



A prototype for real time Gravitational Wave transient signal classifier

EGO - Elena Cuoco (*Scientific Supervisor*)

Trust-IT - Alessandro Staniscia, Emanuel Marzini, Filip Morawski, Alessandro Petrocelli, Silvana Muscella



H2020-ASTERICS project brings together for the first time scientists and communities from astronomy, astrophysics, particle astrophysics & big data. <http://www.asterics2020.eu>

About us



Elena Cuoco
*European Gravitational
Observatory*

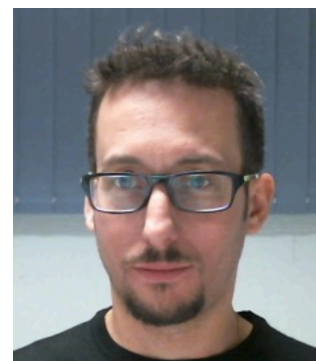
- Ligo/Virgo Co-chair of Machine Learning group

- ESCAPE project scientist

<https://projectescape.eu>



- Chair of CA17137: www.g2net.eu



Alessandro Staniscia
*Temporary Consultant for
Trust-IT Service*

- 10% Freelancer and Temporary Consultant

- 80% Scrum Muster and Team Member on

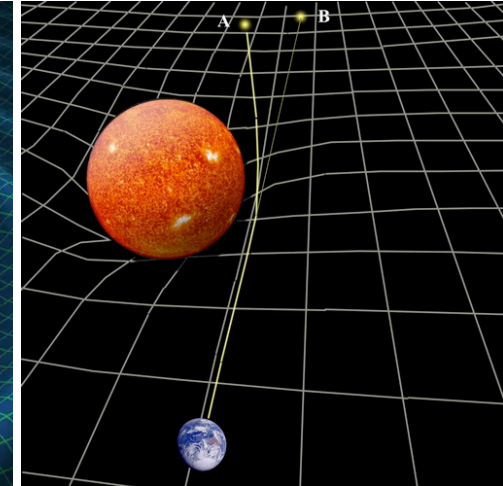
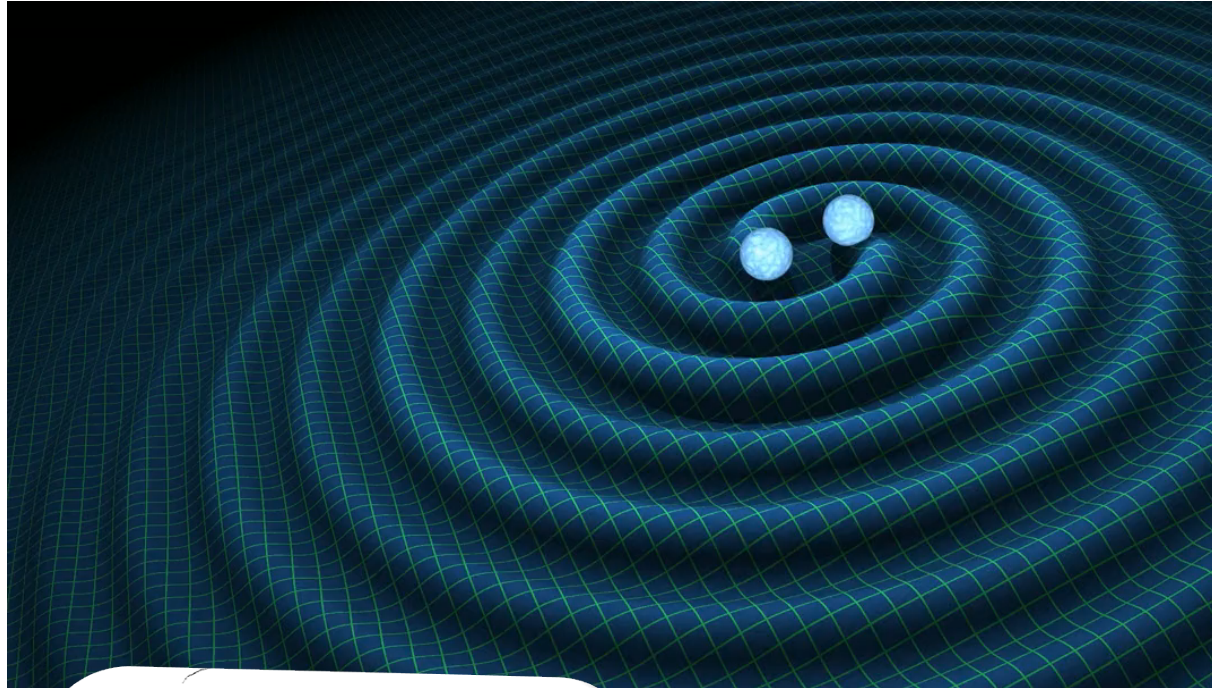
- 10% Dreamer



Outline

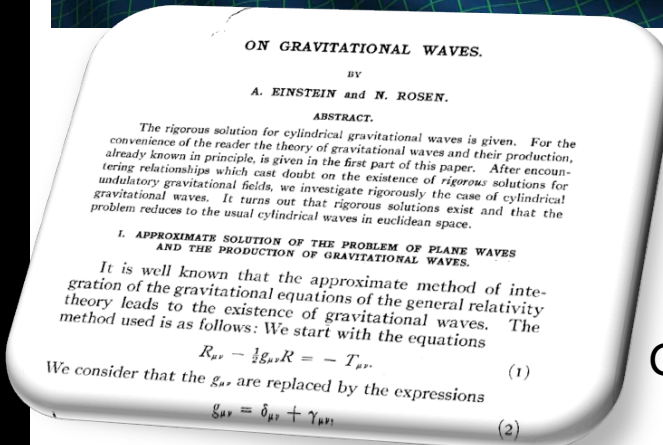
- I. Gravitational waves in a nutshell*
- II. Why real time analysis*
- III. From In-time analysis to Wavefier prototype*
- IV. The technical / implementation aspects*
- V. What's next*

What are Gravitational Waves (GWs)?



General Relativity (1915)

$$G_{\mu\nu} = \frac{8\pi G}{c^4} T_{\mu\nu}$$



Gravitational Waves (1916)

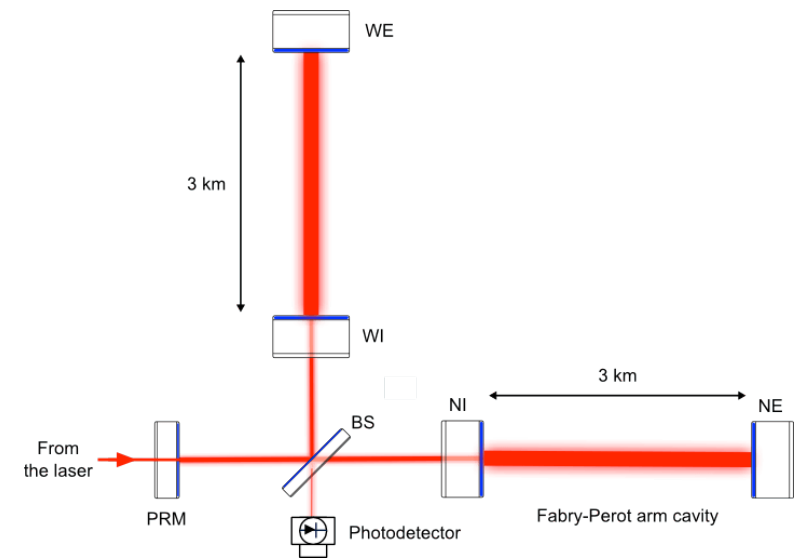
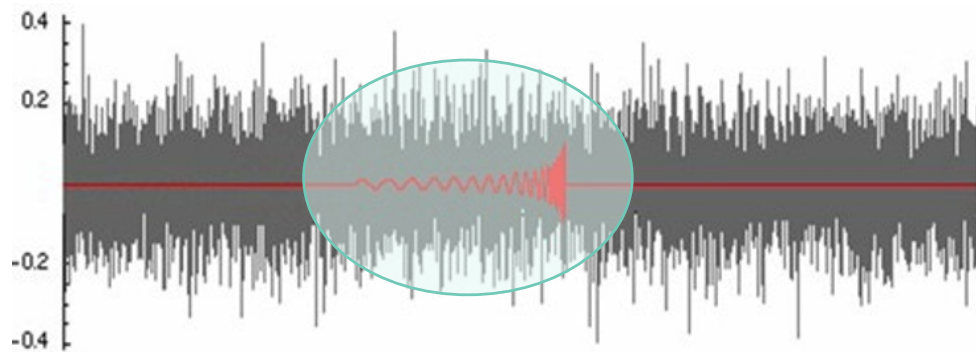
International collaboration



LIGO/Virgo data

Time series sequences... **noisy time series** with low amplitude GW signal buried in

Data flux stream: 50MB/s



I. Gravitational waves in a nutshell

Which kind of astrophysical sources?

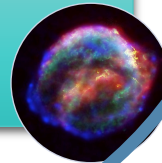
- Known waveform
 - Transient signal

Compact binary coalescence (CBC)



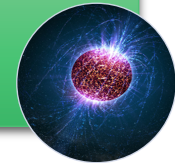
- Unknown waveform
 - Transient signal

Core Collapse Supernovae (CCSN)



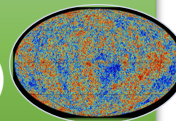
- Known waveform
 - Persistent signal

Continuous Waves (CW)



- Unknown model
 - Persistent signal

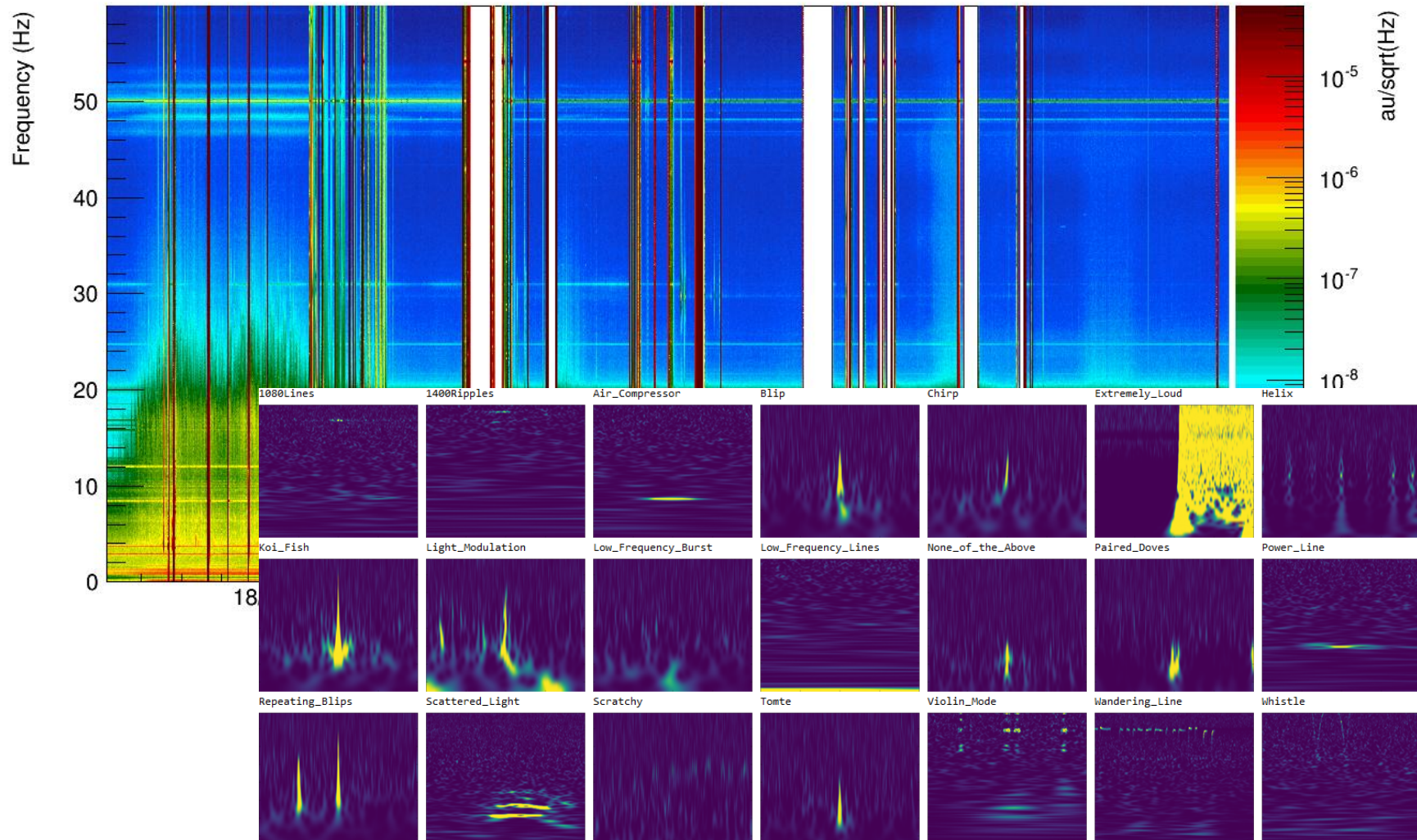
Stochastic Background (SB)



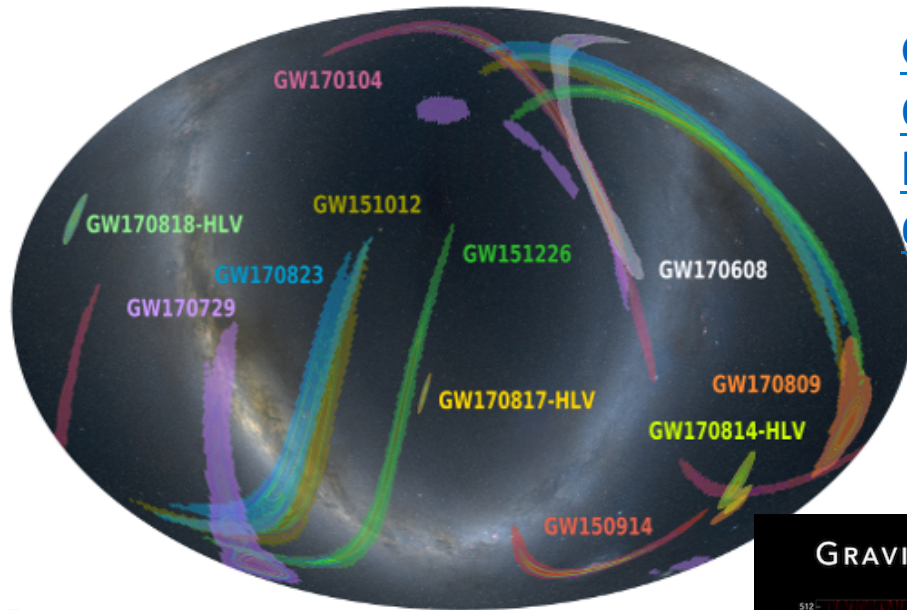
I. Gravitational waves in a nutshell

Which kind of noise?

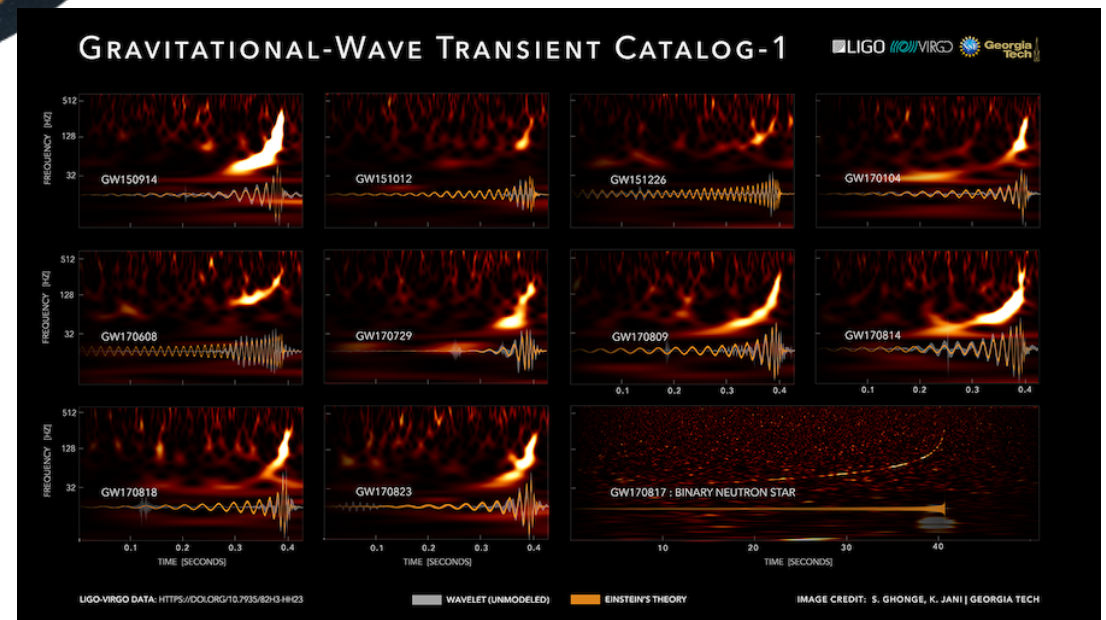
Spectrogram of V1:spectro_LSC_DARM_300_100_0_0 : start=1189644747.000000 (Sun Sep 17 00:52:09 2017 UTC)



The first GW catalog (O1/O2 run)



GWTC-1: A Gravitational-Wave Transient Catalog of Compact Binary Mergers Observed by LIGO and Virgo during the First and Second Observing Runs arxiv.org/abs/1811.12907



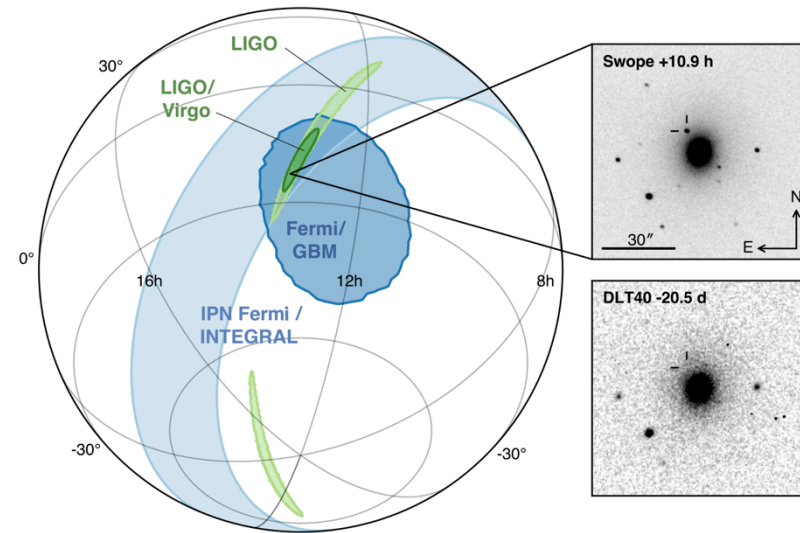
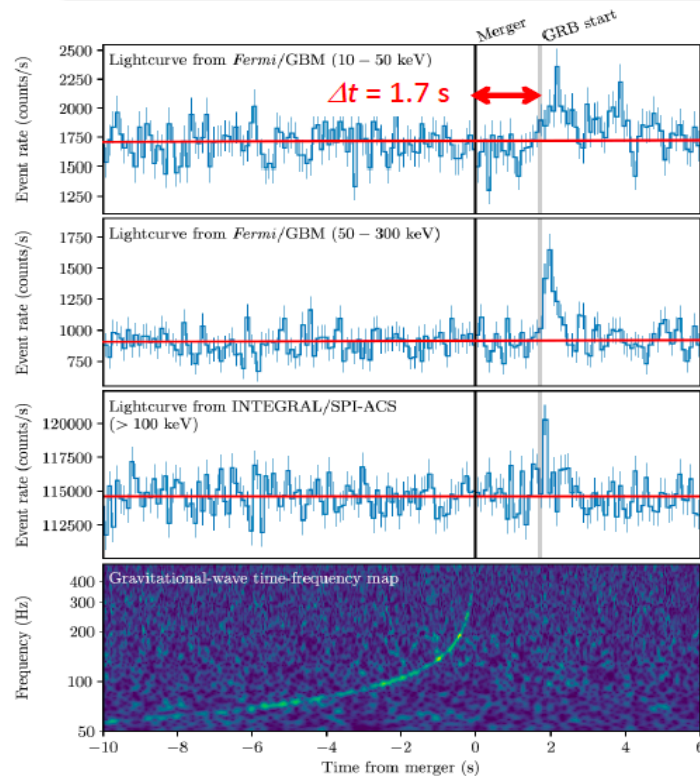
I. Gravitational waves in a nutshell



Why real time analysis

Use of Machine learning

17 August 2017, 12:41:04 UT The MultiMessenger Astronomy



[DOI:10.1103/PhysRevLett.119.161101](https://doi.org/10.1103/PhysRevLett.119.161101)

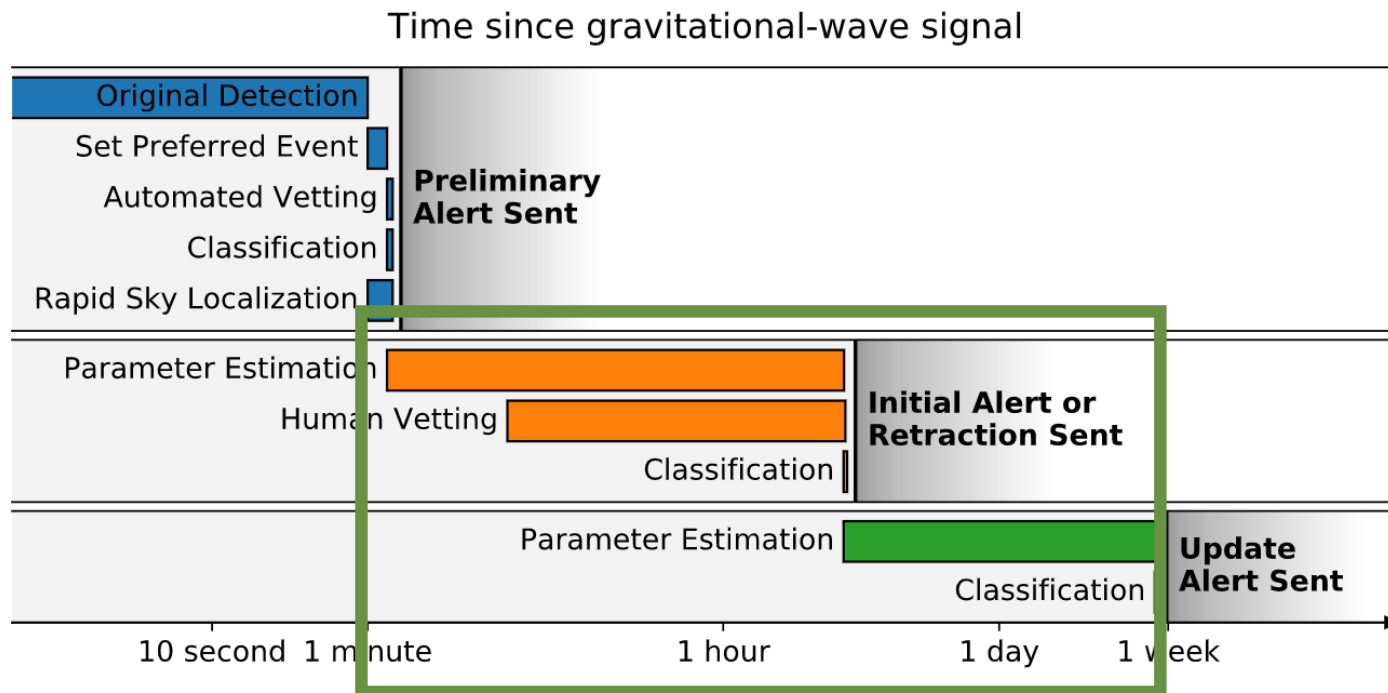
II. Why real time analysis

Low Latency data analysis



II. Why real time analysis

GW alert system



<https://emfollow.docs.ligo.org/userguide/index.html>

O3 event rate ~1/week. Up to now 32 events!

GraceDB — Gravitational Wave Candidate Event Database

HOME	SEARCH	LATEST	DOCUMENTATION	LOGIN
------	--------	--------	---------------	-------

Latest — as of 8 July 2019 13:15:27 UTC

Test and MDC events and superevents are not included in the search results by default; see the [query help](#) for information on how to search for events and superevents in those categories.

Query:

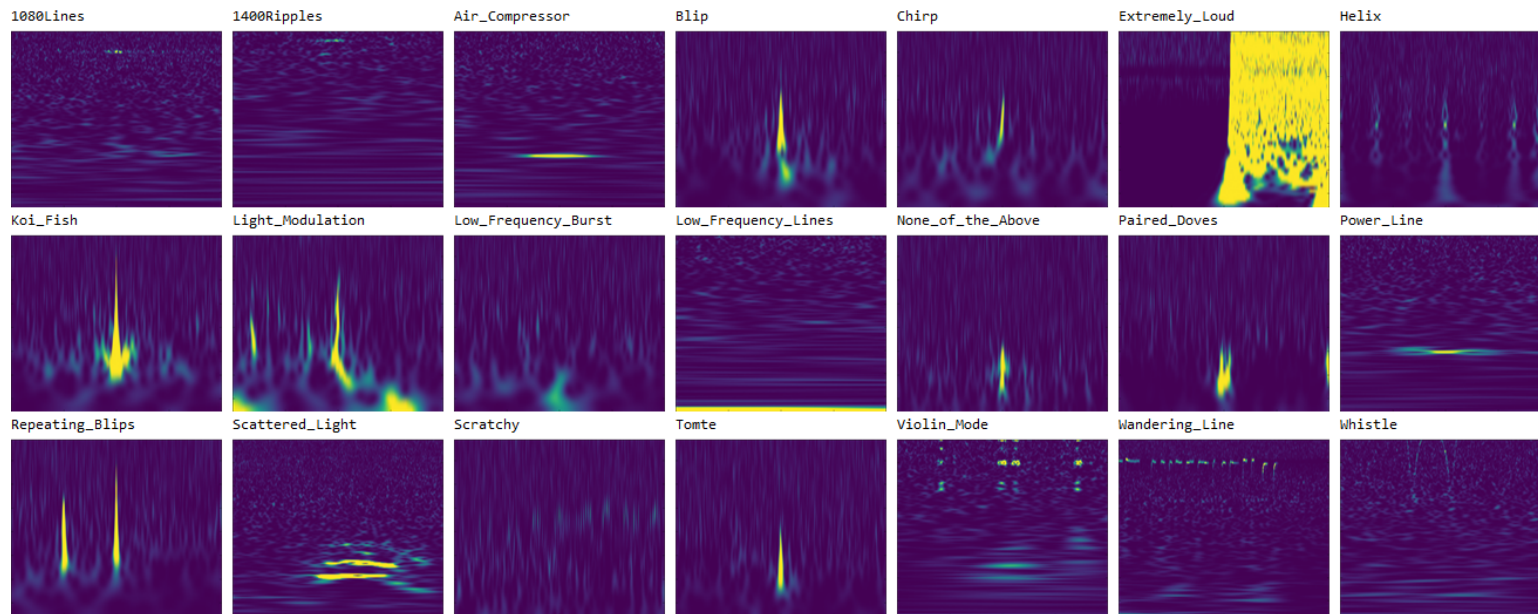
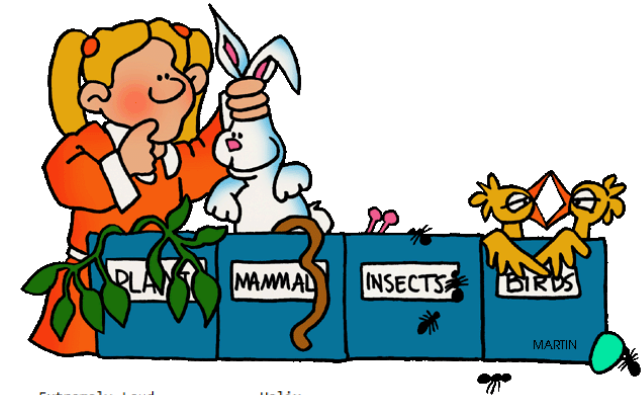
Search for:

UID	Labels	t_start	t_0	t_end	FAR (Hz)	<div>UTC Created</div>
S190707q	ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1246527223.118398	1246527224.181226	1246527225.284180	5.265e-12	2019-07-07 09:33:44 UTC
S190706ai	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1246487218.321541	1246487219.344727	1246487220.585938	1.901e-09	2019-07-06 22:26:57 UTC
S190701ah	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1246048403.576563	1246048404.577637	1246048405.814941	1.916e-08	2019-07-01 20:33:24 UTC
S190630ag	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1245955942.175325	1245955943.179550	1245955944.183184	1.435e-13	2019-06-30 18:52:28 UTC
S190602aq	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1243533584.081266	1243533585.089355	1243533586.346191	1.901e-09	2019-06-02 17:59:51 UTC
S190524q	ADVNO SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1242708743.678669	1242708744.678669	1242708746.133301	6.971e-09	2019-05-24 04:52:30 UTC
S190521r	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1242459856.453418	1242459857.460739	1242459858.642090	3.168e-10	2019-05-21 07:44:22 UTC
S190521g	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1242442966.447266	1242442967.606934	1242442968.888184	3.801e-09	2019-05-21 03:02:49 UTC
S190519bj	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1242315361.378873	1242315362.655762	1242315363.676270	5.702e-09	2019-05-19 15:36:04 UTC
S190518bb	ADVNO SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1242242376.474609	1242242377.474609	1242242380.922655	1.004e-08	2019-05-18 19:19:39 UTC
S190517h	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1242107478.819517	1242107479.994141	1242107480.994141	2.373e-09	2019-05-17 05:51:23 UTC
S190513bm	ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1241816085.736106	1241816086.869141	1241816087.869141	3.734e-13	2019-05-13 20:54:48 UTC
S190512at	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1241719651.411441	1241719652.416286	1241719653.518066	1.901e-09	2019-05-12 18:07:42 UTC
S190510g	ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1241492396.291636	1241492397.291636	1241492398.293185	8.834e-09	2019-05-10 03:00:03 UTC
S190503bf	ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1240944861.288574	1240944862.412598	1240944863.422852	1.636e-09	2019-05-03 18:54:26 UTC
S190426c	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1240327332.331668	1240327333.348145	1240327334.353516	1.947e-08	2019-04-26 15:22:15 UTC
S190425z	ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK	1240215502.011549	1240215503.011549	1240215504.018242	4.538e-13	2019-04-25 08:18:26 UTC
S190421ar	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1239917953.250977	1239917954.409180	1239917955.409180	1.489e-08	2019-04-21 21:39:16 UTC
S190412m	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1239082261.146717	1239082262.222168	1239082263.229492	1.683e-27	2019-04-12 05:31:03 UTC
S190408an	PE_READY ADVOK SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK GCN_PRELIM_SENT	1238782699.268296	1238782700.287958	1238782701.359863	2.811e-18	2019-04-08 18:18:27 UTC
S190405ar	ADVNO SKYMAP_READY EMBRIGHT_READY PASTRO_READY DQOK	1238515307.863646	1238515308.863646	1238515309.863646	2.141e-04	2019-04-05 16:01:56 UTC

Transient Signal Classification using Machine Learning

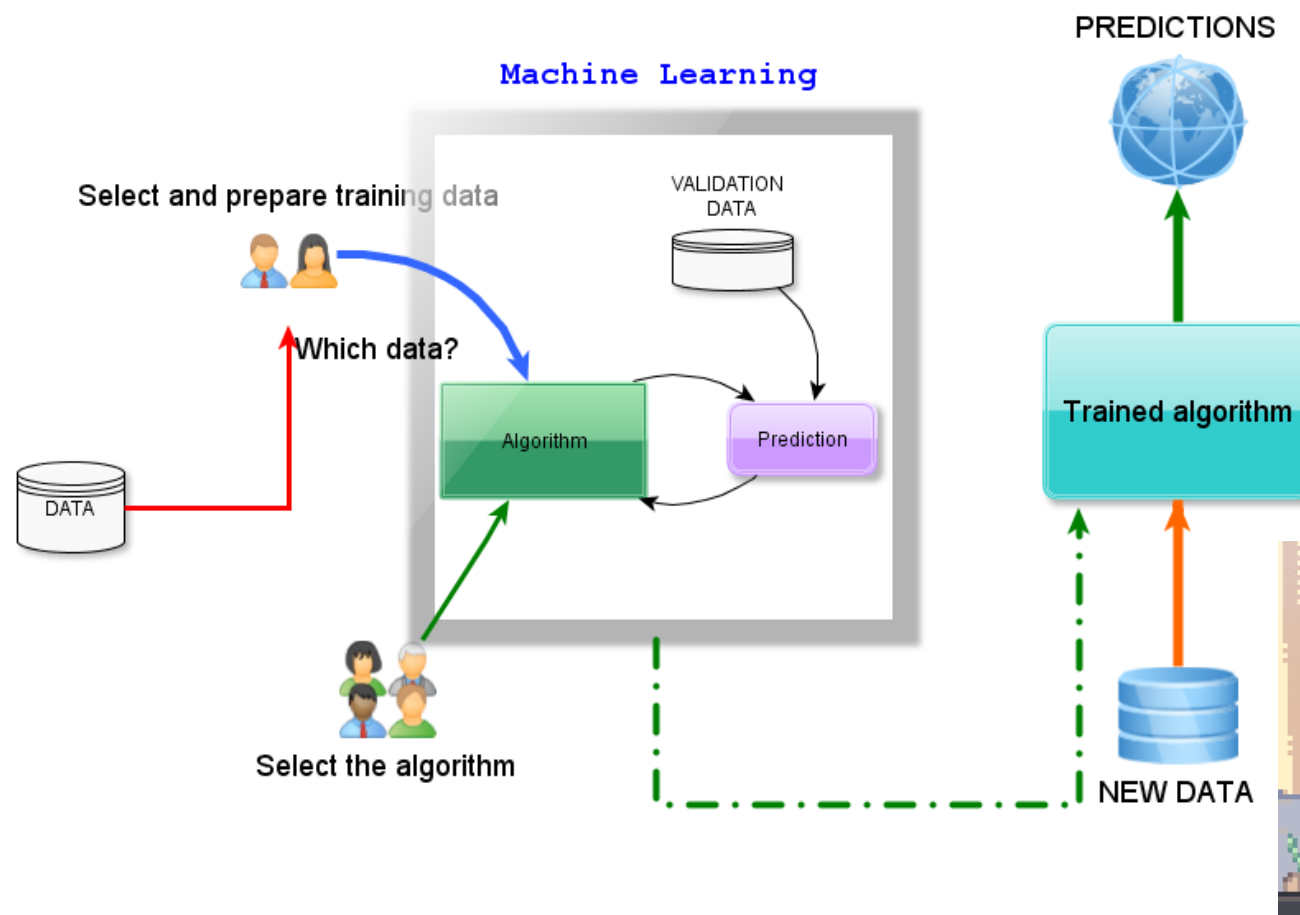
If we are able to classify the noise events,
we can clean the data in a fast way

Gravity Spy, Zevin et al (2017)



<https://www.zooniverse.org/projects/zooniverse/gravity-spy>

Artificial Intelligence workflow



II. Why real time analysis

Machine learning implemented pipeline

Supervised

- We need a **labeled** training data set
- We want to disentangles 'glitches' due to noise from 'transient signals' due to astrophysical signals

Unsupervised

- Extract signal features
- We try to catch main relations among the detected signals to create **cluster** of classes



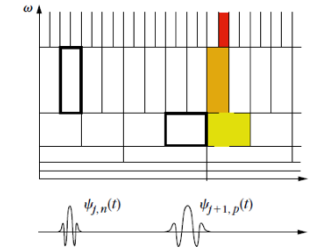
From In-time analysis to Wavefier prototype

ETG*: Wavelet Detection Filter (WDF) workflow

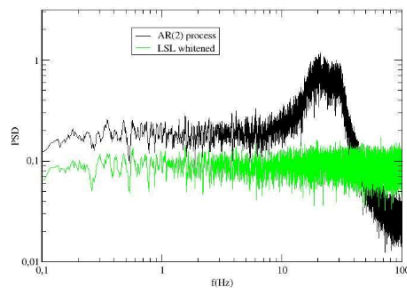
* ETG= Event Trigger Generator

$$x_i = h_i + n_i, \quad i = 0, 1, \dots, N-1,$$

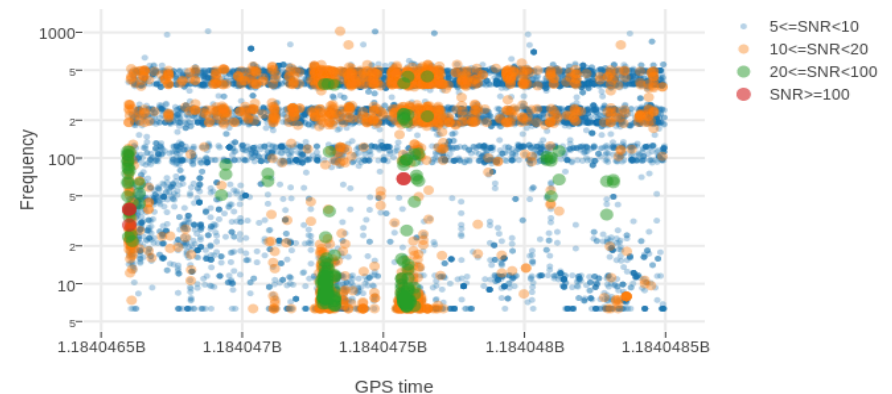
$$Wf(a, b) = \langle f, \psi_{a,b} \rangle = \int_{-\infty}^{+\infty} f(t) \frac{1}{\sqrt{b}} \psi^* \left(\frac{t-a}{b} \right) dt.$$



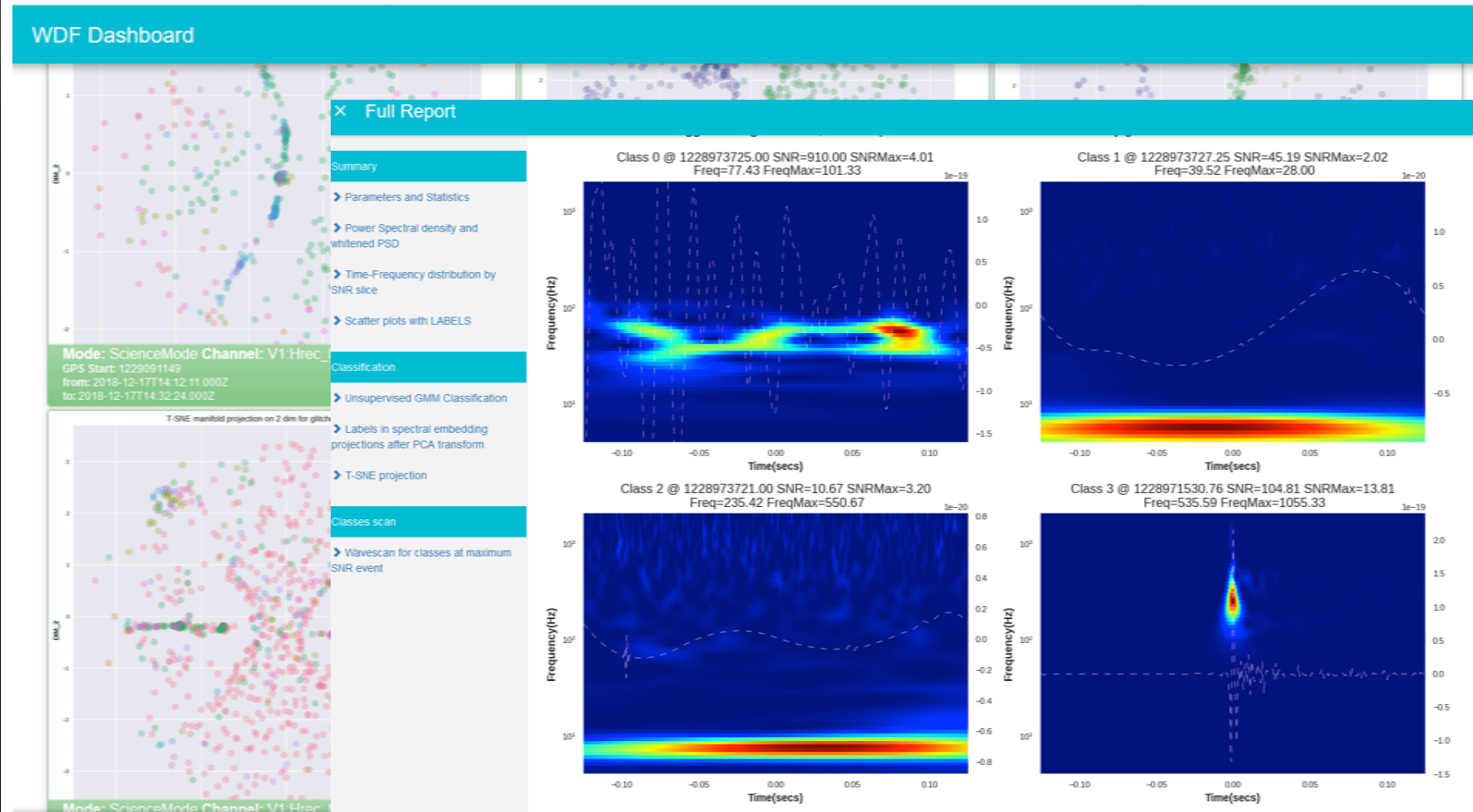
$$\hat{h}_i = W^{-1}(t[W(x_i)]).$$



V1:LSC_DARM: Time frequency glitchgram



In-Time Analysis running @VIRGO



<https://wdf.virgo-gw.eu/>

Why Wavefier and Key Objectives



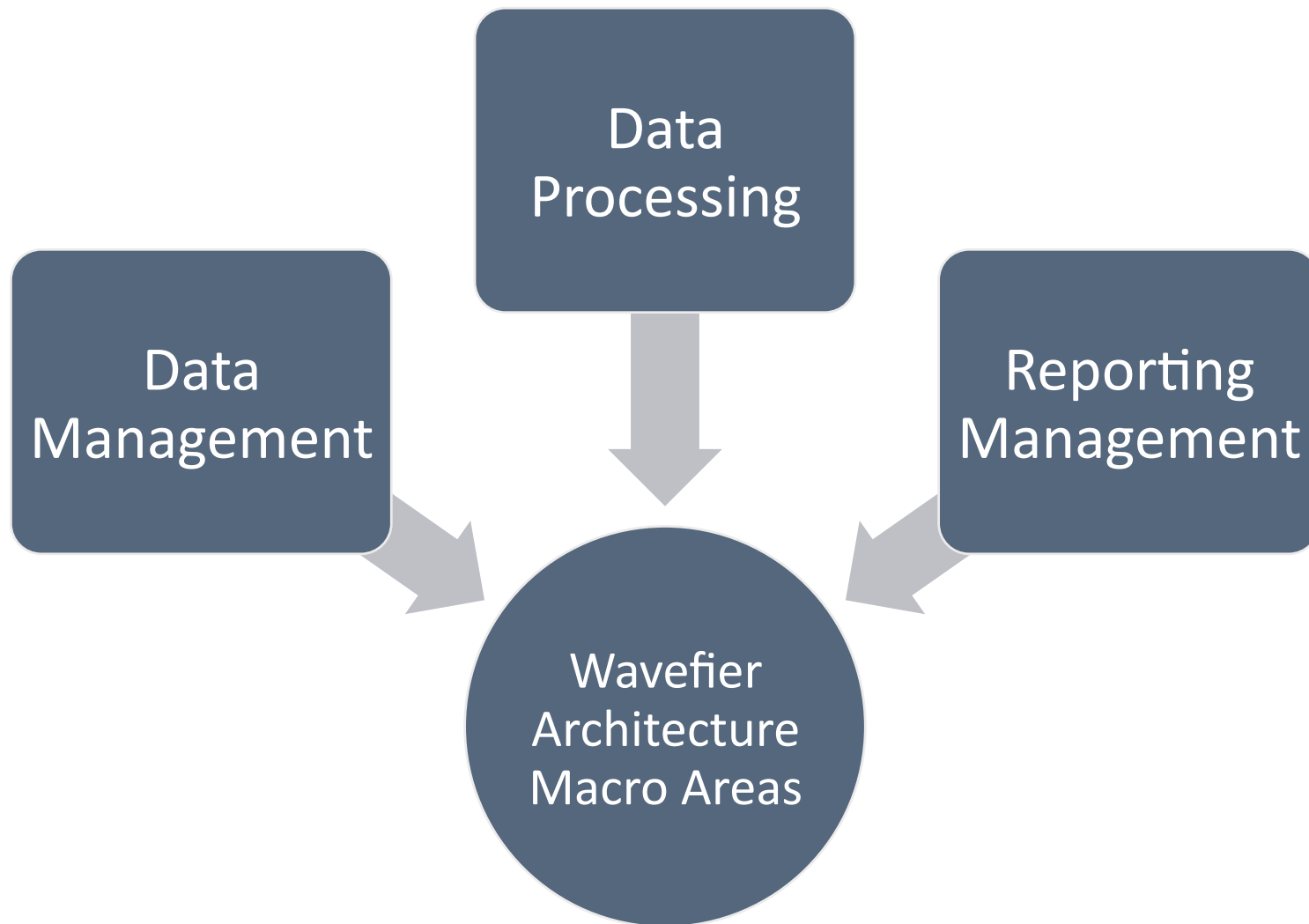
- Setup a **prototype** for a real time pipeline for the **detection of transient signals and their automatic classification**
- Best practices for **software management**
- Test different software architecture solutions to prototype a **scalable pipeline for big data analysis** in GW context.
- **Interoperability** and access to data and services
- **ICT services** supporting research infrastructures
- Use of data in **network infrastructures** and services
- Big data and **Machine Learning** solution easy to **plug-in**
- Test on **cluster** infrastructures



IV

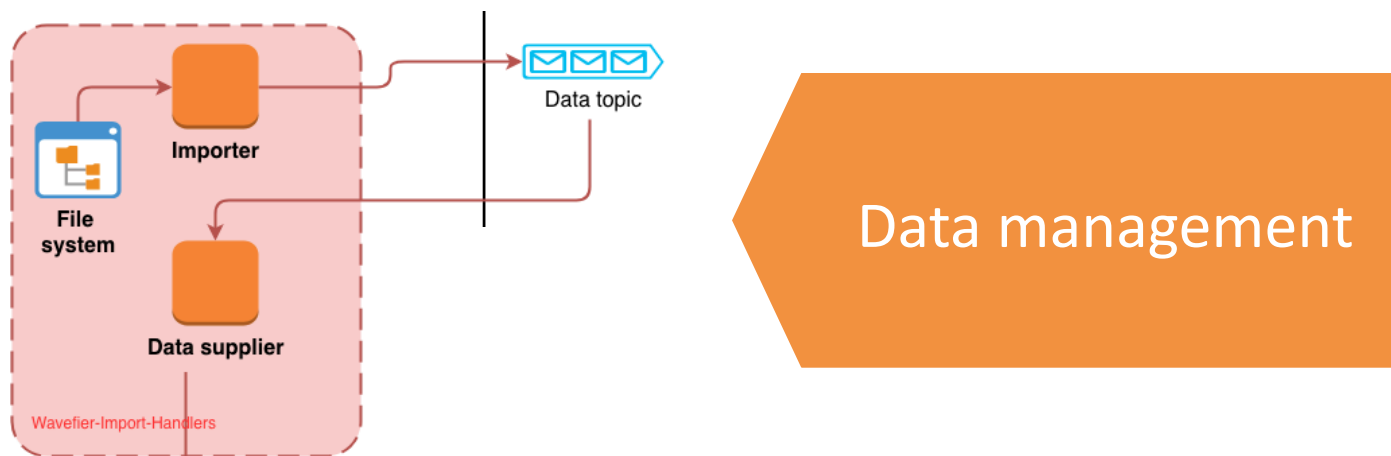
The technical and implementation aspects

Architecture Overview

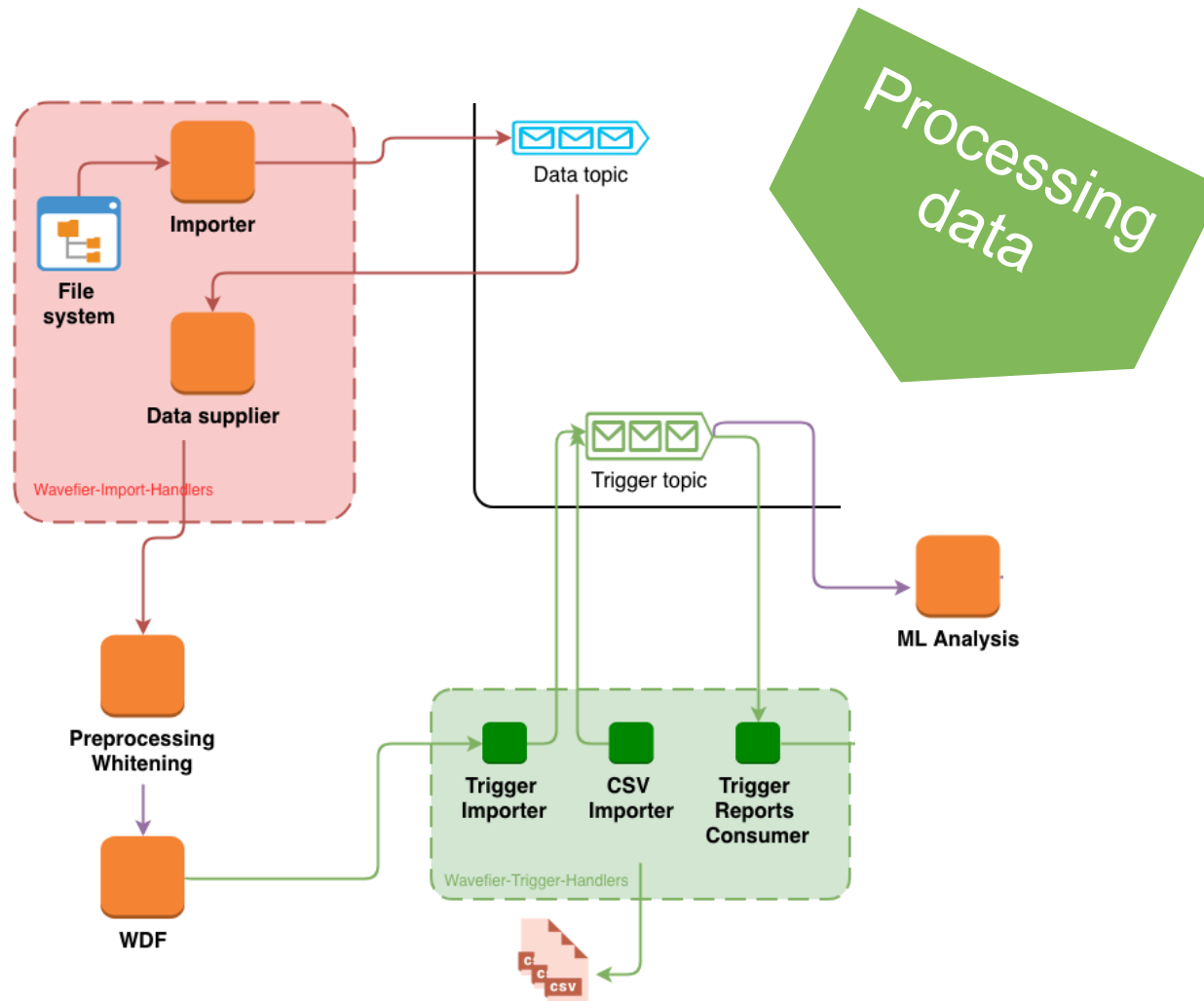


IV. The technical and implementation aspects

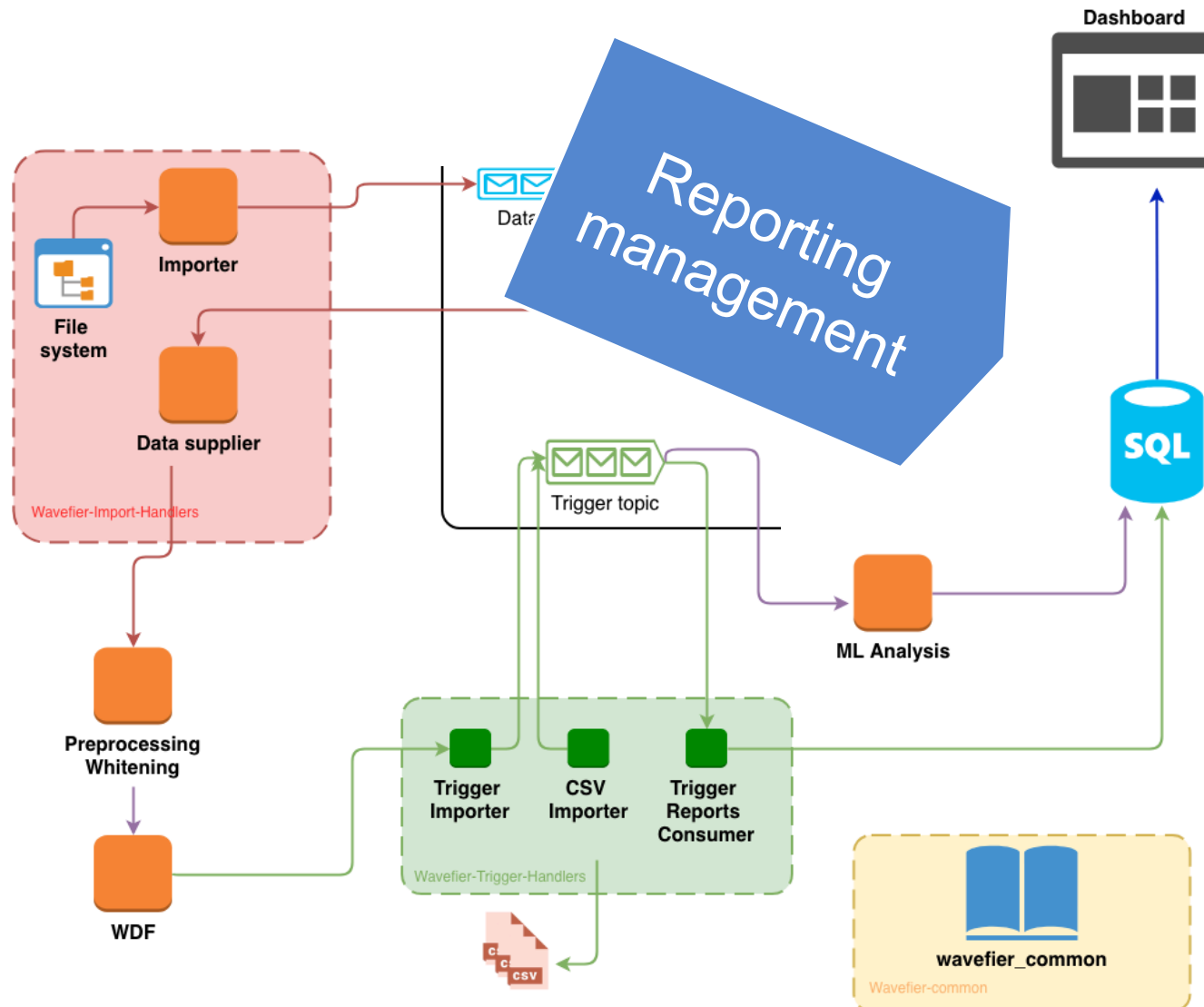
Architecture Overview – Data Injection



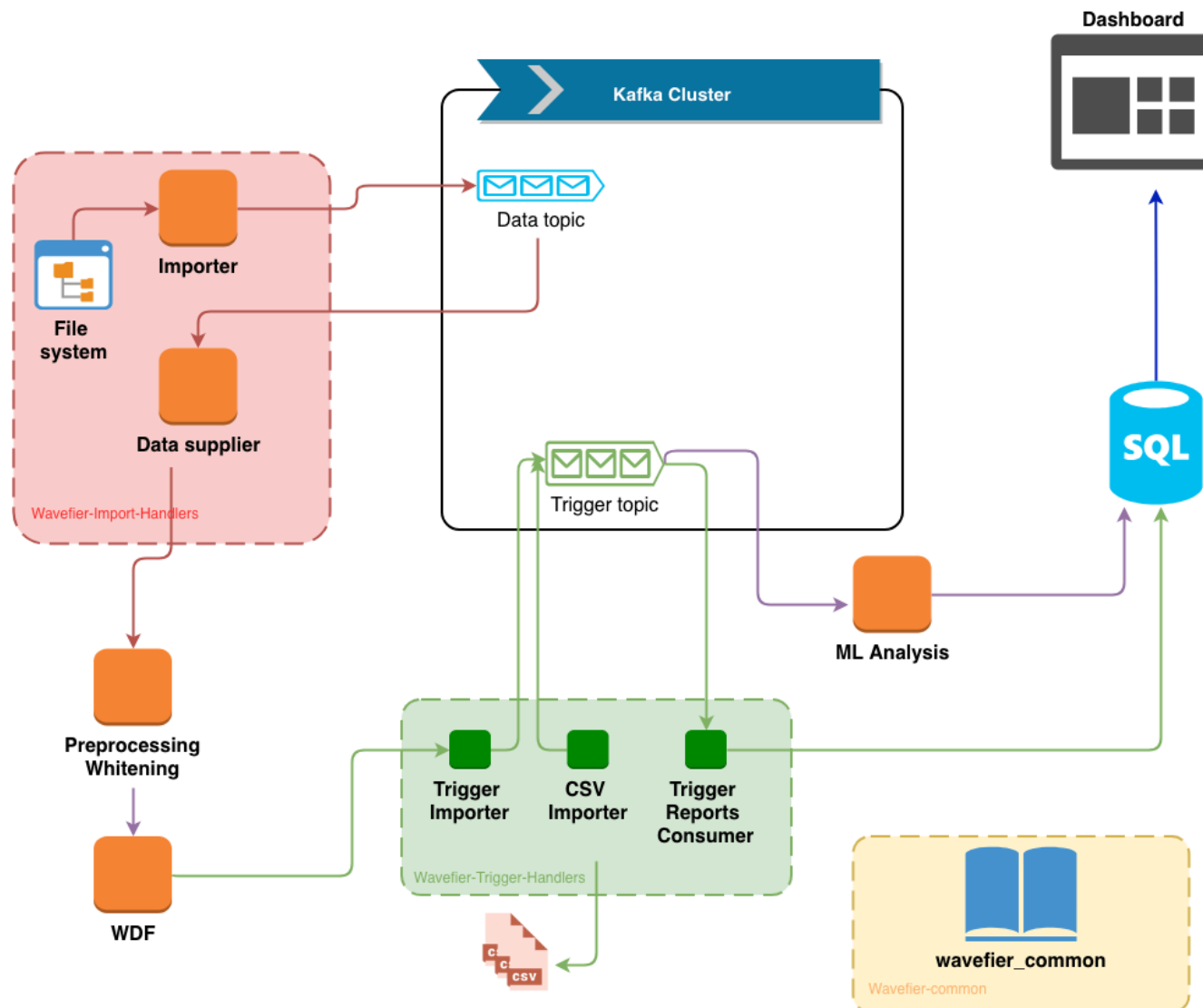
Architecture Overview – Processing Data



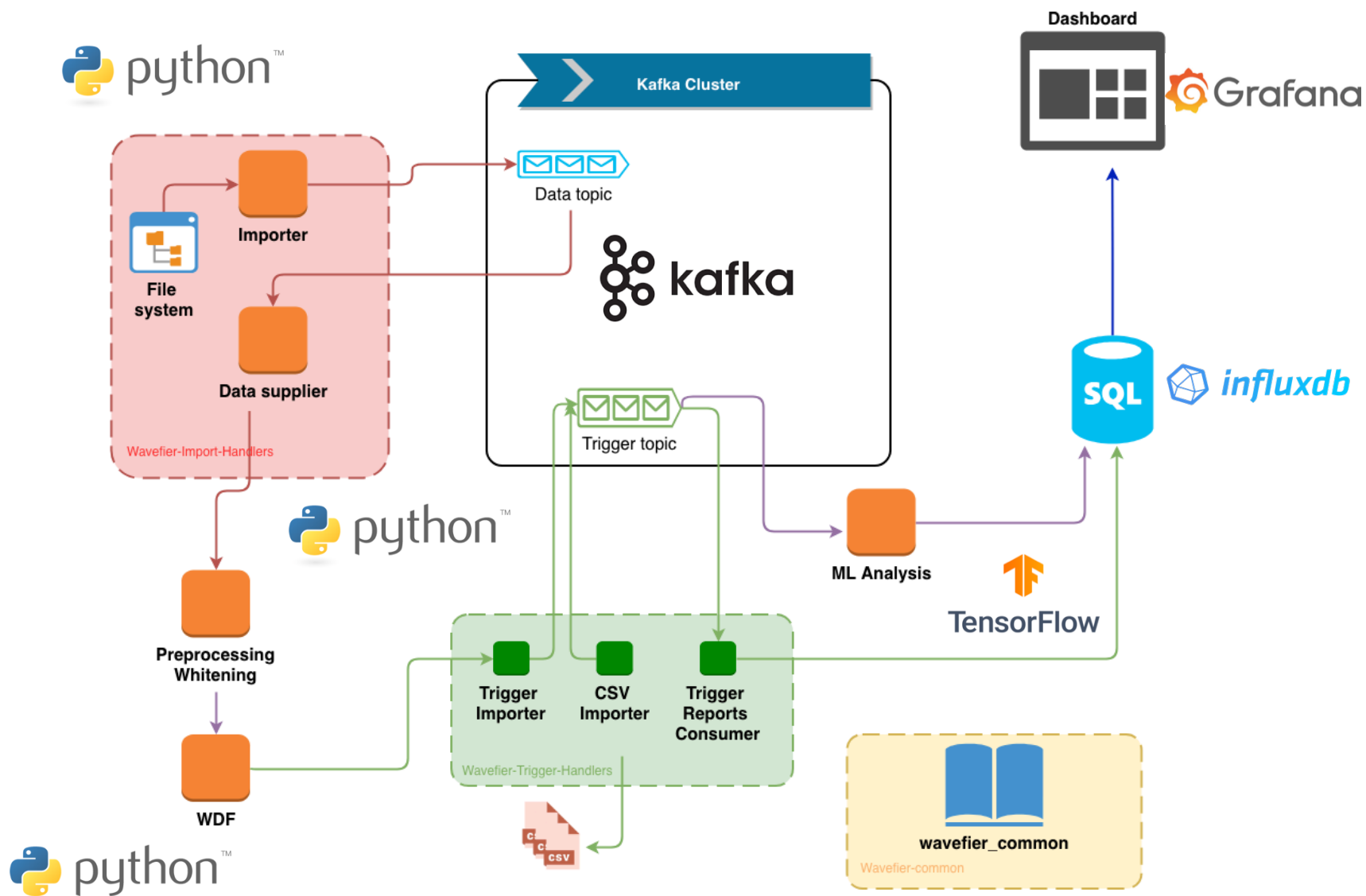
Architecture Overview – Reporting



Architecture Overview



Software Deployment



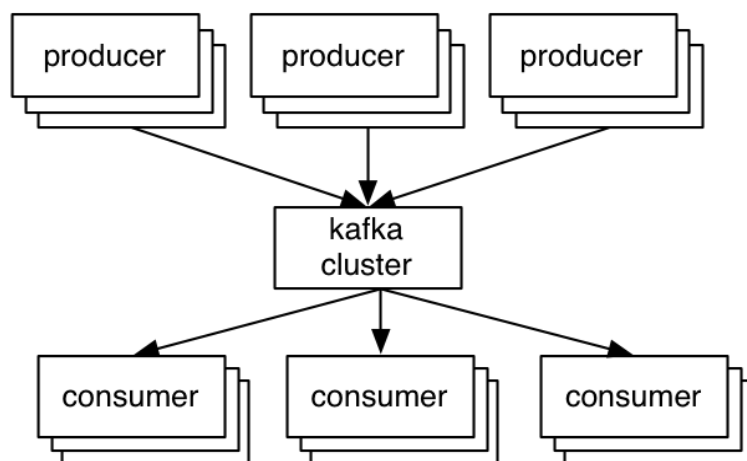


Apache Kafka

Open-source stream-processing software platform developed by LinkedIn and donated to the Apache Software Foundation

Apache Kafka® is a distributed streaming platform.

What exactly does that mean?



More info on:

<https://kafka.apache.org>

A streaming platform has three key capabilities :

- Publish and subscribe to streams of records
- similar to a message queue
- Store streams of records in a fault-tolerant durable way.
- Process streams of records as they occur.

Kafka is generally used for two broad classes of applications:

- Building real-time streaming data pipelines that reliably get data
- Building real-time streaming applications that transform or react to the streams of data

IV. The technical and implementation aspects





Grafana

“Grafana is the open source analytics & monitoring solution for every database “

Why grafana?

Useful build-in features

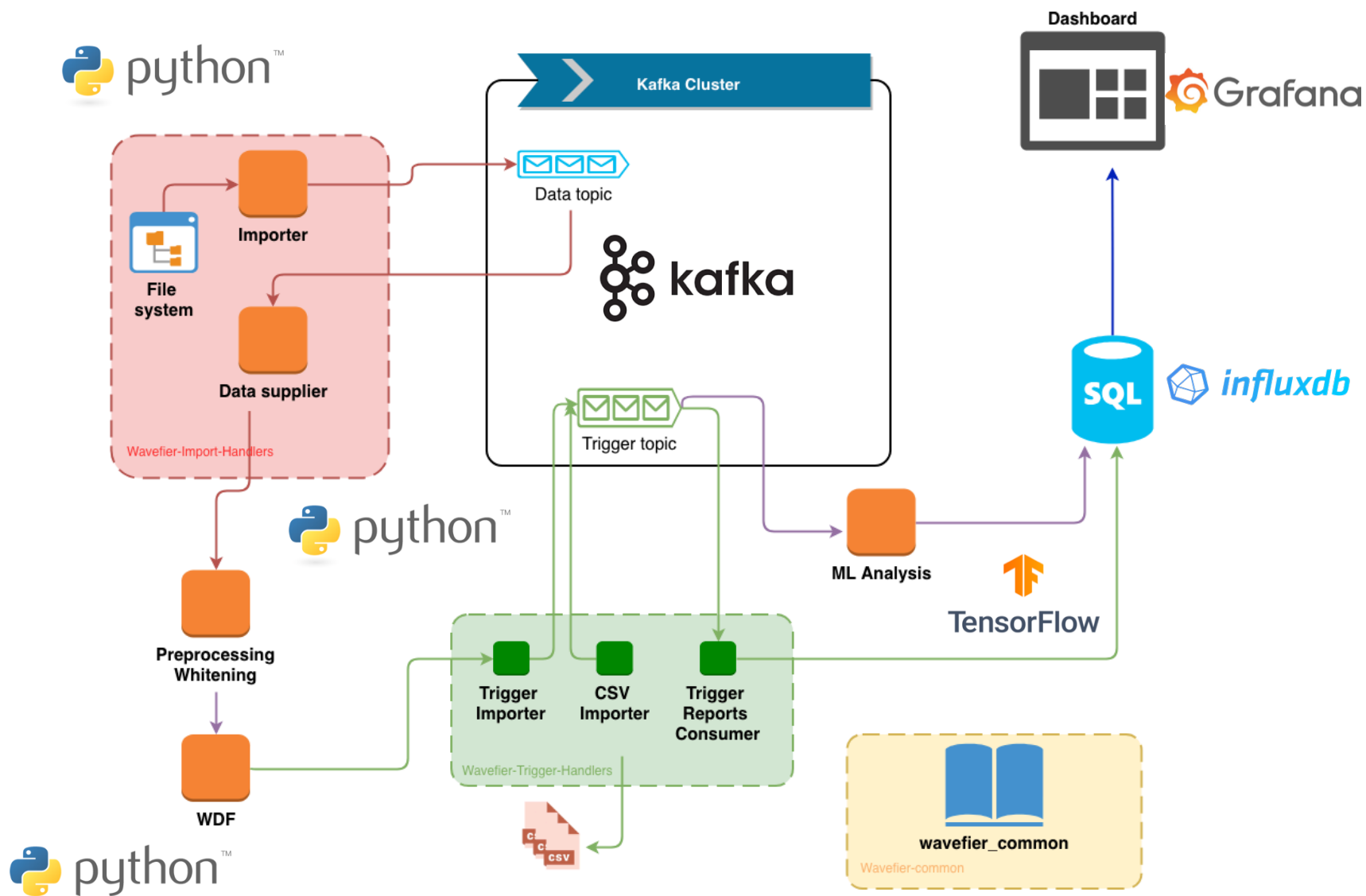
- Authentication, Organization and user settings

Mixed Datasource, Mix different data sources in the same graph

- Grafana supports dozens of databases, natively. Mix them together in the same Dashboard.

Native Notification and Alerting

Software Deployment





InfluxDB

“InfluxDB is a time series database designed to handle high write and query loads.”

Why InfluxDB?

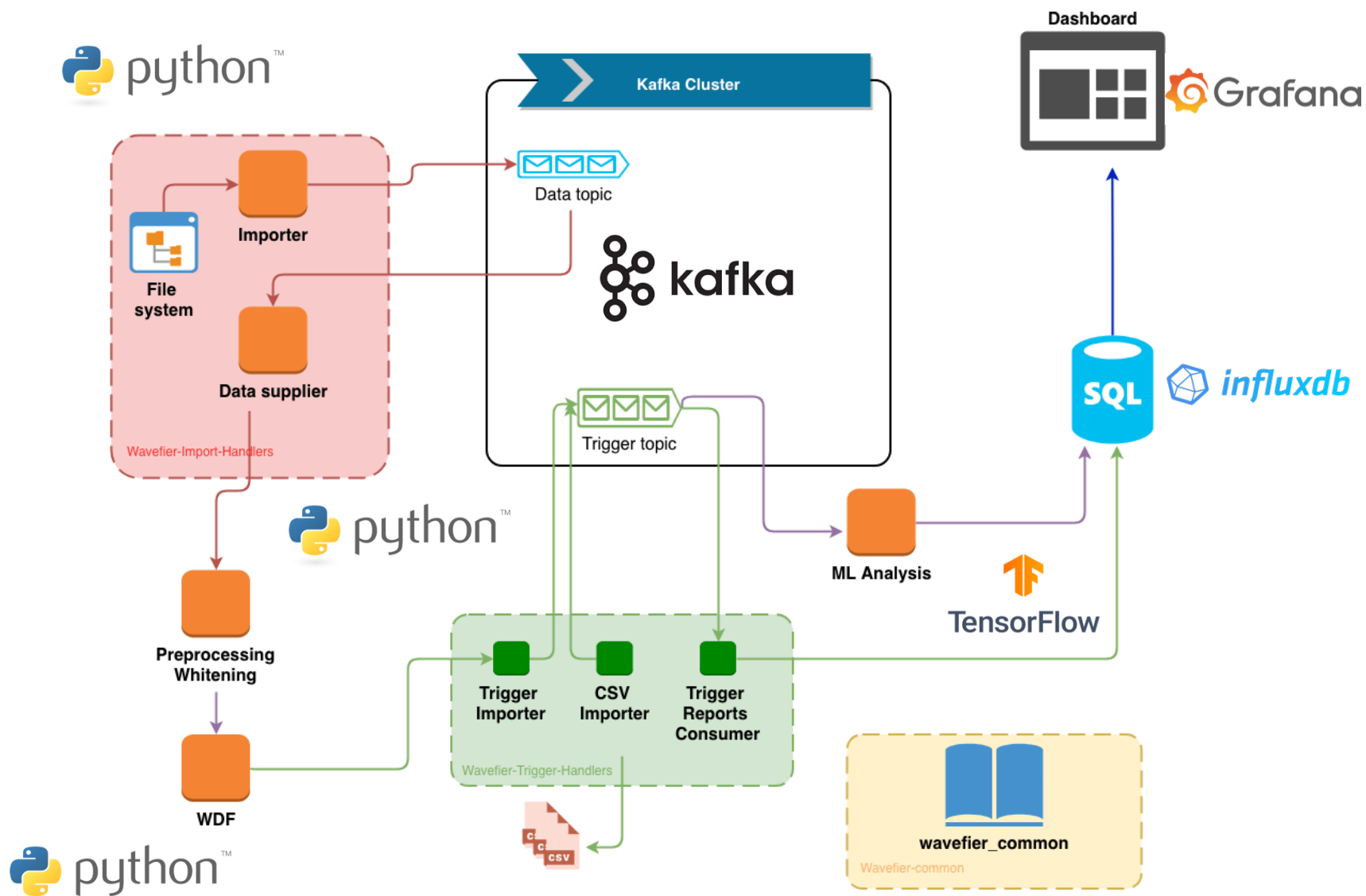
Specific for time series database (TSDB)

All is designed as time series

Friendly because InfluxDB have a SQL-like query language for interacting with it

Grafana has Native support for InfluxDB

Software Deployment





Software Management

- Distribuite Team: Trust-it, EGO
- Different expertises and point of view: Physics, Software Engineer and Computer science
- Different Environment
- One unique target... the prototype

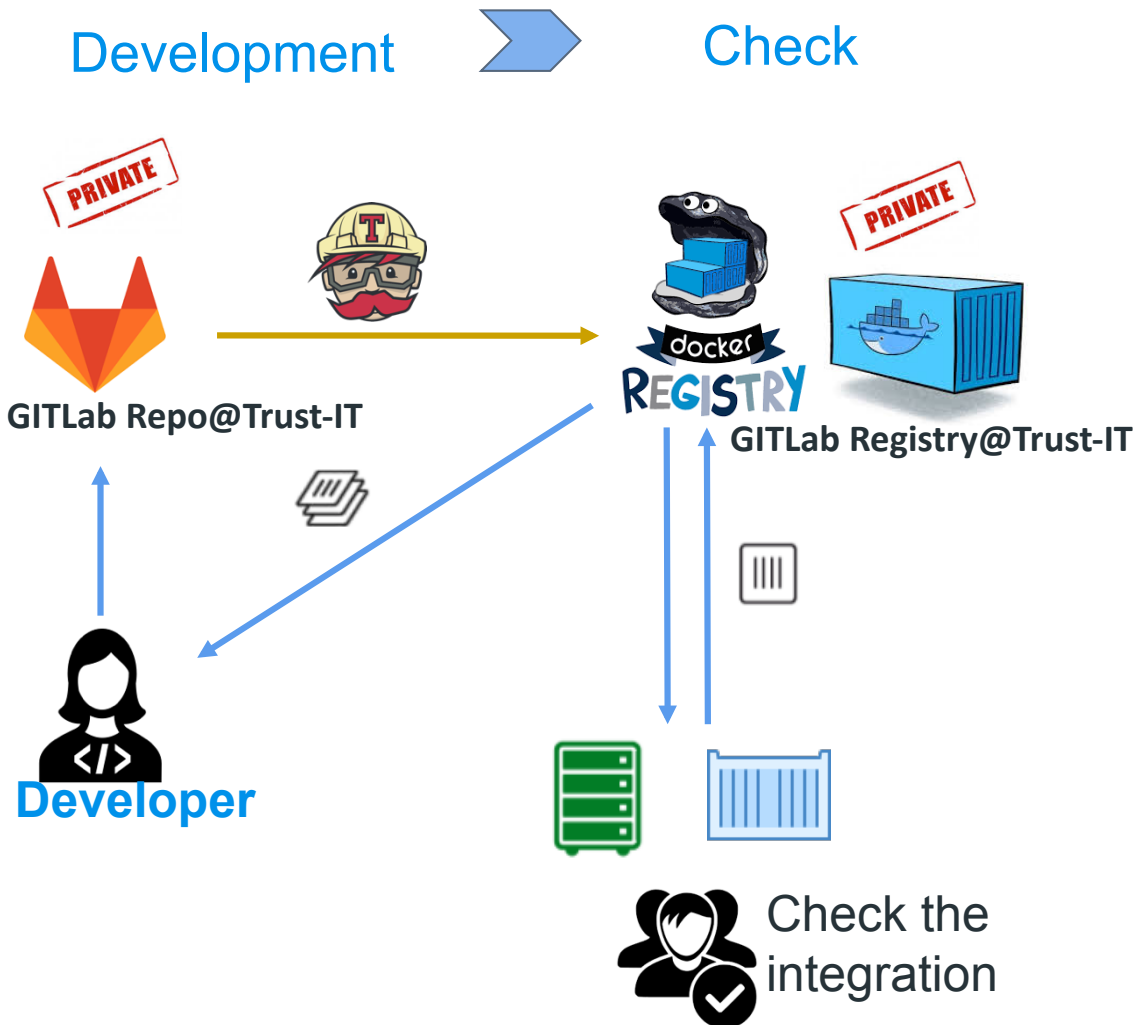
We use Docker!

Docker is an open platform for developers and sysadmins to build, ship, and run distributed applications.



Docker take the concept of container and build an ecosystem around it that would simplify its use

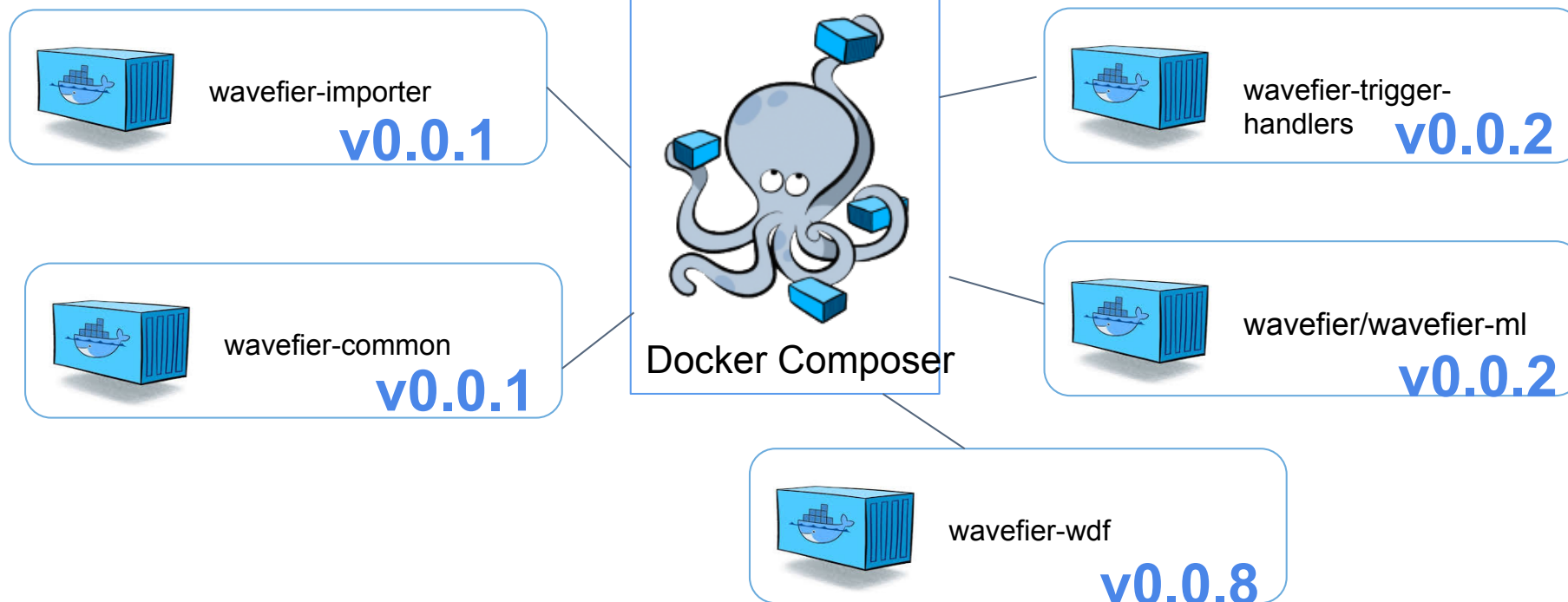
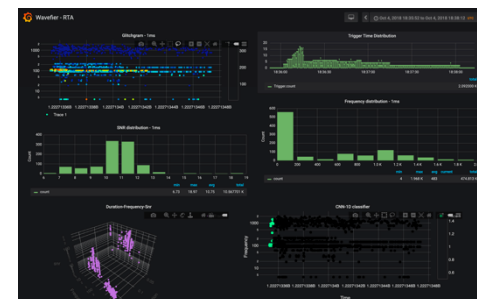
Automation and Continuous Integration



Check the integration - Local Deployment

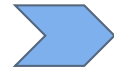


```
#> docker-compose up
wavefier-wdf
```



Automation and Continuous Integration

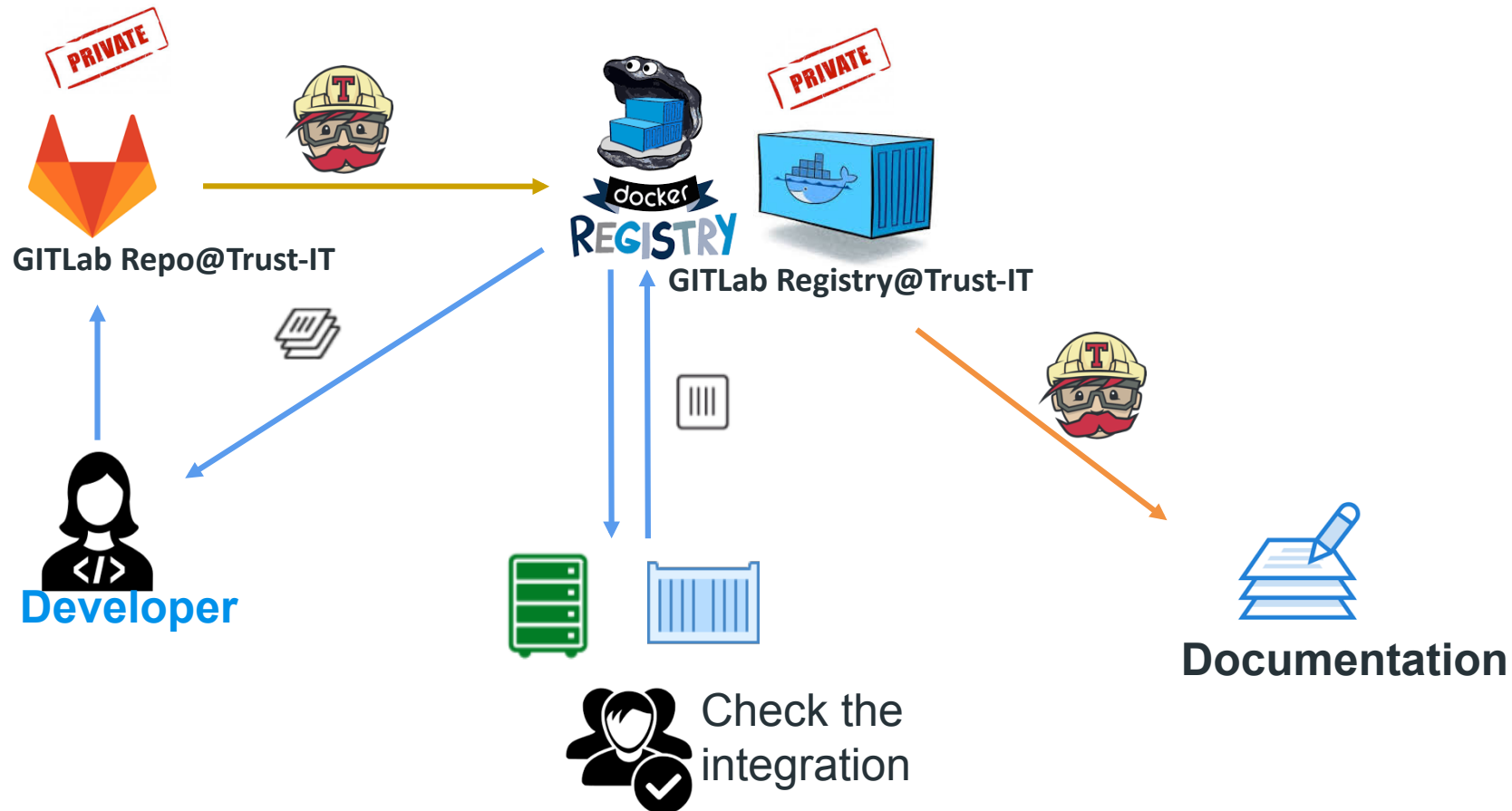
Development



Check

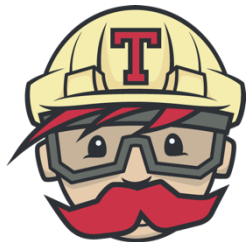


Deployment



IV. The technical and implementation aspects

Automation on Documentation



The Documentation is also generated foreach git Commit

The screenshot displays the Wavefier documentation website. The top page is titled 'TriggerAdapter' and the bottom page is titled 'Introduction'. The 'Introduction' page contains a diagram illustrating the system architecture:

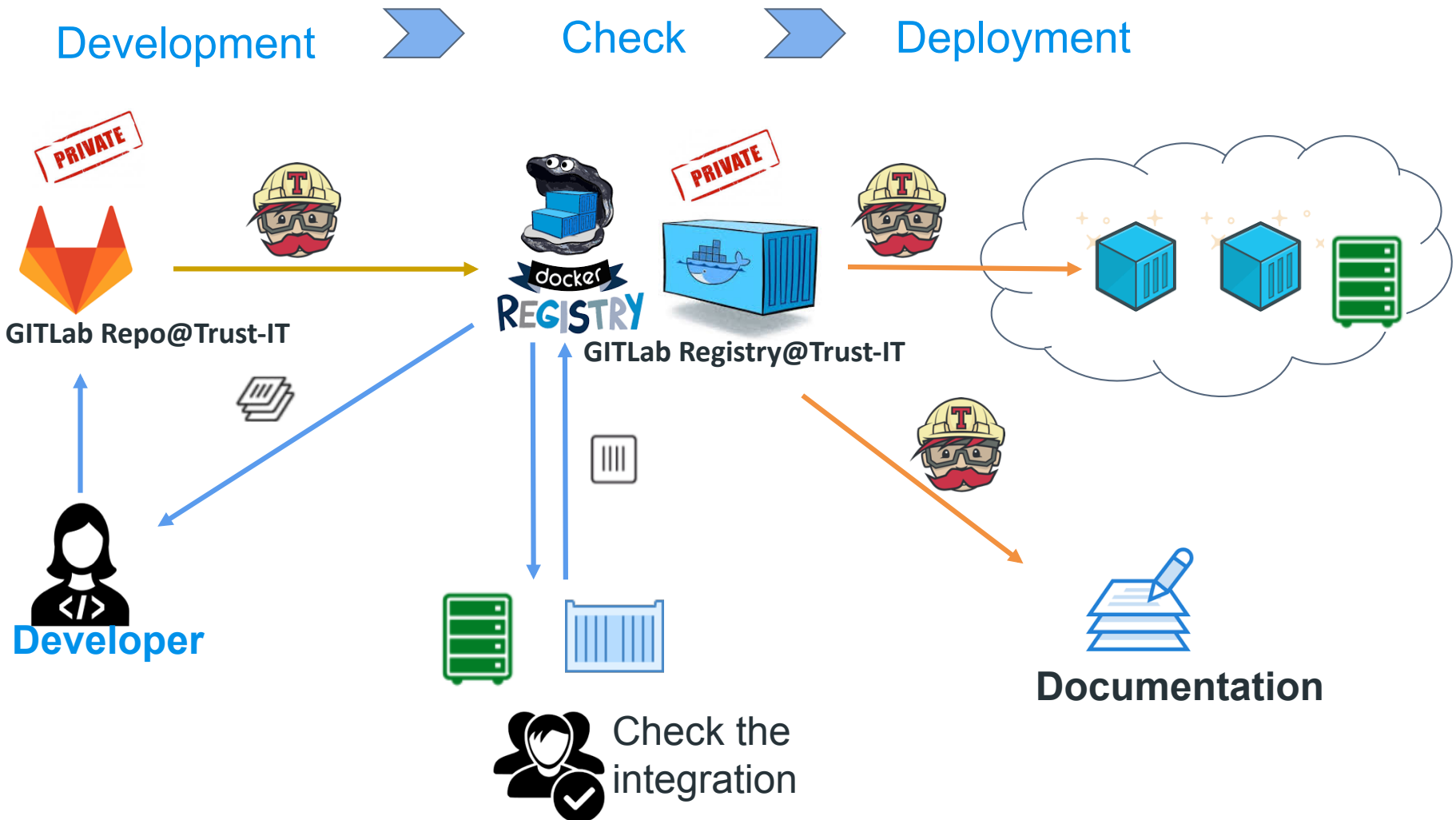
```

graph TD
    CLI["#> CSVImporter  
Command Line Interface"]
    FW["TriggerSupplier.py  
TriggerDistributor.py  
Framework Interface"]
    RI["Grafana  
influxdb  
Report Interface"]
    K["APACHE kafka"]
    ID["InfluxDB Importer Daemon"]

    CLI --> FW
    FW <--> RI
    K --> ID
    ID --> RI
  
```

The diagram shows the flow of data from the Command Line Interface (CLI) through the Framework Interface to the Report Interface (Grafana and influxdb). It also shows the integration with Apache Kafka and the InfluxDB Importer Daemon.

Automation and Continuous Integration





Kubernetes (K8s)

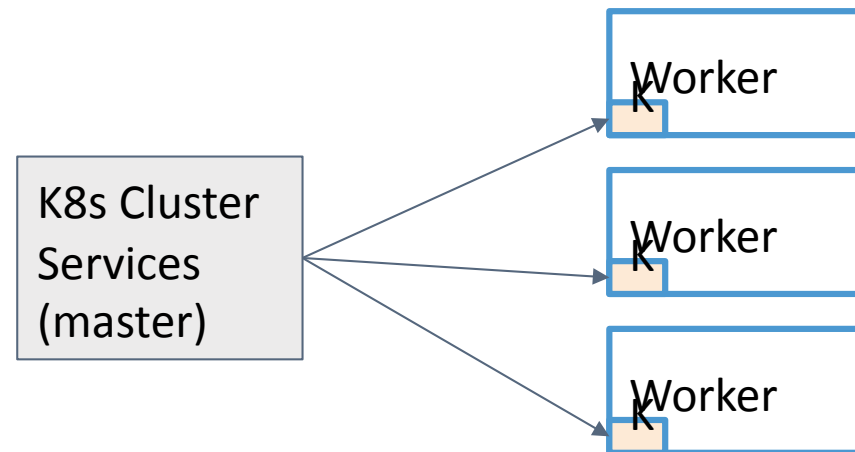
was a project spun out of Google as a open source next-gen container scheduler, designed as a loosely coupled collection of components centered around deploying, maintaining, and scaling applications.

Architecture overview

Kubernetes abstracts away the underlying hardware of the nodes and provides a uniform interface for applications to be both deployed and consume the shared pool of resources.

Masters

are responsible at a minimum for running the API Server, scheduler, and cluster controller. They commonly also manage storing cluster state, cloud-provider specific components and other cluster essential services.



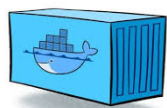
Nodes

Are the 'workers' of a Kubernetes cluster. They run a minimal agent that manages the node itself, and are tasked with executing workloads as designated by the master.

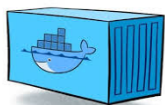
Do you remember local deployment?



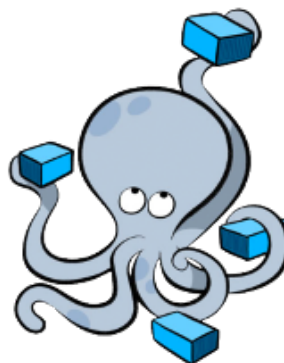
```
#> docker-compose up
wavefier-wdf
```



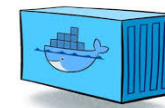
wavefier-importer
v0.0.1



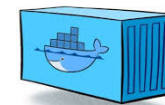
wavefier-common
v0.0.1



Docker Composer



wavefier-trigger-handlers
v0.0.2

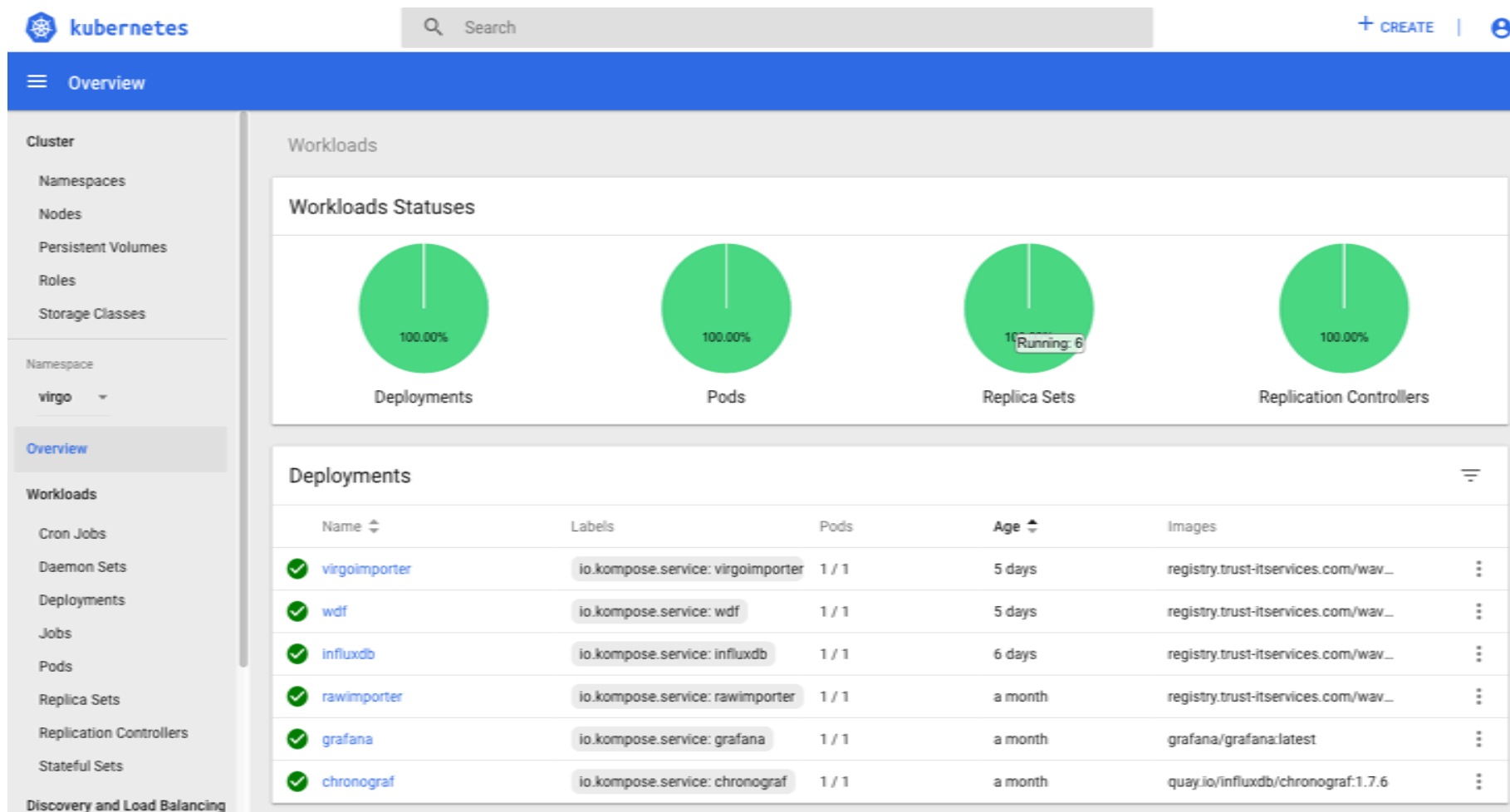


wavefier/wavefier-ml
v0.0.2



wavefier-wdf
v0.0.8

WaveFier running on Kubernetes



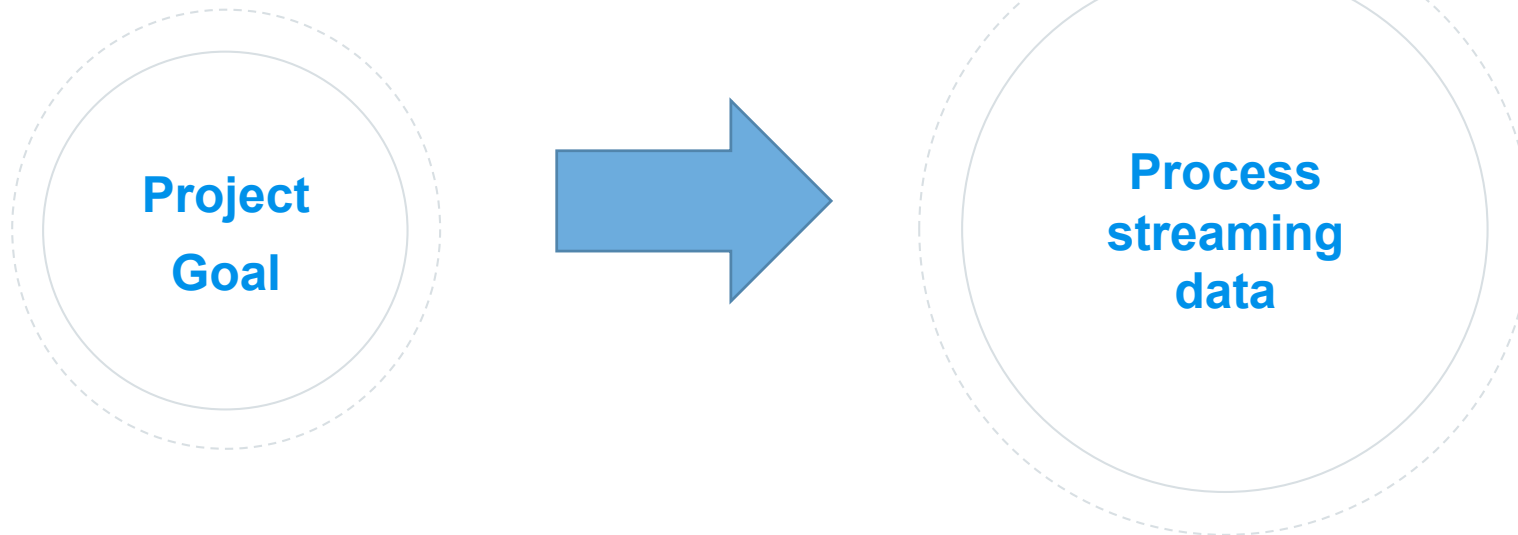
From Offline data to Online data

Offline data

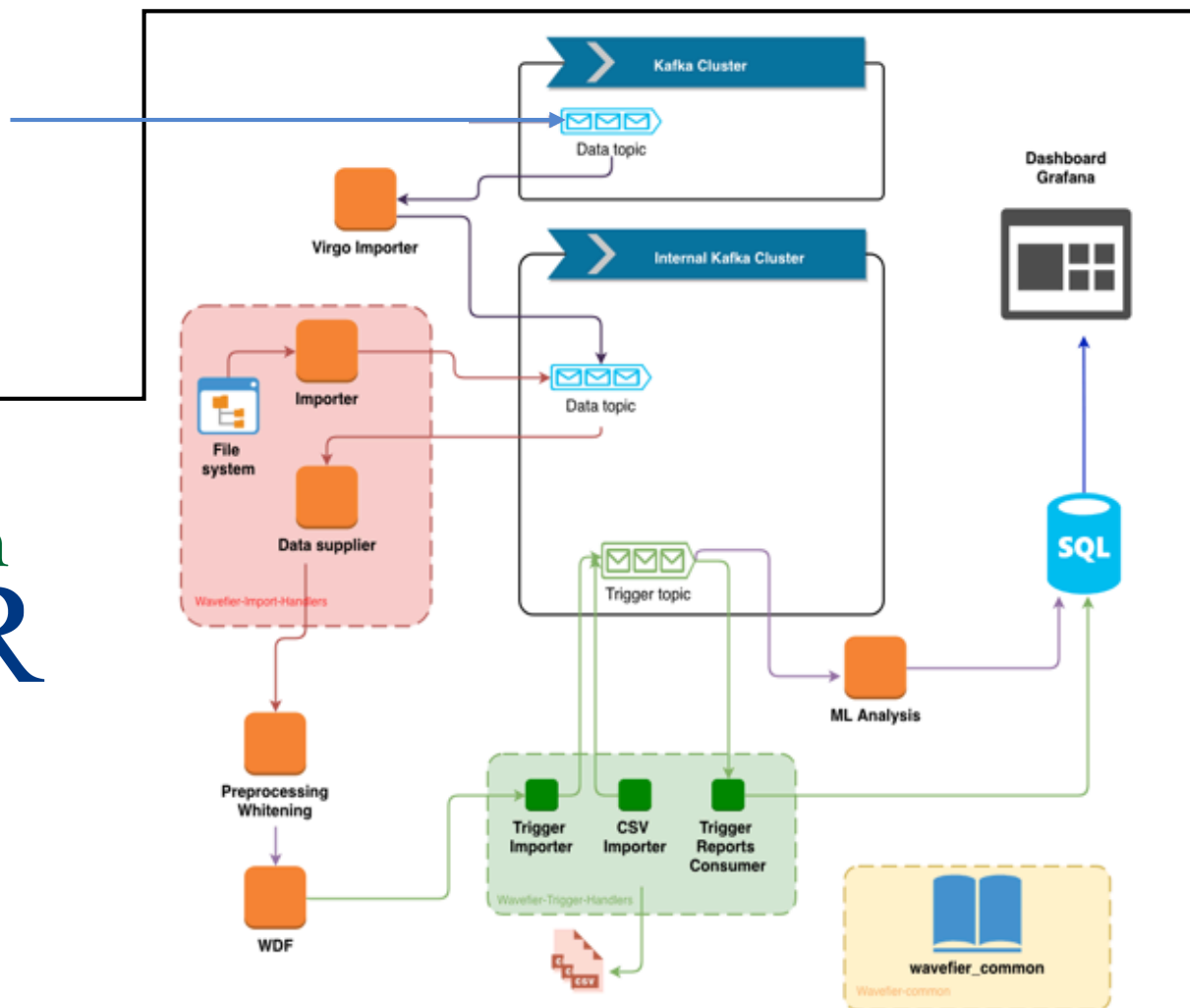
- ⊙ Pick-up interferometer data
- ⊙ Store data in files
- ⊙ Access to cluster
- ⊙ Move data in cluster

Online data

- ⊙ Receiving data from different sources
- ⊙ Receive data from interferometer



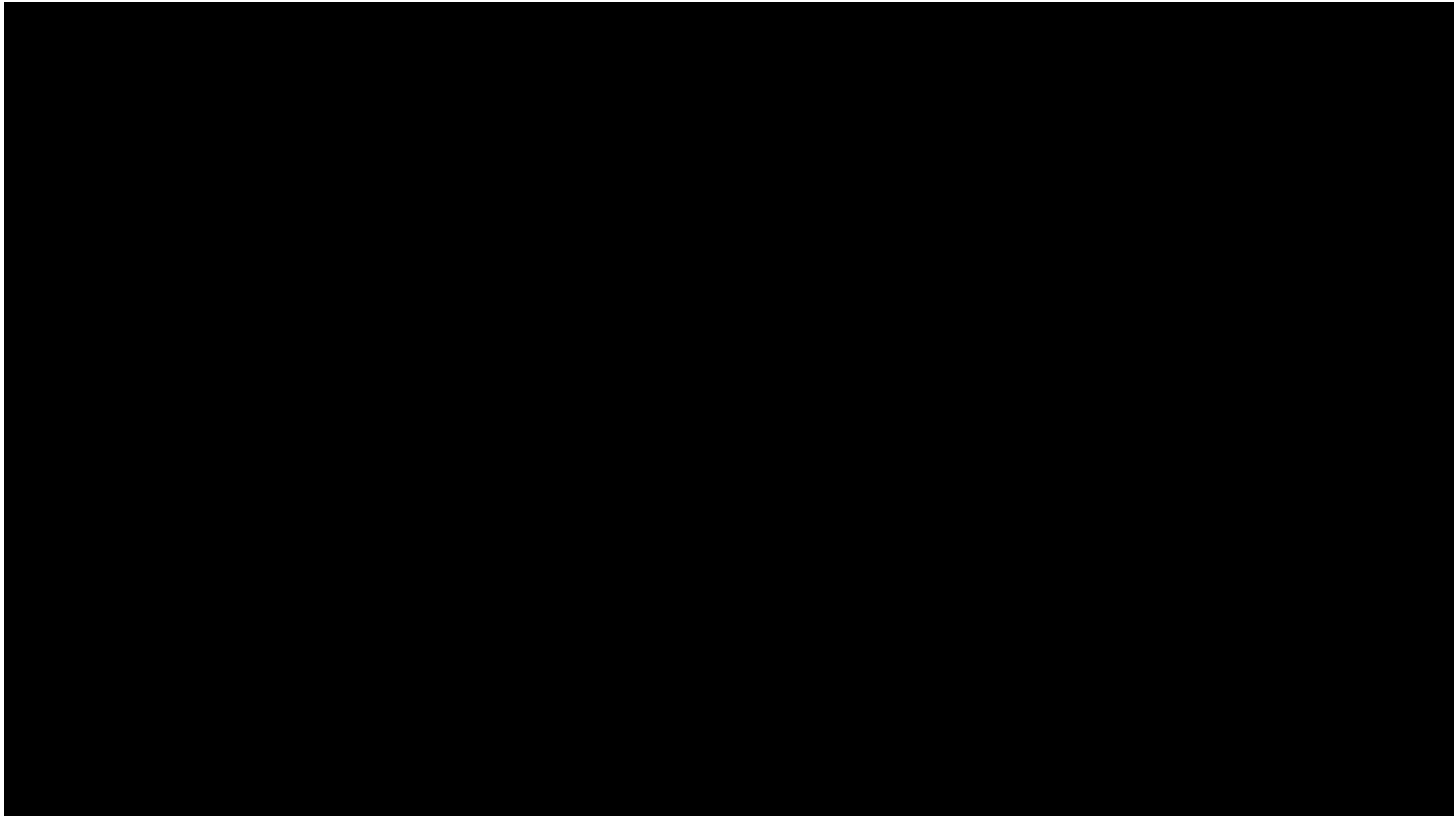
Online data - Deployment



Current Grafana Dashboard with Classification Results

Don

IV. The technical and implementation aspects



**V**

Next step?

Move from prototype/development to production

Investigate the use of much more sophisticated Machine/Deep learning algorithm, using GPU

Wavefier' repos

DOI: 10.5281/zenodo.3356656

August 1, 2019

Software Open Access

A prototype for a real time pipeline for the detection of transient signals and their automatic classification

Elena Cuoco; Emanuel Marzini; Filip Morawski; Alessandro Petrocelli; Alessandro Staniscia

WaveFier is the result of an industrial collaboration project with Trust-IT Services LTD Chase Side Enfield, Middlesex - EN2 6NF - UK and "CNRS - Center National de la Recherche Scientifique in Paris" acting in behalf of the "Laboratoire d'Annecy de physique des particules - LAPP UMR n. 51814" carried out in the context of the H2020 Asterics / Obelics project of the European union's Commission Programme.

The detection of gravitational waves has inaugurated the era of gravitational astronomy and opened new challenges for the multi-messenger study of cosmic sources. Thanks to their sensitivity, the Advanced LIGO and Advanced Virgo interferometers will probe a much larger volume of space and expand the capability of discovering new gravitational wave emitters. However, noise identification and its removal remains one of the most challenging problem in GW data analysis. A single GW detector typically produces data with a rate of 7-8 Tb per day with a flux of 40Mb/s. These data have to be analysed in the faster and most efficient way to increase the detection confidence and to obtain information in real time, about likely noise sources and to help the fast alert system for Electromagnetic Follow Up systems.

Glitches are transient noise events that can impact the data quality of the interferometers and their classification is an important task. Outlier/noise detection has been studied for decades in time-series analysis. ML methods generally employ a semi-supervised approach, with a few others using supervised or unsupervised techniques. Supervised ML techniques require a training phase with labeled data, in order to learn the data model which classifies outliers and inliers. This can be an expensive operation with massive data since generating a labeled training data set can be time consuming, especially if the data need to be labeled manually. Characterizing the glitches is an important task to reduce the impact of transient noise on the detectors. Inspecting glitches manually is a time-consuming and error-prone task. Furthermore, the increase of sensitivity in advanced detectors will lead to more classes of glitches. The use of machine learning looks a promising way to tackle the classification of glitches.

The goal of the project developed is setup a prototype for a real time pipeline for the detection of transient signals and their automatic classification. Moreover, the project has the goal to test different software architecture solutions to prototype a scalable pipeline for big data and deep learning analysis in GW context.

2 views

0 downloads

See more details...

Indexed in

Publication date:
August 1, 2019

DOI:
DOI: 10.5281/zenodo.3356656

Keyword(s):
Machine Learning, Big Data, Apache Kafka, Gravitational Wave

Grants:
European Commission
• ASTERICS - Astronomy ESFRI and Research Infrastructure Cluster (653477)

License (for files):
MIT License

Versions

Version 1.0	Aug 1, 2019
10.5281/zenodo.3356656	

<https://zenodo.org/record/3356656>

<https://repository.asterics2020.eu/content/wavefier>

<https://gitlab.in2p3.fr/escape2020/wavefier/dockers>

V. Next step?

LAPP and CNRS: Giovanni Lamanna, Jayesh Wagh

Trust-IT: Silvana Muscella, Emanuel Marzini, Filip Morawski,
Alessandro Petrocelli, Alessandro Staniscia
And Many thanks to **GARR CSD support**

CREDITS

