# DARE as a platform to support Climate Data Analytics using Cloud Infrastructures

## Christian Pagé

*Research Engineer / Climate Research Domain*

CERFACS  Toulouse, France

Iraklis A. Klampanos, *NCSR "Demokritos", Greece*

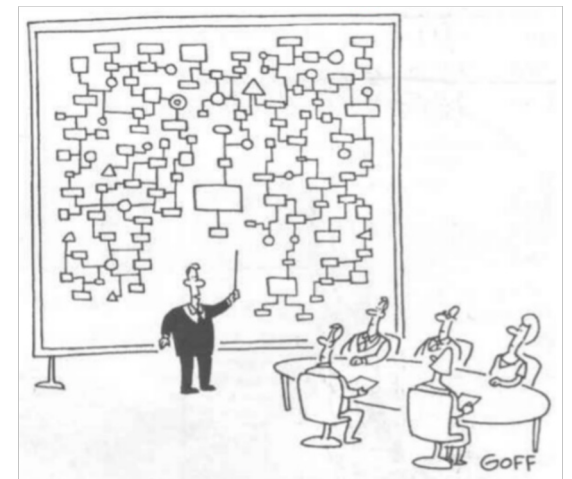Malcolm Atkinson, *The University of Edinburgh, UK*

Alessandro Spinuso, Maarten Plieger, Wim Som de Cerff, *KNMI, Netherlands*

# DARE IS-ENES Climate Use Case Motivations: Scientific, Technical, Societal

- Perform efficient Data Analysis
  - Large number of realizations (ensemble of scenarios)
  - Uncertainties range estimation
  - Process Higher spatial and temporal resolution
  - Easily share intermediate results with collaborators
  - Comparisons when doing numerical model developments

- Achieve a more robust and flexible Data Life Cycle
  - More robust experiments setup
    - Explore several experiment configurations to answer scientific questions
  - Reproducible and traceable experiments
  - **Download locally then Analyze**: a workflow that cannot be sustained
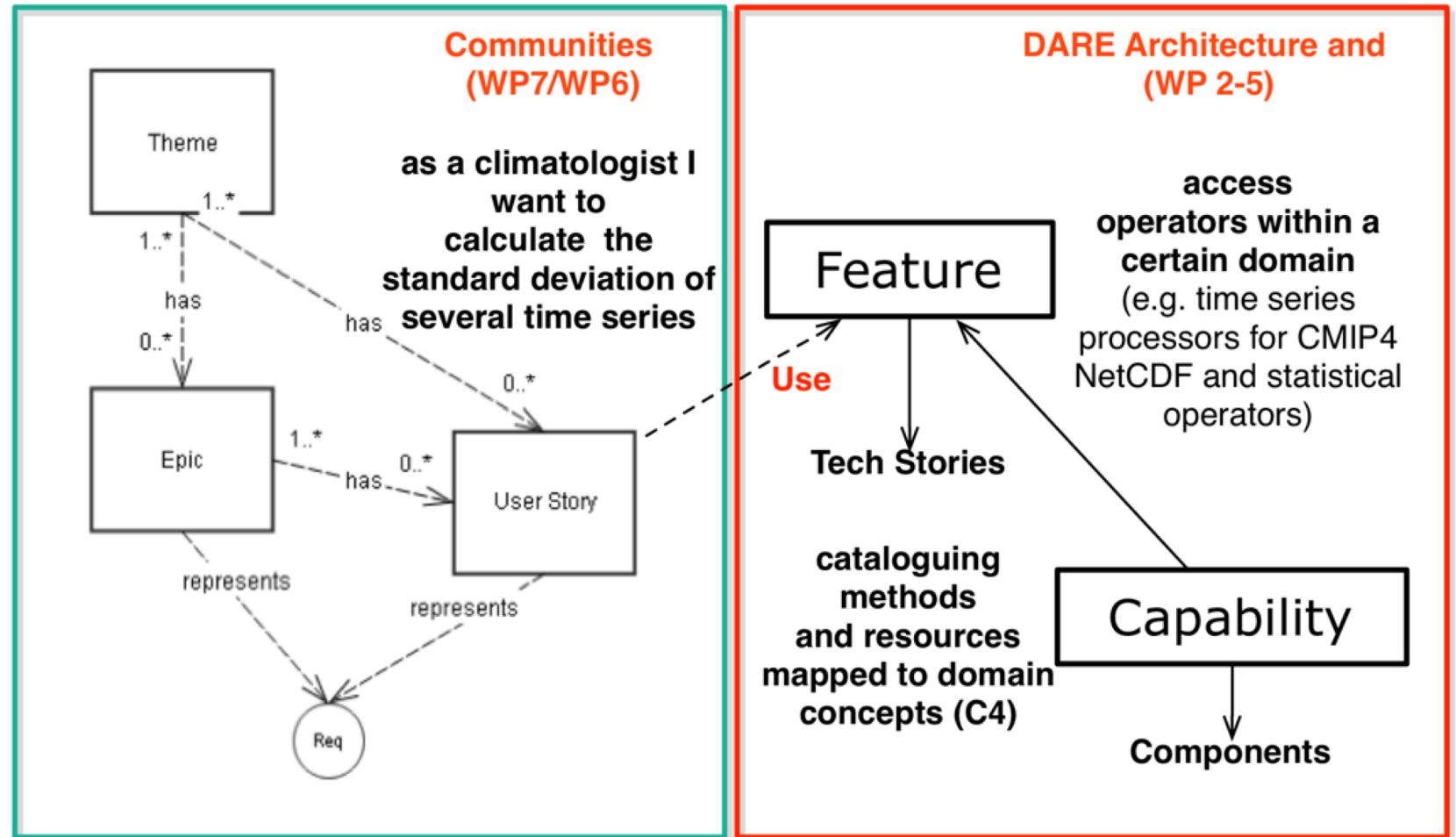
# DARE Mapping Community Needs

**Mapping from Community User Stories to DARE Features and Capabilities**

- The user-story requires to access the right implementation of a component (Feature), which may be implemented through resolving services (Tech Story).

- Cataloguing is a capability of DARE
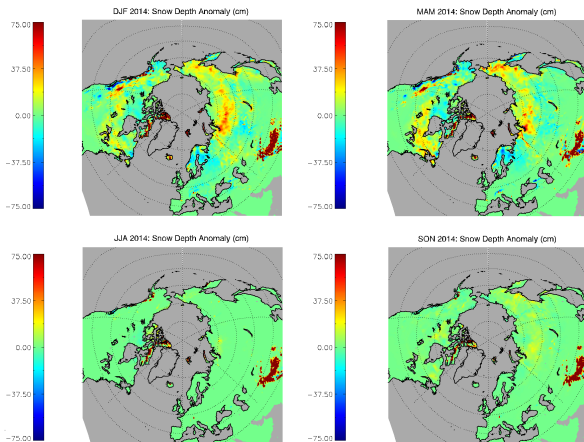
**Requirements Analysis for each Pilot**

- Identify existing architecture, standards and tools components
- Identify communities' needs

# Climate Data Users: Current situation

## Practical Example: A Climate Research PhD Student

- **I want to study how the feedback of the snow cover in Northern Europe and Russia on the weather circulation patterns and temperature extremes over Western Europe is impacted in the future climate**
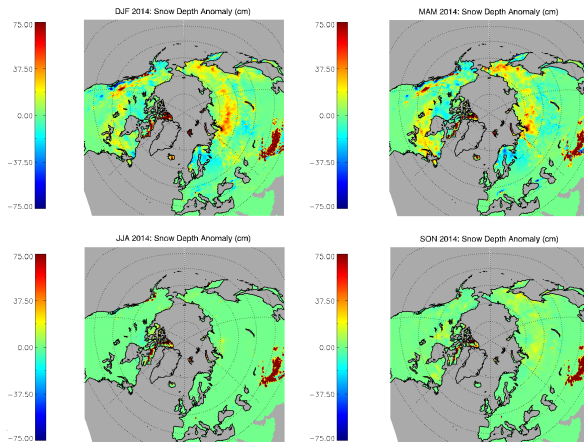


- Surface Temperature (+max/min), Pressure, Humidity, Snow Cover, Precipitation (Solid&Liquid): 8 surface fields

- Historical + All RCPs
- Combination of models an ensemble members
- EUR-44 Euro-Cordex Grid

- ~11 200 files of ~50 Mb each per field

**TOTAL**: **~560 Gb**

# Climate Data Users: Current situation

## Practical Example: A Climate Research PhD Student

- **I want to study how the feedback of the snow cover in Northern Europe and Russia on the weather circulation patterns and temperature extremes over Western Europe is impacted in the future climate**
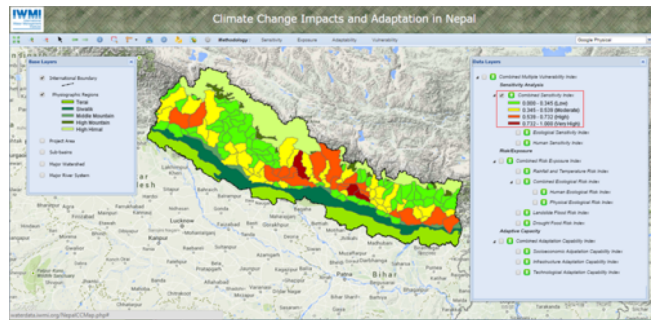


### Needs and questions

- I need to calculate several statistics for analyses
- I need derived quantities (climate indices, indicators)
- I want to assess if higher resolution data is needed or other datasets
- I want to do some Quality Check
- ...

# Climate Data Users: Current situation

## Practical Example: An Impact Engineer

- **My region needs to assess the impact of climate change on how we perform water management. I work with GIS Software to overlay several informational data layers.**



- Surface Temperature (+max/min), Precipitation, Winds : 6 surface fields

- Historical + All RCPs:
- Combination of models an ensemble members
- EUR-11 Euro-Cordex Grid

- 1378 files of ~600 Mb each per field

**TOTAL: ~5 Tb**

# Climate Data Users: Current situation

## Practical Example: An Impact Engineer

# Climate Data Users: Current situation

## Practical Example: An Impact Engineer

- **My region needs to assess the impact of climate change on how we perform water management. I work with GIS Software to overlay several informational data layers.**



<span style="color:red">Needs and questions</span>

- How to reduce the dataset to a representative subset?
  - My client cannot cope with too many realizations
- I need to do the calculations remotely and download the results
- I cannot use NetCDF, I need to import the data into my GIS software

# Climate Data Users: Current situation

## Many common needs

- Guidance and tools for data and scenarios subsetting: selecting a subset of representative scenarios

- Lower significantly the total data size to download
  - Calculate as much as possible remotely

- Reformat/Repackage the data into easier formats and organization/homogenization (implies smaller datasize)

- Full Provenance and Lineage information

- Proper Metadata description, especially for derived data

- Variety of Access Interfaces for adoption: OGC, REST, Jupyter, APIs



GOFF

# Climate Data Distribution: ESGF RI

**IS-ENES CDI C4I**
- Tailored for end-users
- Supports on-demand data processing

**ESGF Data Nodes 2015**
- 40 worldwide
- 18 in Europe (coordinated in IS-ENES)

| | CMIP5 | CMIP6 | CMIP7 |
|---|---|---|---|
| Year | 2012 | 2017 | 2022 |
| Power factor | 1 | 30 | 1000 |
| Npp | 200 | 357 | 647 |
| Resolution [km] | 100 | 56 | 31 |
| Number of mesh points [millions] | 3,2 | 18,1 | 108,4 |
| Ensemble size | 120 | 214 | 388 |
| Number of variables | 800 | 1068 | 1439 |
| Interval of 3-dimensional output (hours) | 6 | 4 | 3 |
| Years simulated | 90000 | 120170 | 161898 |
| Storage density | 0,00002 | 0,00002 | 0,00002 |
| Distributed Archive Size (Pb) | 3,19 | 86,05 | 2260,20 |

# European Landscape & Components
# EUDAT & EGI

## EUDAT CDI B2 Service Suite

▶ Integrated B2 Services

▶ B2ACCESS: Common AAI

▶ Interface between EUDAT B2 Services and Communities infrastructures, such as Climate

▶ Prototype Workflow Service: GEF (Generic Execution Framework)



Registered Users: 21714  VOs: 233
LCPUs: 470,000  Disk: 143PB  Tape: 138PB
Jobs: 1.62 million/day

| Resource Centres | EGI-InSPIRE & EGI Council members | 319 |
|---|---|---|
| | Including integrated RPs | 351 |
| Countries | EGI-InSPIRE & EGI Council members | 42 |
| | Including integrated RPs | 54 |

Integrated EGI-InSPIRE Partners and EGI Council Members

External Resource Providers (integrated)

Internal/External Resource Providers (being integrated)

Peer Resource Providers

▶ Computing Power (FedCloud) resources

# DARE IS-ENES Climate Use Case Draft Architecture

# Monitoring and Exploration of WPS workflows via Provenance
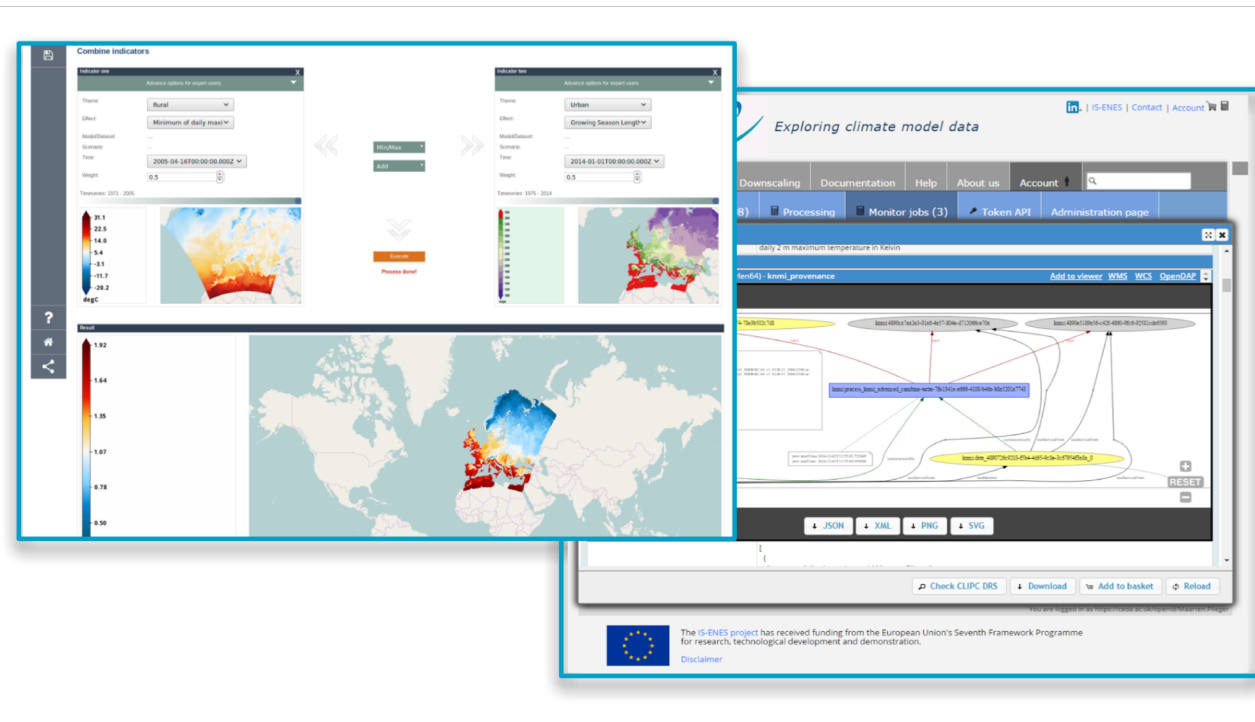
## Standards

- Provenance module: WPS_PROV
- Provenance metadata is stored in NetCDF
- W3C PROV-DM standard



## Visual analytics techniques on provenance

- Highlighting data-reuse, even for cached data
- User interactions
- Exploitation of resources



NetCDF dependencies

Searchable metadata

# DARE IS-ENES Climate Use Case Significant benefits

- Enable the on-demand **delegation of IS-ENES C4I** Platform **Data Analytics** and **Processing** to the **DARE** infrastructure (cloud-ready).
  - Typically, data reduction **on the order of 70-90%** can be achieved, depending on the users' analyses
- Streamline and **ease** the whole **data lifecycle**, with proper data **provenance** and **lineage**

- The DARE Platform will also:
  - **Be interoperable** with **EUDAT CDI** by using its standards and B2 Services
  - **Interface** with **EOSC** and **Copernicus C3S-DIAS**

# Questions & Comments! ☺

http://project-dare.eu



christian.page@cerfacs.fr

# Questions & Comments! ☺

http://project-dare.eu



christian.page@cerfacs.fr

# Open Questions



- Several European platforms will be available: C3S-DIAS, EOSC, ESGF Data/Computing Nodes, IS-ENES CDI & ECAS, EUDAT CDI, EGI, DARE, National Platforms, MAIDK
  - How do we ensure that we do not have duplicate efforts (too much)?
  - Which kind of users do they each address? How users will know which one to use? The ones they can access? With what kind of resources limitations?
  - How do we "educate" different kind of users for wide adoption and usage of those platforms?
  - How can they be interoperable? APIs, AAIs, ...
  - How to ensure that they make available promptly new datasets
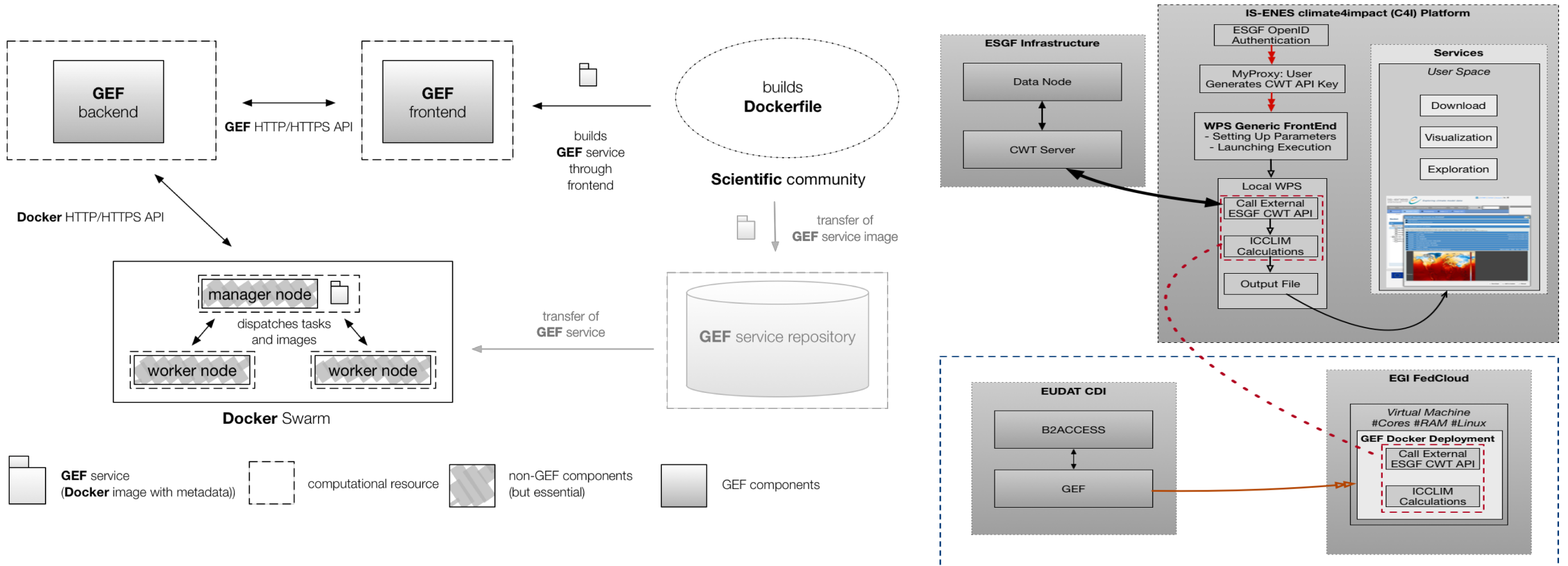  - Will they be scalable enough?

# Open Questions

- How do we deal with non-mature services, changing APIs?

- On-demand remote data processing and data sharing is really needed

- Containerized solutions: distributed processing, orchestration, AAIs...

- What about Data Locality (Distributed Input Data)?

- Metadata Aspects and Reproducibility for the DLC: metadata mappings, full provenance and lineage information, PIDs

# European Landscape & Components
## EUDAT GEF & EGI

# WP7-JRA2 IS-ENES/C4I Pilot
# Generic Use Case

**Objective**: Generate a multi-model multi-scenario time series average of the surface temperature using CMIP5 data

### Scientific Workflow
- Spatially average over Western Europe (continents only)
- Time Period 1950-2100
- RCP 8.5 GES scenario
- All Global Climate Models available
- All members available
- Calculate the average time series
- Calculate the standard deviation
- Extract separate time series of every simulation
- Plot all those time series on a single graph