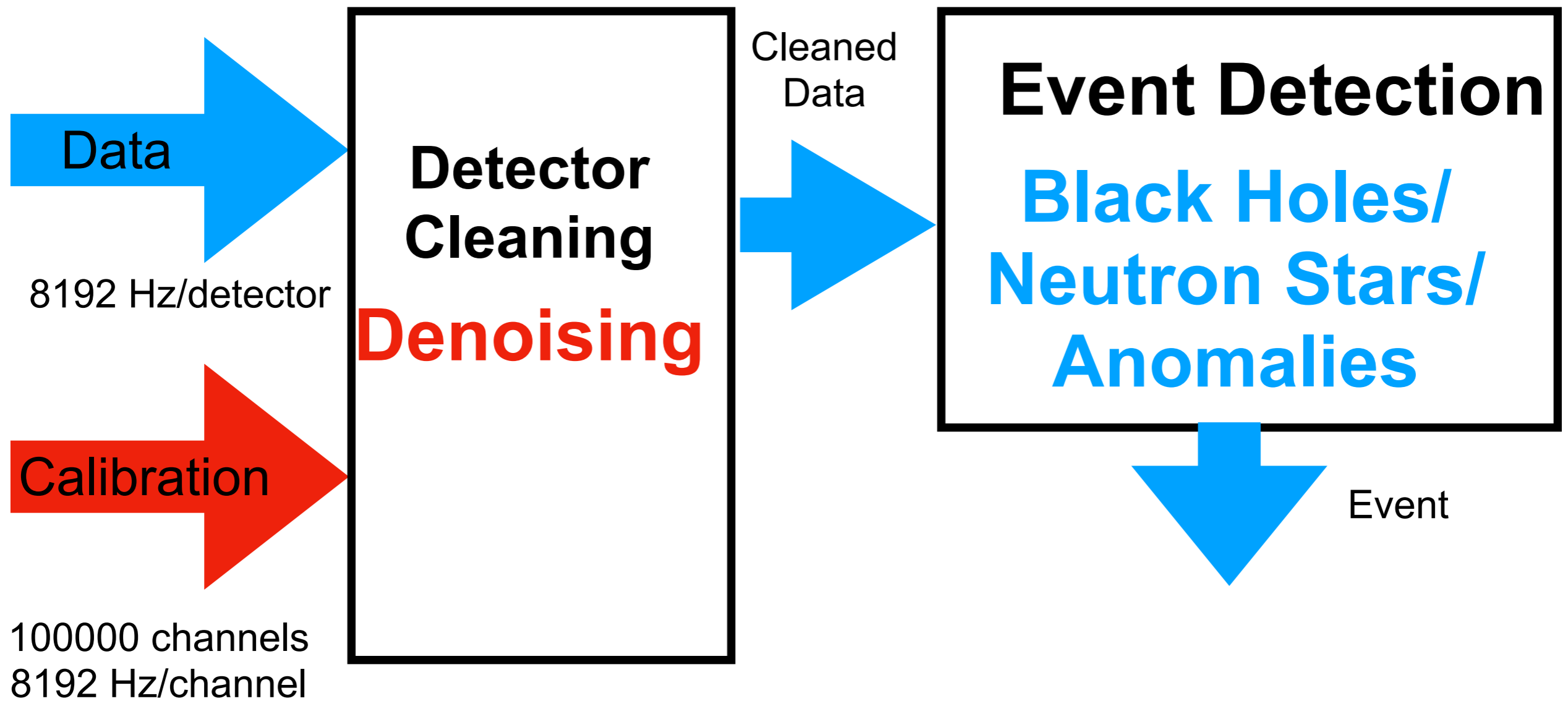
Performing fast identification of GWs critical to alerting world!

LIGO Data Workflow



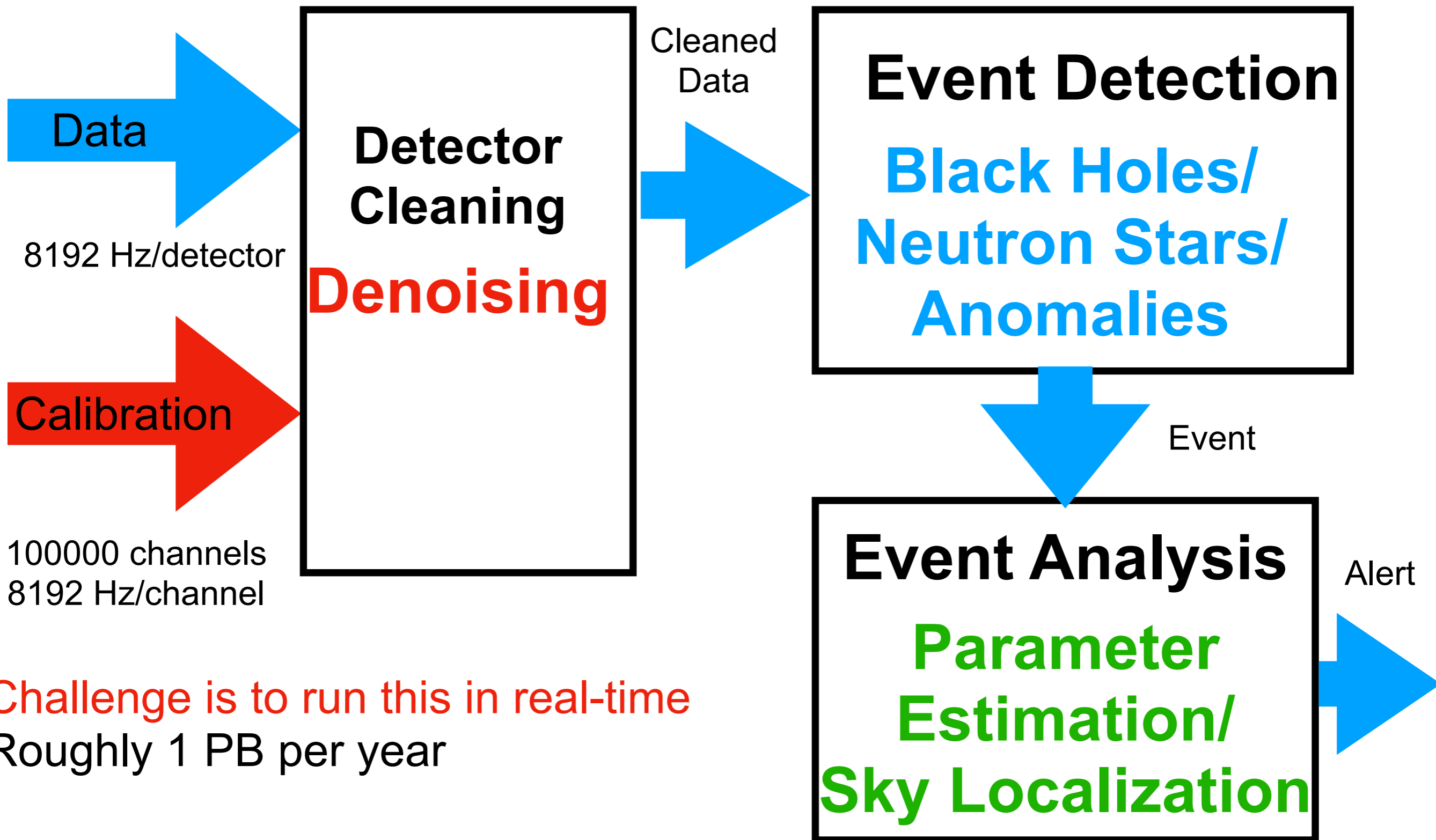
Challenge is to run this in real-time
Roughly 1 PB per year

LIGO Data Workflow



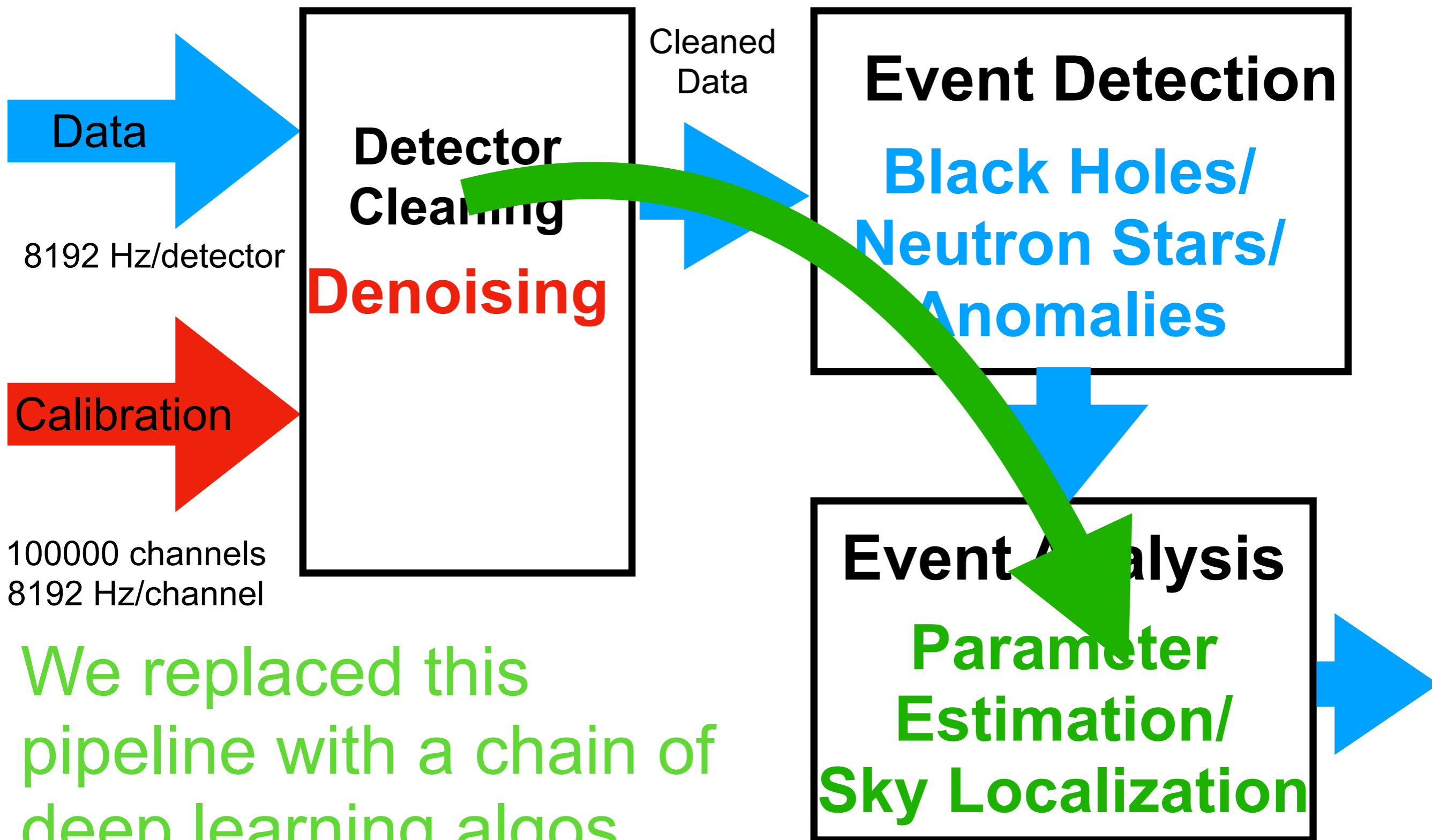
Challenge is to run this in real-time
Roughly 1 PB per year

LIGO Data Workflow

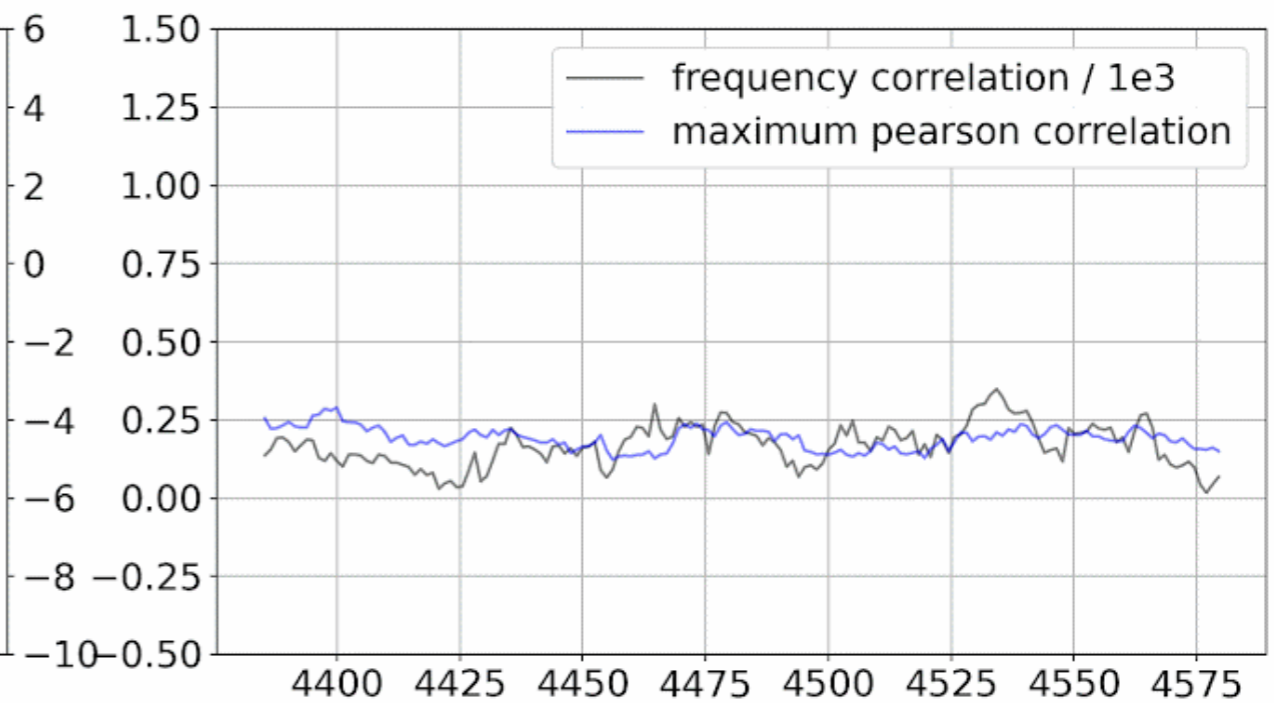
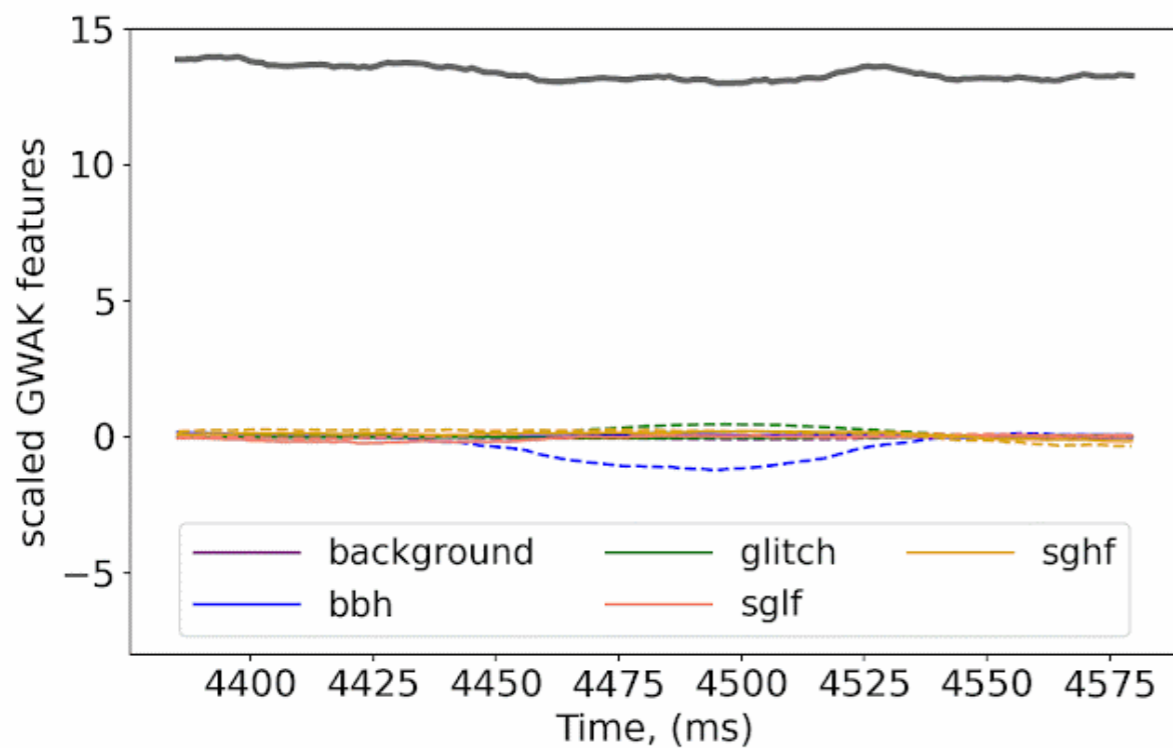
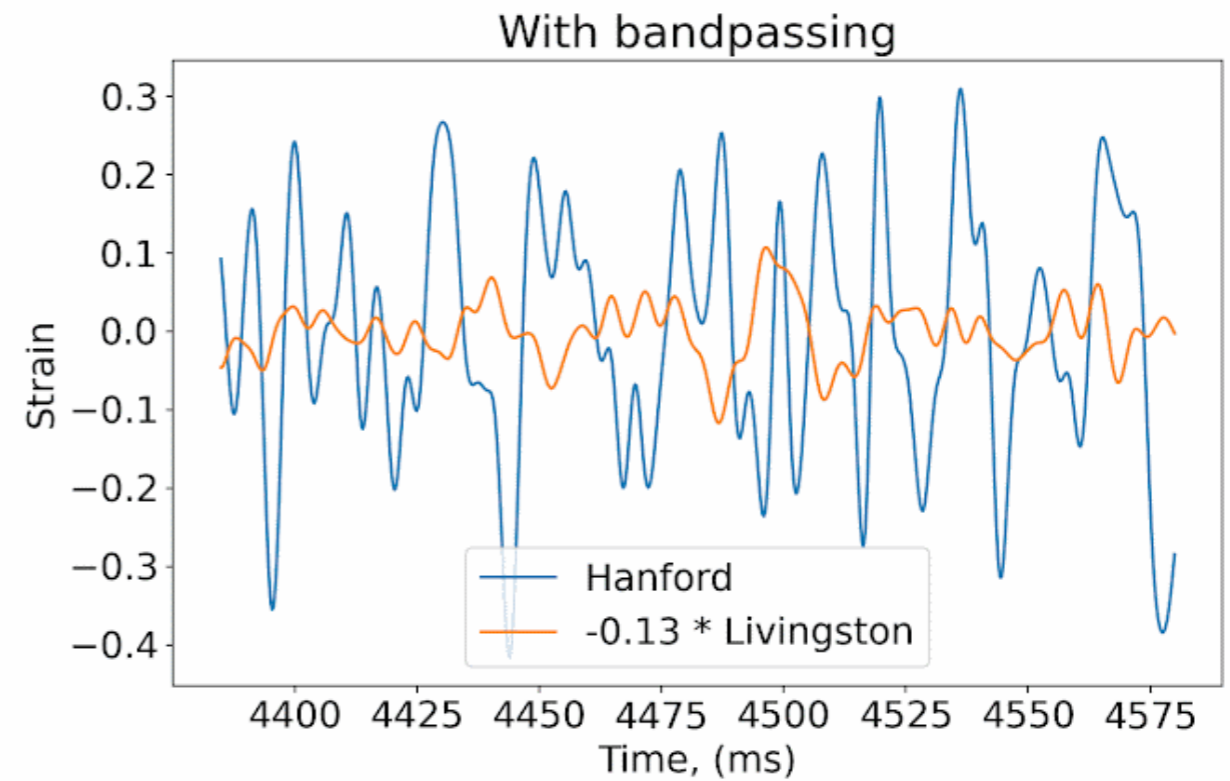
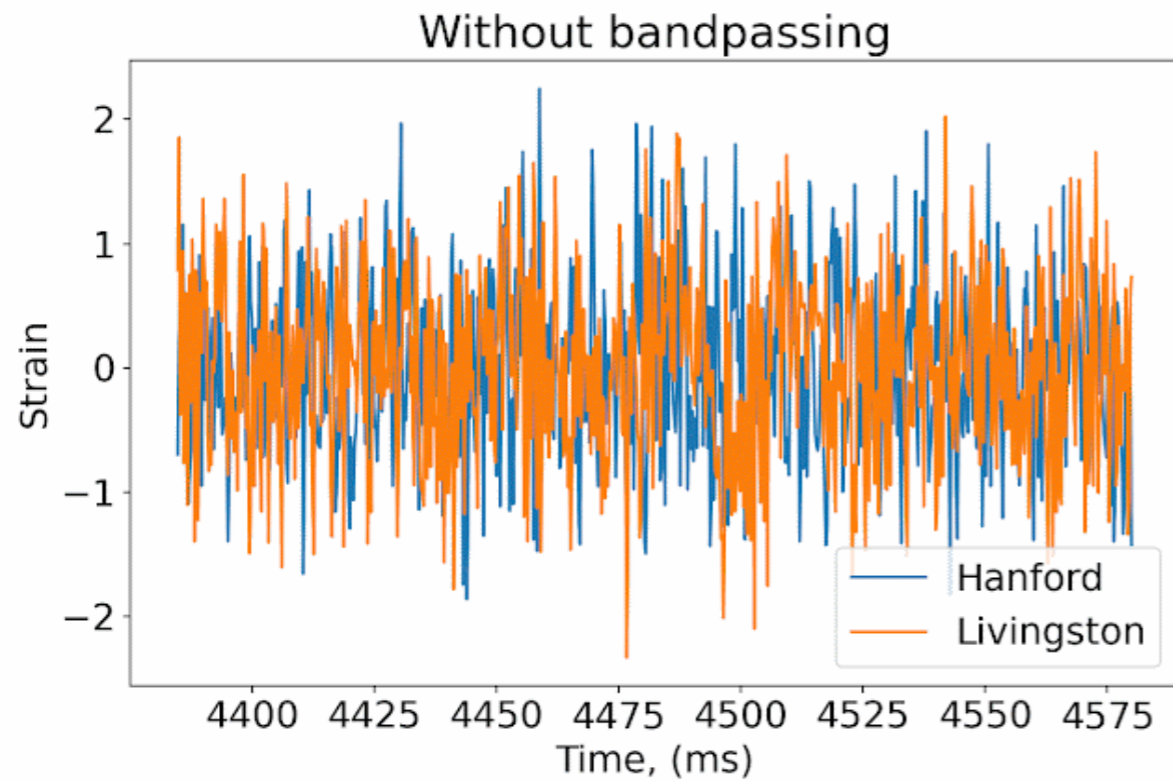


Challenge is to run this in real-time
Roughly 1 PB per year

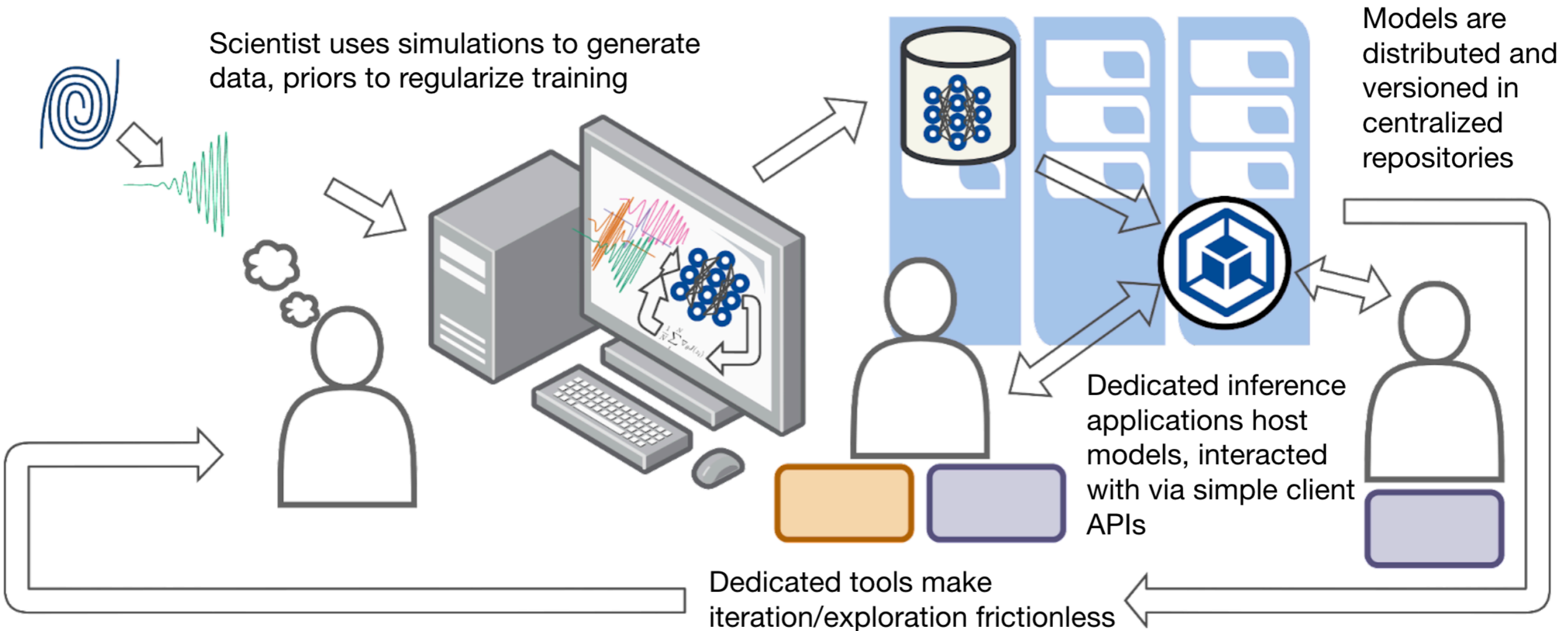
Our Upgraded Workflow



LIGO Data

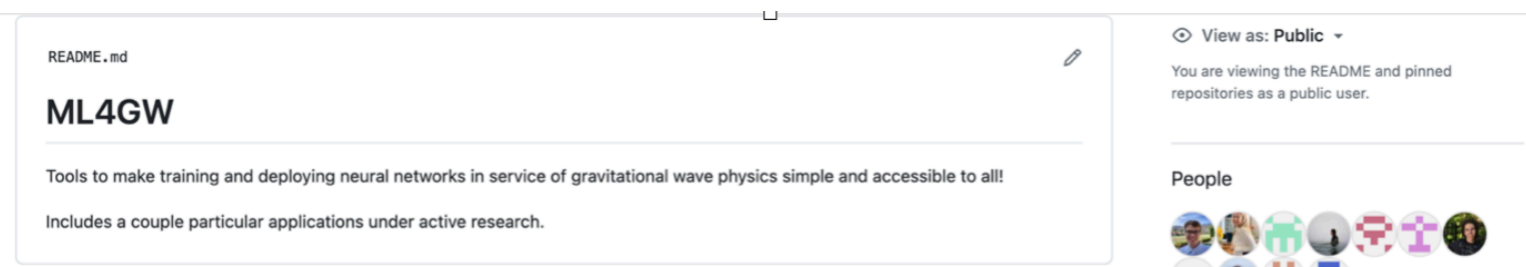


MLOps Toolkit



ML4GW Toolkit

To enable fast deployment + optimized heterogeneity



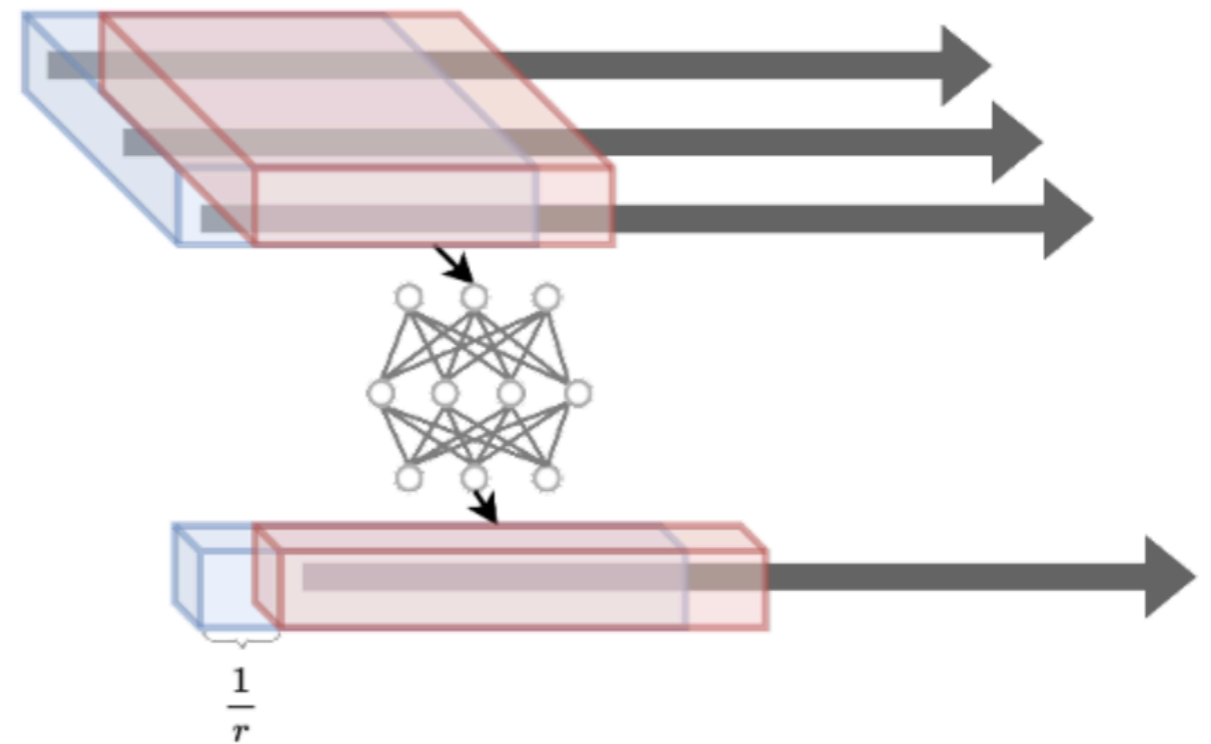
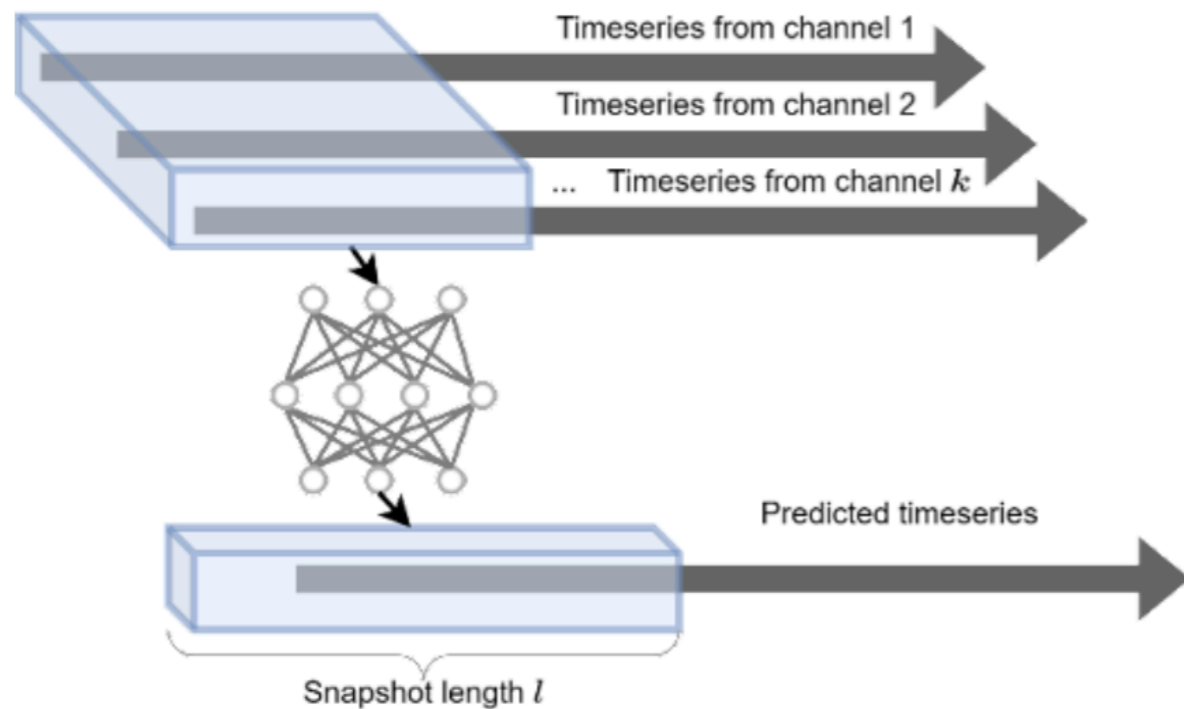
Example: Stateful Caching

Conventional Processing of Time series data

First ML Inference on data

Second ML Algorithm for data

Traditional IaaS



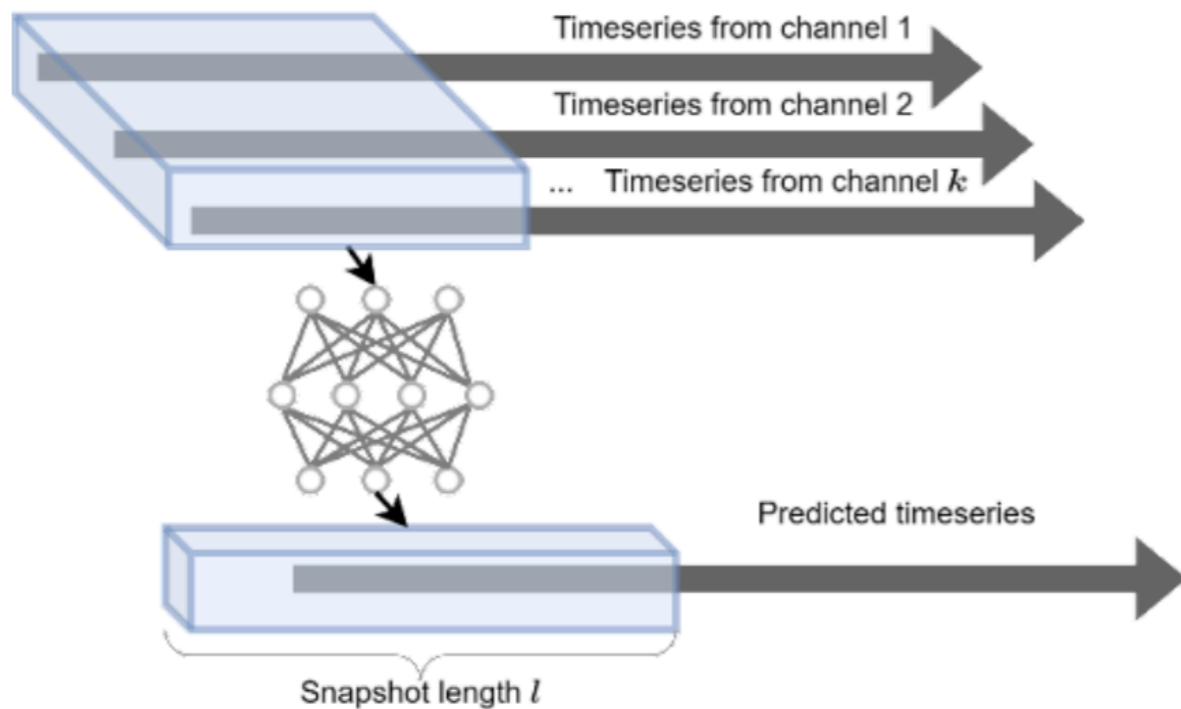
Example: Stateful Caching

Optimized Processing of Time series data

First ML Inference on data

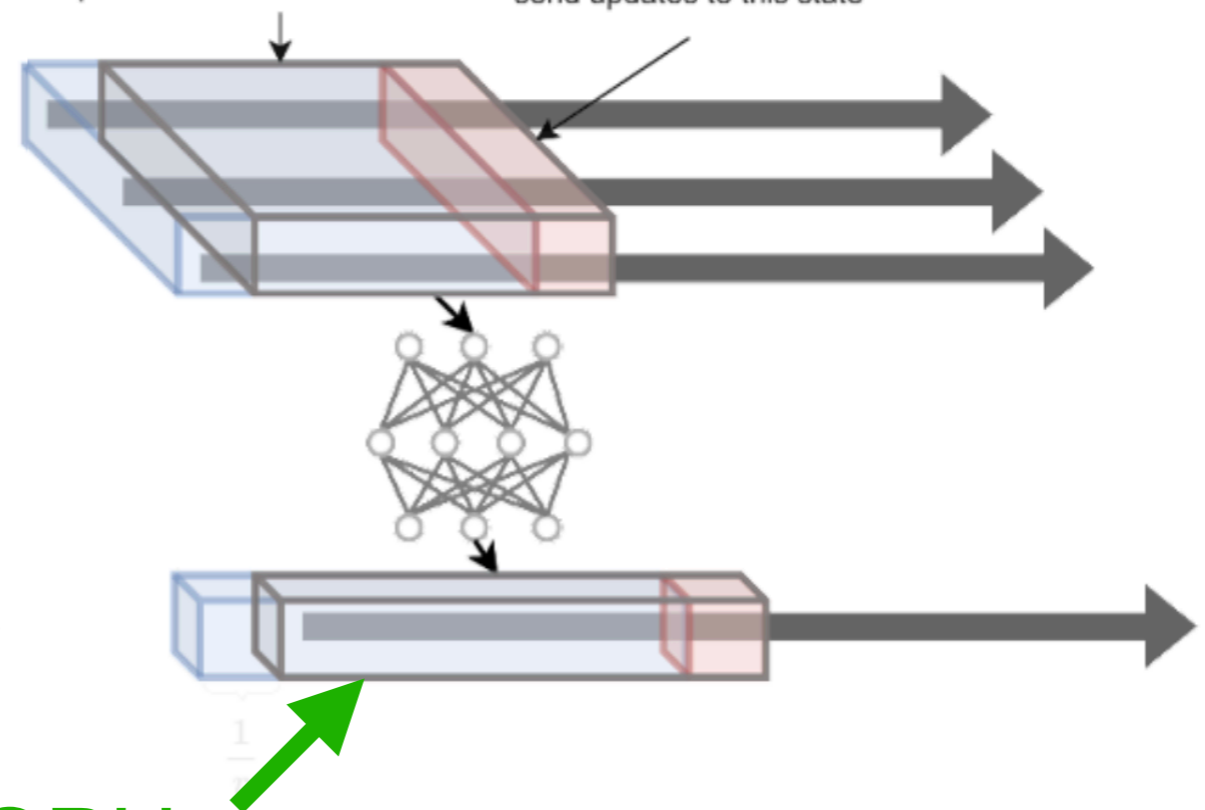
Second ML Algorithm for data

Traditional IaaS



Snapshotter model on inference service maintains most recent input to model as a state

Client only needs to send updates to this state



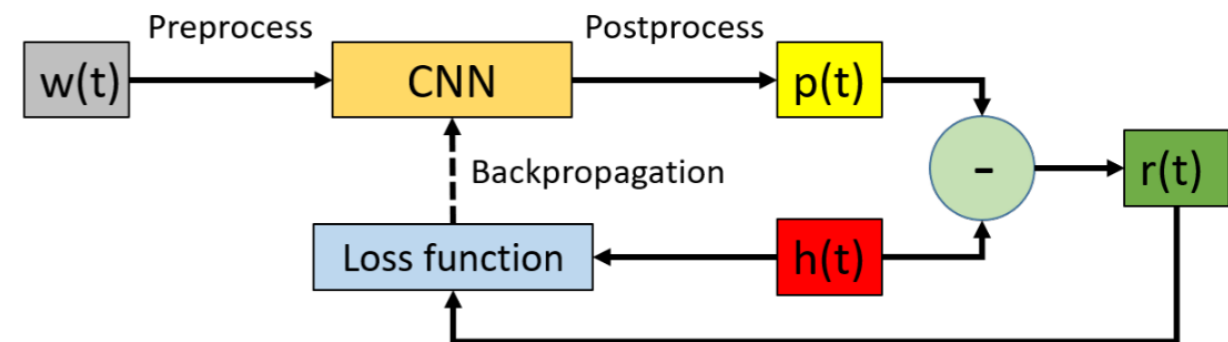
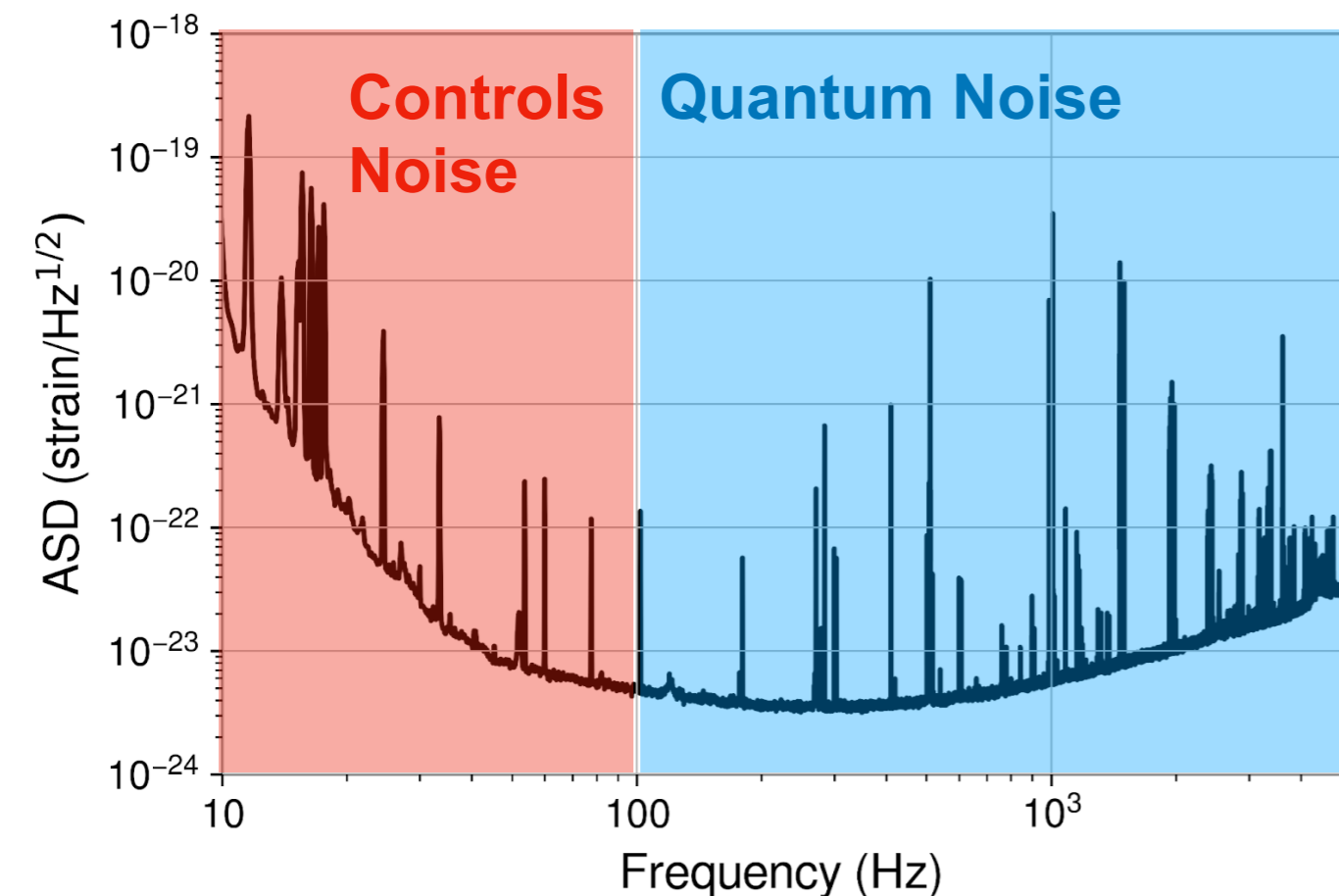
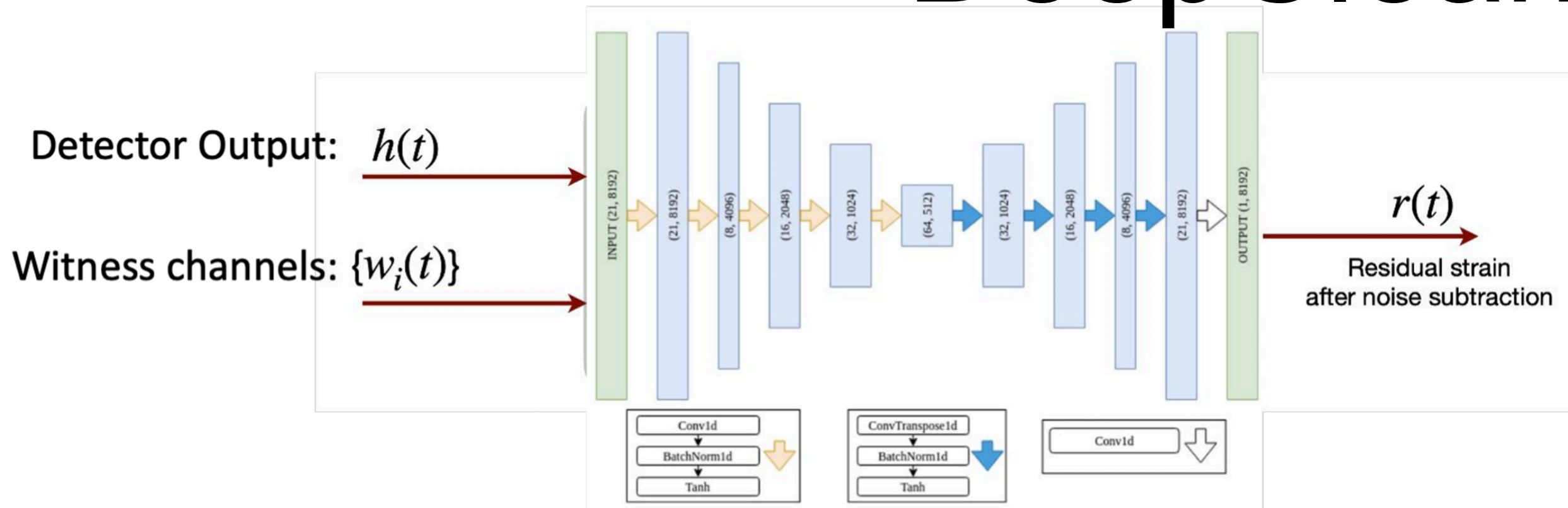
Cache on the GPU

Cacheing+optimizations speeds up ML by 10x over vanilla pytorch

The Algos

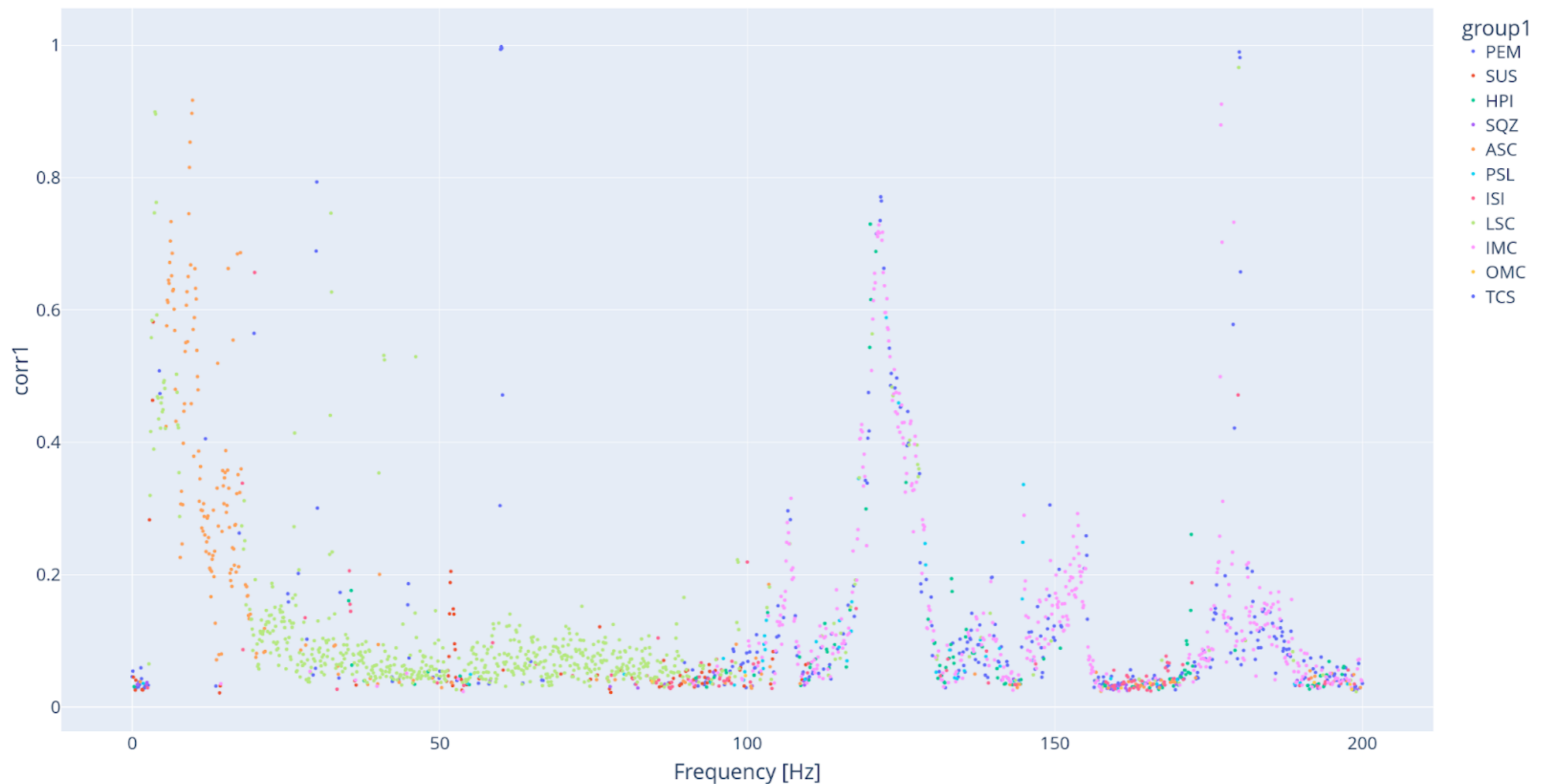
- Step 1: Denoising
 - Deepclean: Cleaning up detector in real-time
- Step 2: Transient detection
 - A-frame: black hole merger detection in real-time
 - GWAK : unmodeled (anomalous) GW detection
- Step 3:
 - AMPLIFI: real-time parameter estimation using LFI
- Step 4: Alert

DeepClean



1000 Monitored Witness channels
 100k total Witness channels
 Noise reduced through cross correlation w/witness channels

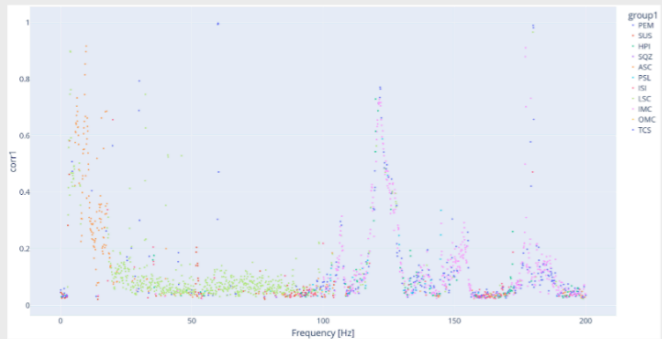
Real time Monitoring



- Daily monitor noise and its correlation w/1000 channels
 - Find channels on a daily+frequency basis w/highest correlation
 - Feed them into the NN for training and decorrelation

Real time Monitoring

Coherence Monitor



Requirements:

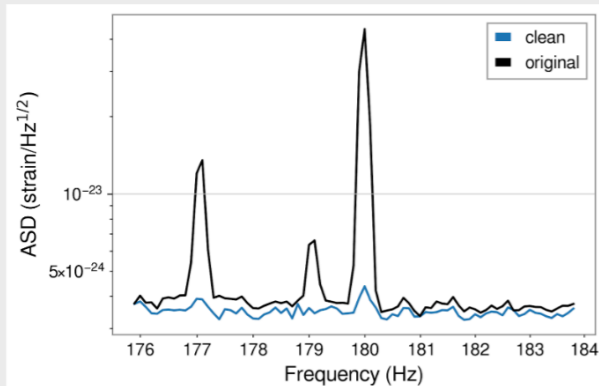
- 1st/2nd highest coherence
- coherence > 0.2

Witness channels

```
H1:ASC-CSOFT_P_OUT_DQ
H1:ASC-DHARD_Y_OUT_DQ
H1:PEM-CS_MAG_LVEA_OUTPUTOPTICS_X_DQ
H1:ASC-CHARD_Y_OUT_DQ
H1:ASC-DSOFT_P_OUT_DQ
H1:PEM-CS_MAG_EBAY_LSCRAK_X_DQ
H1:ASC-REFL_A_RF45_I_YAW_OUT_DQ
H1:LSC-POP_A_RF45_I_ERR_DQ
H1:ASC-DHARD_P_OUT_DQ
H1:SUS-ETMX_L3_OPLEV_PIT_OUT_DQ
H1:ASC-AS_A_DC_NSUM_OUT_DQ
```

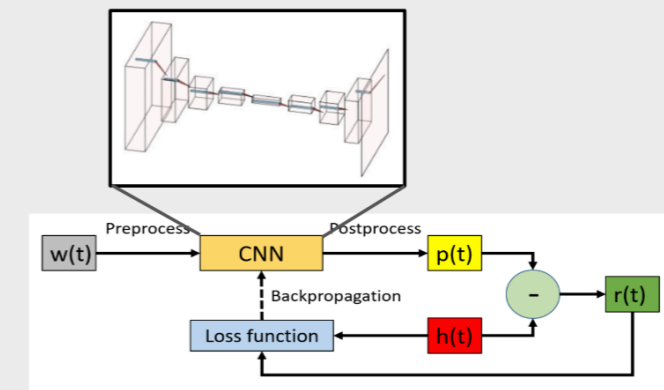
Fully automated training
(hyperparameters, ...)

Cleaned Strain



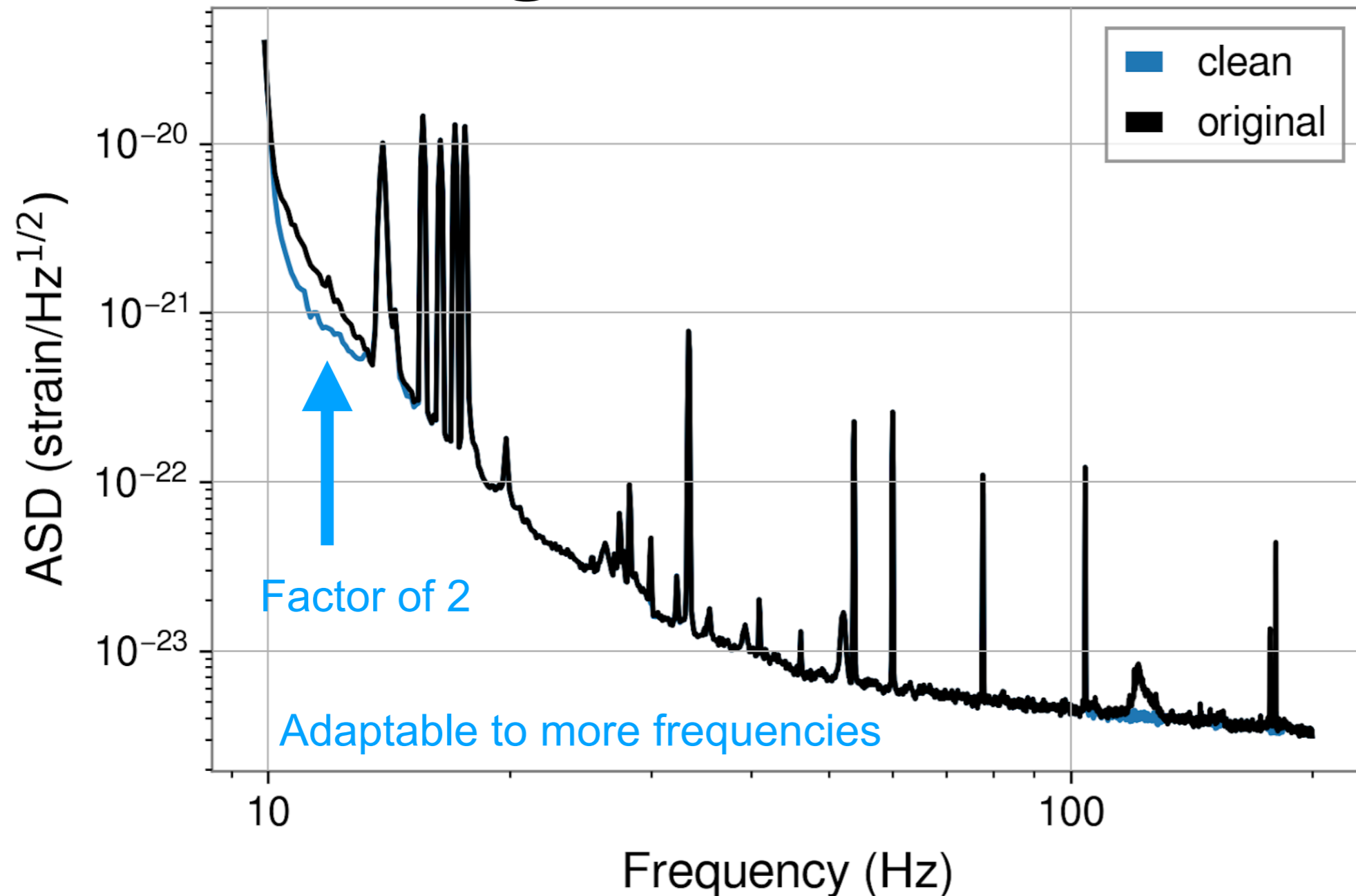
Apply cleaning on
independent data

DeepClean



- Aim to run end-to-end turnaround in < 1h of arrival

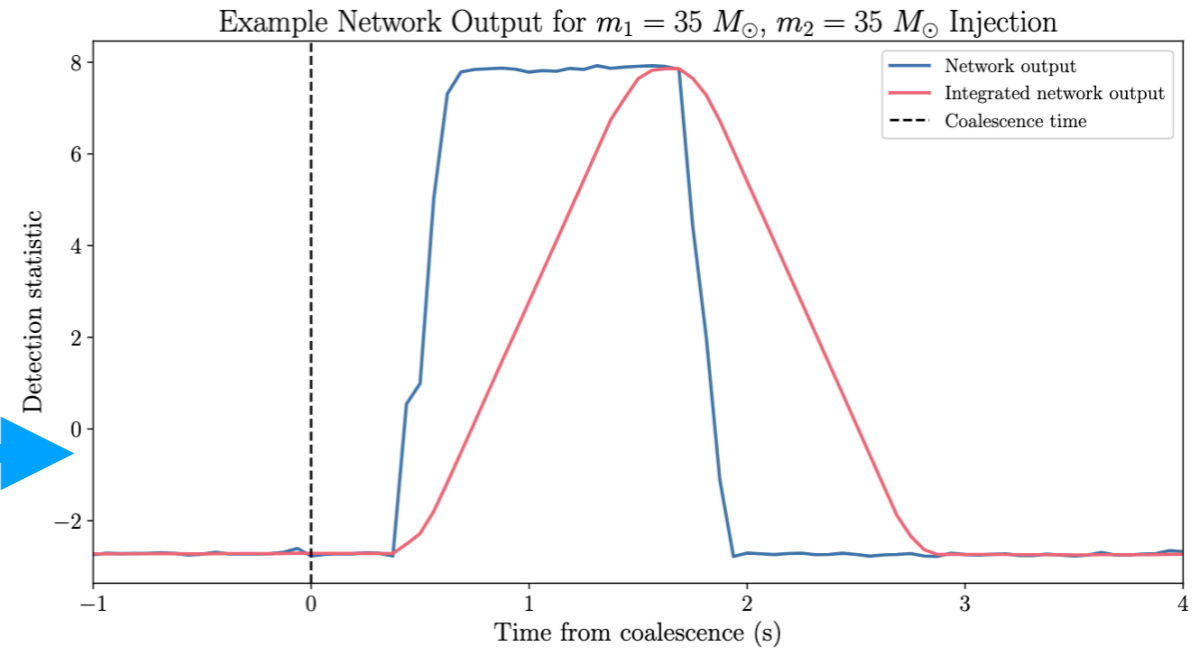
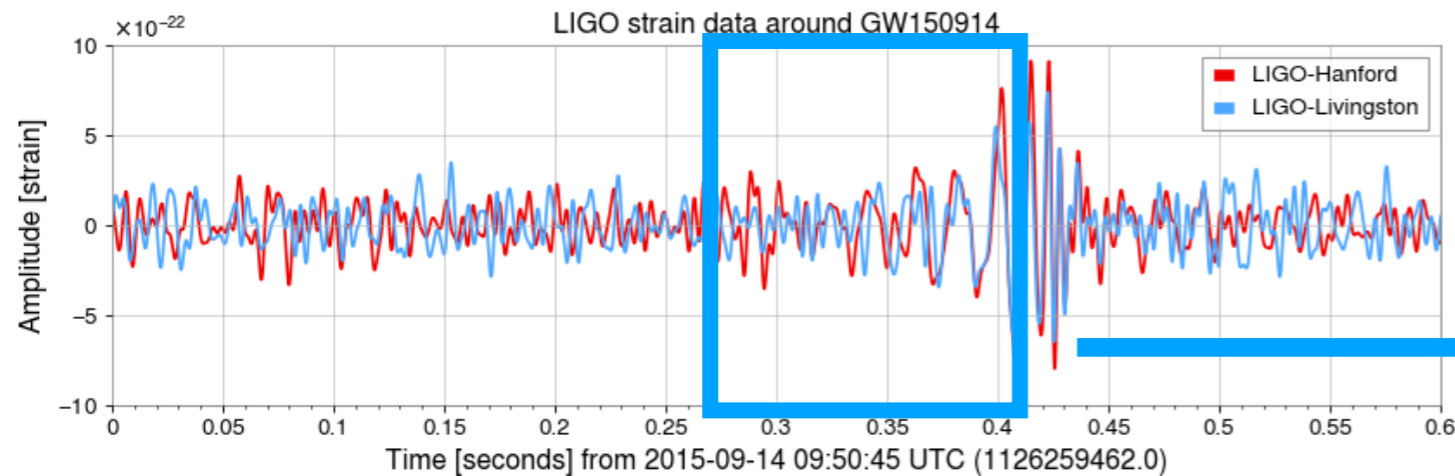
Significant Reduction



- Toolkit is adaptable for many different frequency problems
 - Denoising in frequency space is a common problem

A-Frame

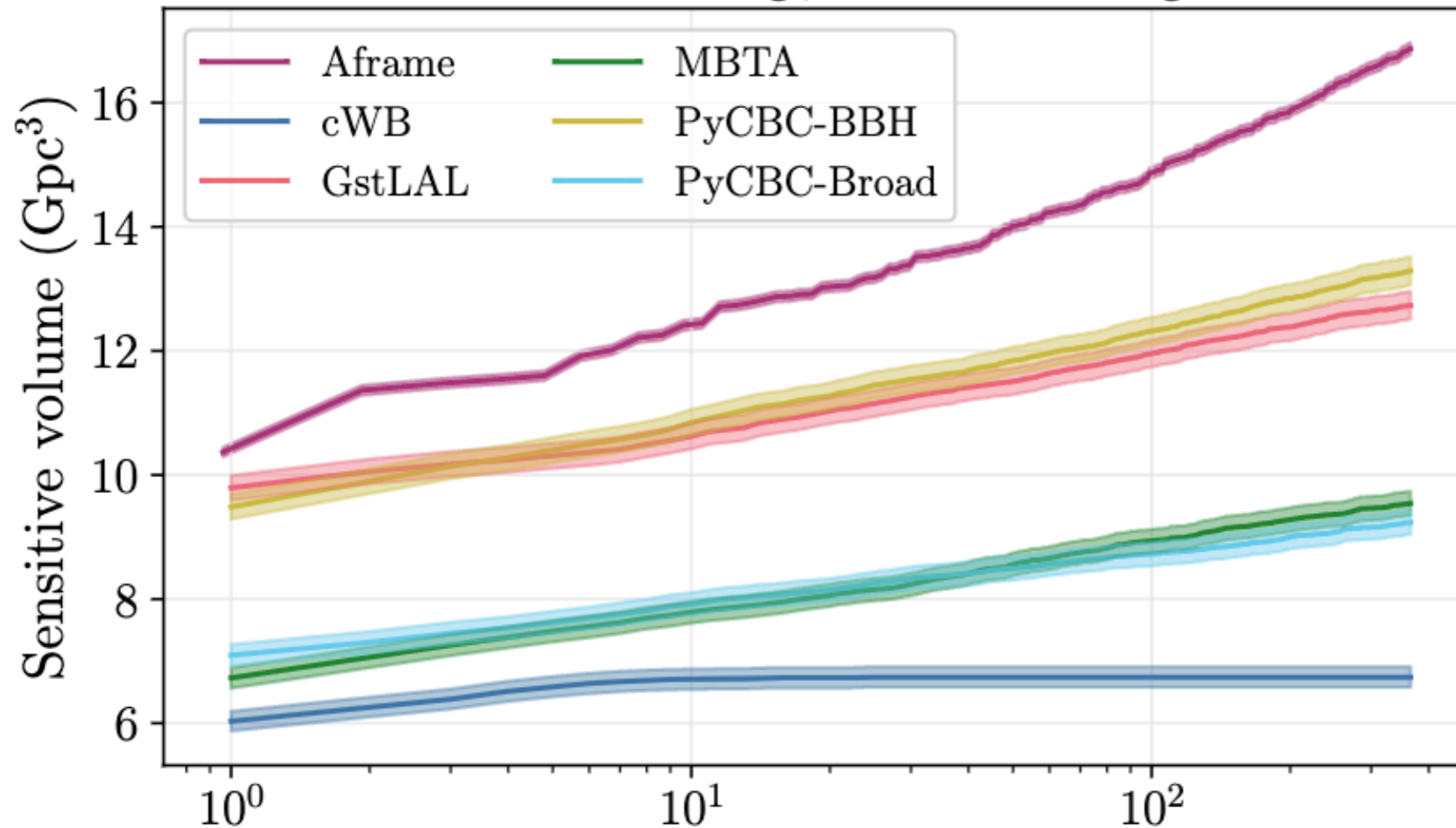
Sliding NN Score



- Neural network targeting Compact Binary Coalescence
 - Aka Black Hole Mergers
 - Working to adapt this to other signatures
 - ▶ Neutron star mergers
 - Curriculum learning scheme and Glitch mitigation essential
 - ▶ Sophisticated on GPU mixing tools for this

A-Frame

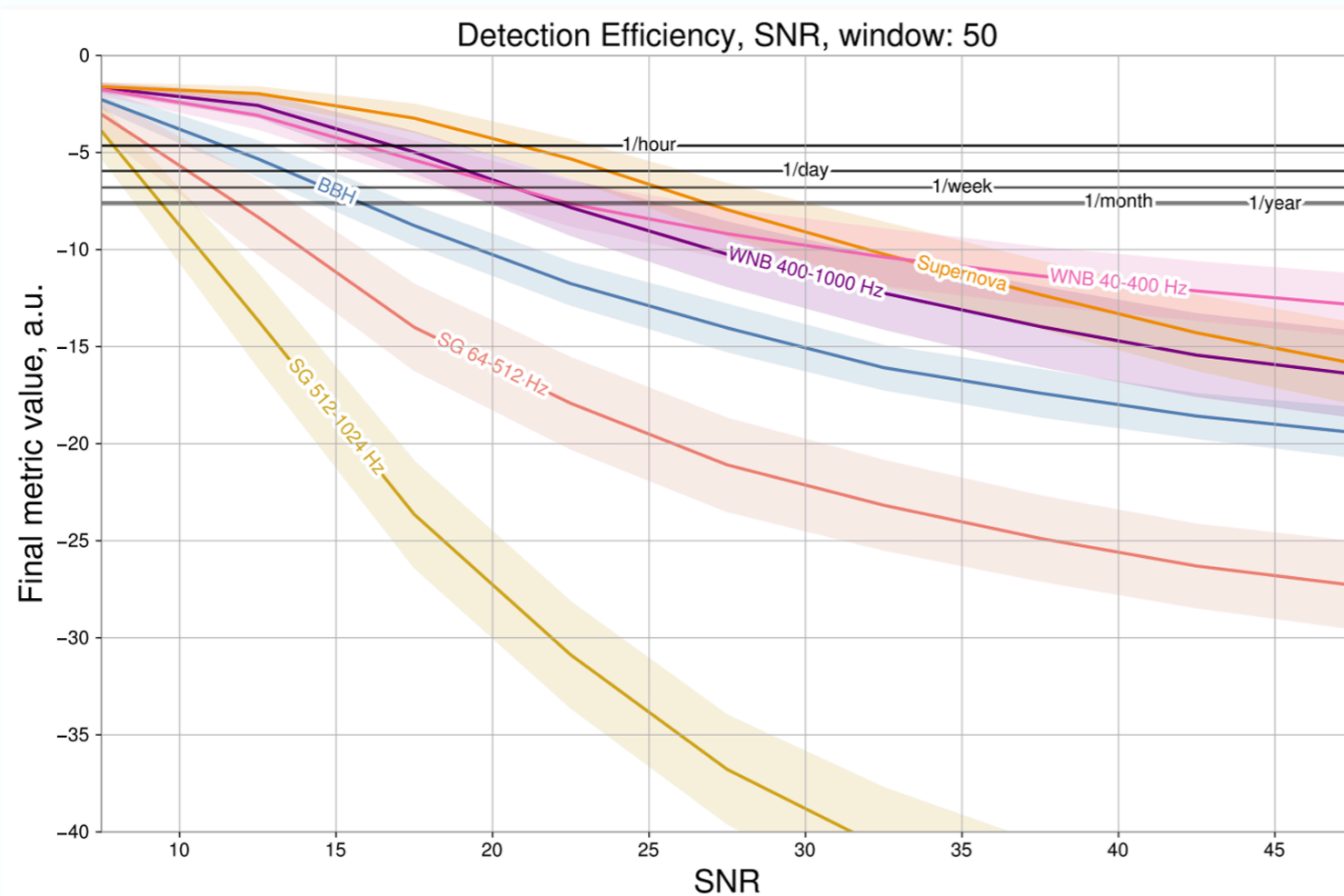
$$m_1 = 35 M_{\odot}, m_2 = 35 M_{\odot}$$



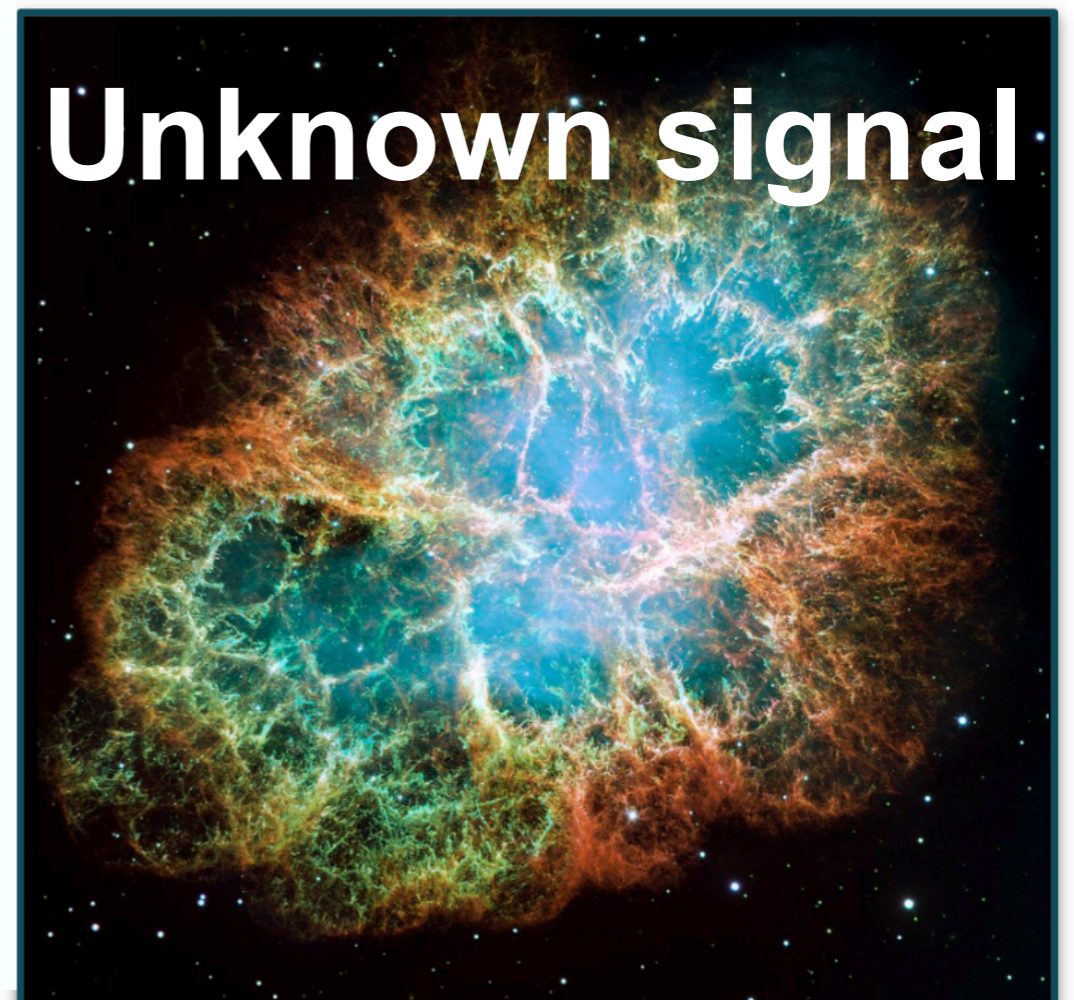
- Best sensitivity over the highest mass mergers
 - Consistent observations with running data
 - In the final stages of review to become a public merger

GWAK: Anomaly Detection

- Gravitational wave signatures can have unknown shapes
 - We can use AI based anomaly detection to identify these
 - AI based algorithm can runs on the fly (real-time)



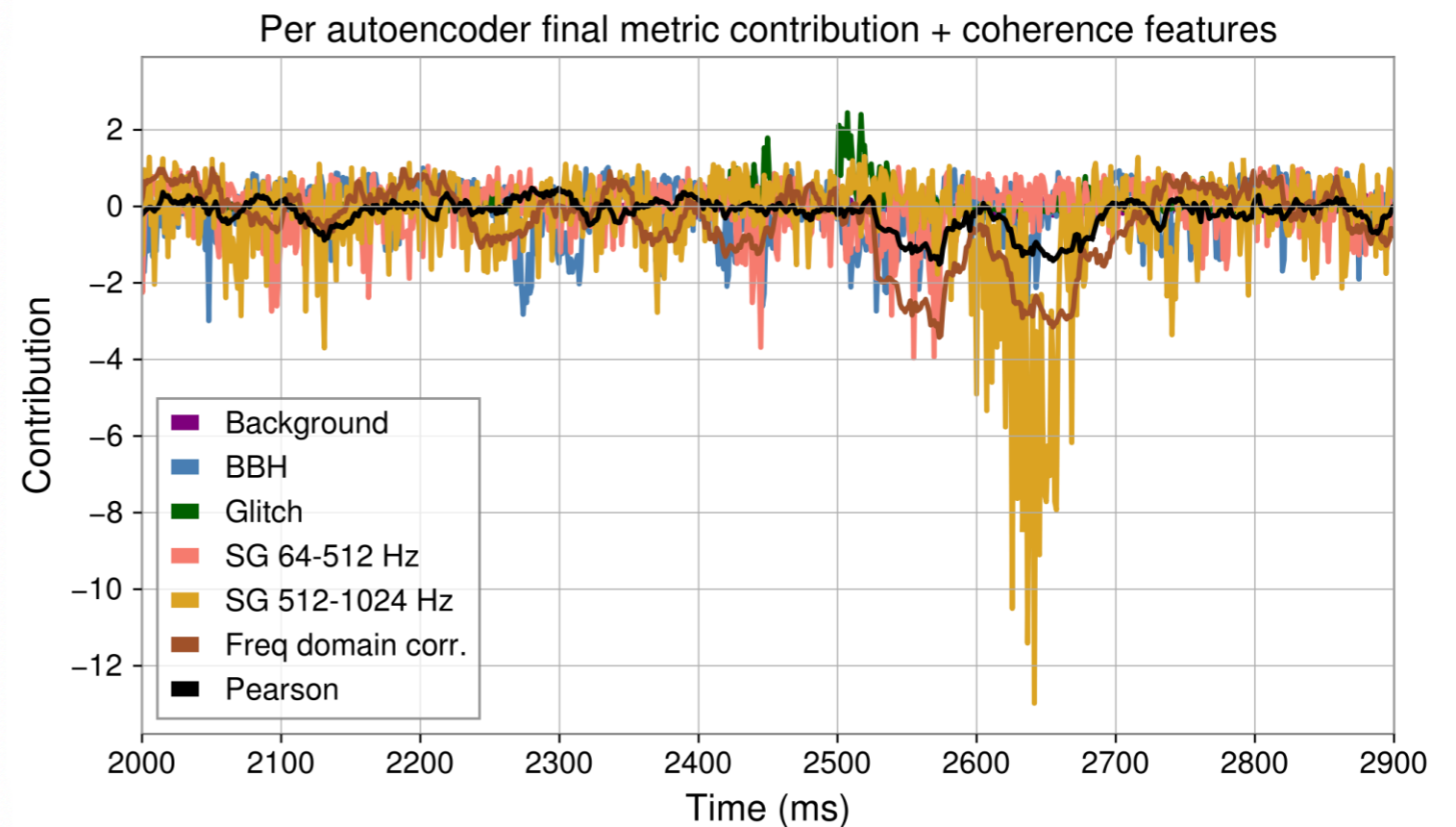
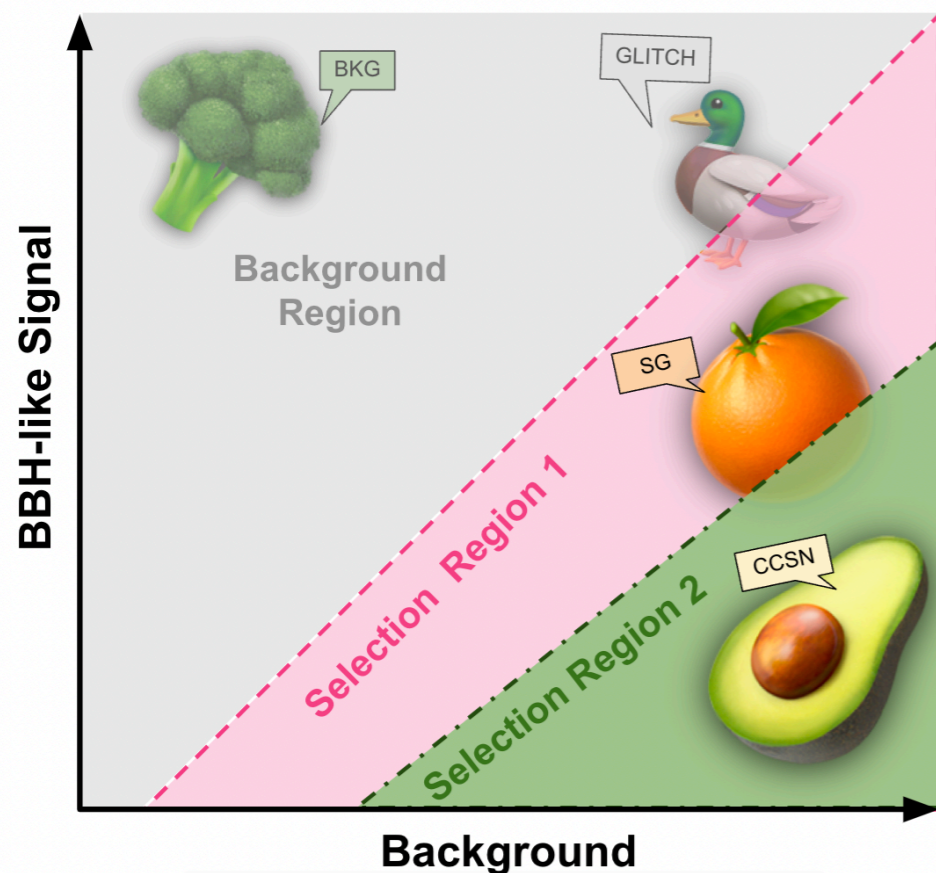
Core-collapse supernova (CCSN)



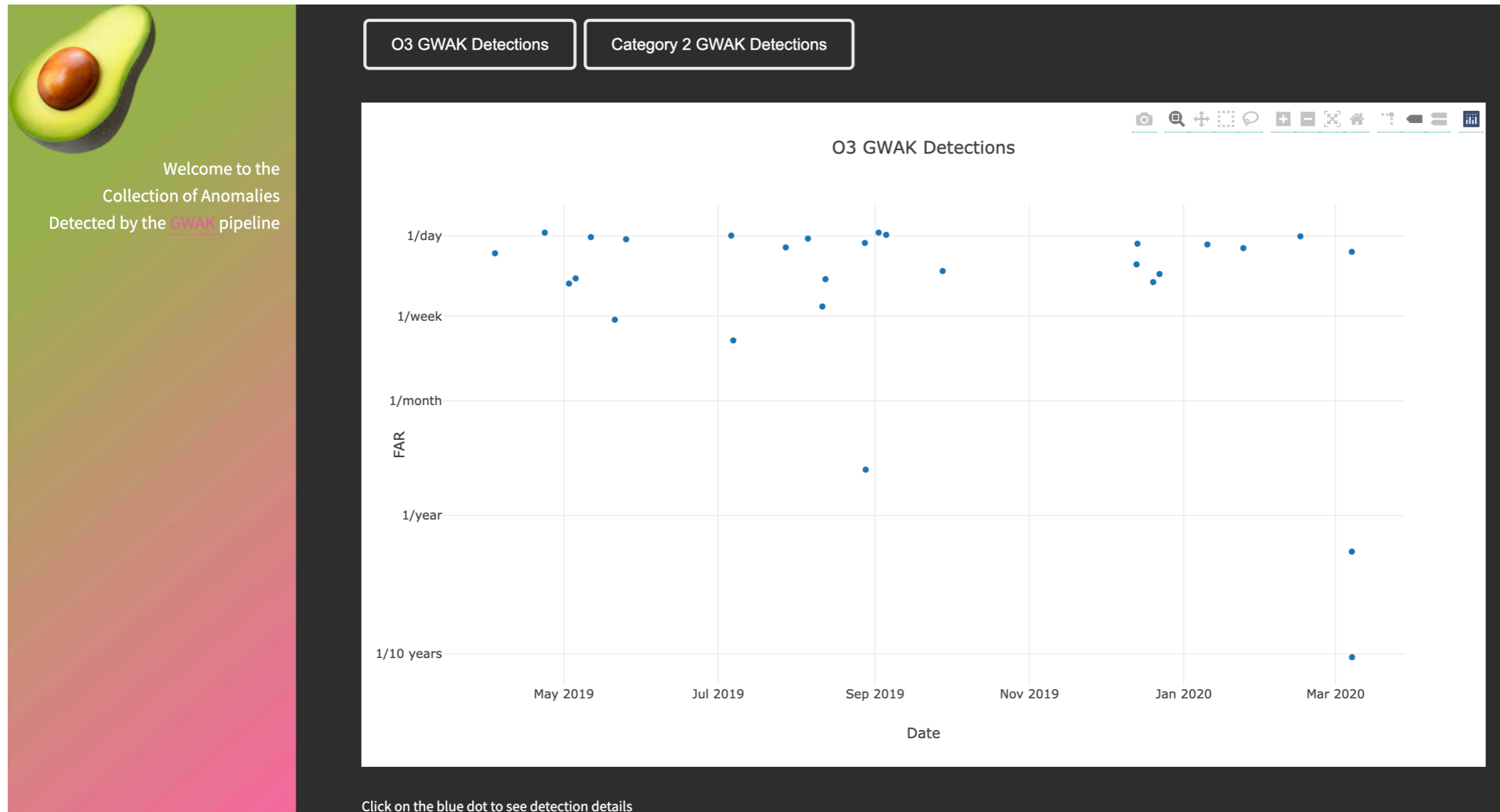
Embedded Space

- GWAK works by constructing a 11 dimensional space
 - Space presents a likelihood of a signal in a specific region
 - ▶ Likelihood on typical signals (bbh/sine gaussian)
 - Metric is constructed by a hyperplane in the space

2D GWAK Space

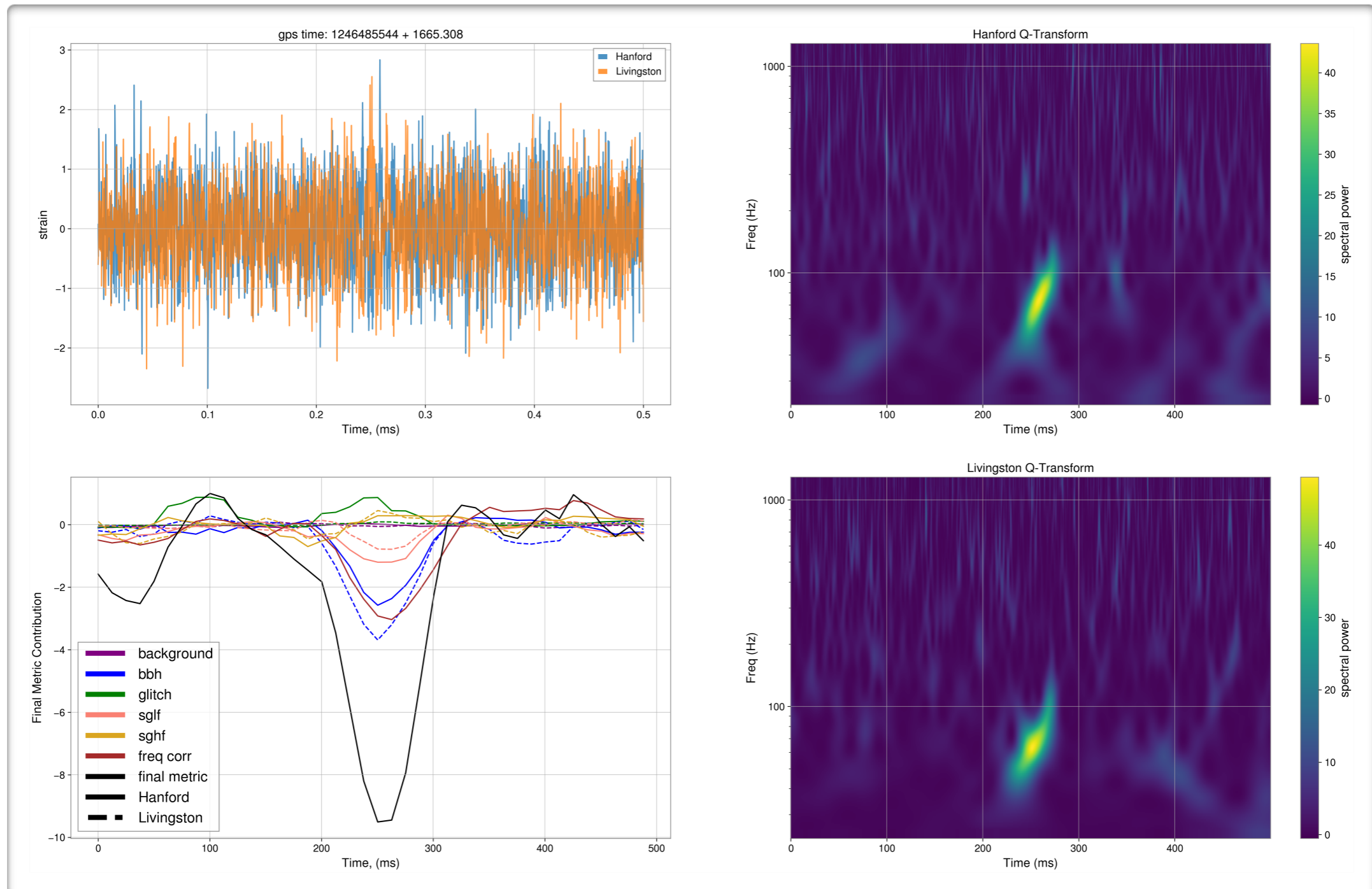


GWAK on Data



- Working to add this to real-time alerts
 - Already running internally in LIGO

GWAK on Data



- Working to add this to real-time alerts
 - Already running internally in LIGO

Anomaly detection ML²⁵ challenge



NSF HDR A3D3: DETECTING ANOMALOUS GRAVITATIONAL WAVE SIGNALS

48 PARTICIPANTS

110 SUBMISSIONS

Edit

Participants

Submissions

Dumps

Migrate

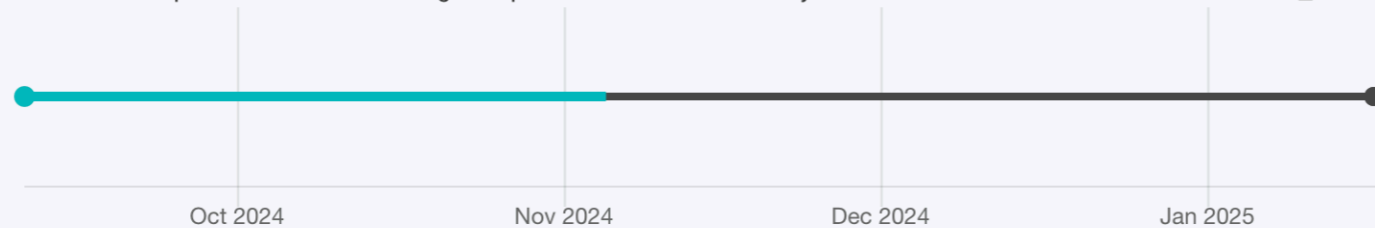
ORGANIZED BY: A3d3hdr

CURRENT PHASE ENDS: January 16, 2025 At 7:00 PM EST

CURRENT SERVER TIME: November 4, 2024 At 9:56 AM EST

Docker image: ghcr.io/a3d3-institute/hdr-image:latest

Secret url: https://www.codabench.org/competitions/2626/?secret_key=c95f242c-3a1d-4965-8a3a-b5b6bf51089d



Get Started

Phases

My Submissions

Results

Forum

?

Challenge Overview

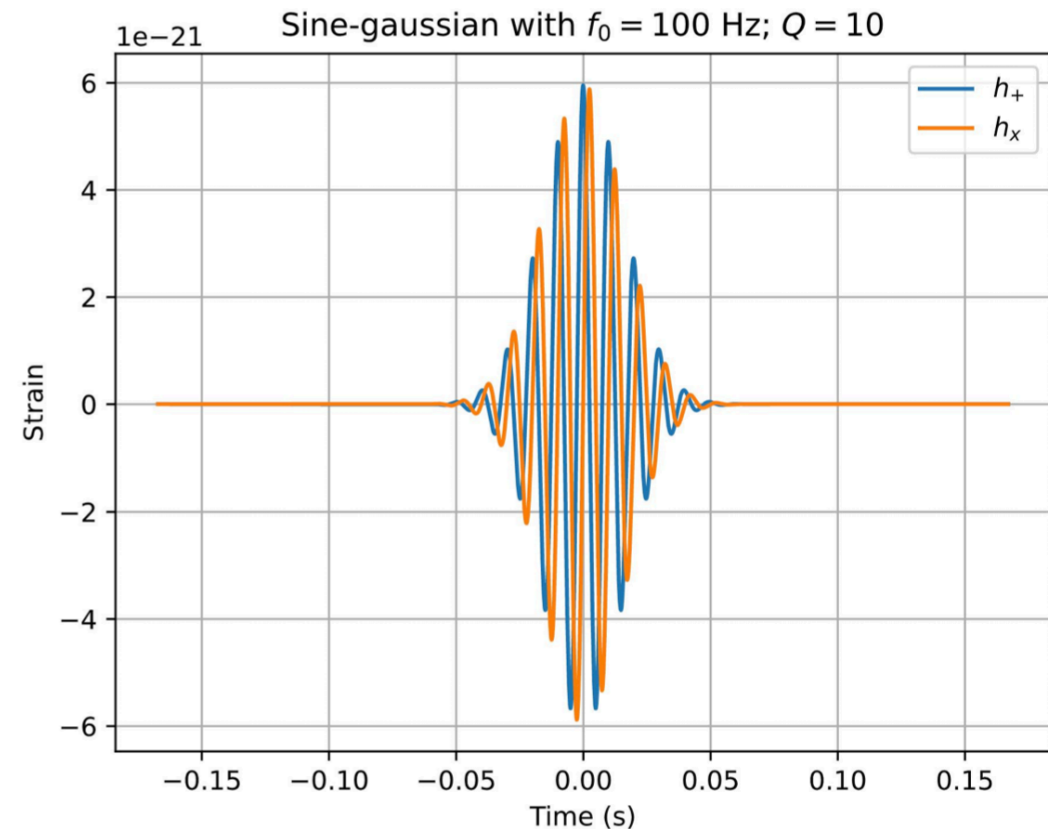
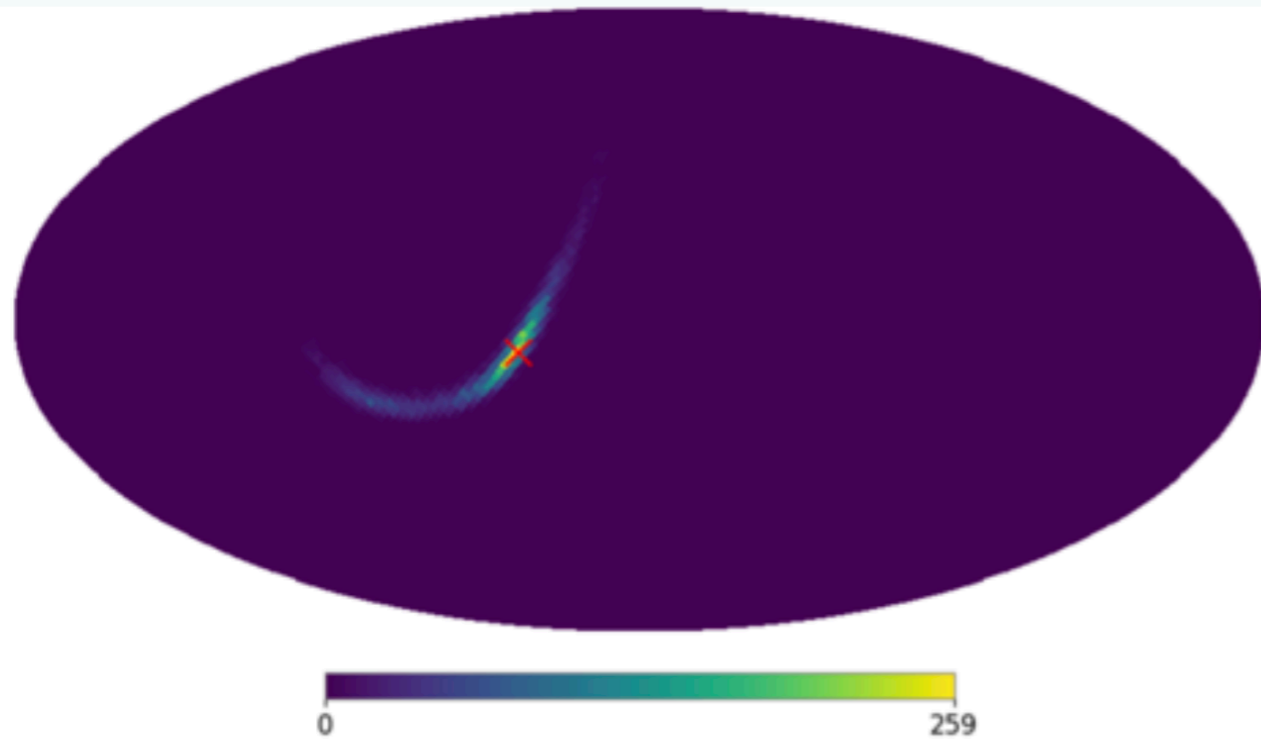
Datasets

Starting kit and sample

Overview

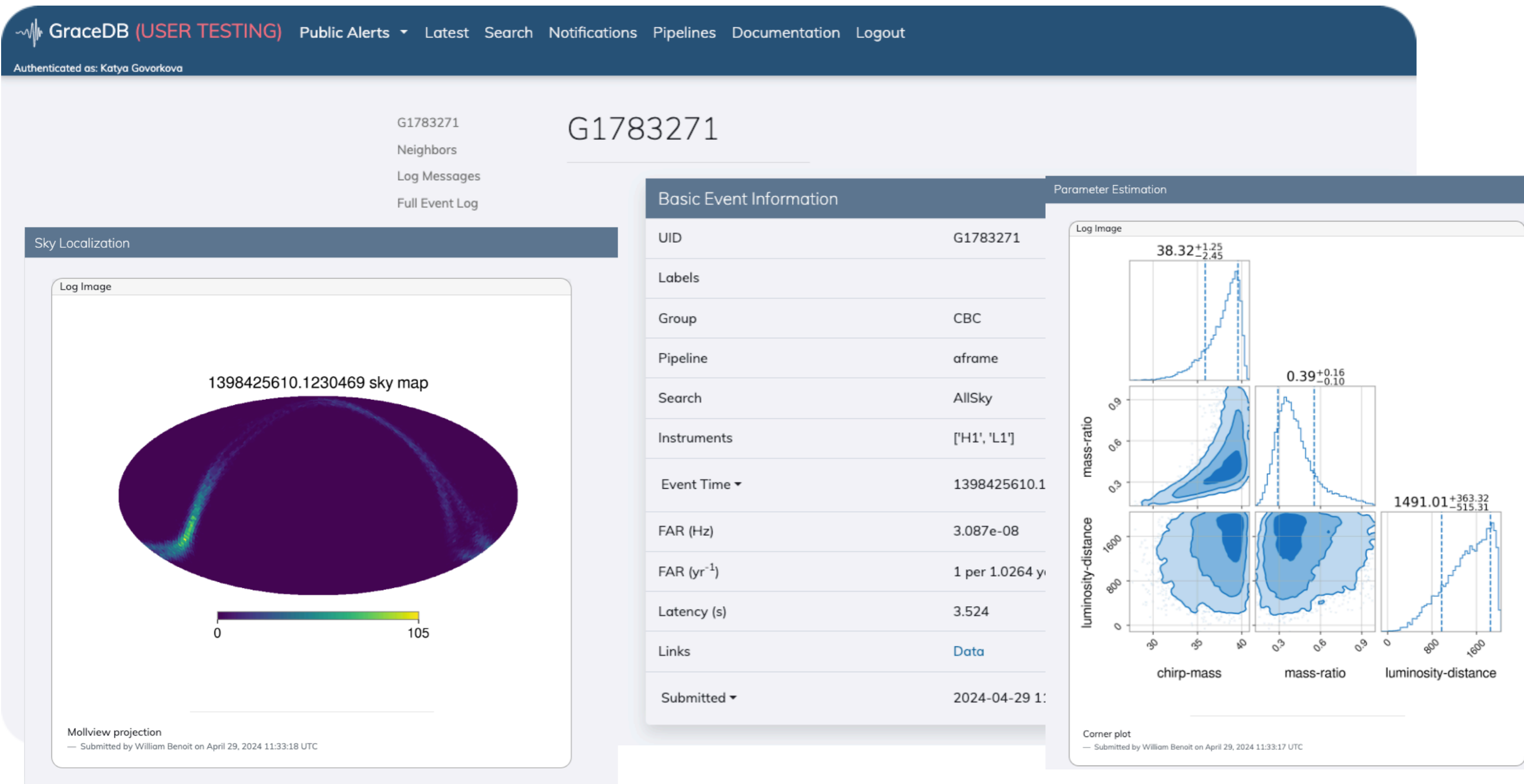
<https://www.nsfhdr.org/mlchallenge>

Parameter Estimation



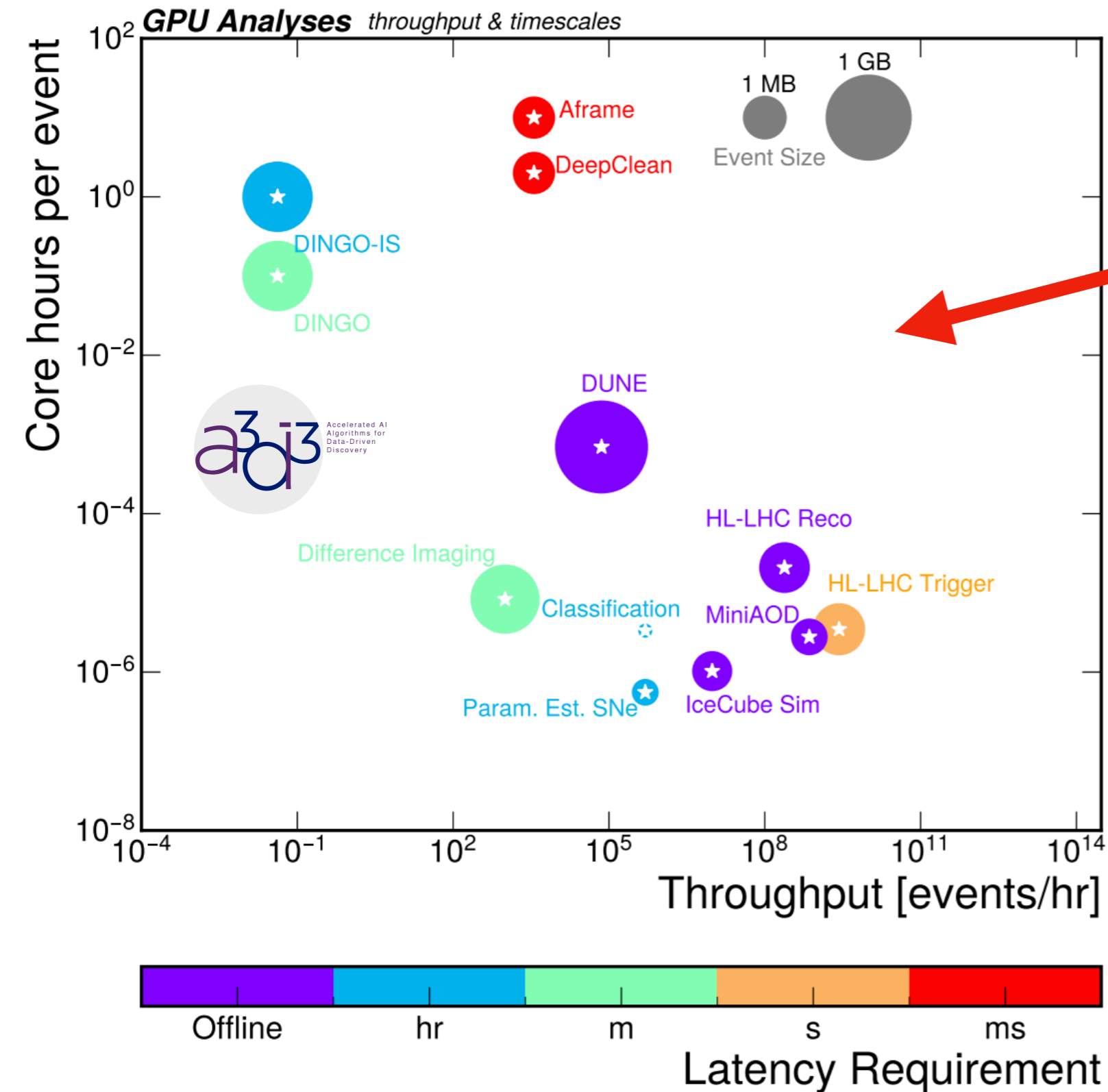
- Perform sky localization and gravitational wave parameters
- Perform fast parameter estimation using likelihood free inference
 - Normalizing flows embed broad knowledge of waveforms
 - Customized embedding to ensure compressed info
 - Parameter estimation done within seconds (or potentially less!)

Building a Pipeline



Pipeline is up and running

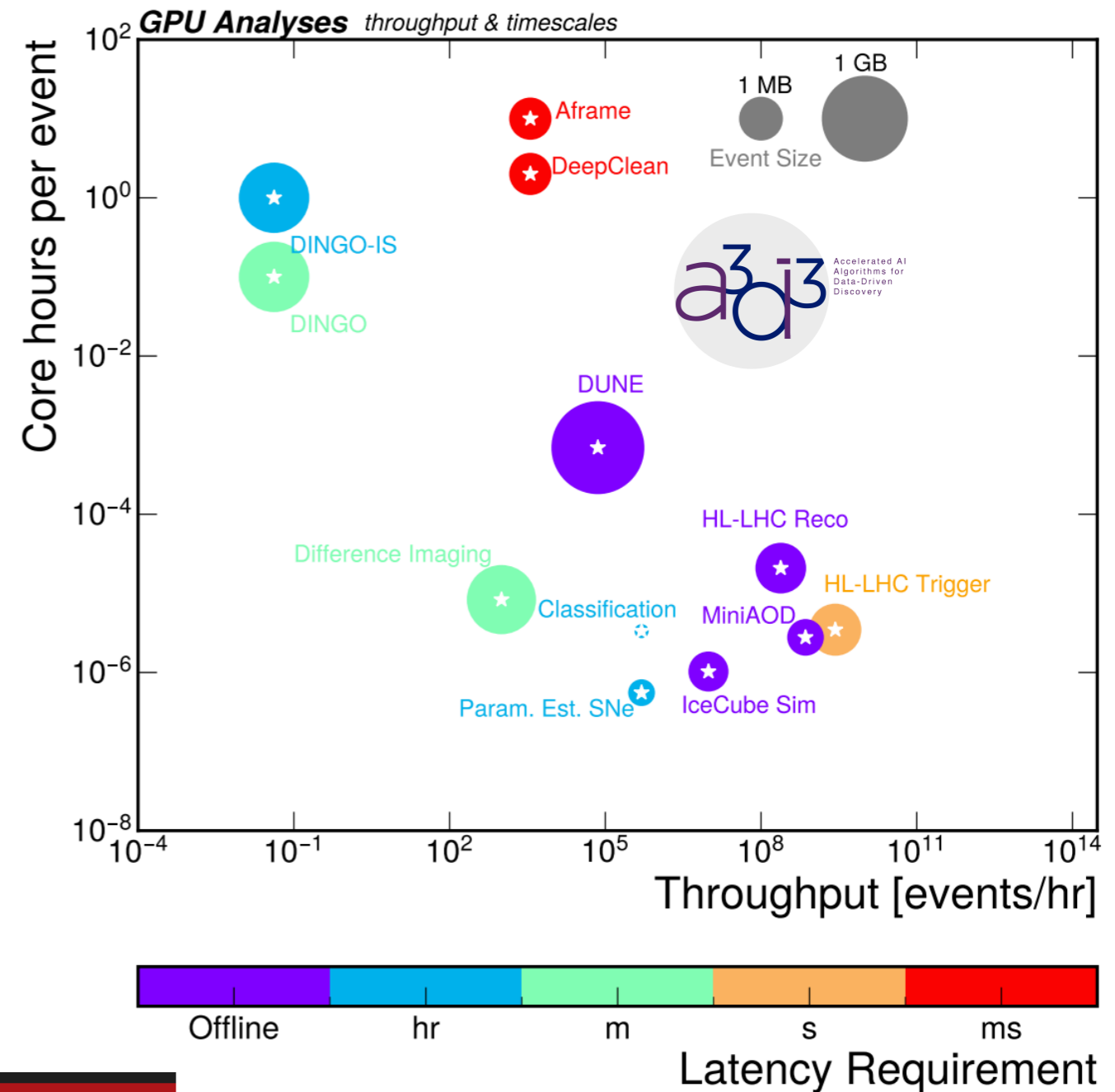
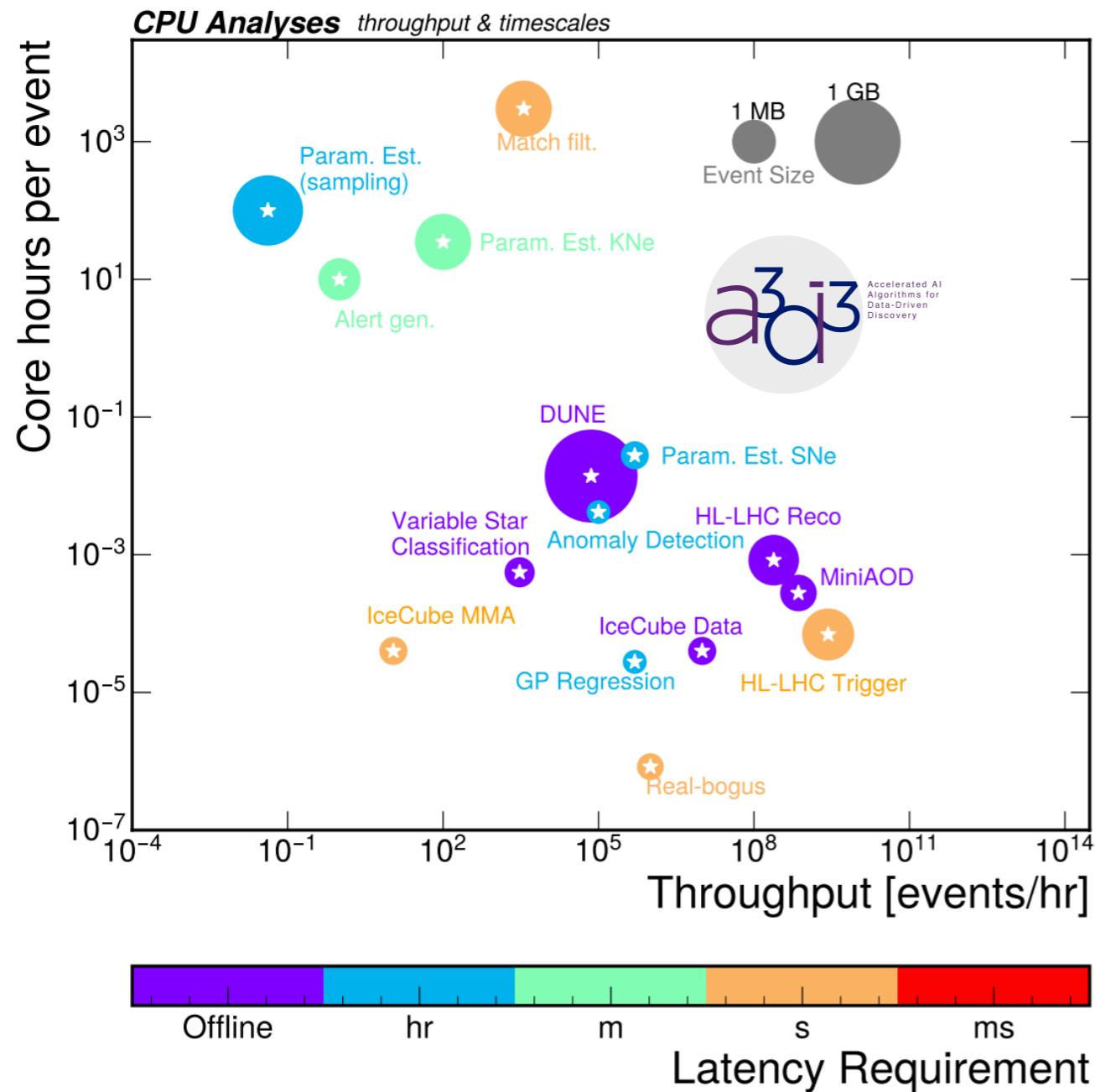
GPUs for ML



For large algorithms parallelizability and shear compute is unprecedented

Major demand emerging in many domains

Computing Demands



arXiv > hep-ex > arXiv:2306.08106

High Energy Physics – Experiment

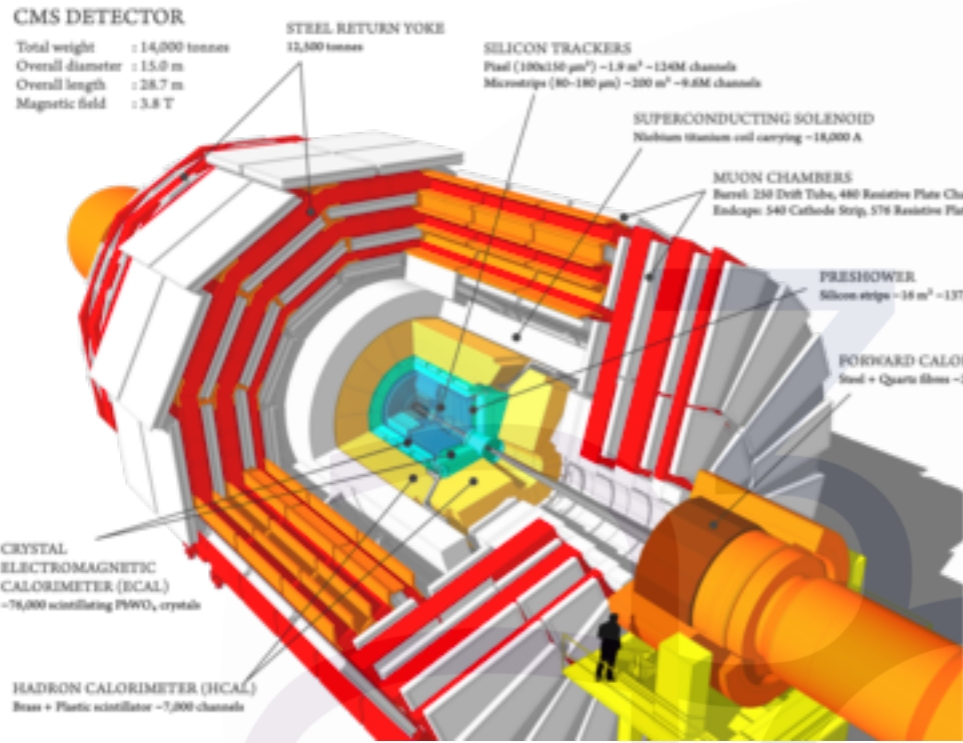
[Submitted on 13 Jun 2023]

Applications of Deep Learning to physics workflows

Arxiv: 2306.08106

Have a whitepaper outlining
Inference Workflows Demands

Building an Ecosystem



DEEP UNDERGROUND
NEUTRINO EXPERIMENT
 DUNE inference as-a-service
[arxiv:2301.04633](https://arxiv.org/abs/2301.04633)

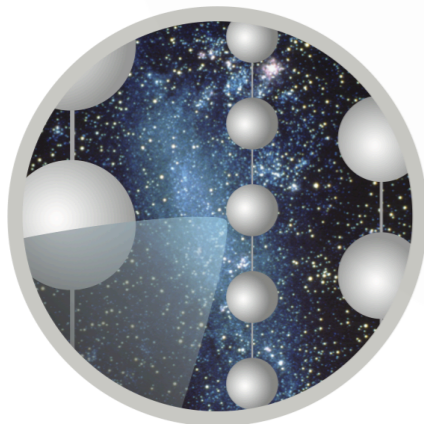
CMS Inference-as-a-service
 (arxiv:2402.15366v1)



Tracking-as-a-service



LIGO inference-as-a service
[arxiv:2108.12430](https://arxiv.org/abs/2108.12430)



ICECUBE

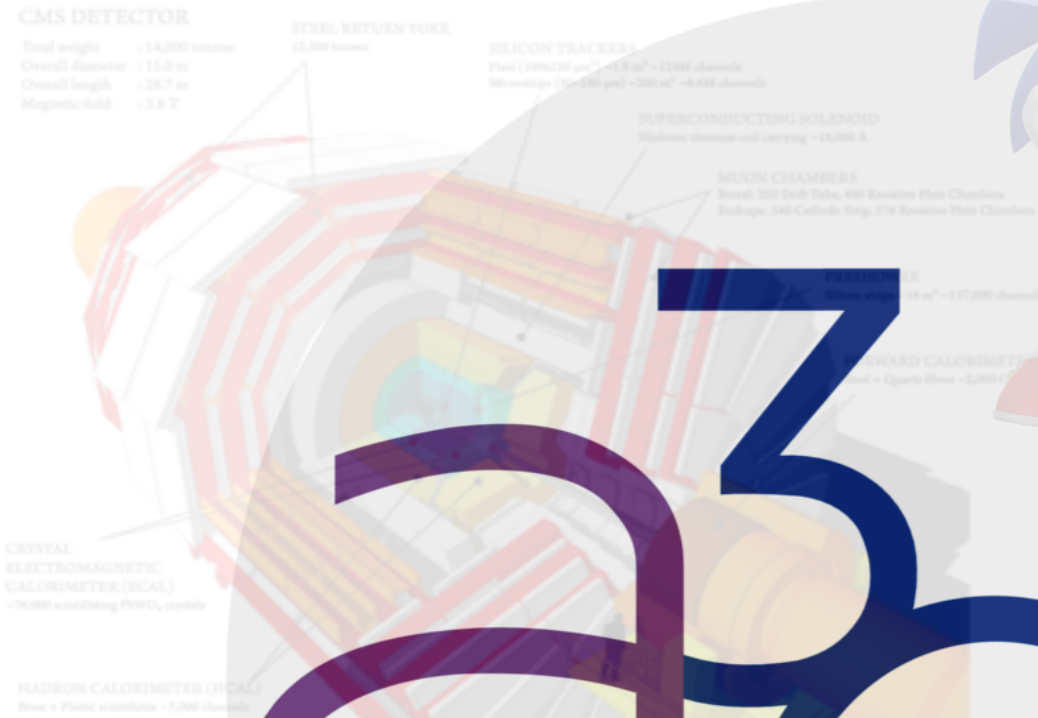
Starting to Investigate

- Super-SONIC

- Landing page of SW for Iaas for Physics

shared toolkit across many experiments

Building an Ecosystem



Accelerated AI Algorithms for Data-Driven Discovery
[DUNE Inference-as-a-service arxiv:2301.04033](https://arxiv.org/abs/2301.04033)

[CMS Inference-as-a-service \(arxiv:2402.15366v1\)](https://arxiv.org/abs/2402.15366v1)



[LIGO inference-as-a service arxiv:2108.12430](https://arxiv.org/abs/2108.12430)



ICECUBE
 Starting to Investigate

- Super-SONIC

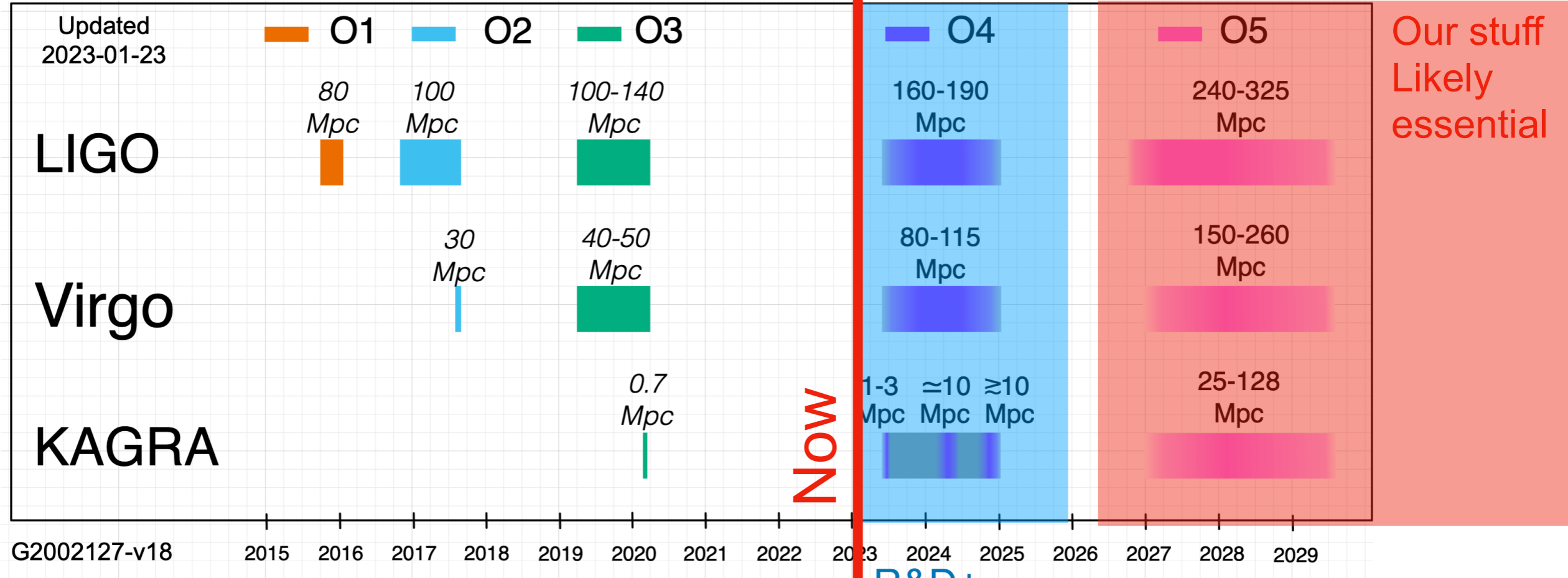
- Landing page of SW for iaas for Physics

shared toolkit across many experiments

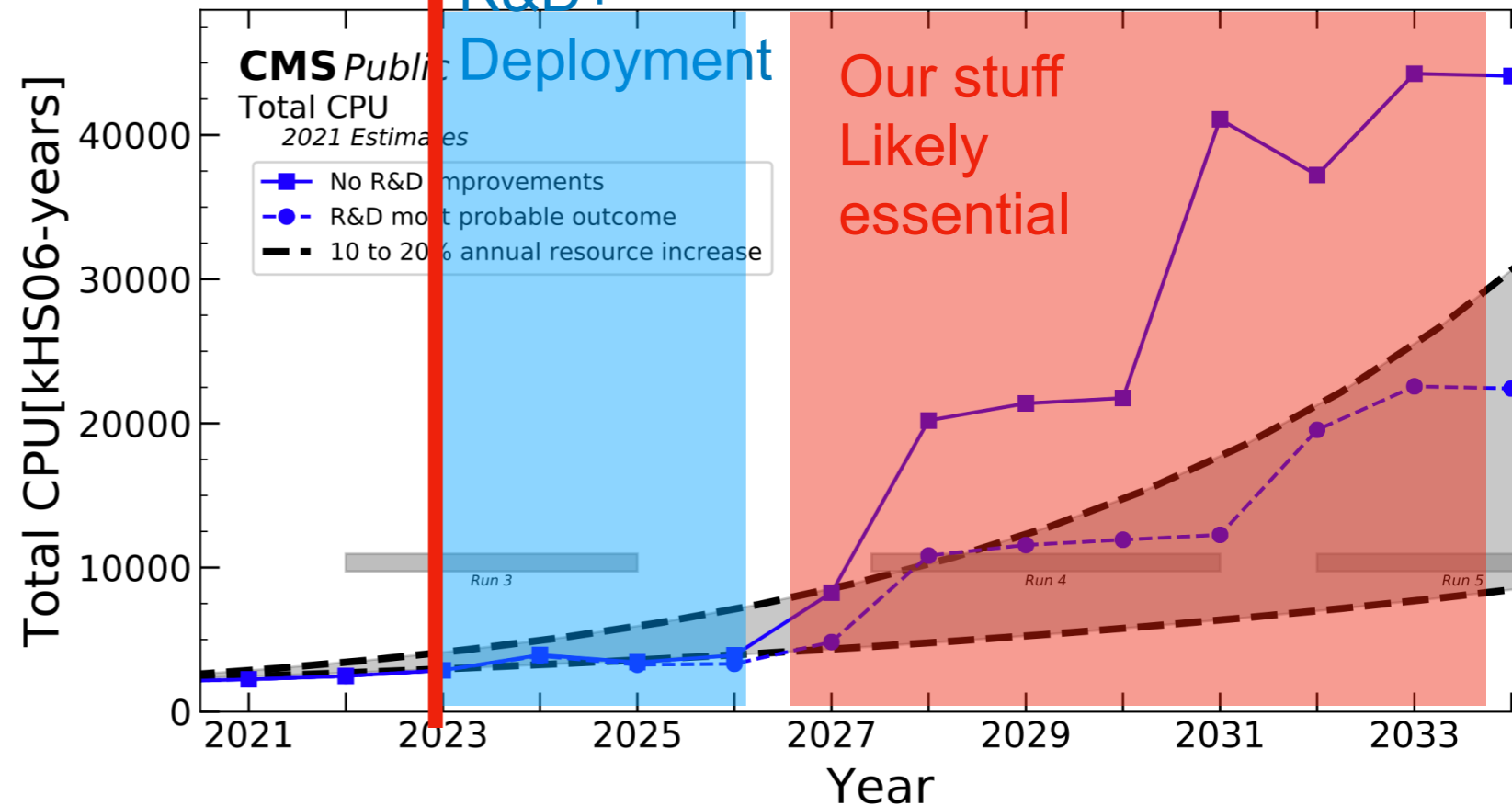
What is Critical?

- We would like to highlight commonalities across domains
 - **Computing demands**
 - ▶ Critically connected infrastructure for ML science deployment
 - ▶ Inference differs from training → **Efficiency is Key**
 - **Software Stack**
 - ▶ With all ML algorithms aim for a set of core software tools
 - ▶ Containerization: Apptainer/Kubernetes/...
 - **ML Problems**
 - ▶ Awareness of the diversity of problems is critical (Not just LLM)
 - ▶ **Highlighting the similarity across scientific domains is critical**

Timelines ³³



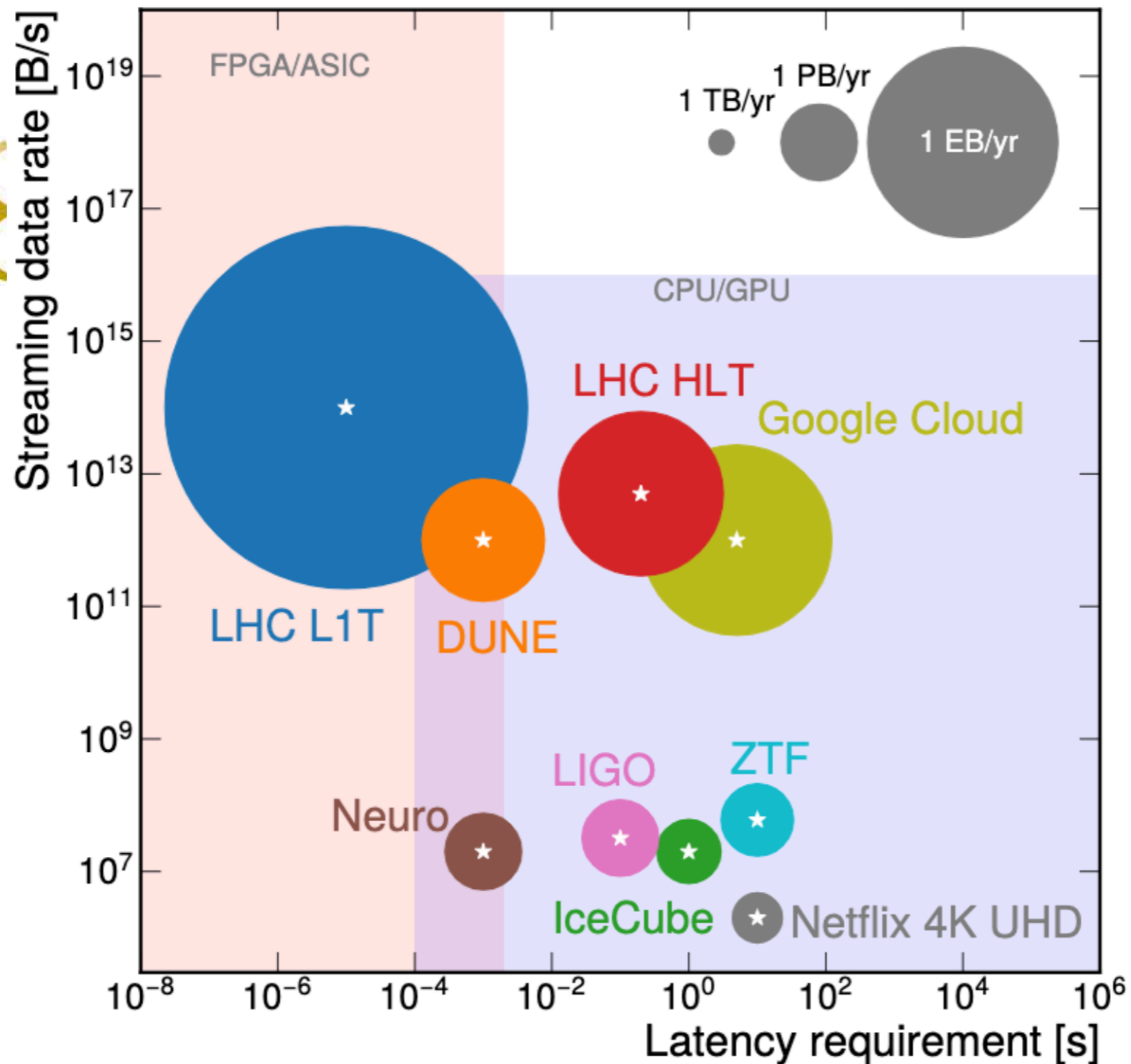
DUNE timeline and various astro timelines (Rubin/LSST) Should also figure in our overall schedule



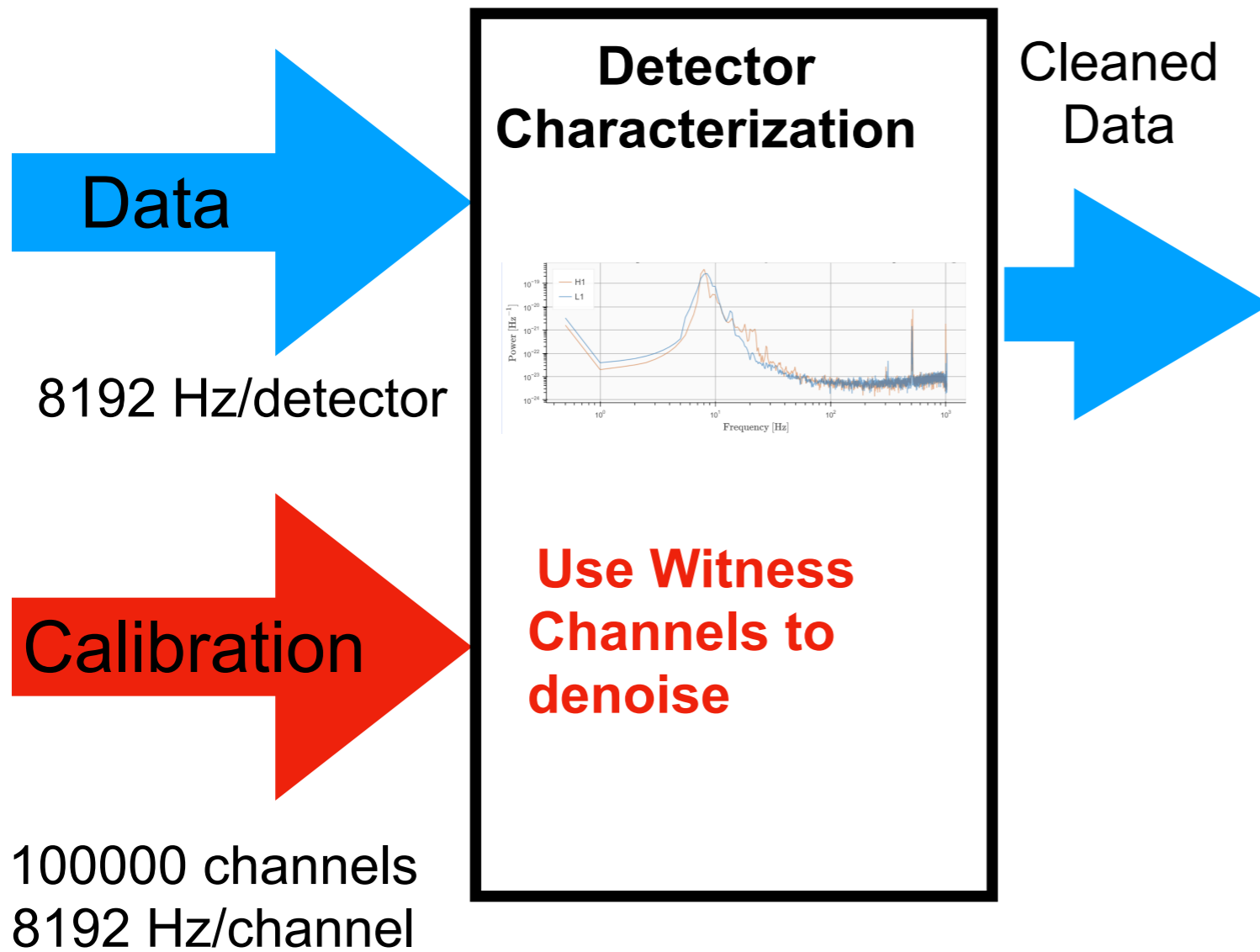
Recap

- We are building a real-time GW pipeline using ML
 - Latency and throughput are critical element in the design
 - Effective integration of heterogeneous compute is critical
- Our pipelines are up and running internal and soon externally
 - We see this as a path to be a main pipeline for future (O4) running
 - We encourage many others to build on our toolkit
- While our focus is on GW the tools here are broadly applicable
 - SBI toolkit has already been adapted for Kilonovae
 - In discussions for adapting parts of the toolkit for other time series
- Looking to expand the scope of our work under ML4GW toolkit
 - More collaborations within GW and beyond

Thanks!



LIGO Data Workflow



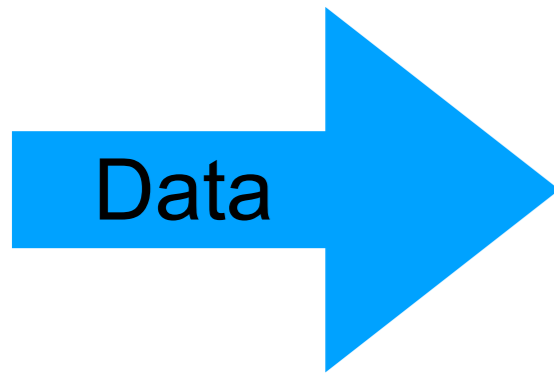
Challenge is to run this in real-time
Roughly 1 PB per year

LIGO Data Workflow



Challenge is to run this in real-time
 Roughly 1 PB per year

LIGO Data Workflow



8192 Hz/detector

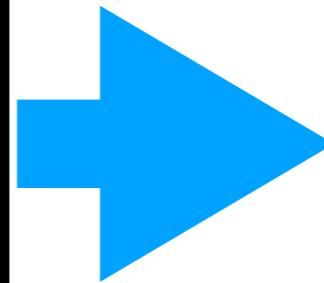


100000 channels
8192 Hz/channel

Detector Characterization

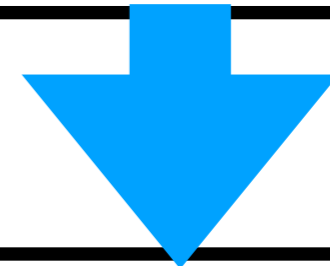
Use Witness Channels to denoise

Cleaned Data



Event Detection

Black Holes/Neutron Stars/???

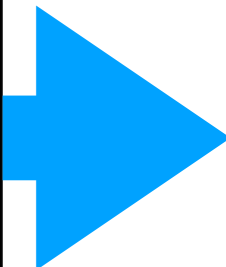


Event

Event Characterization

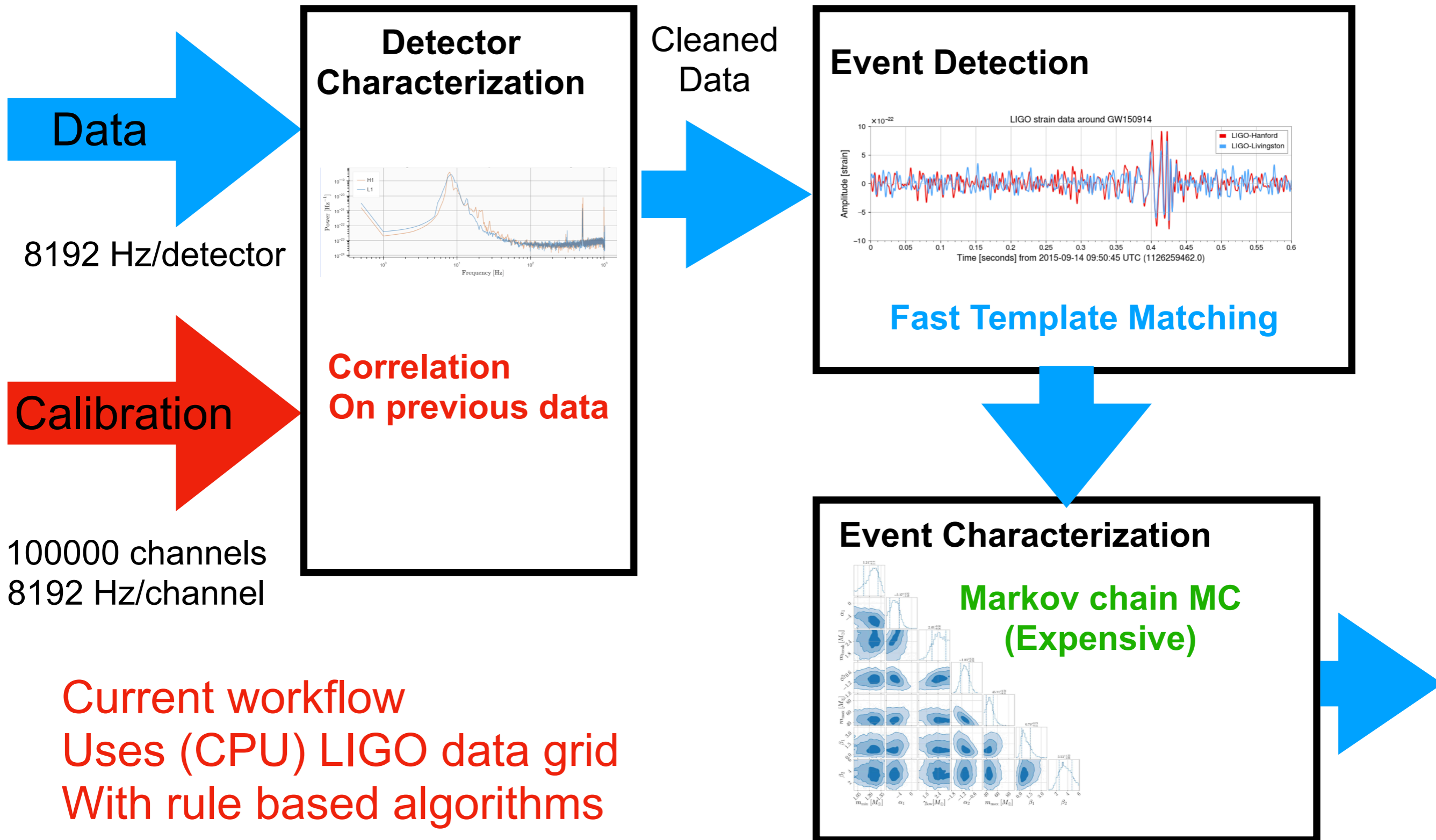
**Parameter Estimation/
Sky Localization**

Alert



Challenge is to run this in real-time
Roughly 1 PB per year

Current Workflow



Our Upgraded Workflow

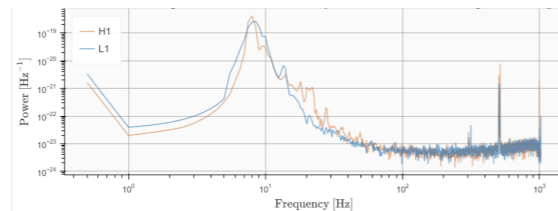
Data

8192 Hz/detector

Calibration

100000 channels
8192 Hz/channel

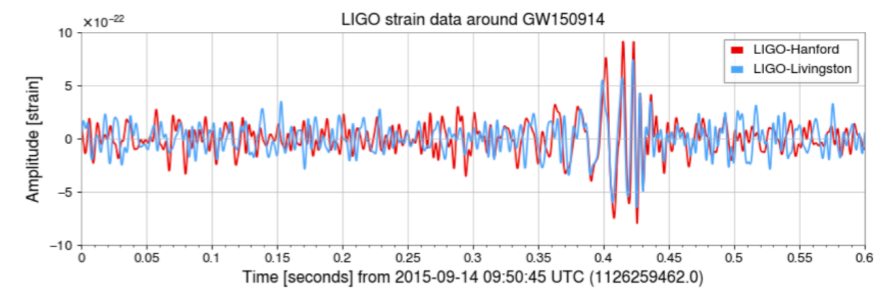
Detector Characterization



DeepClean
NN based AE

Cleaned
Data

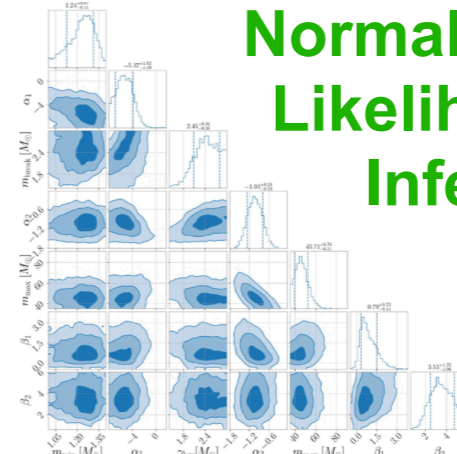
Event Detection



NN-based Algos
A-Frame/Neutron Star/GWAK

Event Characterization

Normalizing Flow
Likelihood Free
Inference



ML Based Workflow

Almost 100% GPU based

Only a handful of GPUs for real-time

Computing Resources

LIGO Data Grid

- LIGOs computing ecosystem of mostly CPU resources
- Limited GPUs, workloads not scalable
- GPUs resources are not sufficient to support large ML workflows

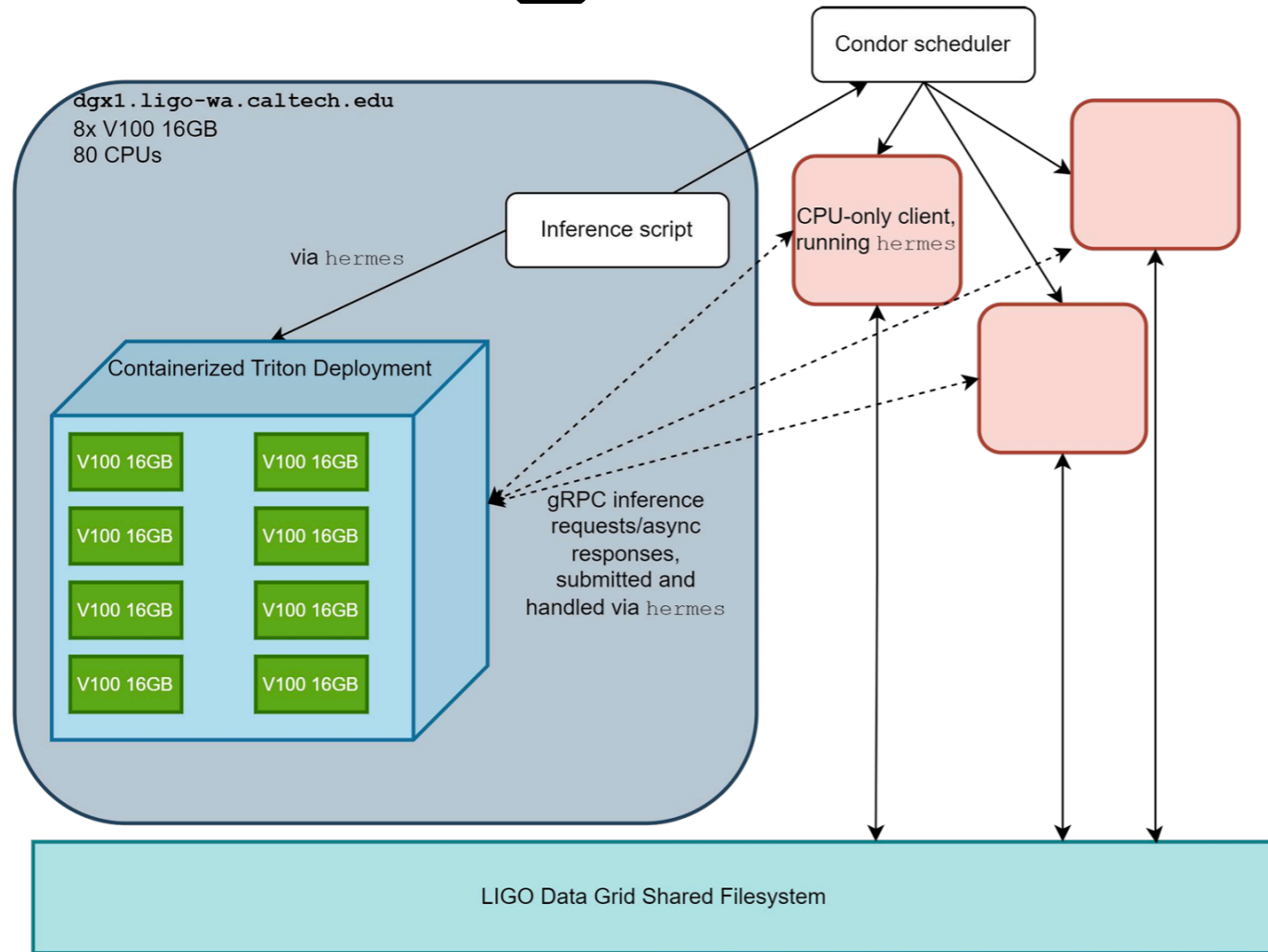
Nautilus HyperCluster

- Collection of computing clusters containing 1000s of GPUs
- Containerized workloads
- Easily scalable with Kubernetes
- With Kubernetes infra, can easily migrate to other cloud resources



NRPs GPU resources makes it possible for us to scale to full analyses

Large Scale deployment



Throughput of 3800 s'/s achieved

100 years of data in 2 weeks!

10x speed up from conventional GPU + More possible (FP16)

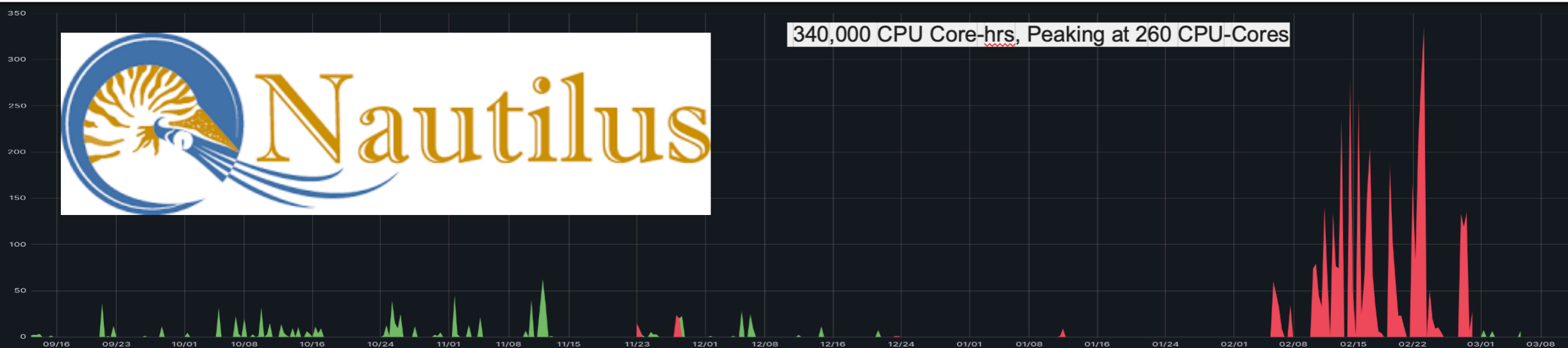
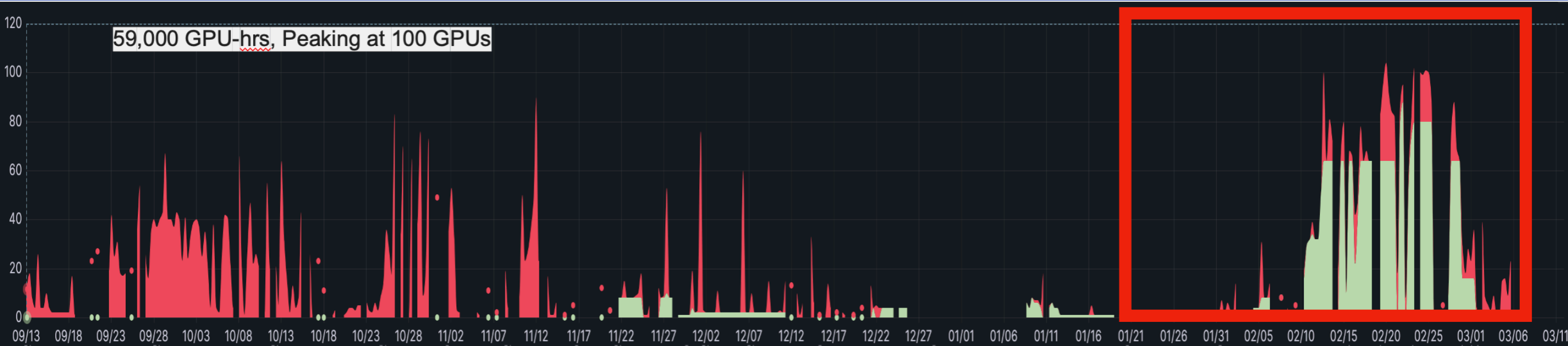


Ability to scale processing to large clusters with k8s

Current Usage

Phil Harris, MIT

LIGO GPU/CPU Usage Per Day, Last 6 Months



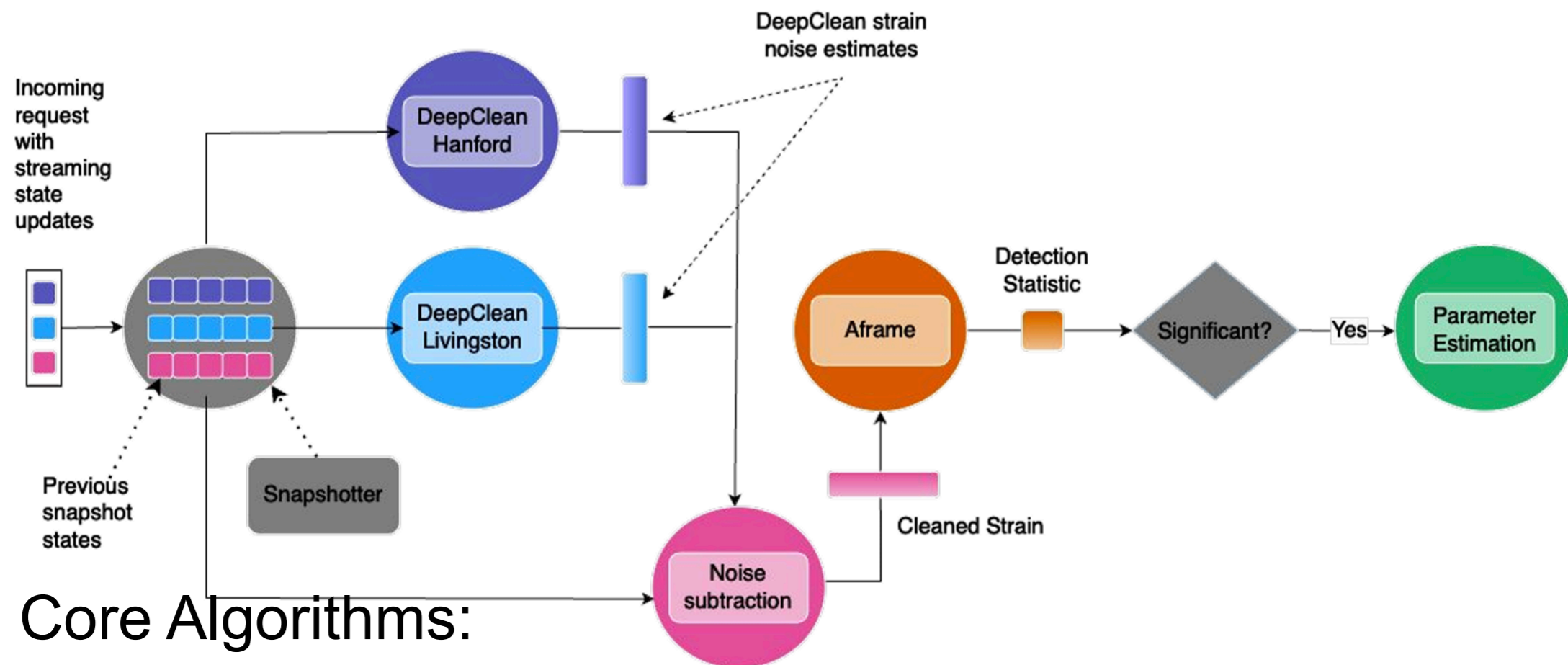
NRP NATIONAL RESEARCH PLATFORM

Namespace [osg-ligo](#), [bbhnet](#)

NRP Nautilus

Algorithm Training of a black hole merger algorithm

What Algorithms exist?

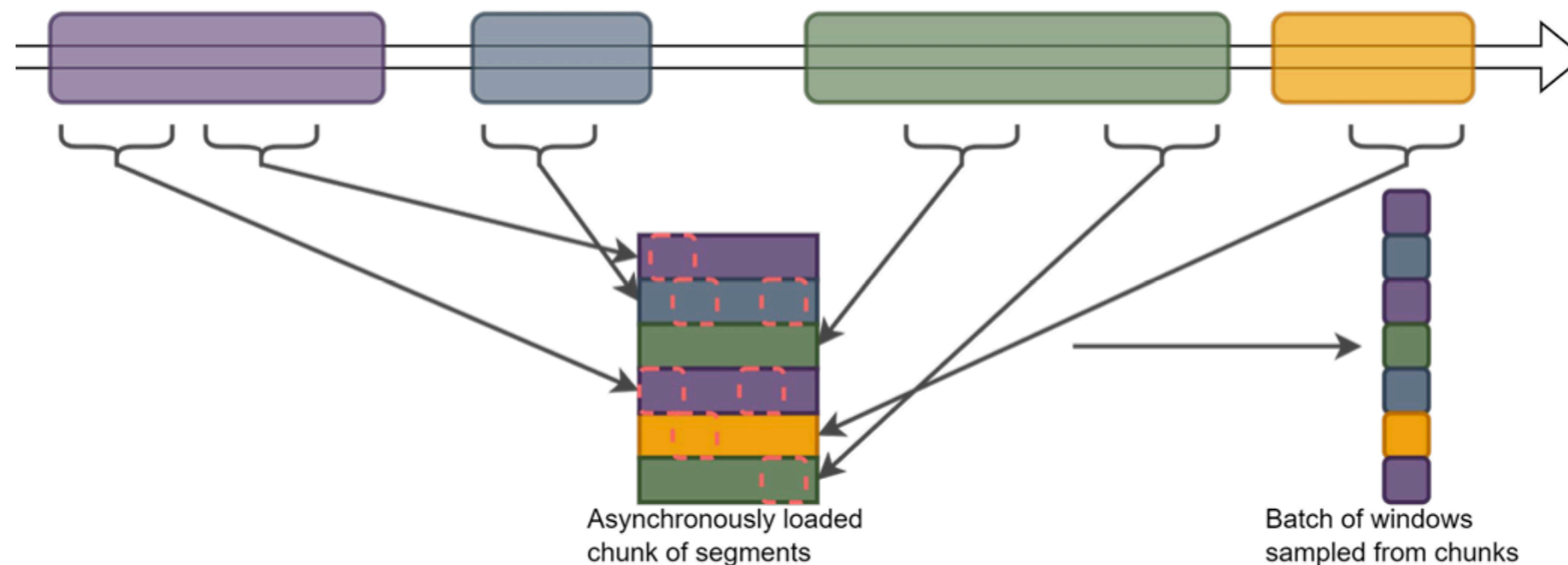


- Core Algorithms:

- Detector Denoising => Deep Clean
- Black Hole Merger detection => A-frame
- AI based anomaly detection => GWAK
- Neutron Star Merger detection
- Parameter Estimation

Time Series Caching

Transitioning to larger datasets



Chunked loading of background data

GPU batching can be enabled through chunked data loading
Parallel time series processing allows fast training

Computing Challenge

- GW data is time series data
 - Our toolkit targets critical time series setup
- Vanilla ML processing workflow
 - Load ML model and shove data through it :
 - ▶ 512s of data per second (s'/s) on 16 GB V100 GPU
 - ▶ 1 year of data is 17hrs of computing
 - ▶ 100 years (needed for analysis) is 70 days
 - Too slow to iterate on ideas
- Our workflow utilizes optimized schemes to avoid these issues

ML Challenges

- Aiming to build a website hosting Scientific ML Challenges

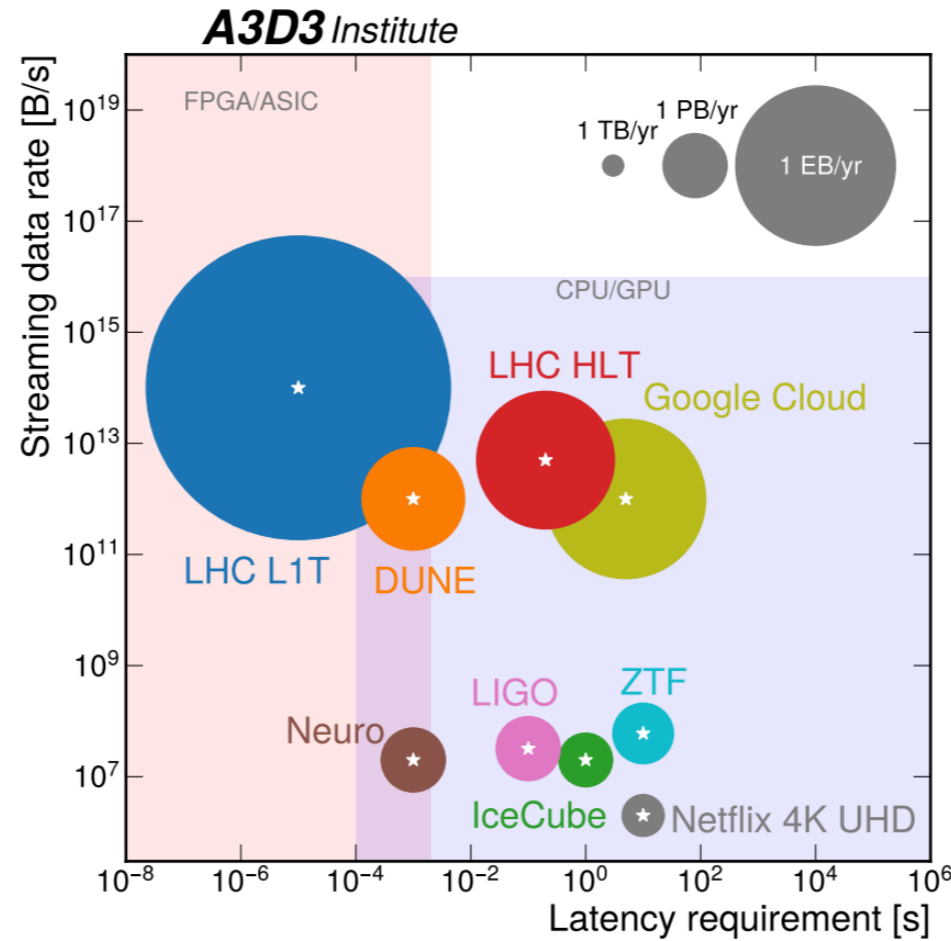
Connecting with ML Commons



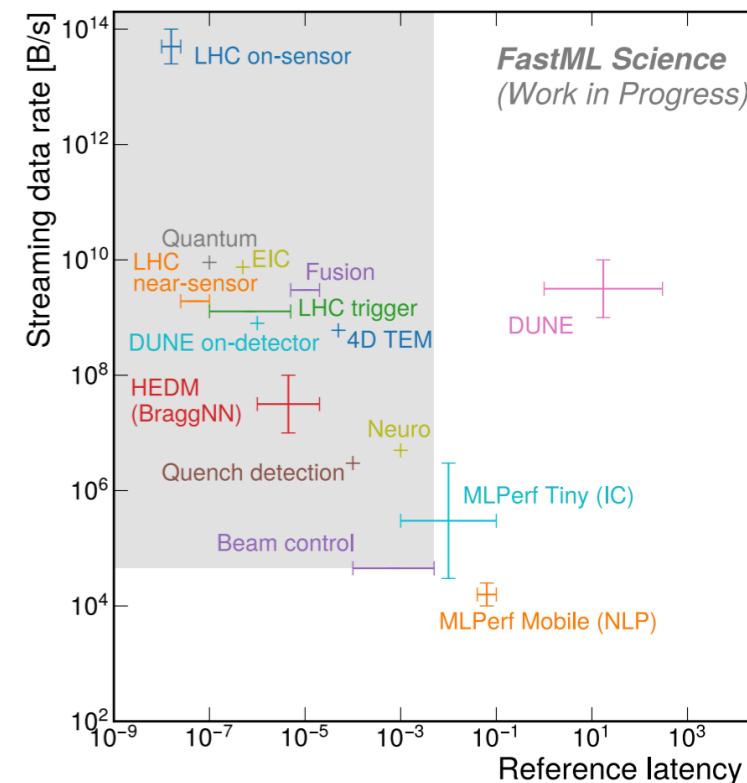
ML Commons

Machine learning innovation to benefit everyone.

MLPerf Tiny Inference
A benchmark suite for ultra-low-power tinyML systems



Connecting With Hardware



Would like to highlight Criticality of Scientific Problems

Support from NSERC FAIRUniverse

A Vision

- Can we align science across ML Challenges?
 - Details [here](#) following C. Herwig, N. Tran (Fermilab)

		Scientific Moonshots		
		Domain A	...	Domain N
AI thrusts	AI - 1: Real-time	Benchmark 1A		Benchmark 1N
	AI - 2: Control			
	AI - 3: Autonomous			
	AI - 4: Foundation			
	AI - 5: Generative	Benchmark 5A		Benchmark 5N

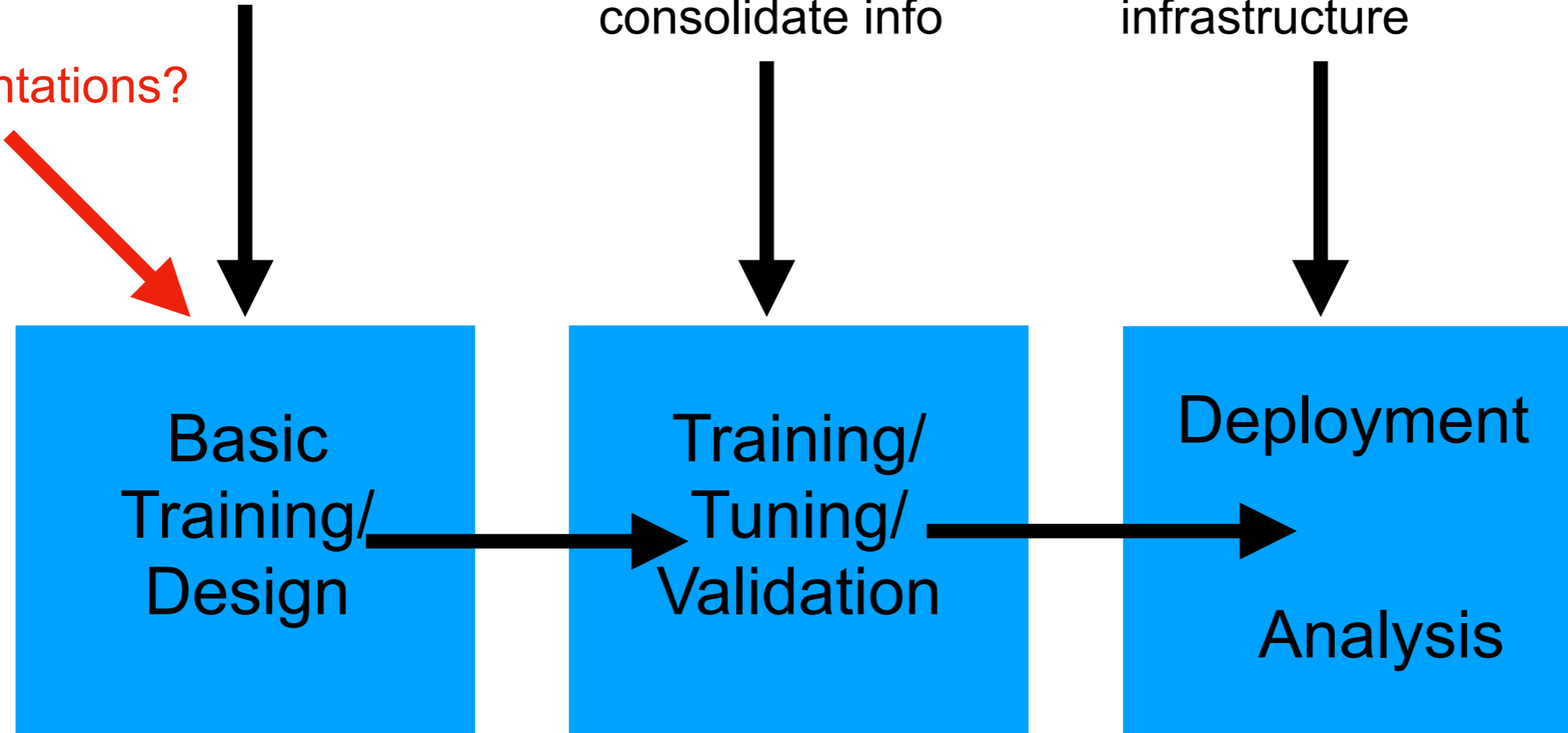
Anatomy of an Algo

Good Data/Simulation
For training

Critical software
tools that
consolidate info

Software/hardware
deployment
infrastructure

Augmentations?



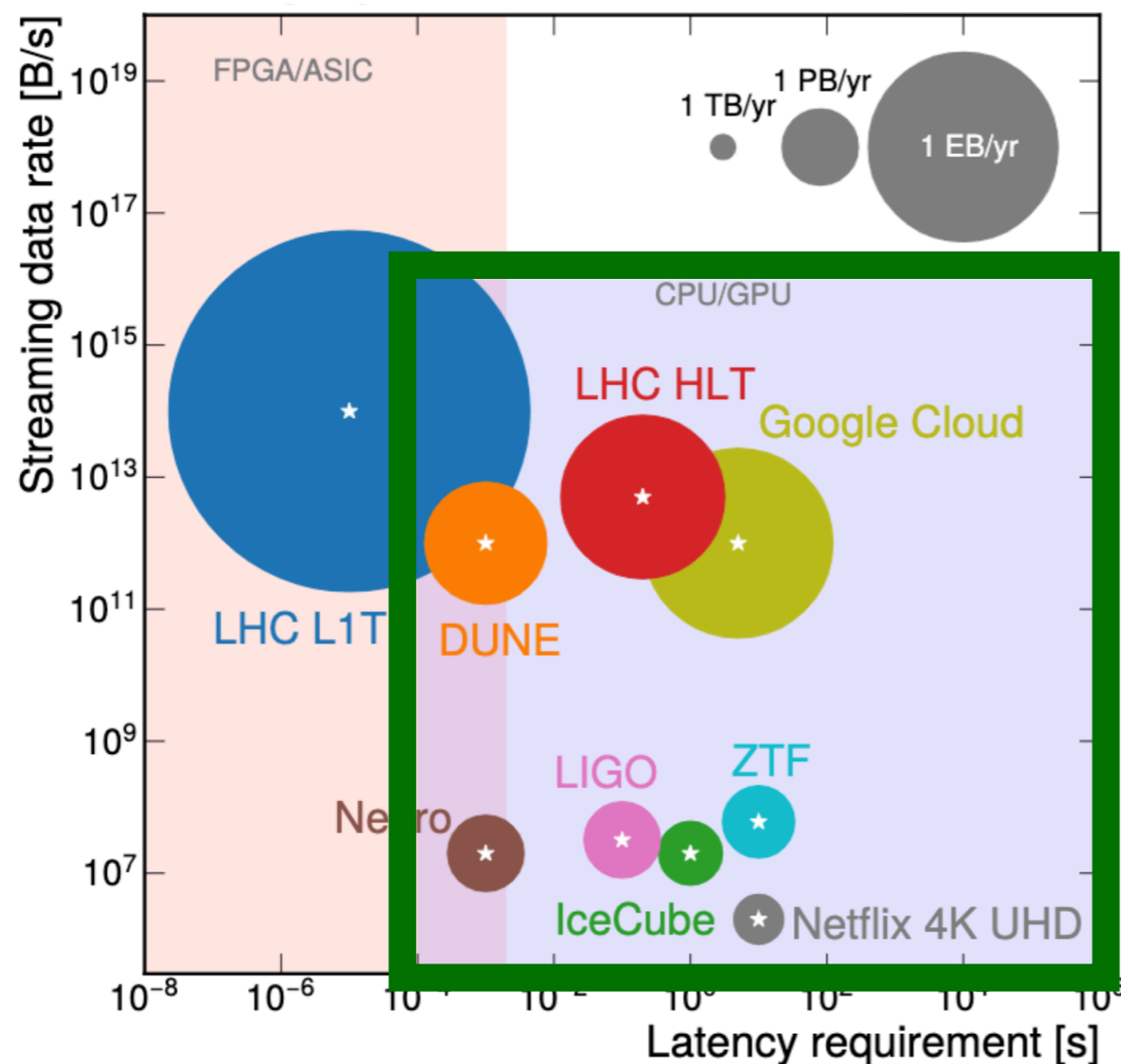
Local
GPU

NRP

NRP

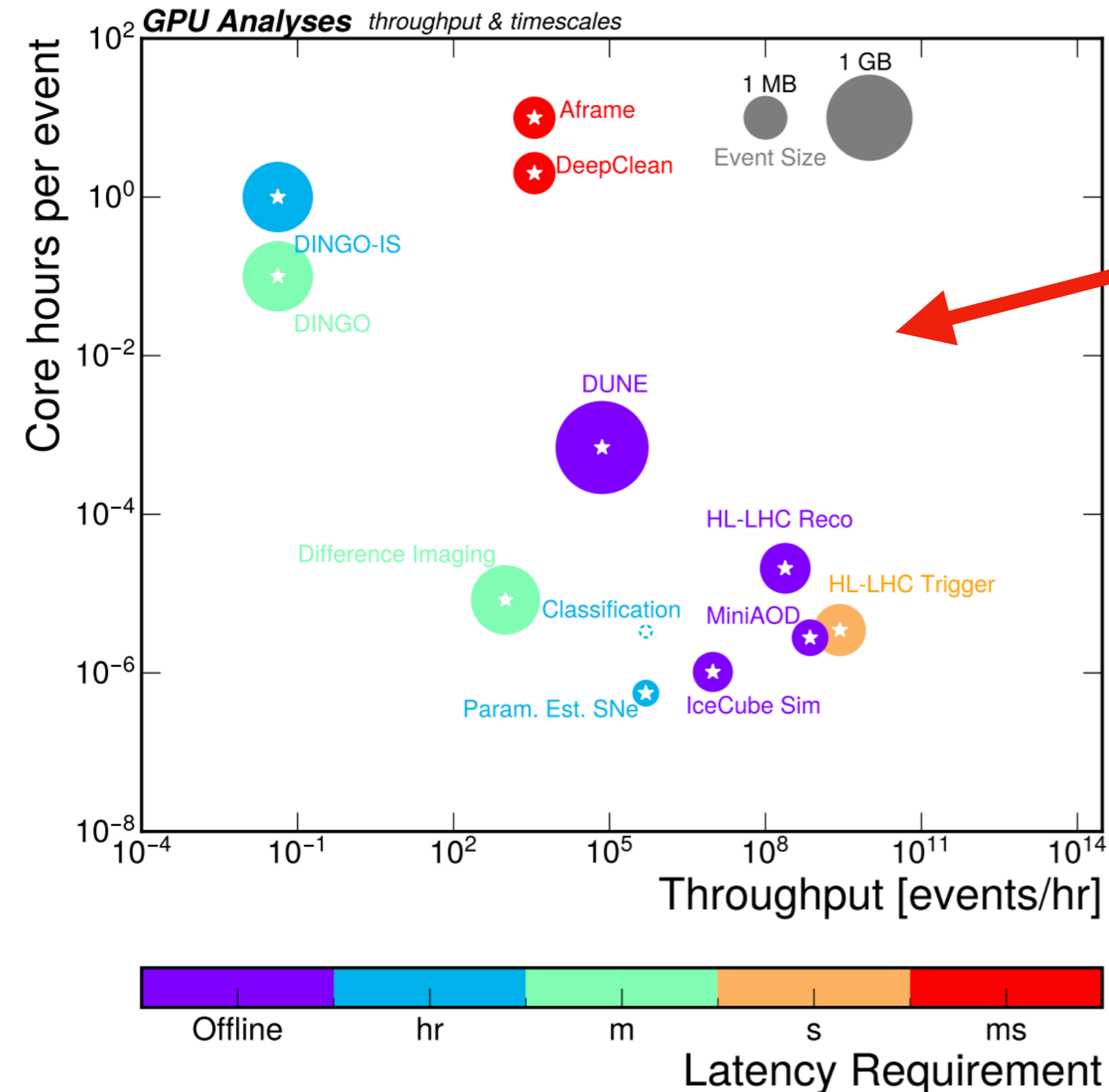
What computes are here?

- Within the FastML Community there is a broad range
 - We often try to characterize this range by customization
 - Low Latency and Low Power need more customization



This is our focus here
 We want to understand the high throughput component

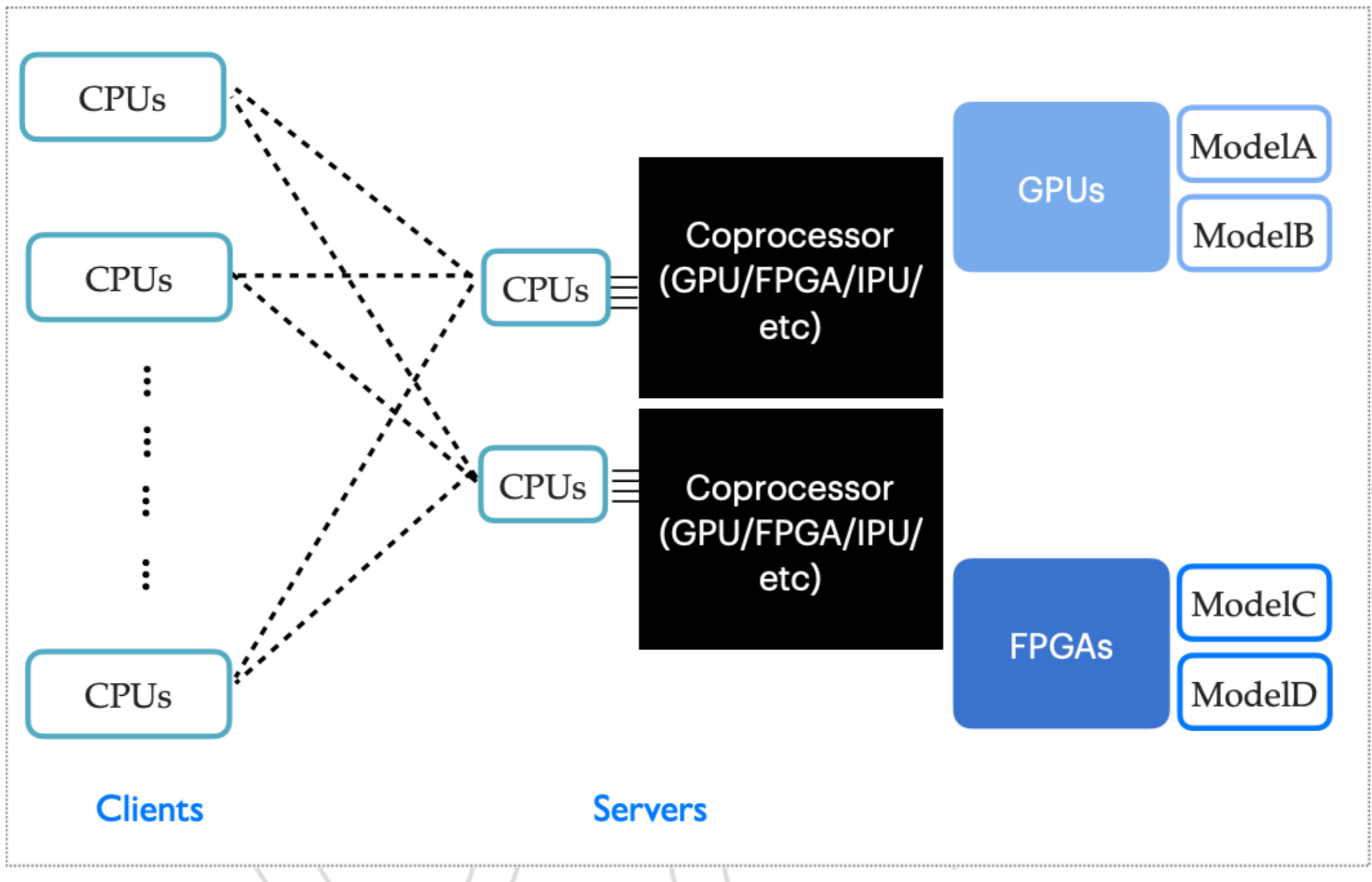
GPUs for ML



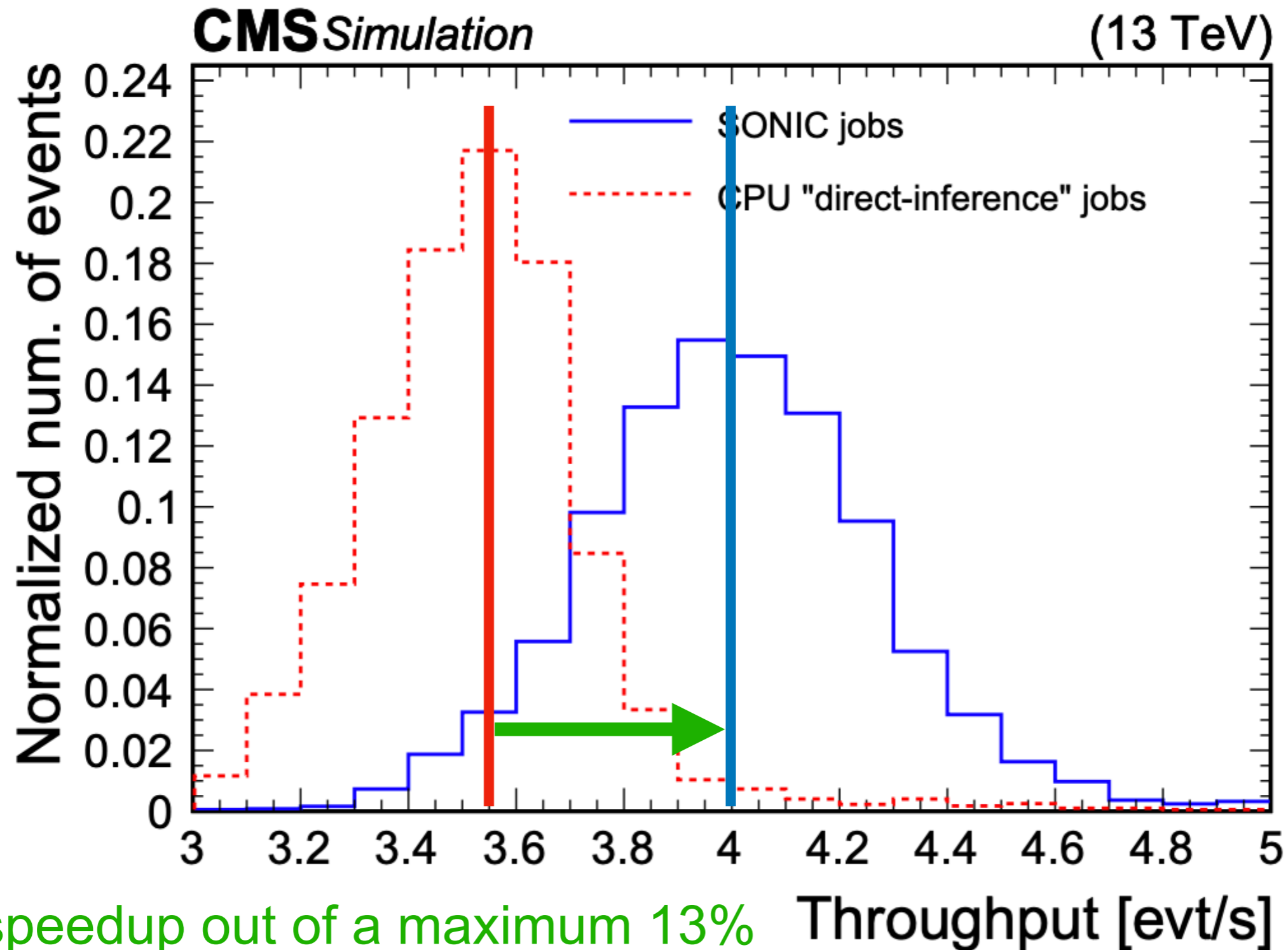
For large algorithms parallelizability and shear comput is unprecedented

Major demand emerging in many domains

Deploying to Scale



Proof we can 10k CPUs & 200 GPUs



A3D3

- An institute to unite real-time AI
 - Quickly looking for people to be part of extended team

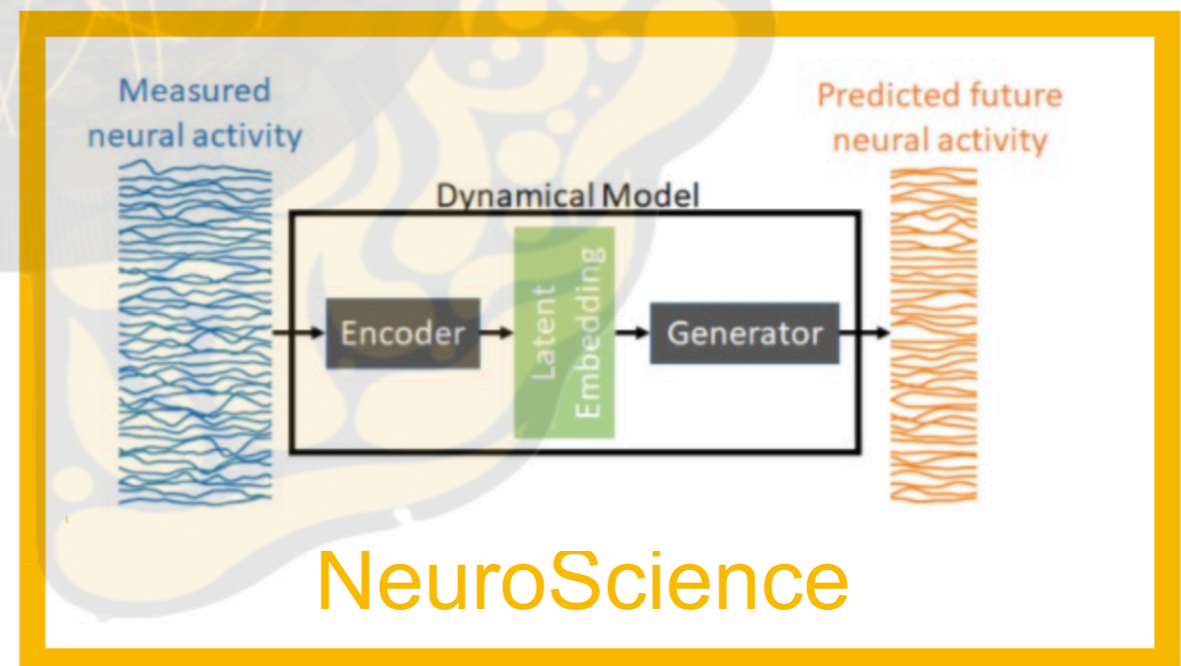
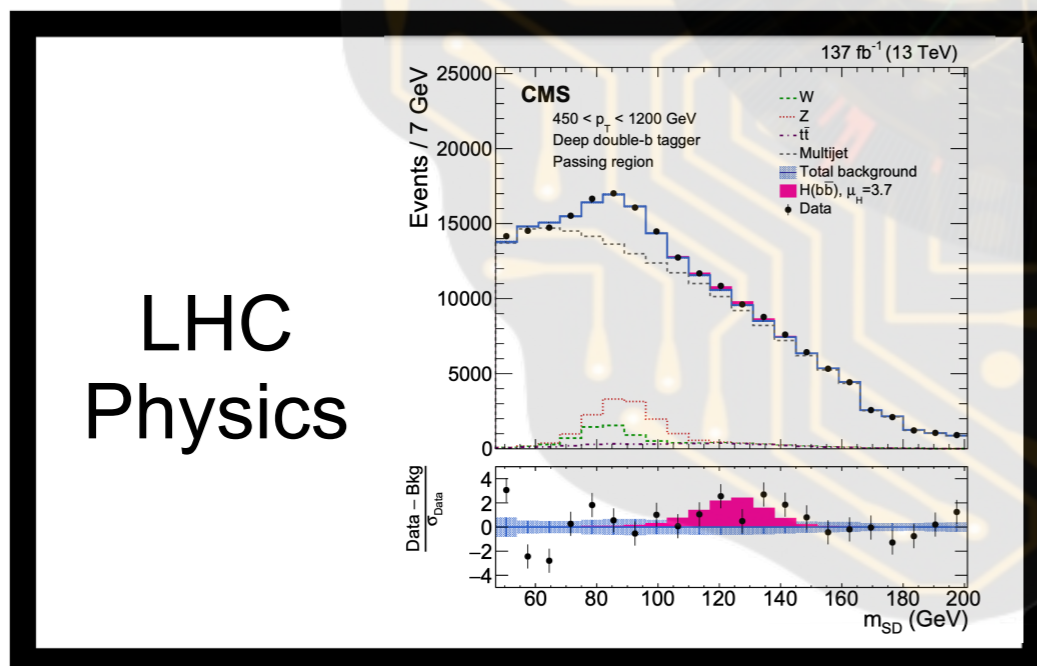
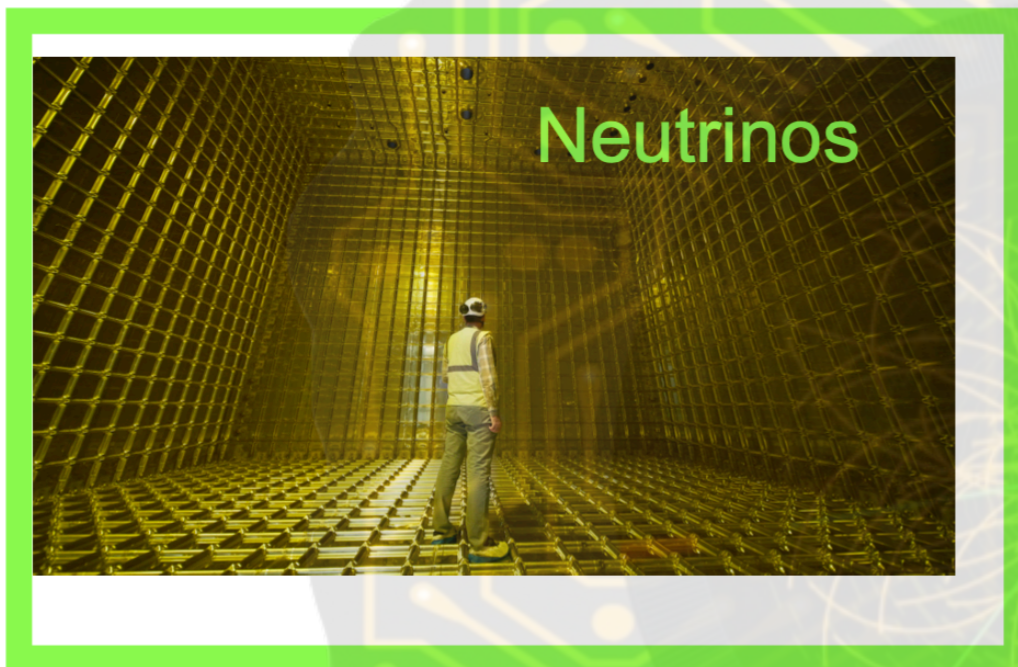


Accelerated AI
Algorithms for
Data-Driven
Discovery

An Institute: A3D3

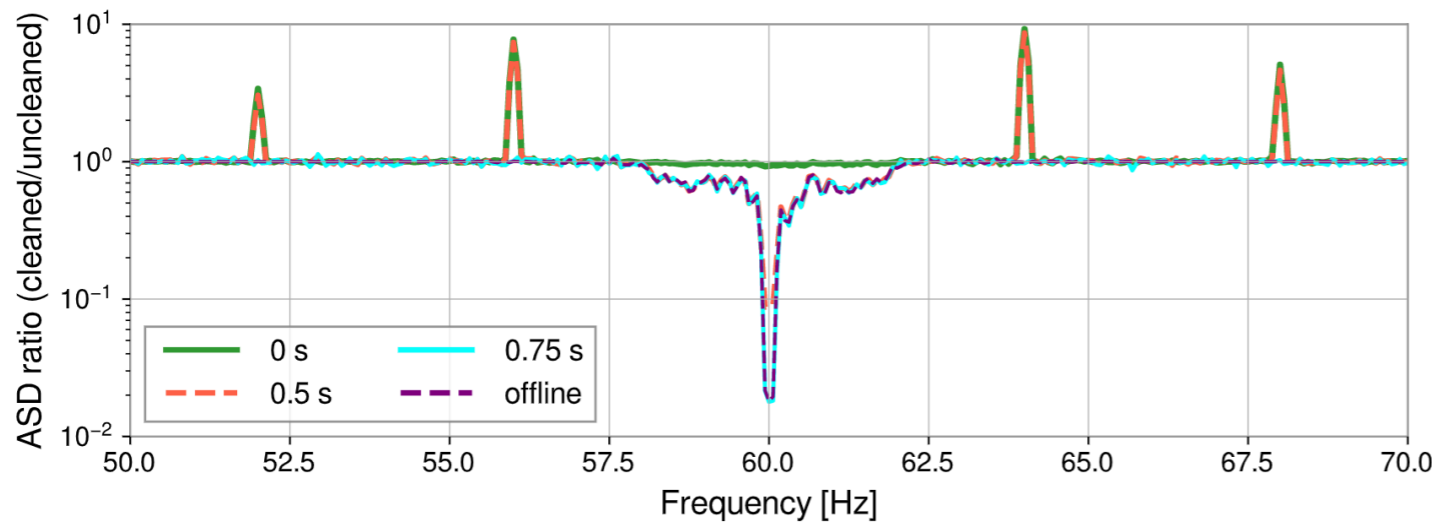
- We have been awarded a new institute to explore real-time AI
 - Accelerated AI Algorithms for Data Driven Discovery (A3D3)

New Types of Computing

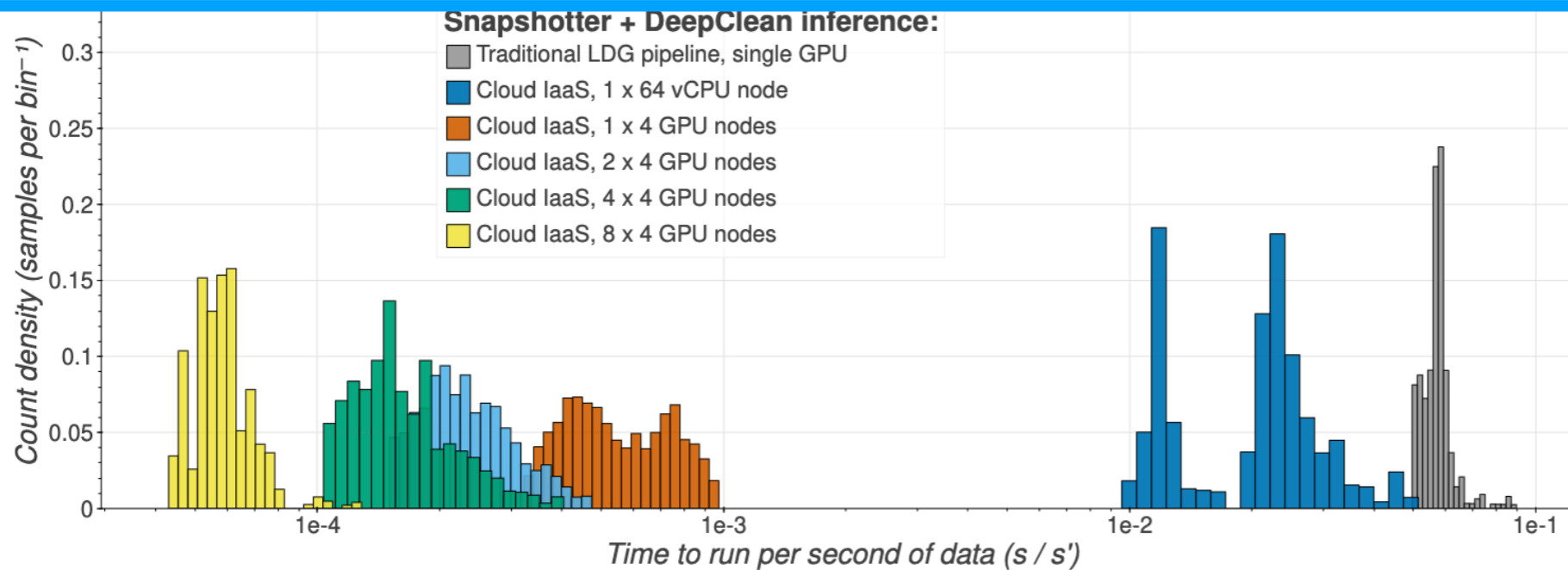


Gravitational Waves

- Actively building an AI alert system to be deployed at LIGO



Developed AI-based Denoising and BBH detection



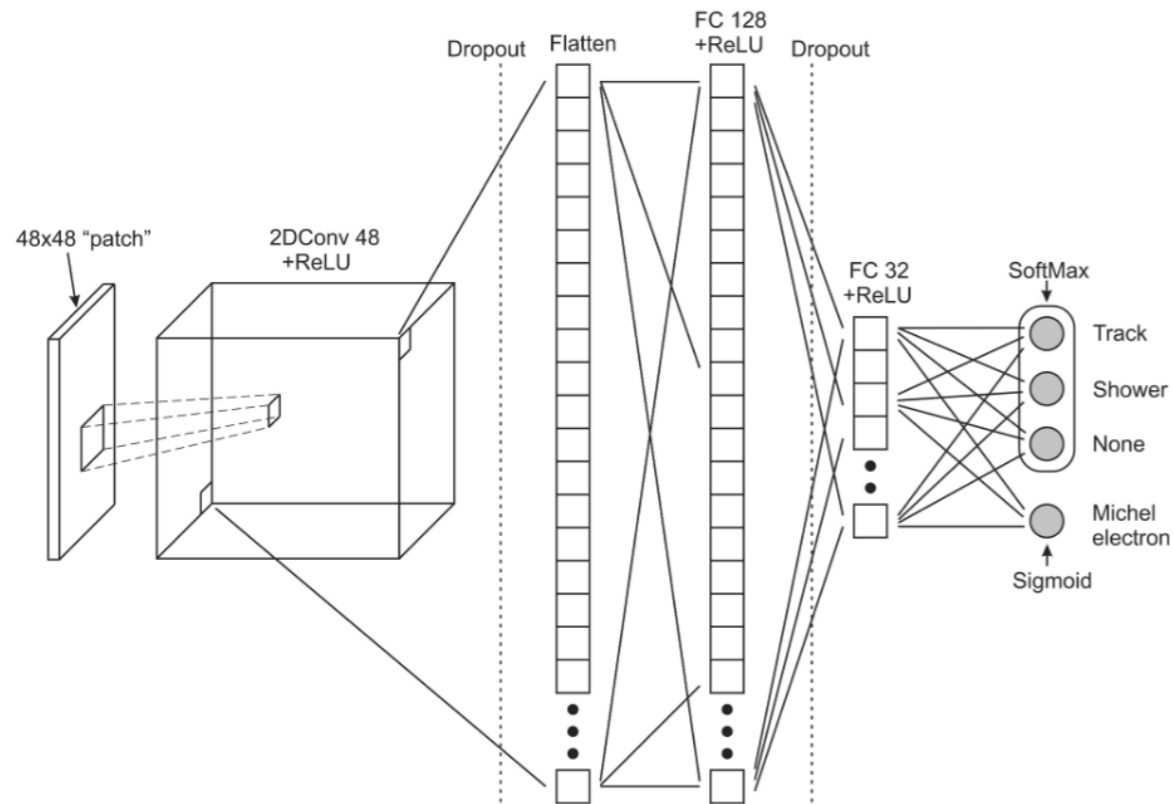
Constructed a GPU-as-a-service integration for GW low latency alerts

x1000 reduction in overall throughput

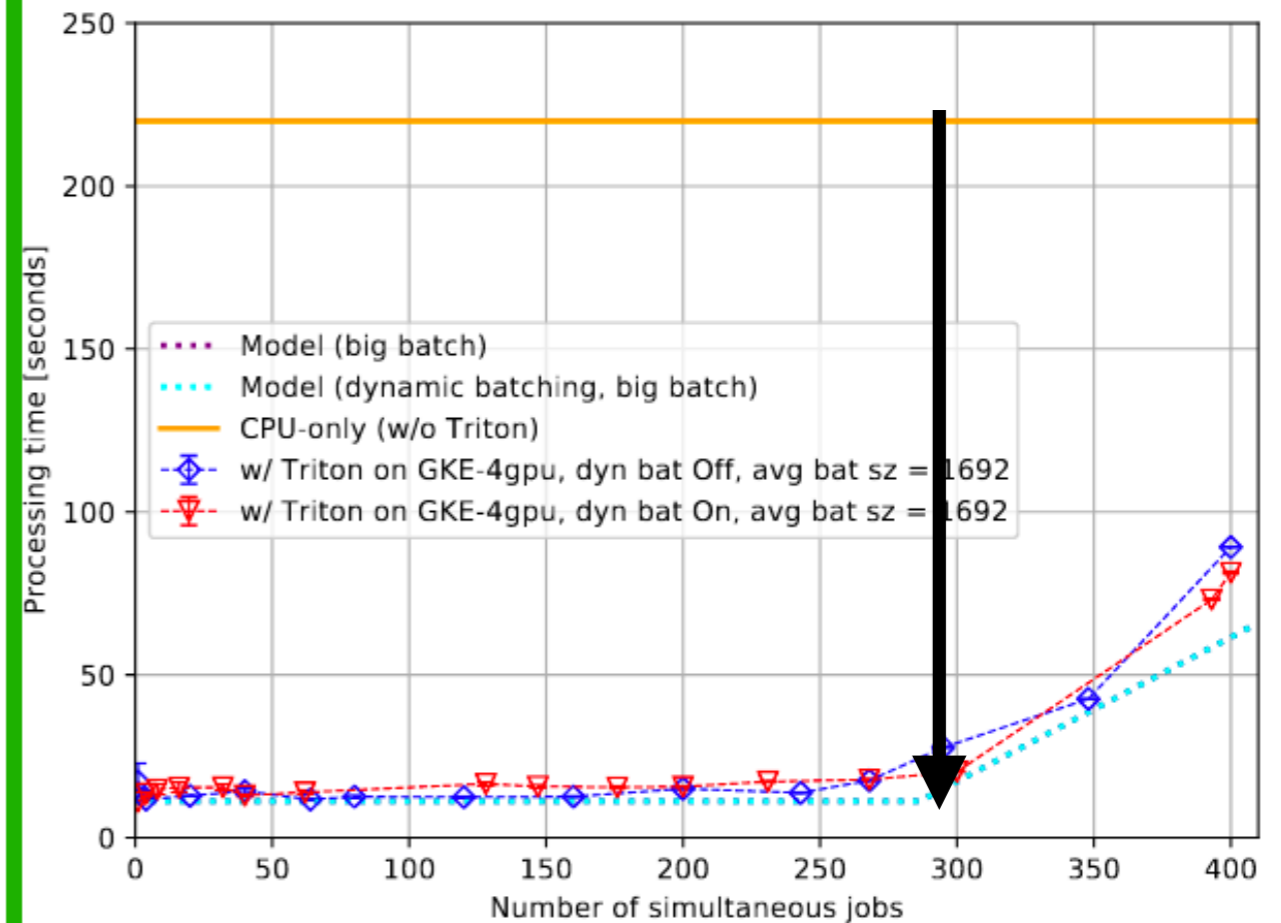


Neutrino Physics

- We are pursuing the same idea in Neutrino physics

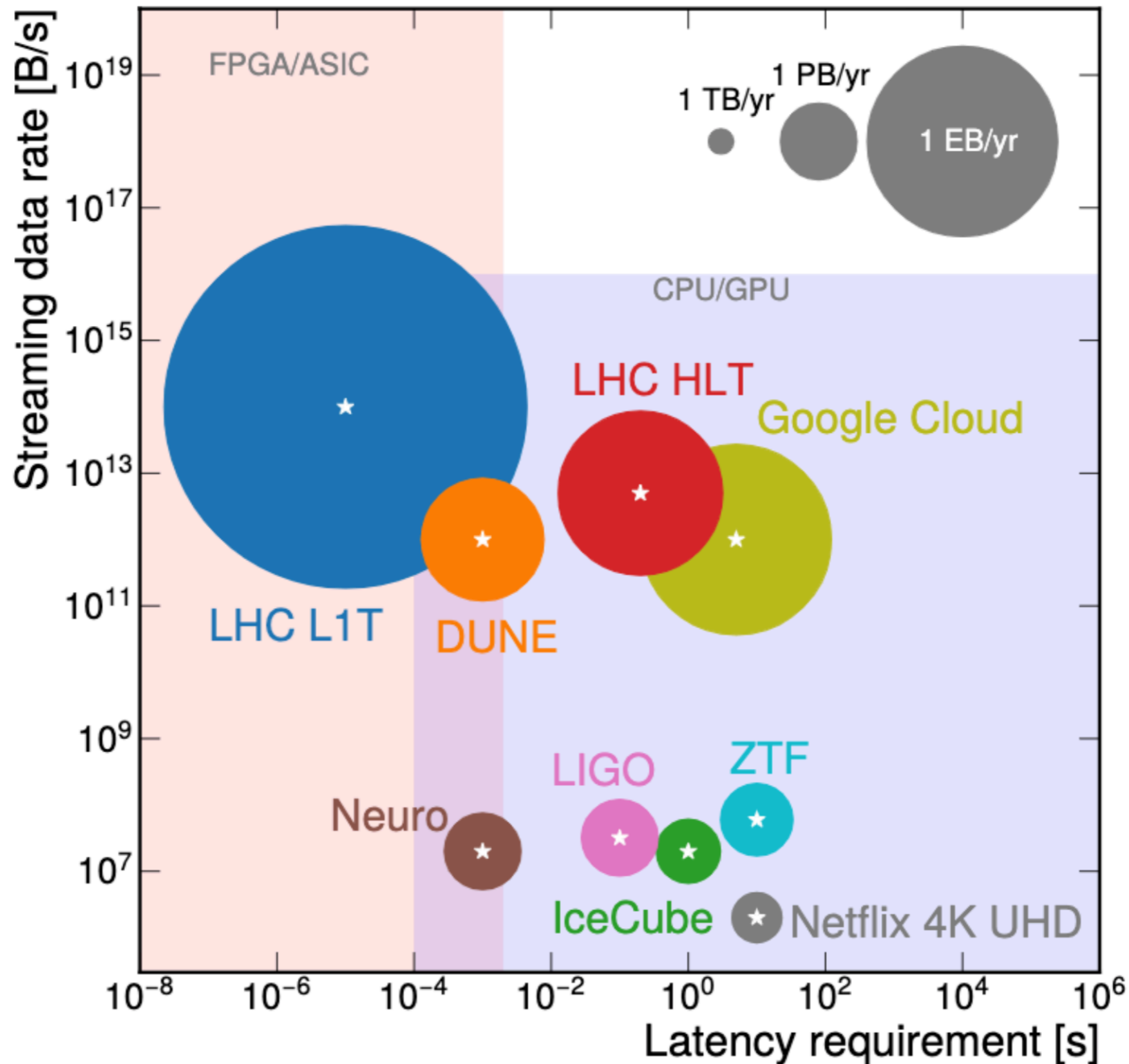


Michel Electron Id NN

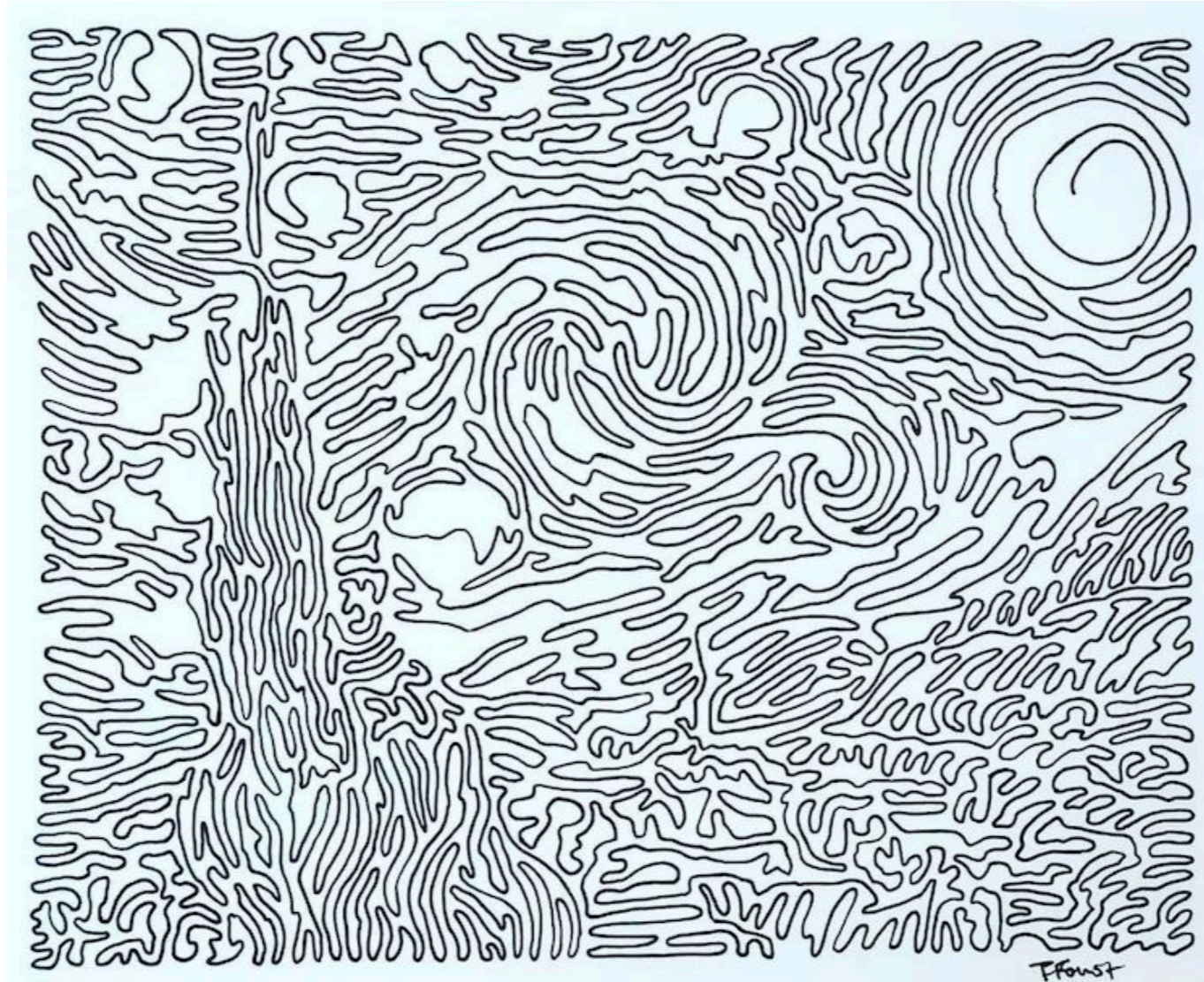


Large Factor in speed up

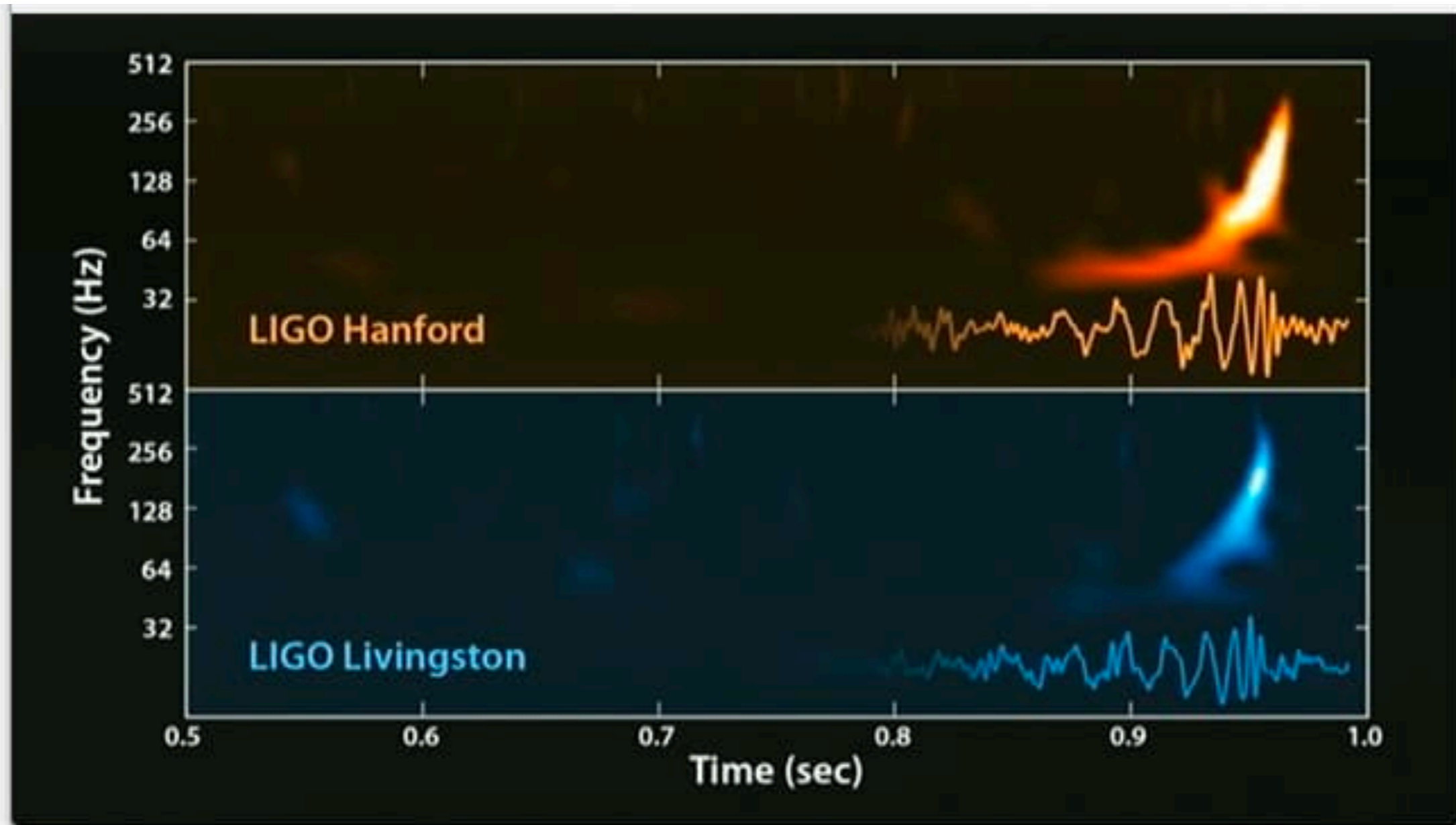
Preparing for the future



Thanks!



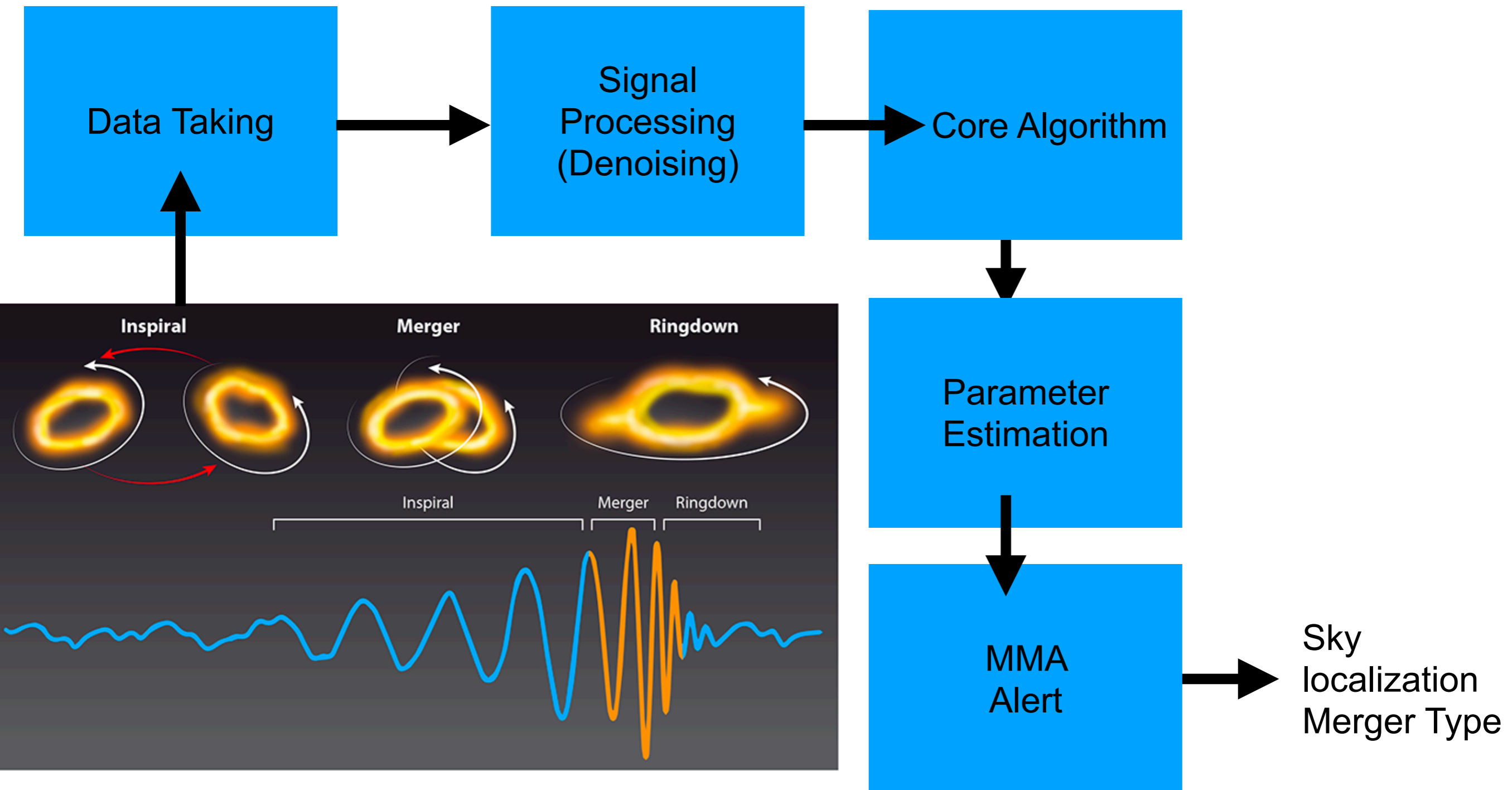
LIGO Challenge:⁶¹ Can we find all mergers



- LIGO has 10^5 channels at 1024 Hertz
- Looking for subtle signals hidden in the noise

Real-time Detailed (10k core) analysis every millisecond

LIGO Path to success



ML4GW toolkit

- Enable a complete AI pipeline we have developed ML4GW
 - Comprehensive toolkit for ML pipeline in Gravitational Waves

README.md

ML4GW

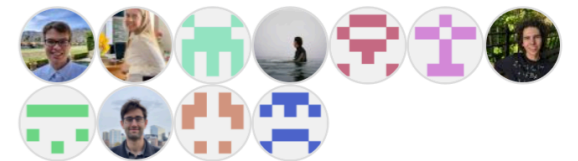
Tools to make training and deploying neural networks in service of gravitational wave physics simple and accessible to all!

Includes a couple particular applications under active research.

View as: Public

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

People



Top languages

Python Jupyter Notebook
Makefile

Most used topics

Manage

deep-learning gravitational-waves
mlops python

Pinned

DeepClean Public

Nonlinear noise subtraction from gravitational wave strain data

Python 3 6

aframe Public

Detecting binary black hole mergers in LIGO with neural networks

Jupyter Notebook 13 16

ml4gw Public

Torch utilities for doing machine learning in gravitational wave physics

Python 11 7

hermes Public

Inference-as-a-Service deployment made simple

Python 2 4

pinto Public

Job environment management and execution tool

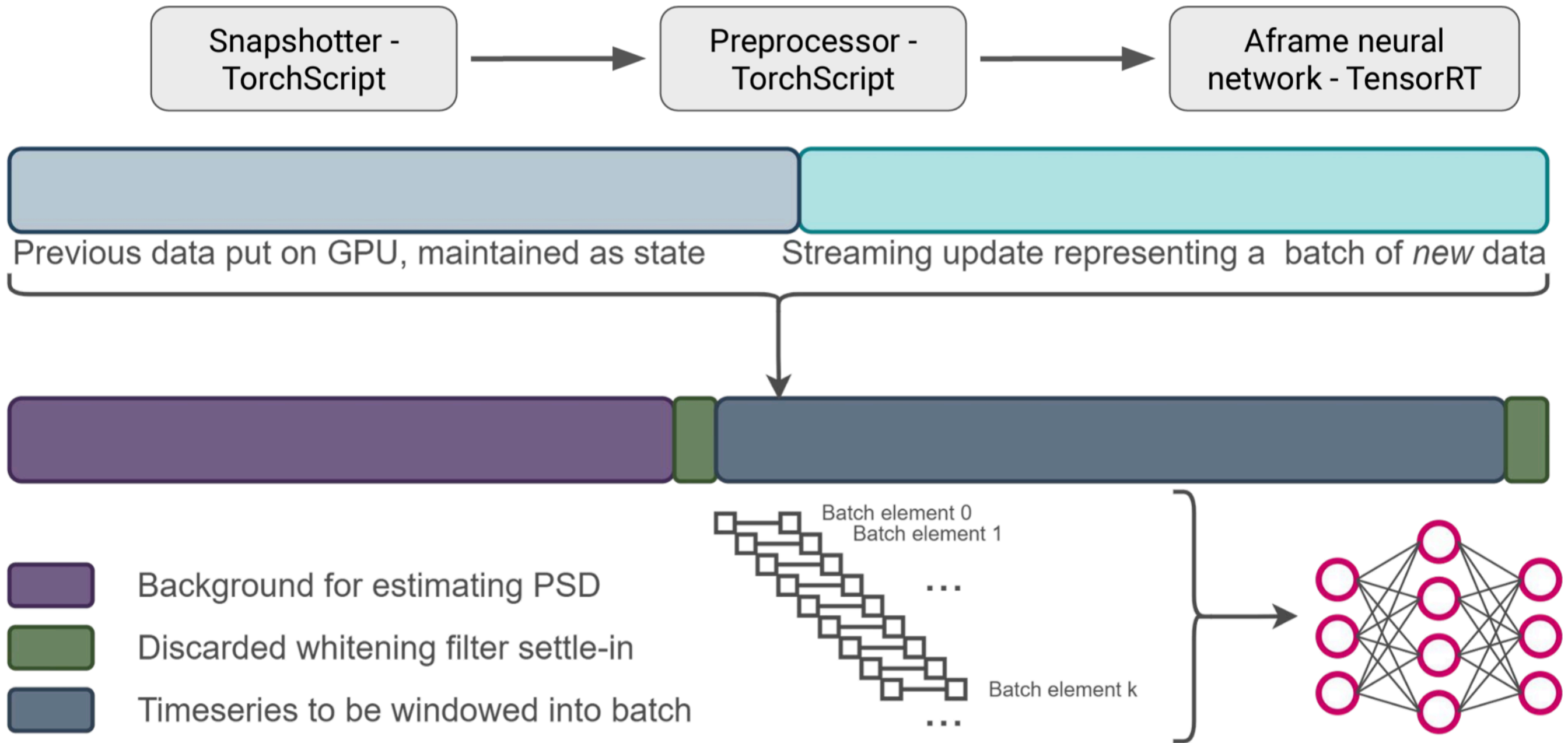
Python 3

typeo Public

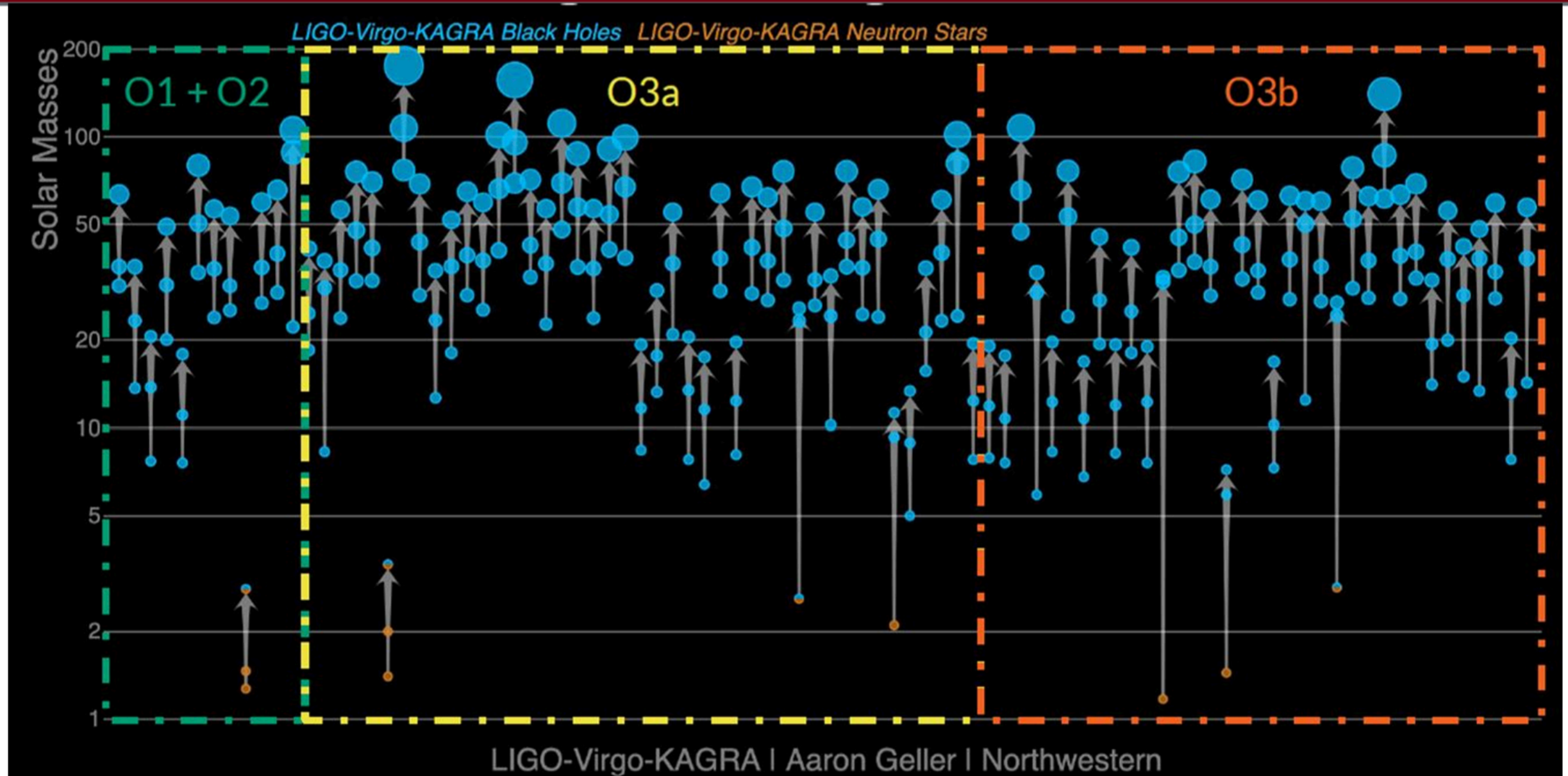
Functions as scripts as functions

Python 2

Hermes: Inference -as-a-⁶⁴ service deployment



Third transient event catalog: GWTC-3



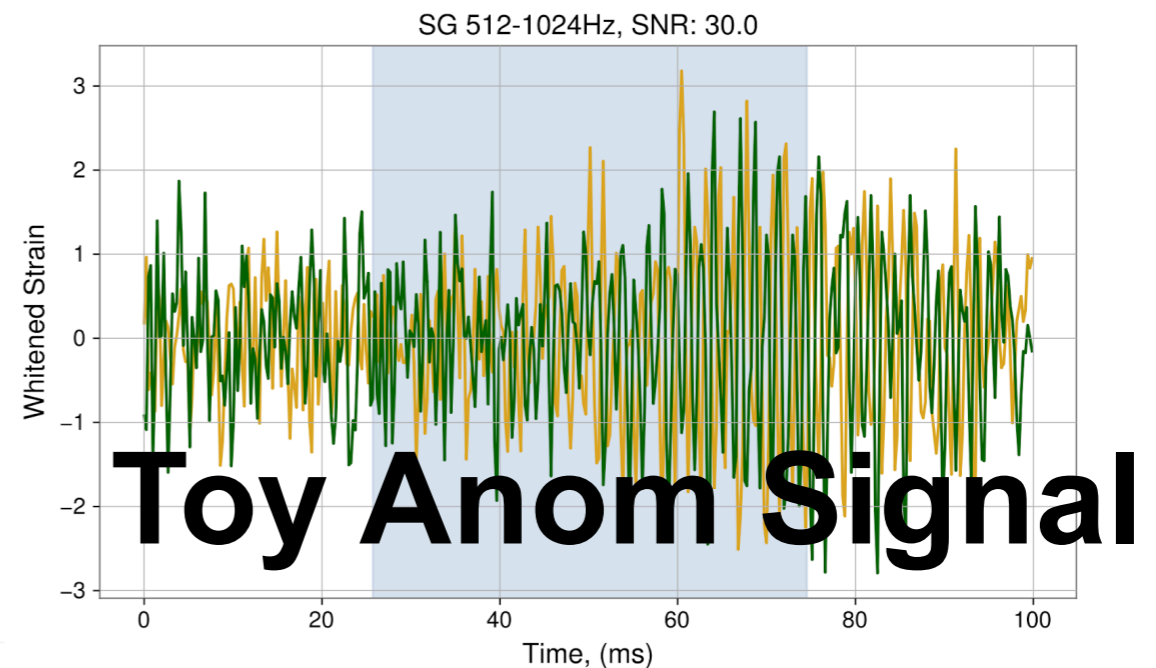
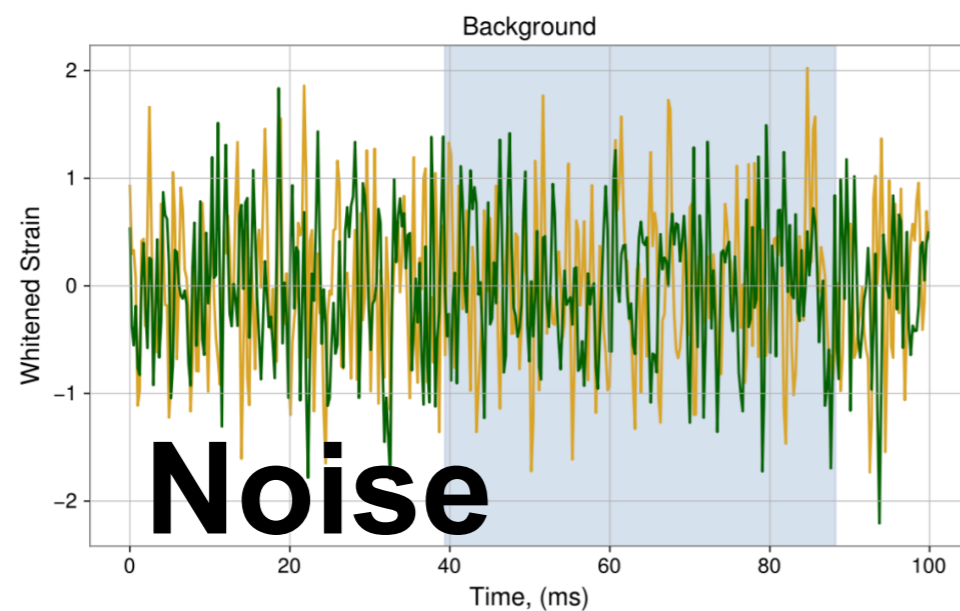
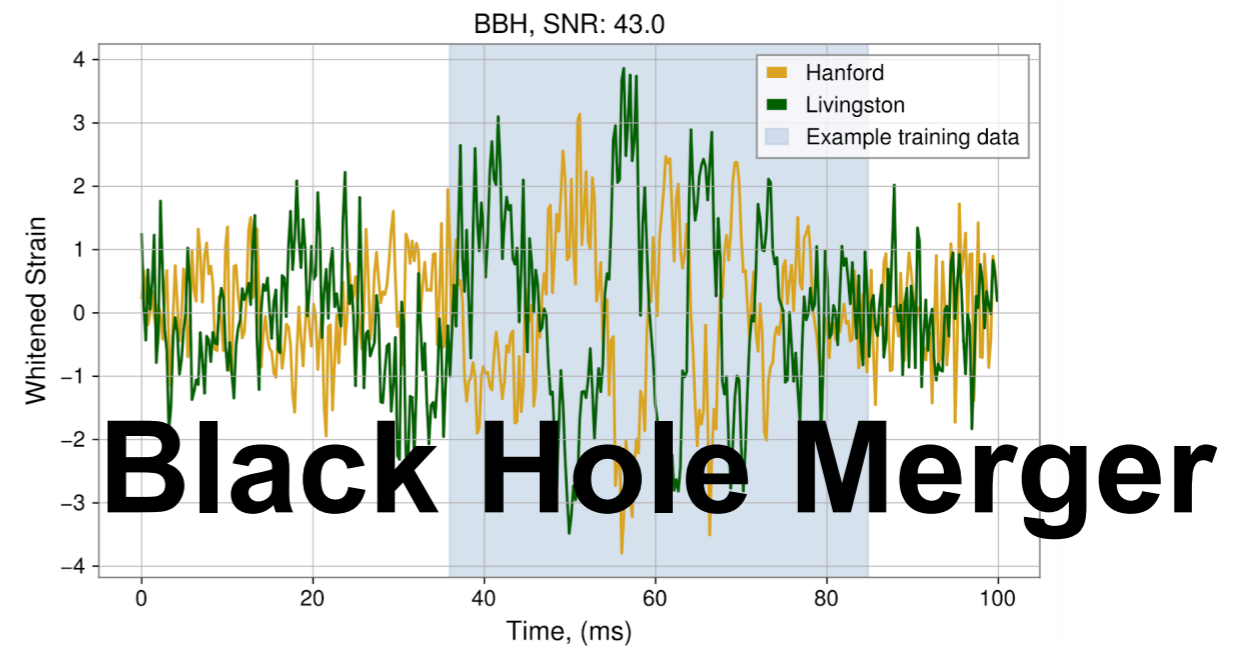
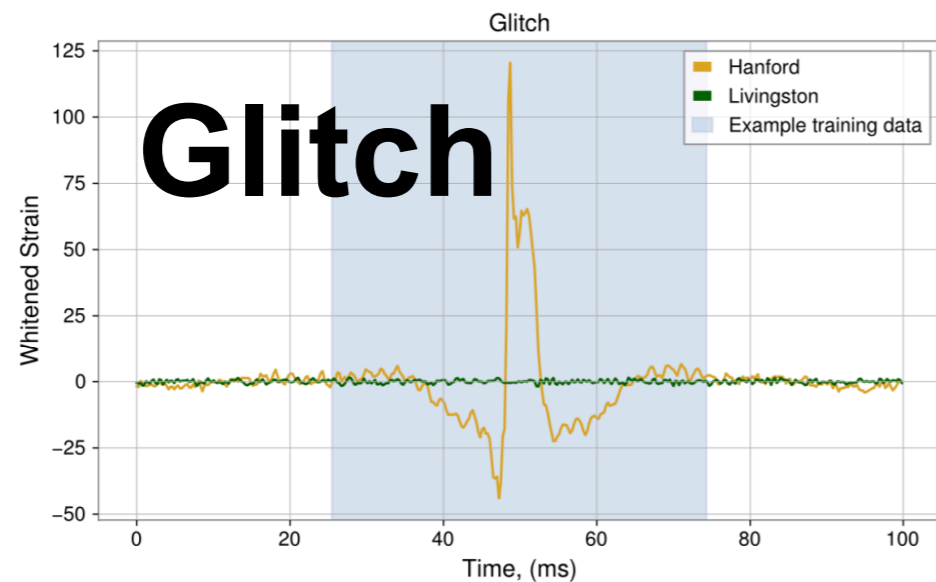
11 events
from O1+O2

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

35 events in O3b, 90 total
(catalogs are cumulative)

GWAK Space

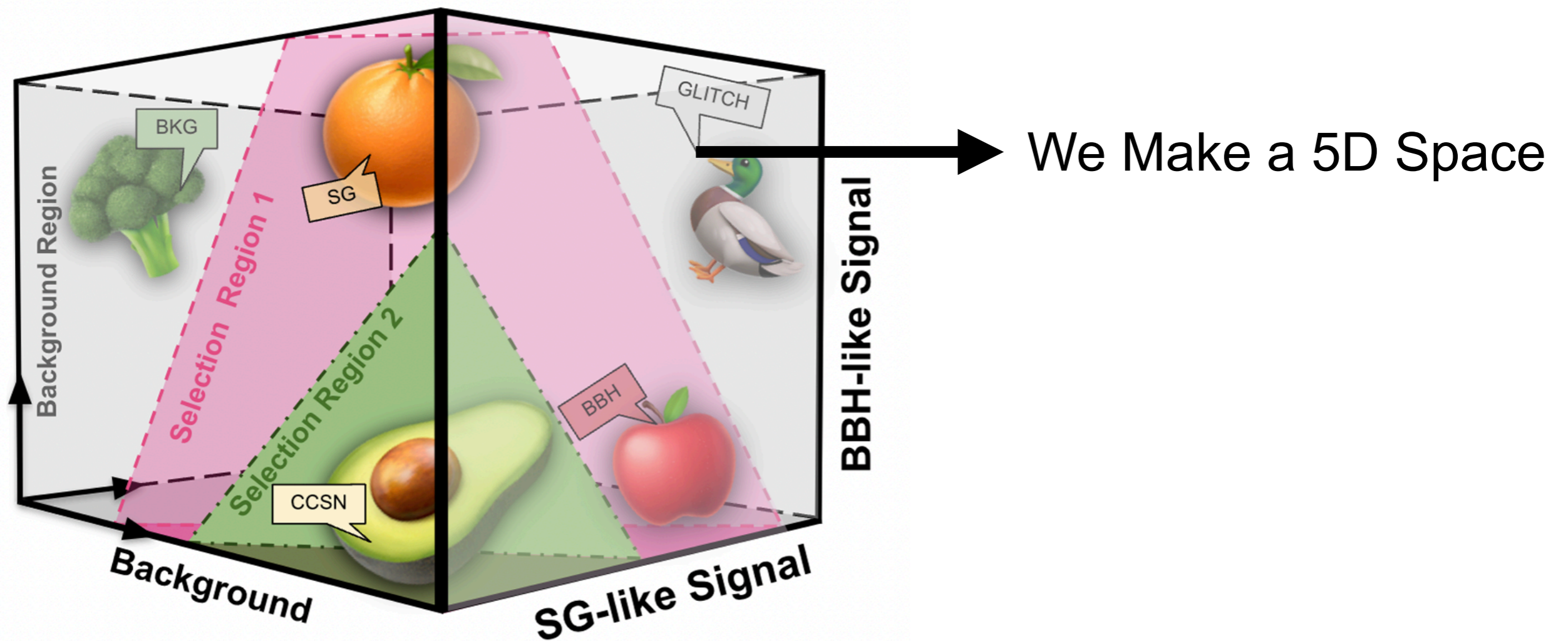
- GWAK stands for GW (QU)AK like guacamole



GWAK Space

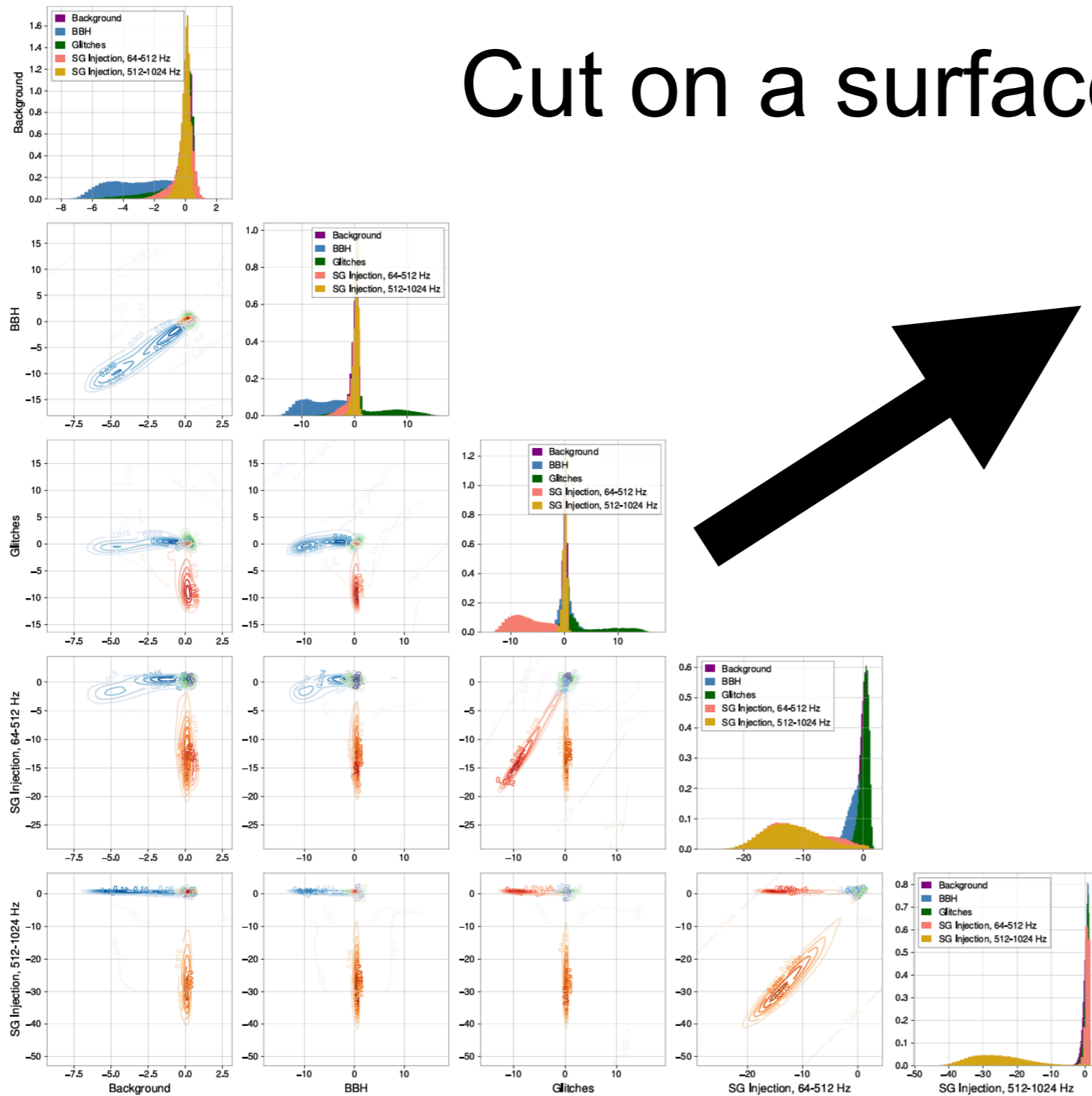
- Using autoencoders we construct a similar space

3D GWAK Space

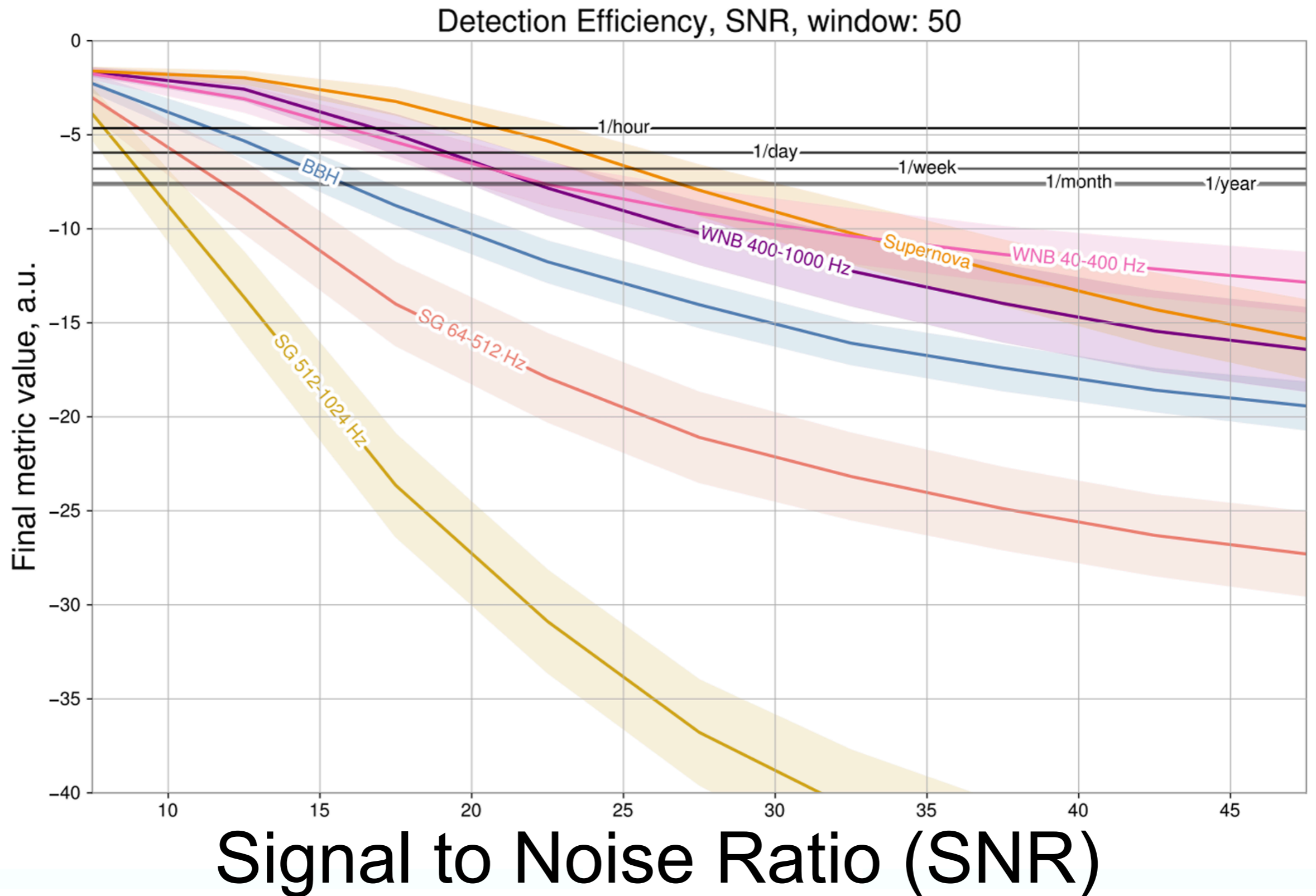


GWAK Space

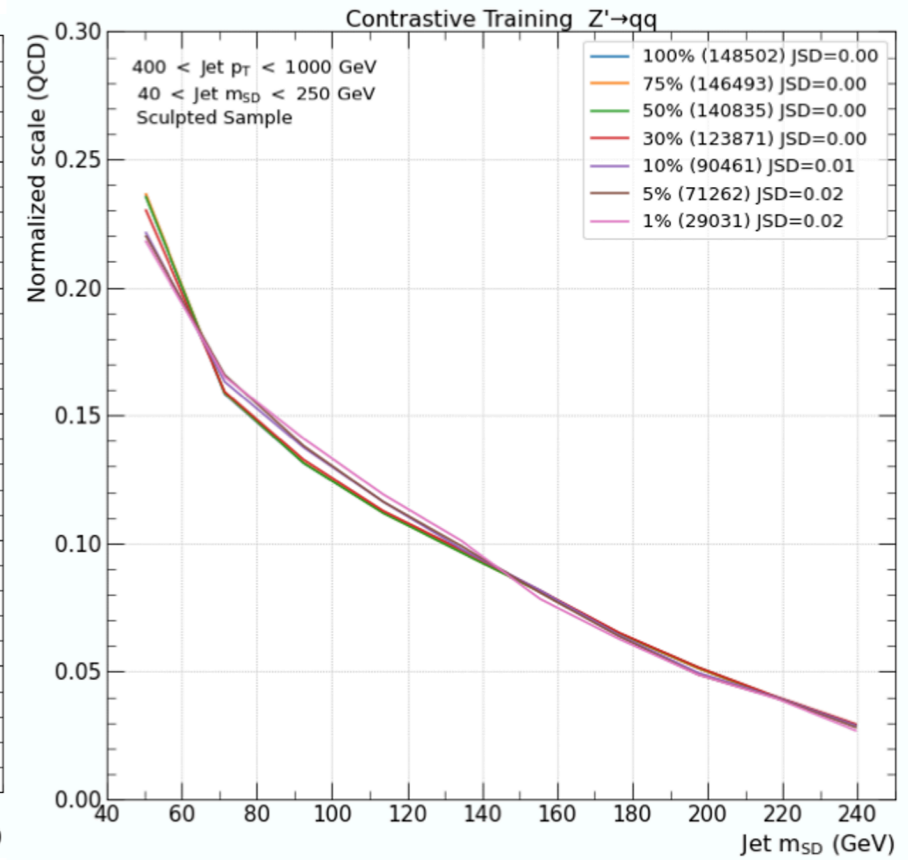
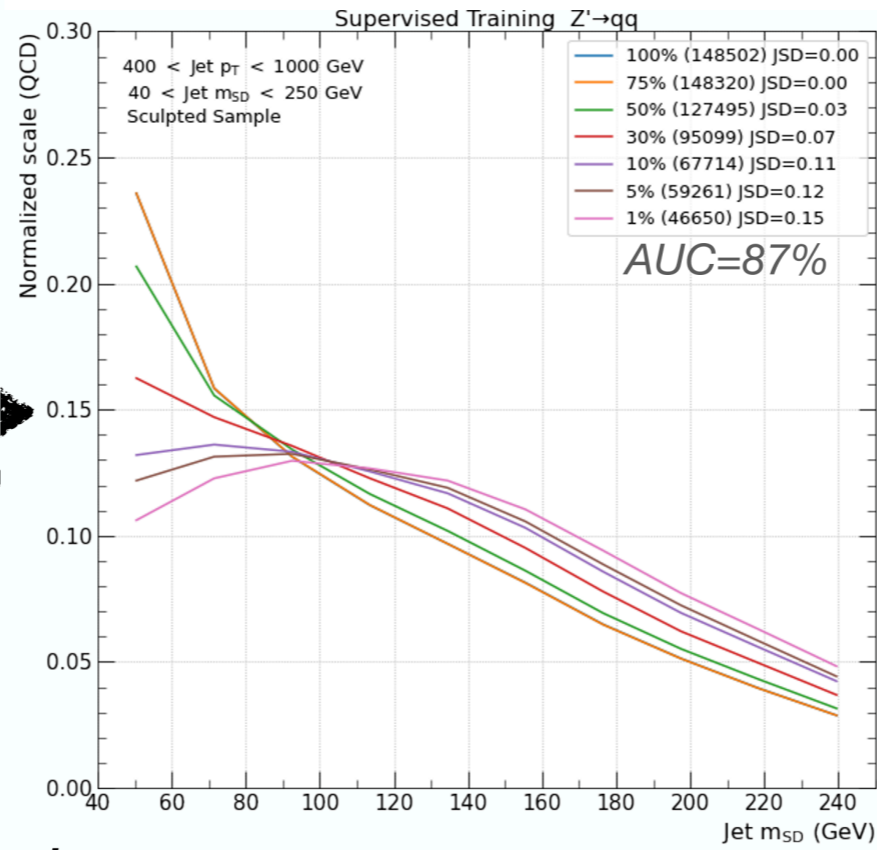
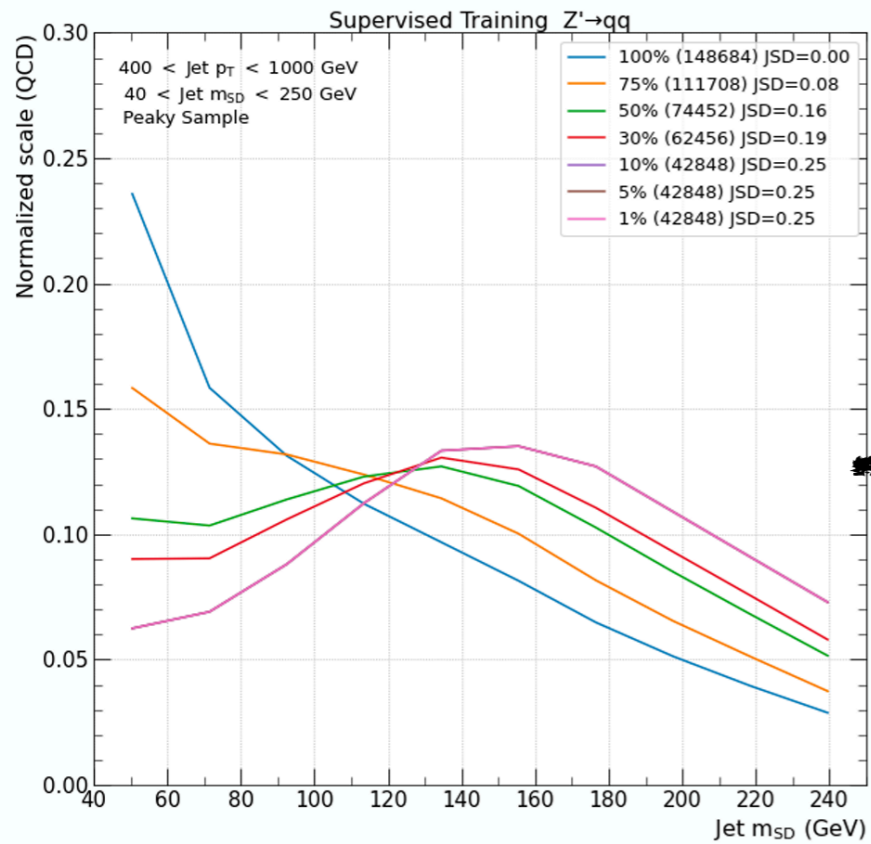
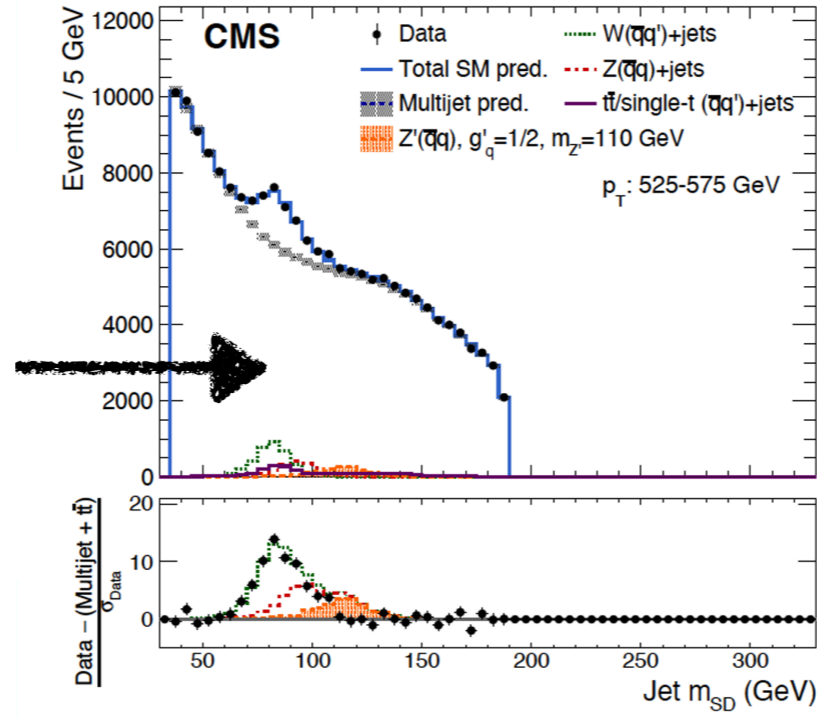
Cut on a surface in the 5D Space

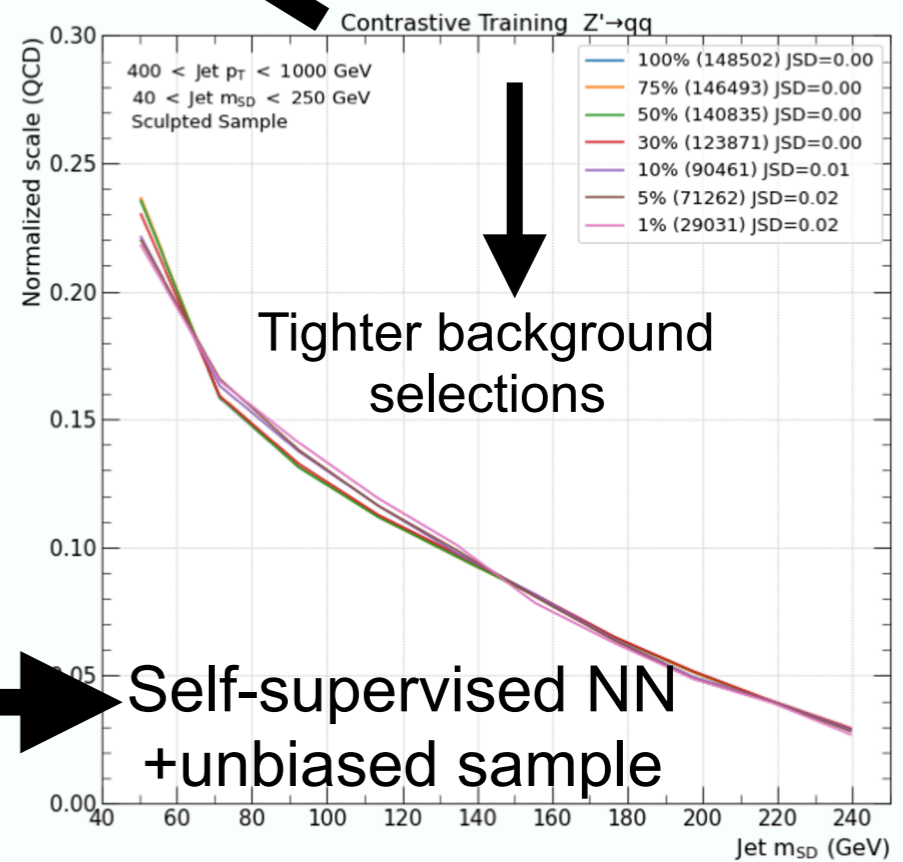
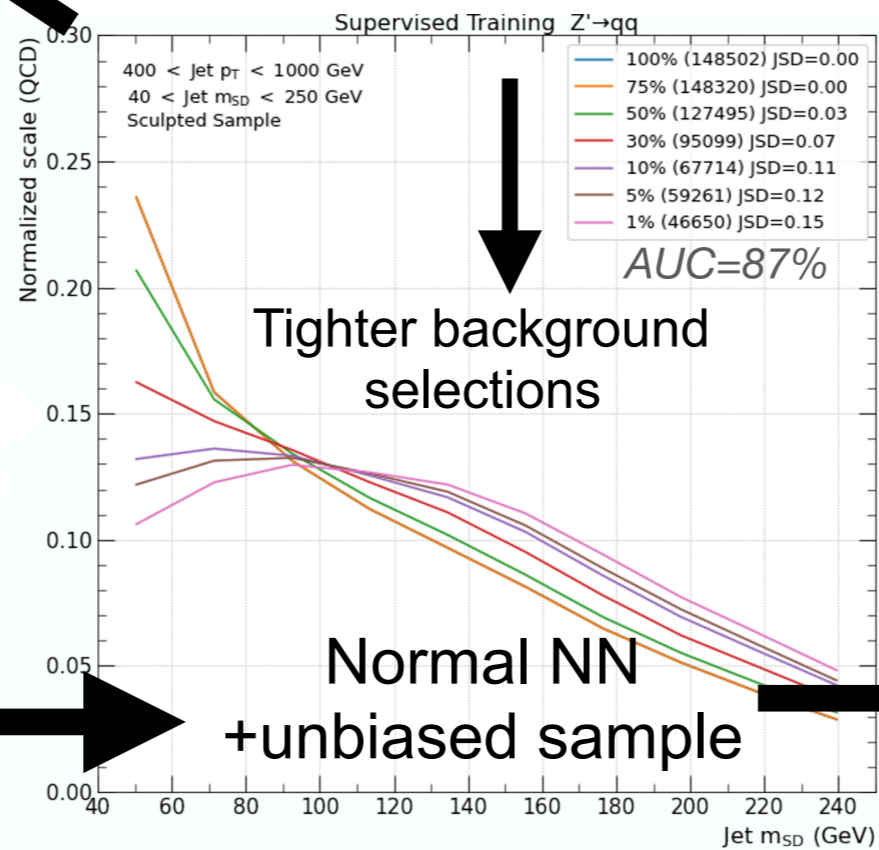
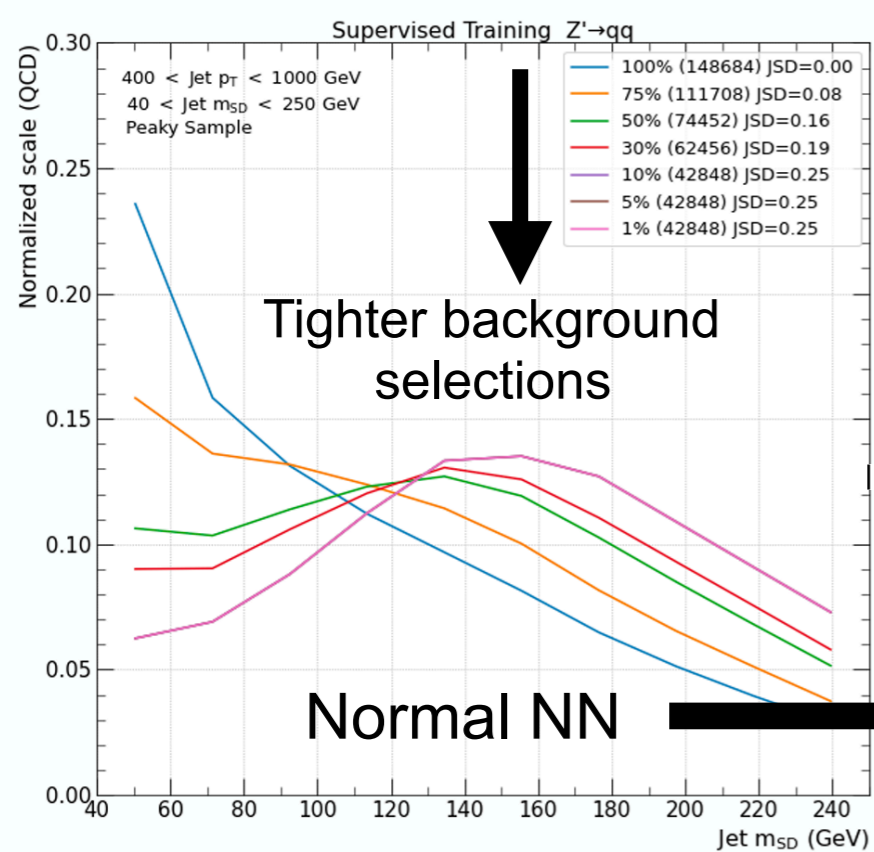
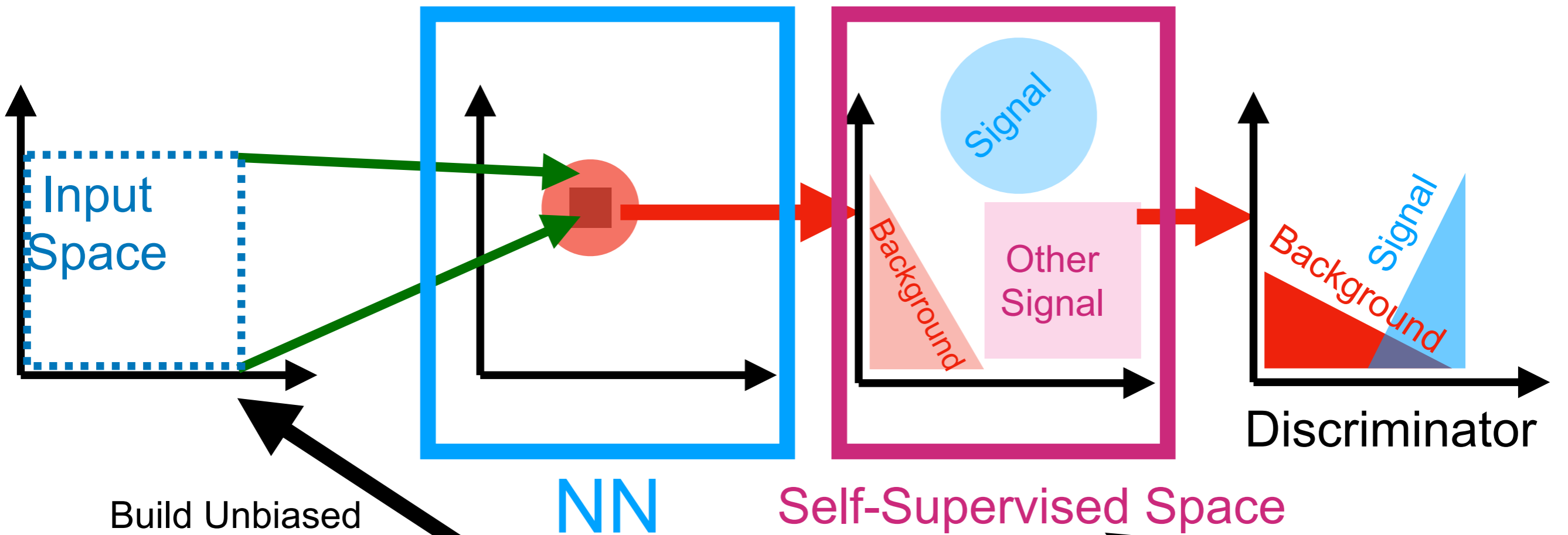


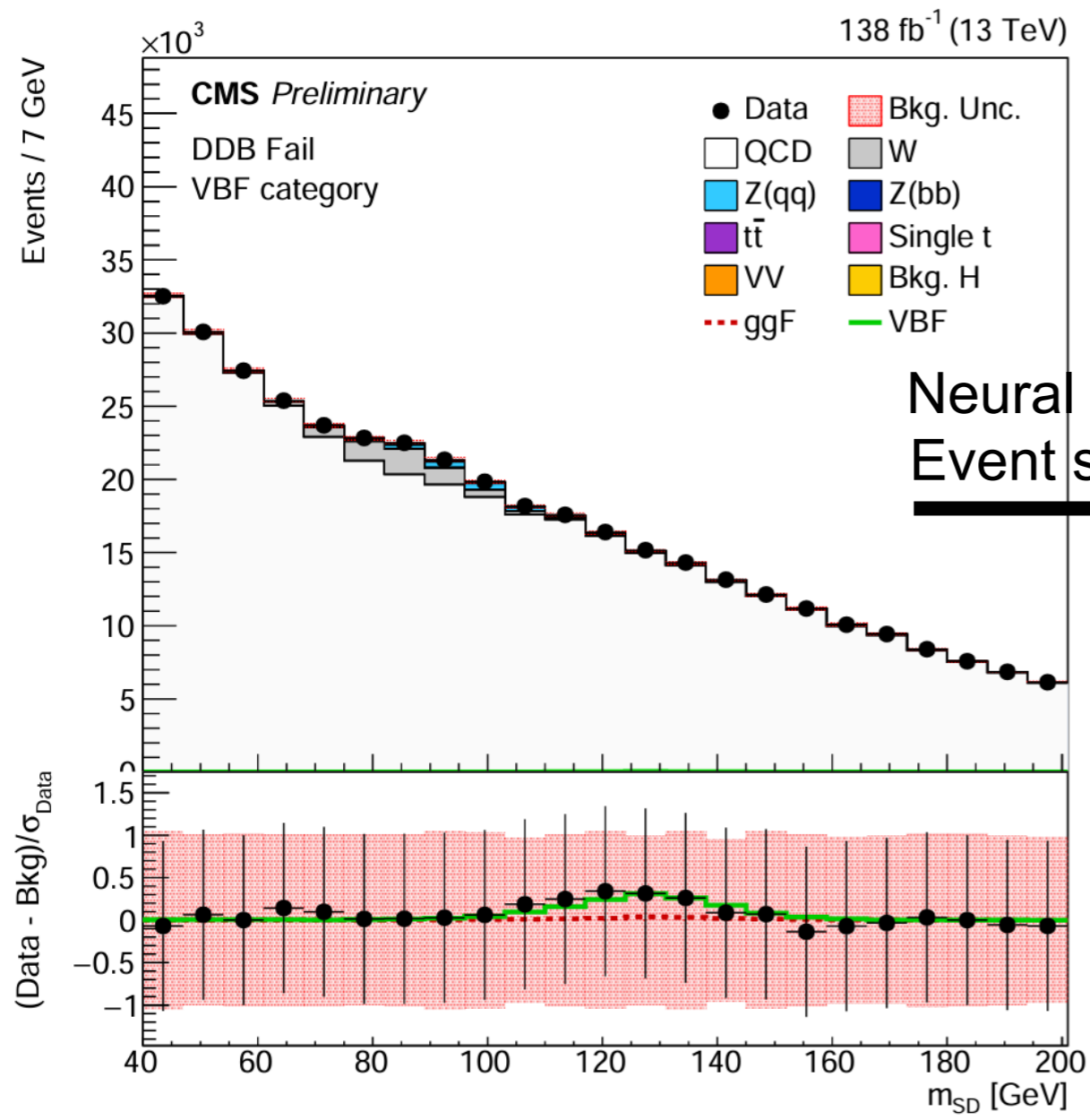
GWAK Algorithm



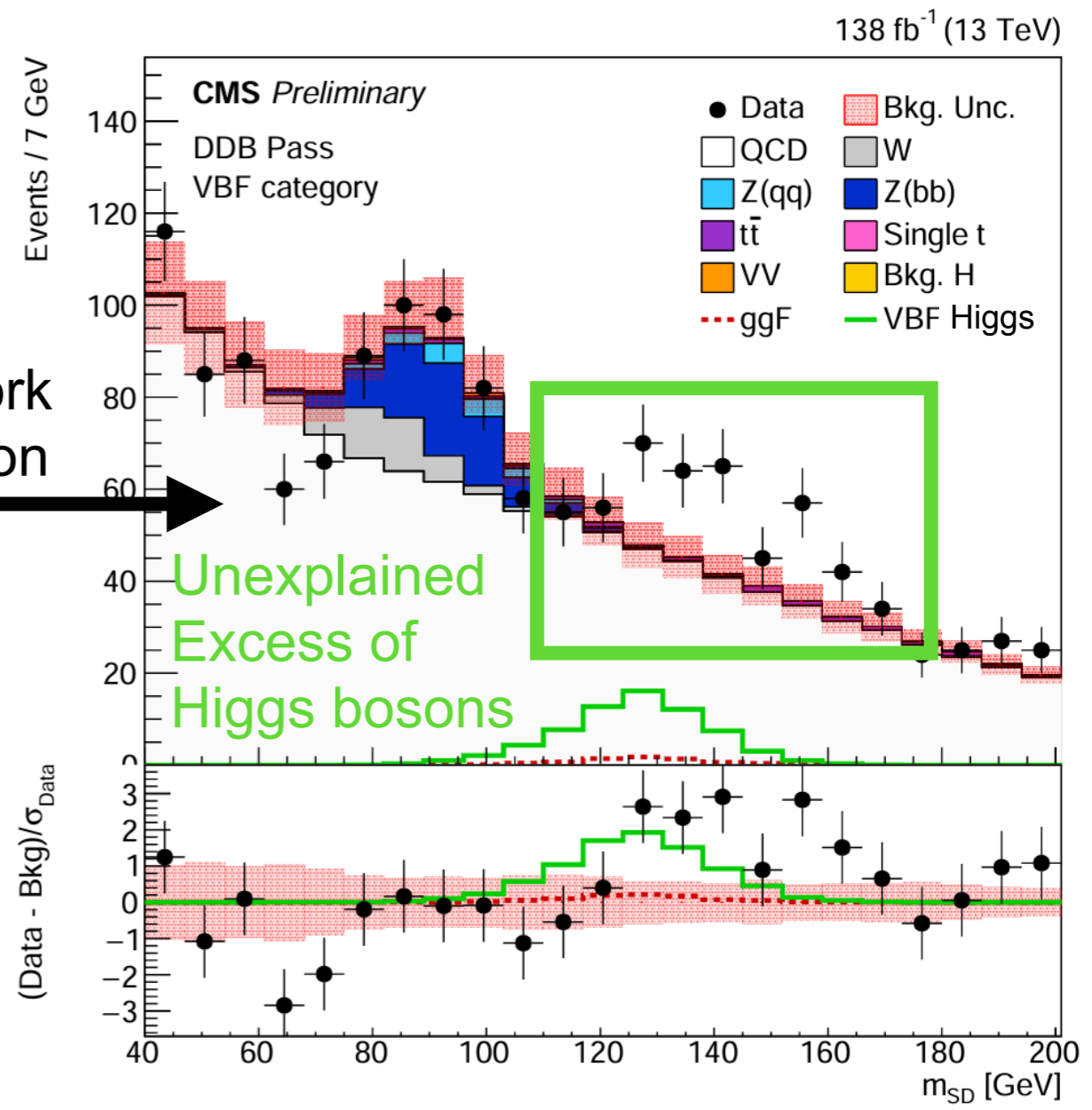
<https://arxiv.org/abs/1603.00027> 41.1 fb⁻¹ (2017) (13 TeV)





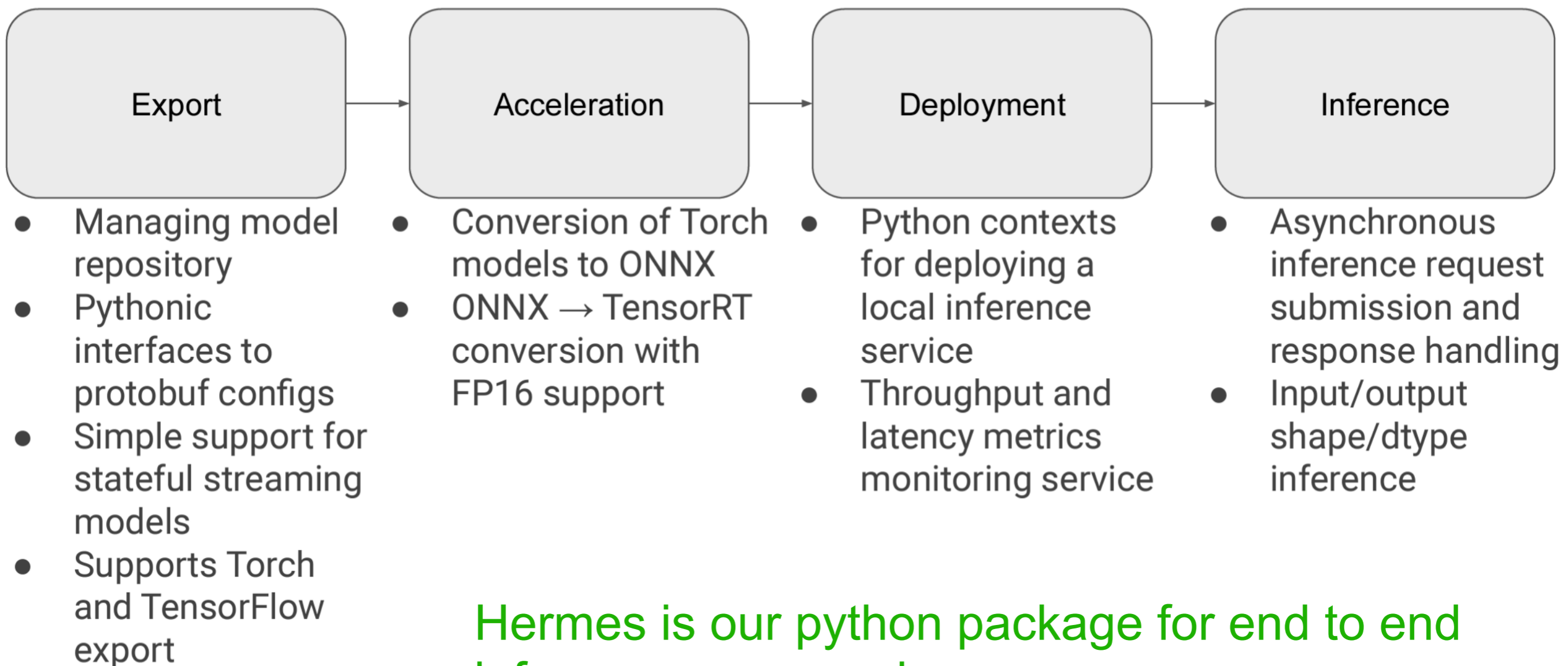


Neural Network
Event selection



Hermes: Inference -as-a⁷³ service deployment

<https://github.com/ML4GW/hermes>



Hermes is our python package for end to end inference-as-a service

Computing Ecosystem



PyTorch Lightning



- Ray: Handle distribution of CPU brokers for data
 - PyTorch Lightning: Optimized training and model development
 - **Apptainer: containerization to enable maximum flexibility**
 - Luigi: lightweight task execution, Kubernetes support contrib (Spotify)
 - LAW: Wrapper for Luigi to enable Condor/Slurm support (HEP)
 - **Kubernetes for distributed server balance**
- Open source tools some from industry**

ML in GW Processing

Online

Real-time analysis with goal of alerting electromagnetic astronomers (MMA) of significant events

Detect events → Localize on Sky
→ Send public alerts

Main consideration is *latency*

Potential Future Use Case

Offline

Large scale analysis of archival data for

- End to end searches
- Validating new methods, performing new research

Main consideration is *throughput*

What we want now w/Nautilus

