# Data Management Plan all'ONT Gruppo EIDA-UF3

Carluccio I., Daneck P., Della Bina E., Fares M., Franceschi D., Mandiello A. , Maniscalco M., Mazza S.,Pintore
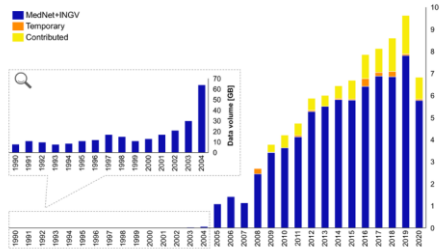
# Keywords

Policy dati INGV
Open Data
Open Access
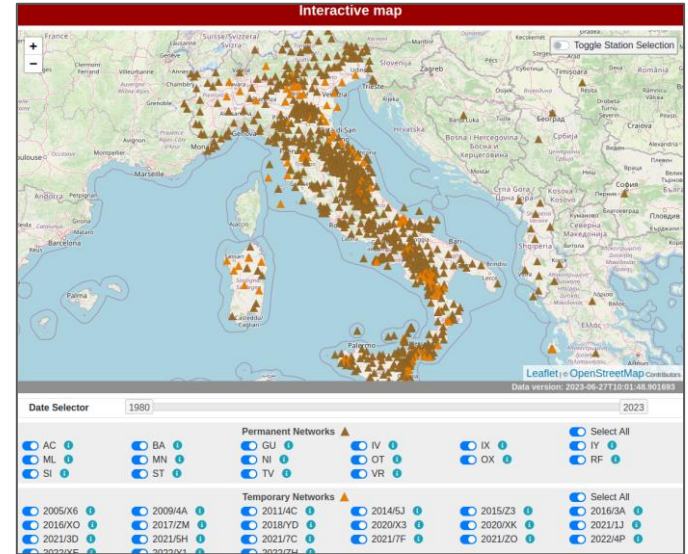Trusted-Archive
FAIR data

# EIDA Italian Node

The EIDA Italia node is the INGV managed datacenter that stores and distributes seismic waveforms collected by the INGV seismic networks, and other Italian and some foreign data providers. The list of institutions that contribute with their seismic networks data to EIDA Italia can be found at www.eida.ingv.it site. Data from other countries in the Euro-Mediterranean area are collected directly by INGV and partners by the MedNet network.

EIDA-ITALIA:

- 15 data providers
- 16 permanent networks
- 21 temporary deployments
- 1000+ stations
- 11k channel-epoch
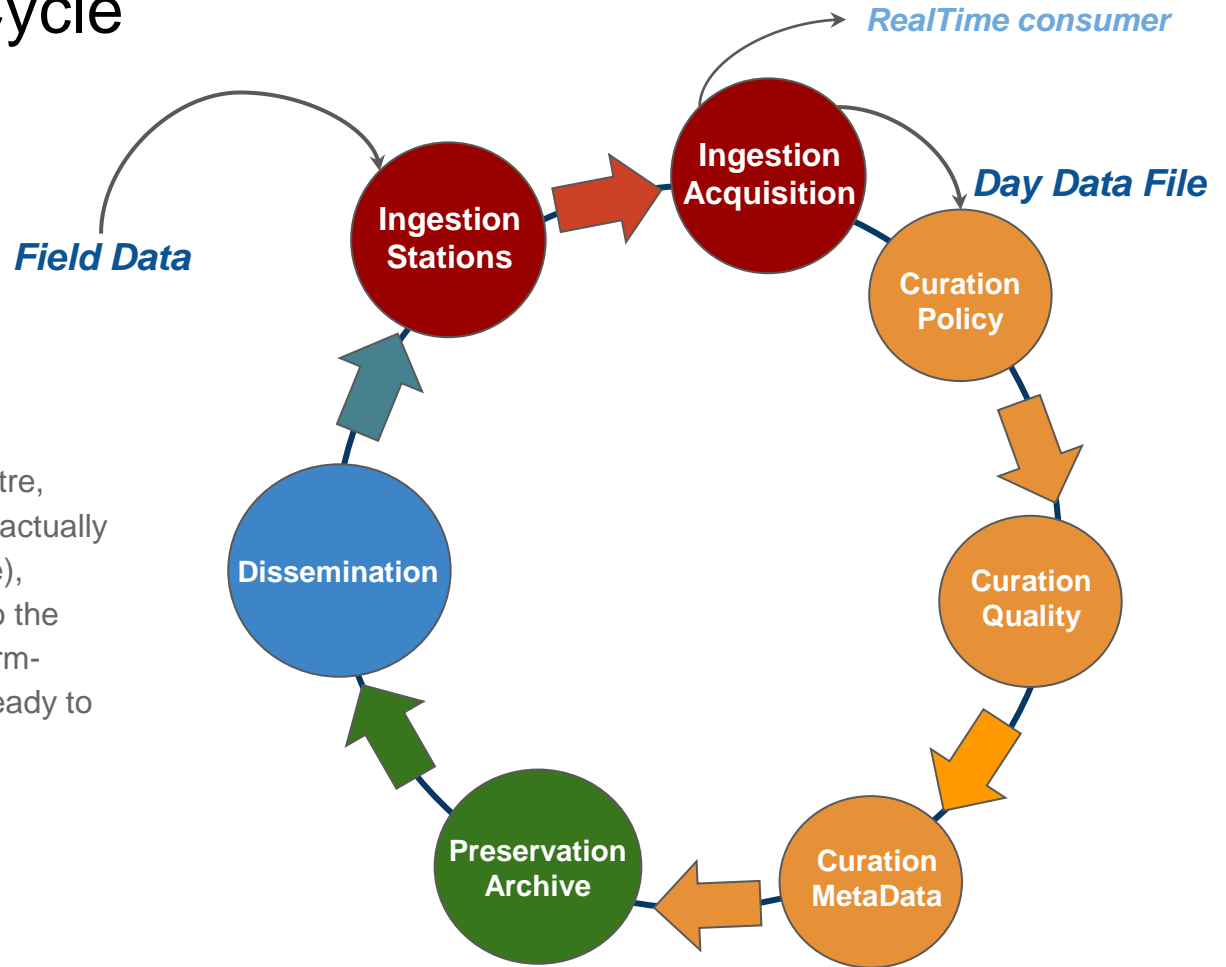- 120+ TB waveform data files

# Seismic Data LifeCycle

Data Centre point of view
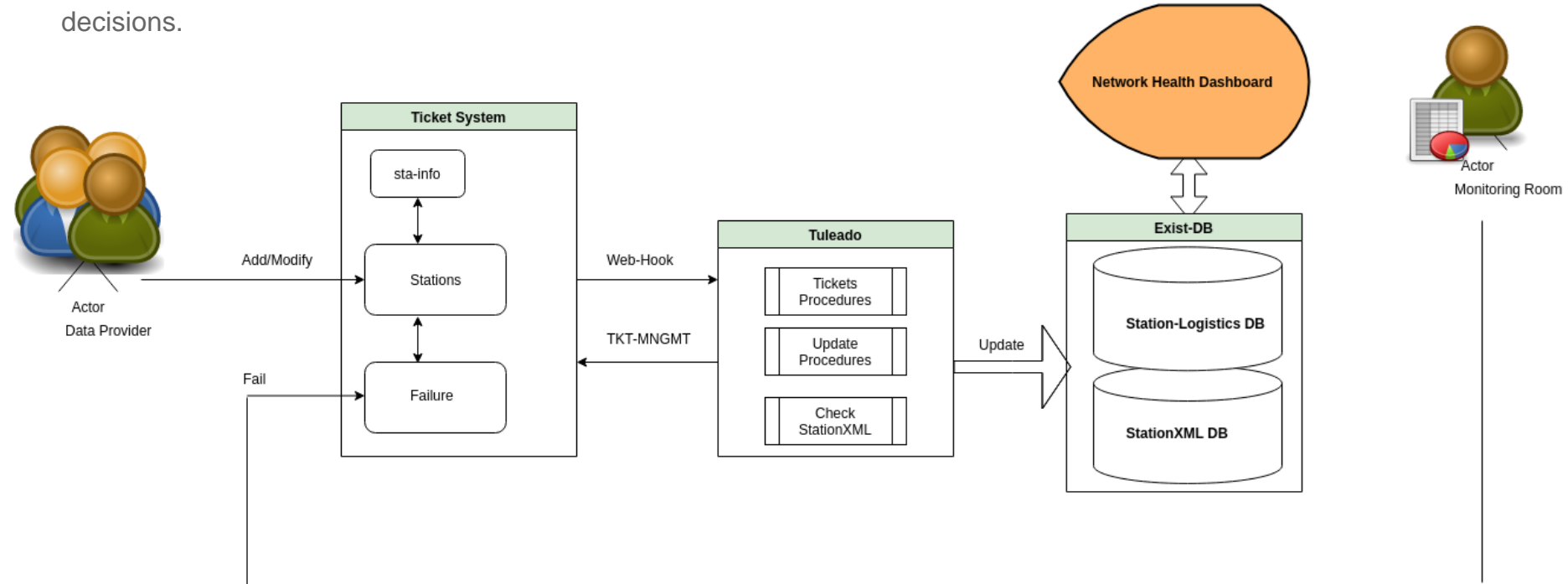
## Seismic Data LifeCycle



Our data life cycle, as Seismic data centre, start from initial generation: the station (actually we consider stations as our data-source), hence via the ingestion procedure, go to the curation phase, then will go into long-term-preservation phase; there data will be ready to be disseminate and re-use.
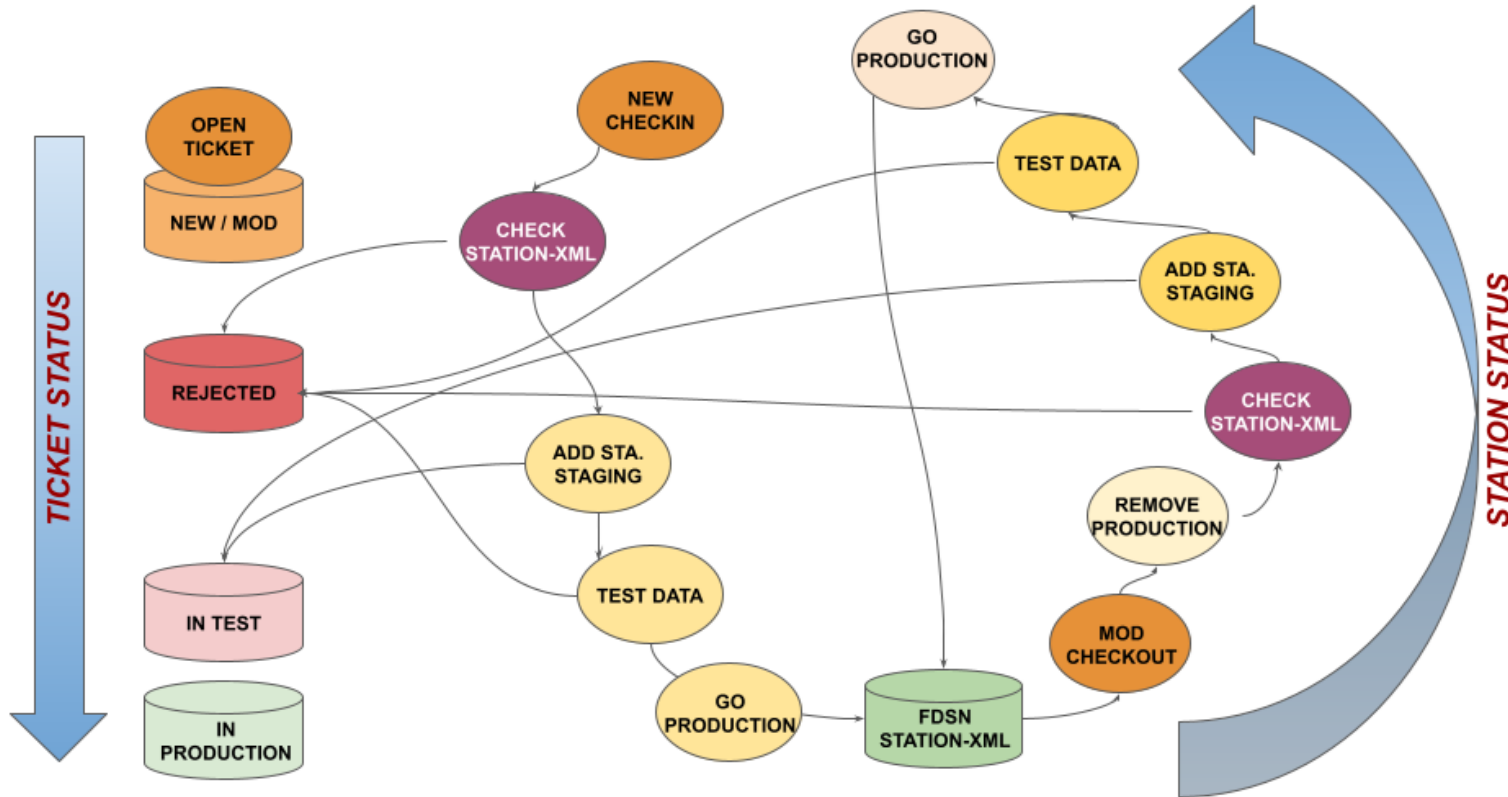
# Data Ingestion Station management system

Our system has a modular design, by leveraging existing open-source solutions i.e. as a customized implementation of a ticketing system, a GUI leverages Grafana creating a dashboard to present an overview of the system status of all components. Ticket status drives the station process through all of our systems (from acquisition to distribution). All these steps are managed through *Tuleado* who is able to follow some policies and make some decisions.
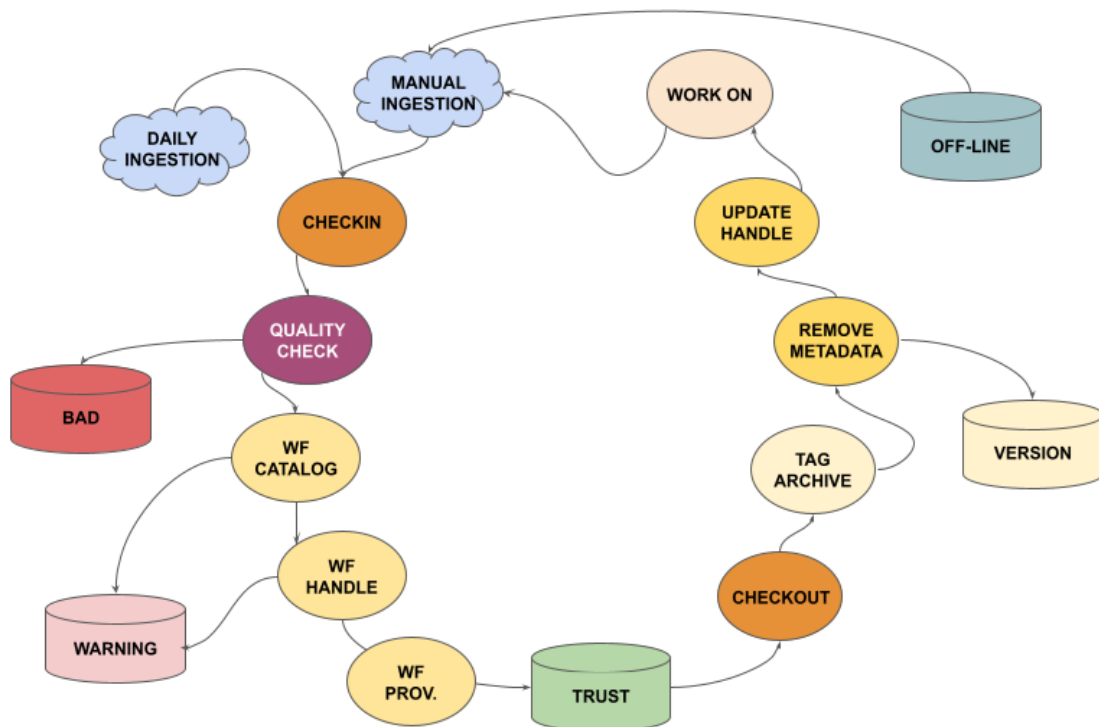
# Data Ingestion Station management model

Our station management model is based on a tight coupling between a ticket system and ours acquisition/dissemination systems; avoiding as much as possible human intervention in back office. To do this a Business Automation is implemented via an in house developed software agent called "*Tuleado*".

# Data Curation Policies & Quality

We have deployed an "**Policy Enforcement Point**" in order to execute our policies, expandable as needed, on data file.
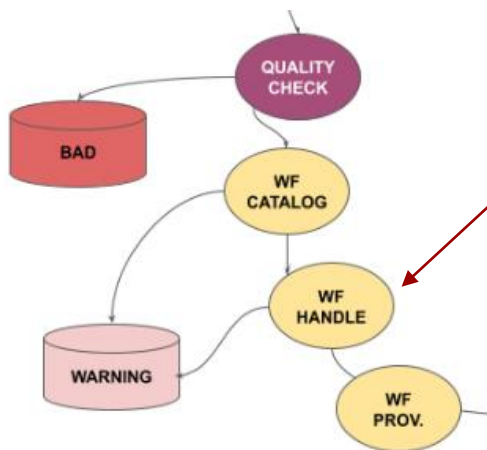
E.g. our 'checkin' policy do this: after pass sanity check and quality assessment, data file is directed to extract community metadata and metrics, then is turned from data file into digital object, lastly others metadata are applied.



- Check-in
- Quality Checks
- Extract Meta
- Mint a PID
- PUT IN Trust Archive

- Check-out
- Remove Meta
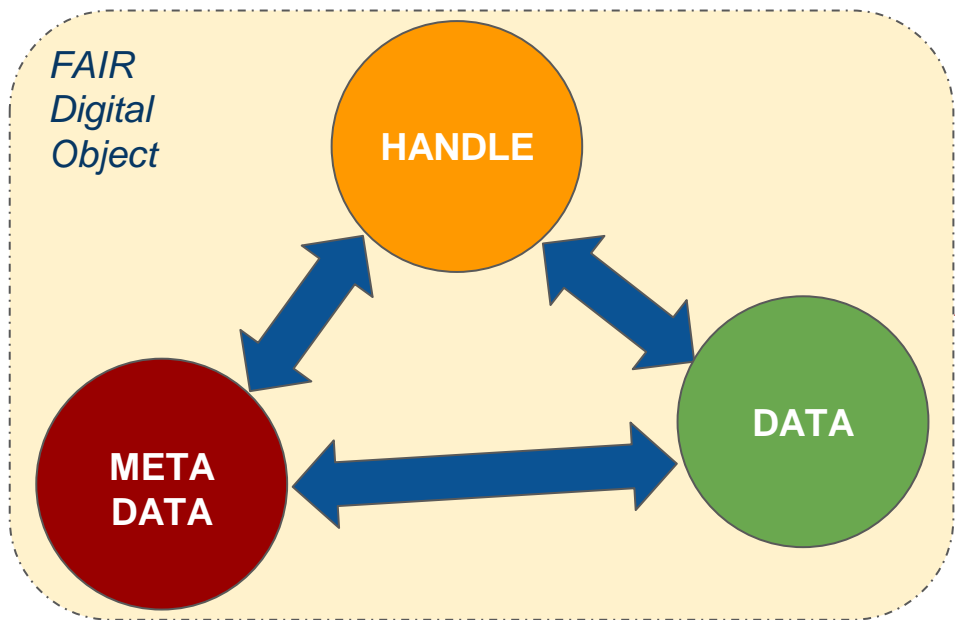- PULL OUT Trust Archive
- Version Update

# Data Curation Data File + unique identifier = Digital Object

*We add a Unique Persistent Identifier to all Data Files that have passed quality check & extracted Metadata:*
***Data File become Digital Object***



*FAIR Digital Object*

HANDLE

DATA

META DATA

# Data Dissemination
Handle resolver: a Digital Object approach

- *Provides uniform 'information' interface*
- *Easy to use*
- *DO Citation (ePIC-PID)*
- *Availability ensured to decades*
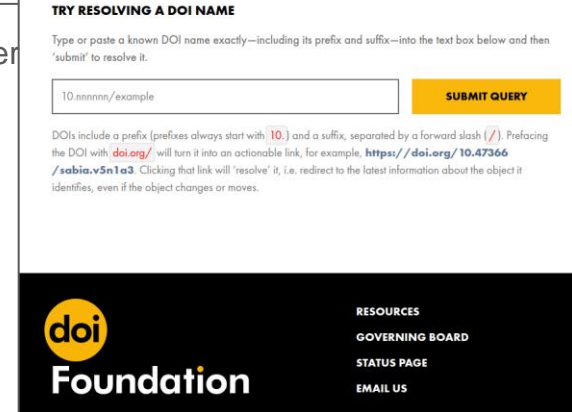- *Hides system details*



ePIC
Persistent Identifiers for eResearch

## Welcome

Resolve your PID

[Submit]

ePIC was founded in 2009 by a consortium of European partners in order to provide PID services for the European Research Community, based on the handle system (TM, https://www.handle.net/ ), for the allocation and resolution of persiste...

**TRY RESOLVING A DOI NAME**

Type or paste a known DOI name exactly—including its prefix and suffix—into the text box below and then 'submit' to resolve it.

10.nnnnnn/example    [SUBMIT QUERY]

DOIs include a prefix (prefixes always start with **10.**) and a suffix, separated by a forward slash ( **/** ). Prefacing the DOI with **doi.org/** will turn it into an actionable link, for example, **https://doi.org/10.47366 /sabia.v5n1a3.** Clicking that link will 'resolve' it, i.e. redirect to the latest information about the object it identifies, even if the object changes or moves.

**https://www.doi.org**/11099/11ed-996f-0242ac120005          :: latest ve...

**http://www.handle.net**/**11099/11ed-996f-0242ac120005**          :: latest version

11099/11ed-996f-0242ac120005**#version=1**     :: specific version

11099/11ed-996f-0242ac120005**#metadata**     :: DC metadata (WF-HANDLE)

11099/11ed-996f-0242ac120005**#provenance**   :: provenance (WF-PROV)

11099/11ed-996f-0242ac120005**#document**     :: doc of file (Human readable)

**11099/wf-search#lat; lon; rad; dstart; dend;**   :: list of PID in that area at given time-window

**doi Foundation**
RESOURCES
GOVERNING BOARD
STATUS PAGE
EMAIL US

# Data Preservation Long term preservation: Storage - Tape backup

Our policy for long-term-preservation is to have 3 copies: in local storage, then make a replication at a new storage archive in Naples and a backup in a tape library, compliant to 3-2-1 backup policy.
A test suite on data backupped is performed on regular base

***Long Term**: Hardware and file systems are upgraded on regularly base in order to allow readable over time (by now last 3 decades are available).*

*Follow **3-2-1 Backup rules** coined by Peter Krogh*

- **Have at least 3 copies of your data**
- **Store the copies on 2 different media**
- **Keep it safe with 1 backup copy off-site**

***Tape** Backup are **tested cyclically** on monthly base*




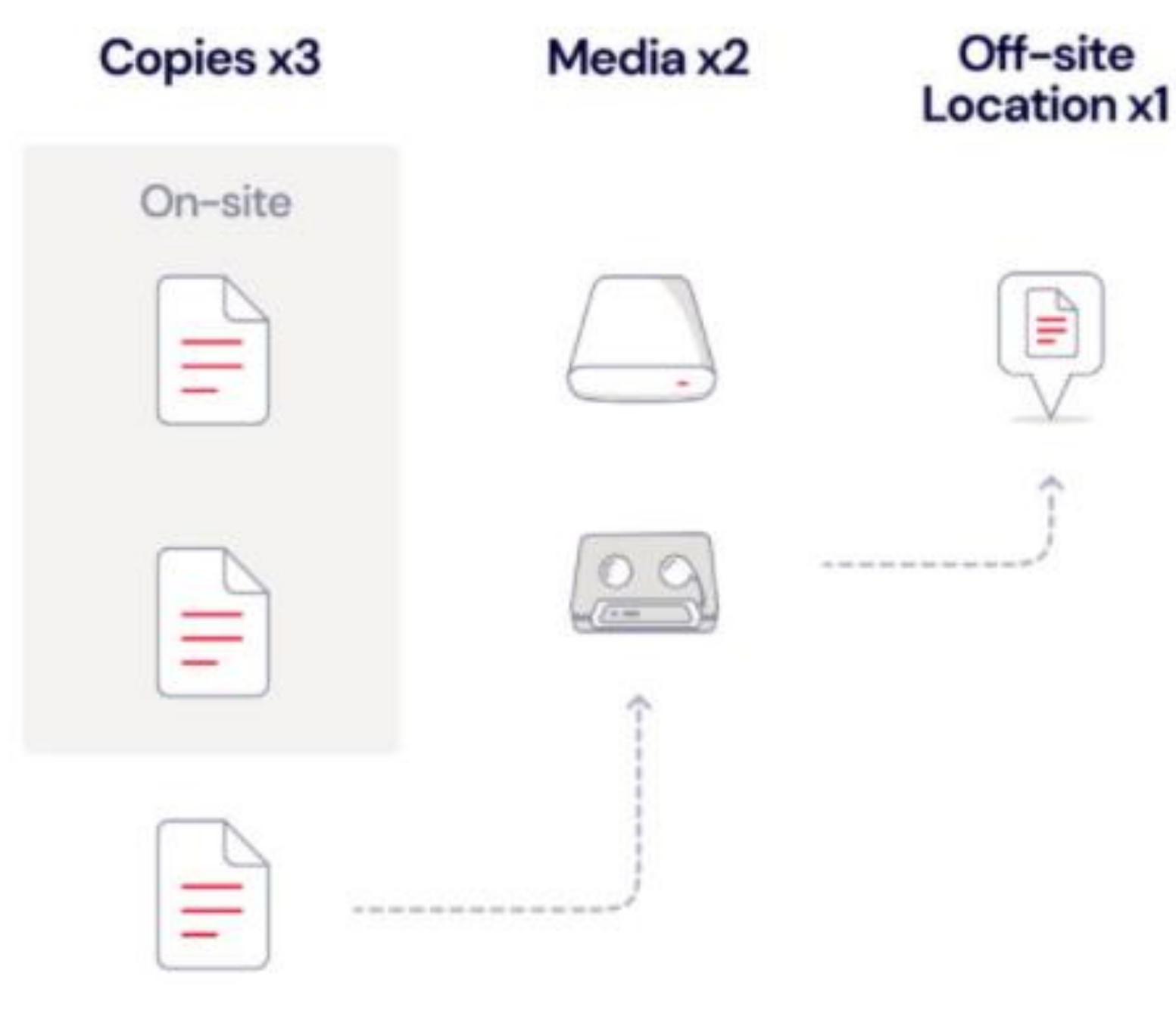

# Data Dissemination WEBsite, Web-services and so on..

After all these previous steps we are ready to disseminate data through web services and website.

Web site:
- EIDA Italia (https://eida.ingv.it/): frontend for browsing and downloading Italian Stations data and metadata belonging to Italian EIDA node.

Web Services:
- FDSN StationXML: nosql xml based system (eXist-db) provides APIs for download stations metadata based on FDSN standard.
- FDSN Dataselect: SeisComP based system provides APIs for download waveforms data in mseed format.
- FDSN availability based on mongo database.
- EIDA WS
  - Routing: find correct url to download across EIDA federation
  - WF-Catalog: data metrics; data quality; data availability

# Data Dissemination
Computational Archive: processing as a service

To serve users we are adding significant computational resources and an adequate processing and analysis framework combining Apache Spark, and ObsPy, creating a "computational archive" where storage resources and computational resources converge.
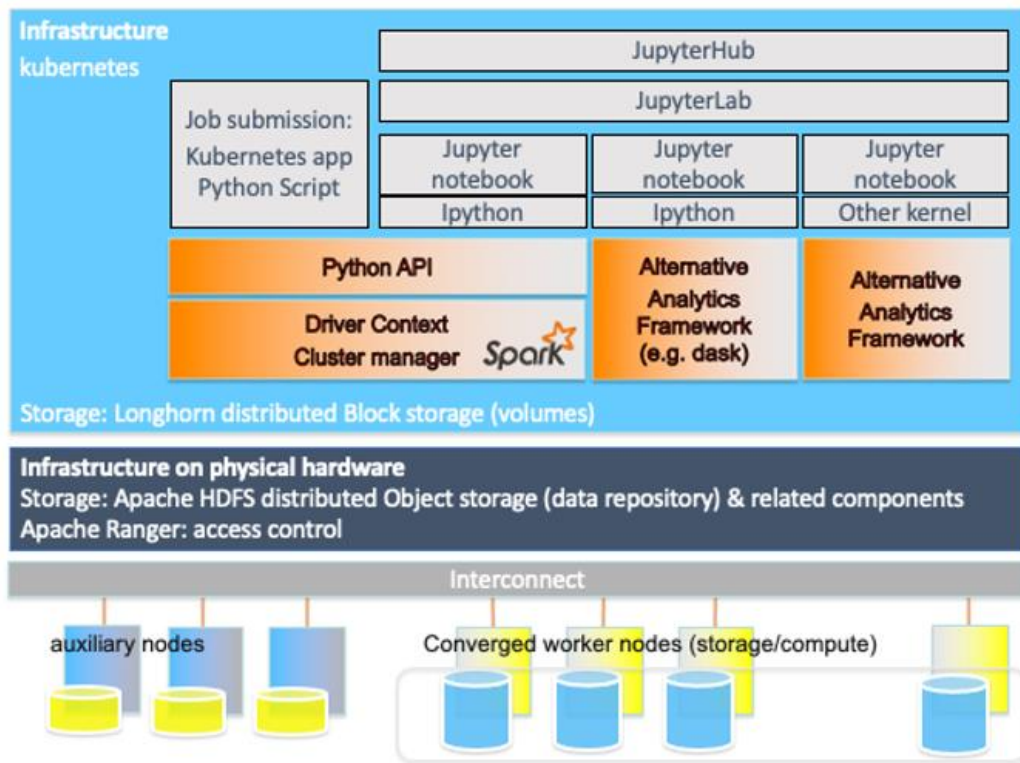


*Worker resource (aggregated)*
CPU: 384 cores (8x2x24); AMD Epyc 7352, 2.3 GHz base clock
RAM: 8 TB; DDR4
Mem bandwidth: 204.8 GB/s /CPU

*Storage resources (aggregated)*
SAS HDD: 1440 TiB (8x10x18 GiB)
SAS SSD: 61,440 GiB (8x2x3,840 GiB)

# Sostenibilità

**Fondi Istituzionali**

**Allegato A DPC WP 3**

**JRU EPOS-Italia** OS1

**Progetti:** **Centro Italia DL-50 OR AA**

# Cosa manca

Una politica dei dati in streaming
Repository INGV dei Dati repo.data.ingv.it

# Grazie