



Koninklijk Nederlands  
Meteorologisch Instituut  
*Ministerie van Infrastructuur en Milieu*

# **User Centered Provenance Management for Data-Intensive Platforms**

***Alessandro Spinuso***  
***[spinuso@knmi.nl](mailto:spinuso@knmi.nl)***



## **Respect des Fonds - Historical Facts.**

**France 1790, Establishment of “Les Archives Nationales”  
trying to merge private and public archives throughout the country**

*... 50 years of archivists headache trying to regroup and classify records ...*

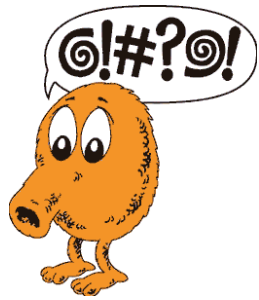
**France 1841, formulation of the **principle of provenance:****

**The unity of the archive took precedence over the material objects.**

**Classification based on administration, organisation, individual, or  
entity by which they were created.**

<http://ingmarbergman.se/en/category/tags/principle-provenance>

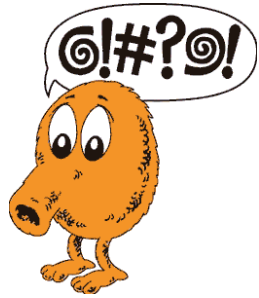
[http://en.wikipedia.org/wiki/Respect\\_des\\_fonds](http://en.wikipedia.org/wiki/Respect_des_fonds)



## *The Ambiguous Origins of the Archival Principle of “Provenance”*

*Shelley Sweeney*

***“Broadly speaking, the word “provenance,” whether used by a rare book librarian, an archaeologist, an art curator, or an archivist, refers to the origins of an information-bearing entity or artifact. But there the consensus ends.”***



**Archeology:** “the place where an object was found or recovered in modern times; the findspot”

**Librarians:** “Information concerning the transmission or ownership, as of a book”. (Never used for retrieval)

**Geology:** “The reconstruction of the history of sediments movements over time” (Many types of detrital records to unveil regional tectonic history)

**Wines:** “A documented history of wine cellar conditions”. (transactions of old wine with the potential of improving with age)





How does this apply to the curation of Scientific Data?  
**How do we use it and which are the Costs?**

- **Who are the archivists ?**
- Which percentage of the actual provenance recordings is relevant ?
- **Unified models (W3C PROV-O) are useful, how about the content ?**
- How do we produce provenance data (Automated / Manual)
- **Now that CPU time and data bandwidth is a business model, shall we benefit from what we use and pay ?**

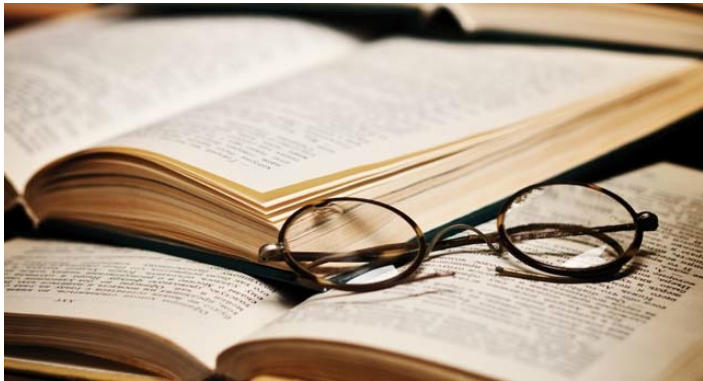


## Computer Science

*“The lineage of data or processes, as per data provenance”*

*“Data sets are reliable when the process used to create them are reproducible and analysable for defects”*

*“Scientific workflows assist scientists and programmers with tracking their data through all transformations, analyses, and interpretations.”*



## Researchers are part of the archiving process.

They know what is relevant to understand their results.

Automated system should provide support for a consistent and effective acquisition of provenance metadata - **Selective and extensible Provenance**.

[A. Misra] [I. Foster.]

**Reuse of workflow settings and data** in *Collaborative Environments* [P. Missier and B. Ludascher].

**Data stream processing engines** for Data Intensive computation, present expensive requirements for provenance collection, either in terms of size or I/O [W. D. Pauw]



## *Virtual Earthquake and Seismology Research Community in Europe*

**Virtual Environment** for of **Earthquakes Simulations** and  
evaluation of **Earth Models**

<http://portal.verce.eu>

Combined access to **computing infrastructures** (EGI, PRACE, Local Clusters), for  
development and execution of large **HPC** computations

Access and use of **European data archives** and services adopting International  
standards (FDSN, GCMT, OneGeology, EFEHR, QuakeML)

**Adoption of Workflow Technologies, Data Management and Provenance  
System**



Setup Results IRods

max extent Help Layers info

Solver Earthquakes Stations Submit Control

File FDSN

Open

Name: DefaultName

File: Select a file Browse... Upload

	0/15	Desc	Date	Depth	Latitude	Longitude	Magnitude	MT	
	<input type="checkbox"/>	LUCCA	2013-06-30T1...	9800.0	44.171	10.2047	4.5		
	<input type="checkbox"/>					10.2108	4.4		
	<input type="checkbox"/>					10.135	5.1		
	<input type="checkbox"/>					10.4543	4.8		
	<input type="checkbox"/>					9.6703	4.3		
	<input type="checkbox"/>					10.9502	4.7		
	<input type="checkbox"/>					10.9795	5.3		
	<input type="checkbox"/>					10.9663	4.1		
	<input type="checkbox"/>					11.0657	5.6		
	<input type="checkbox"/>					11.305	4.2		
	<input type="checkbox"/>					11.4407	4.9		
	<input type="checkbox"/>					11.2635	5.8		
	<input type="checkbox"/>					9.354	4.0		
	<input type="checkbox"/>					10.009	4.9		

VERCEDEMO01402044603337.mp4

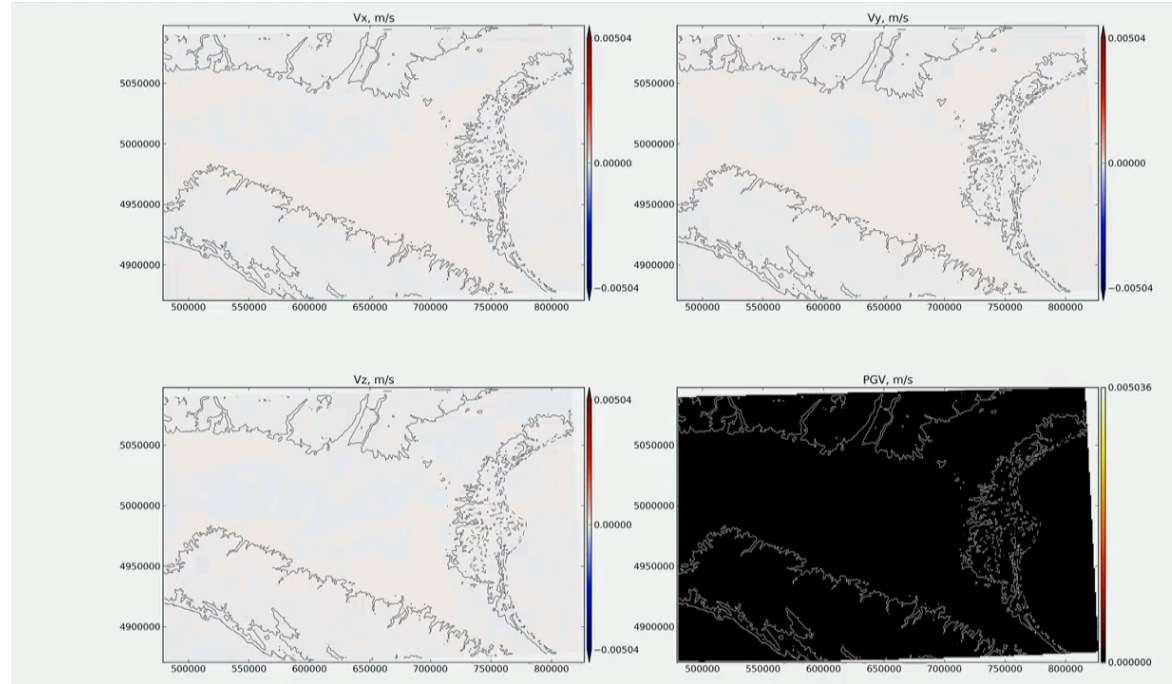
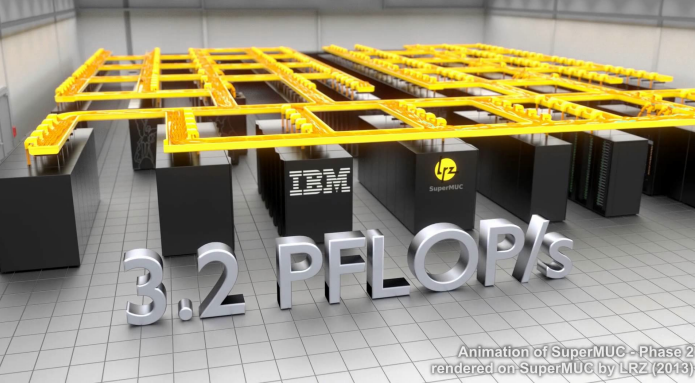
00:30 -00:29



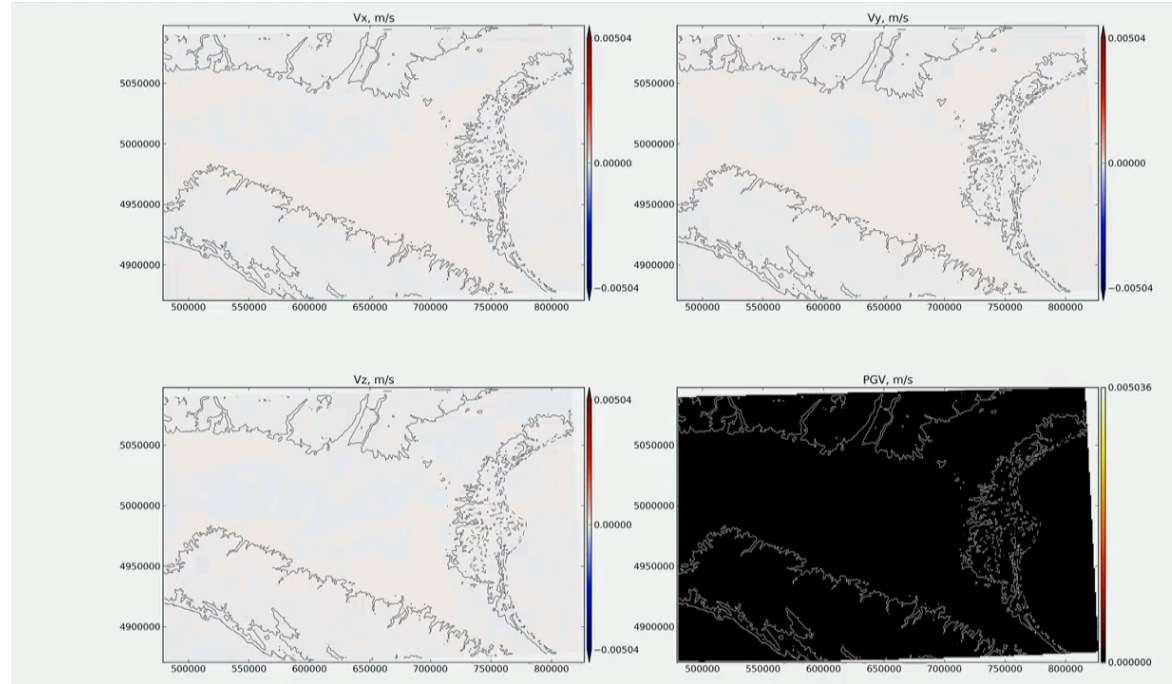
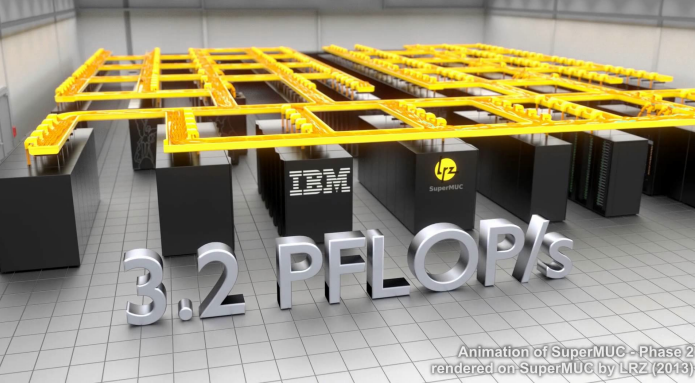
## Forward Modelling Use Case

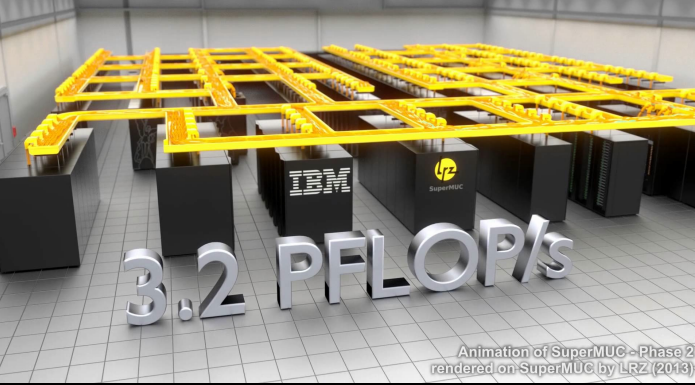
1. Production of **synthetic seismograms** for various **Earth models and earthquakes** on a continental scale Requires the execution of **HPC simulation codes** called solvers (**Simulation**).
2. The synthetic data may be **compared with real observations** (**Raw Data Acquisition, MISFIT**)
3. **Model updates and improvement** (**Inversion**)



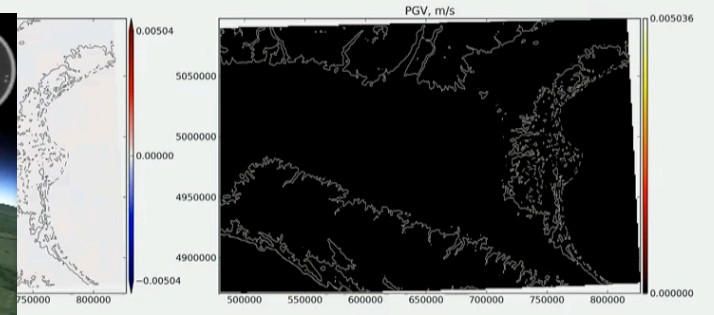
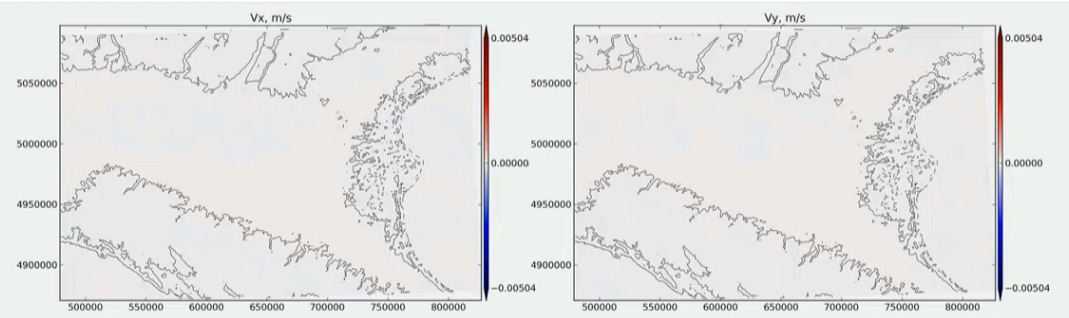




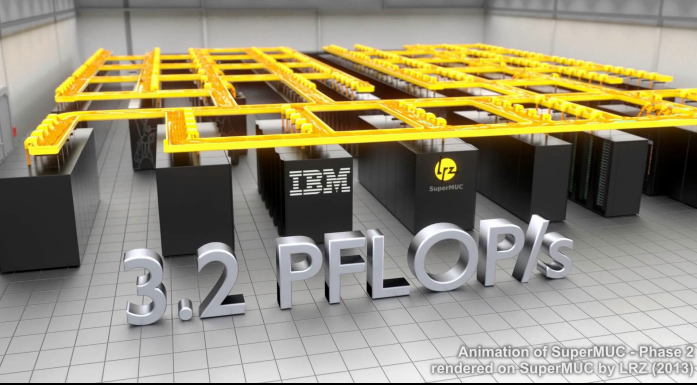




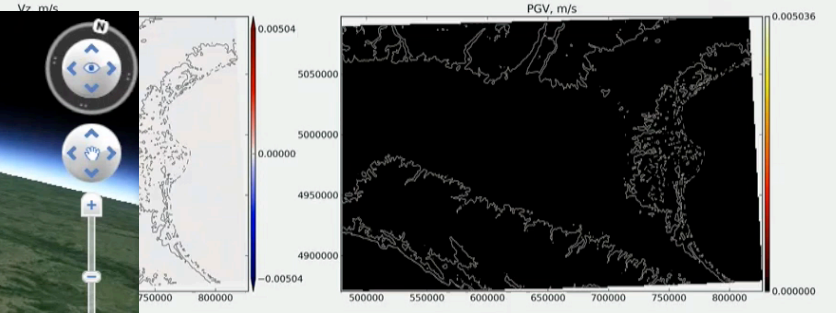
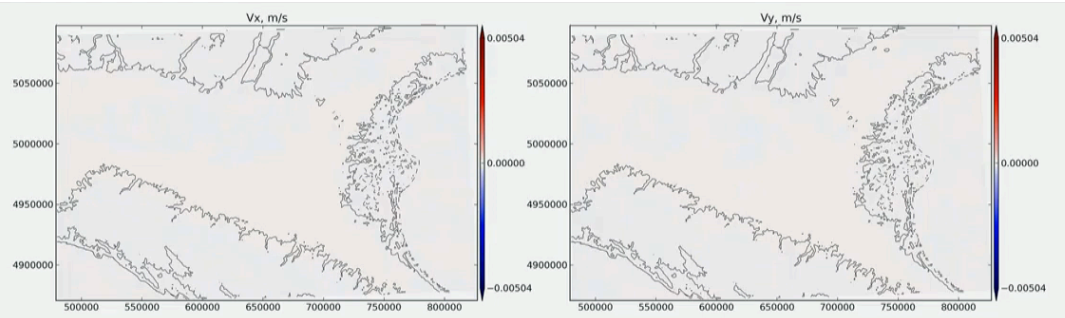
Animation of SuperMUC-Phase 2 rendered on SuperMUC by LRZ (2018)



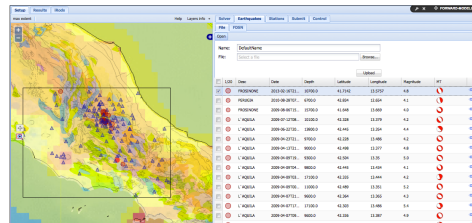




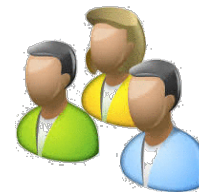
Animation of SuperMUC-Phase 2 rendered on SuperMUC by LRZ (2018)



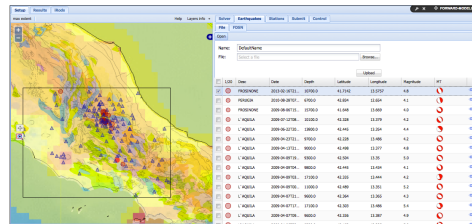
# VERCE Platform, Components Interaction



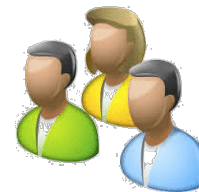
Science Gateway  
Community Applications



# VERCE Platform, Components Interaction



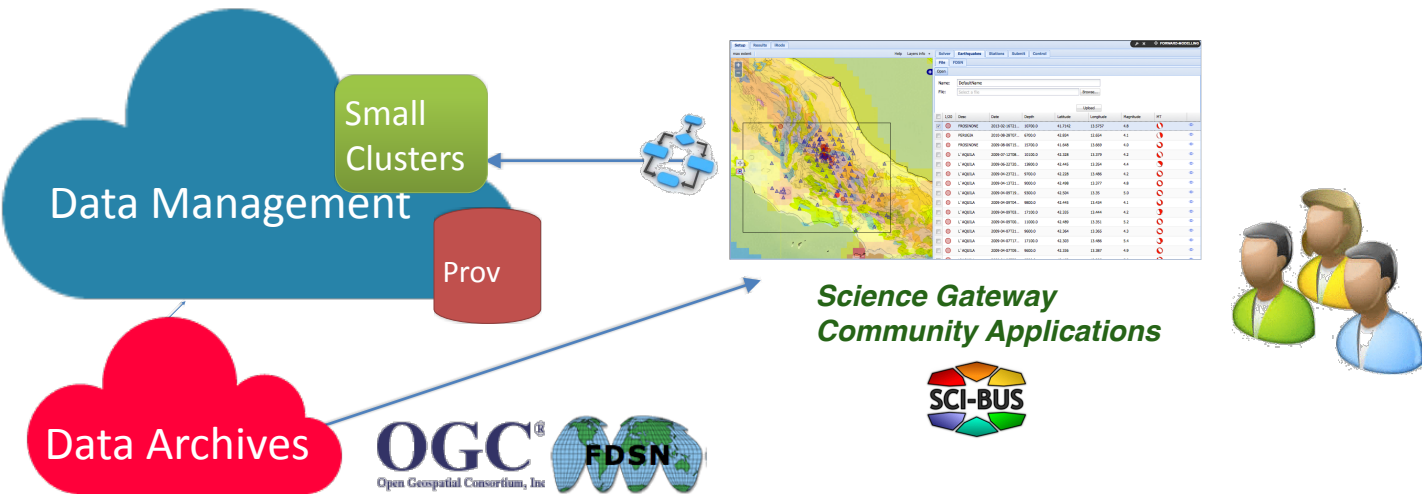
Science Gateway  
Community Applications



# VERCE Platform, Components Interaction

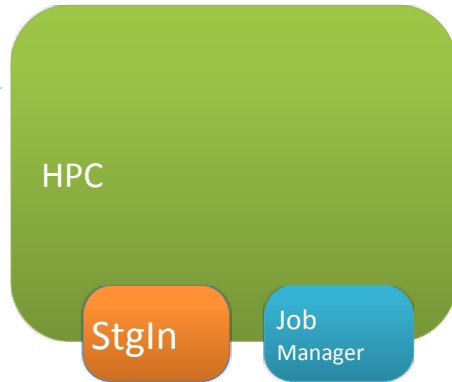


- 1 - Raw data acquisition
- 3 - MISFIT



# VERCE Platform, Components Interaction

2 - HPC Simulation  
(model stage-in)

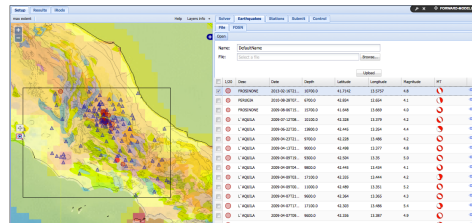


1 - Raw data acquisition  
3 - MISFIT

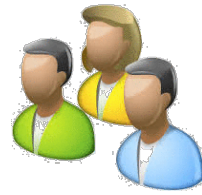


Small Clusters

Prov



Science Gateway  
Community Applications

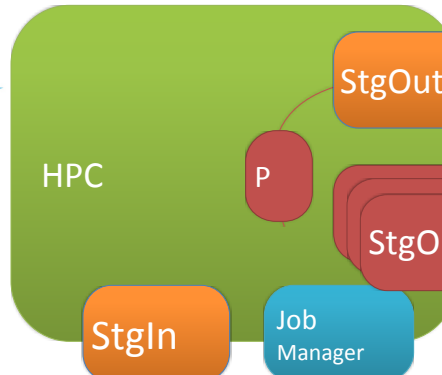




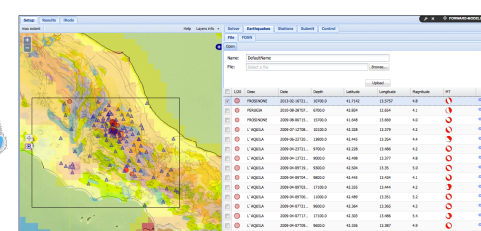
# VERCE Platform, Components Interaction

2 - HPC Simulation  
(model stage-in)

Results and provenance management



1 - Raw data acquisition  
3 - MISFIT



Science Gateway  
Community Applications



Runtime  
Provenance  
messaging

Metadata and Provenance  
Archive and Services

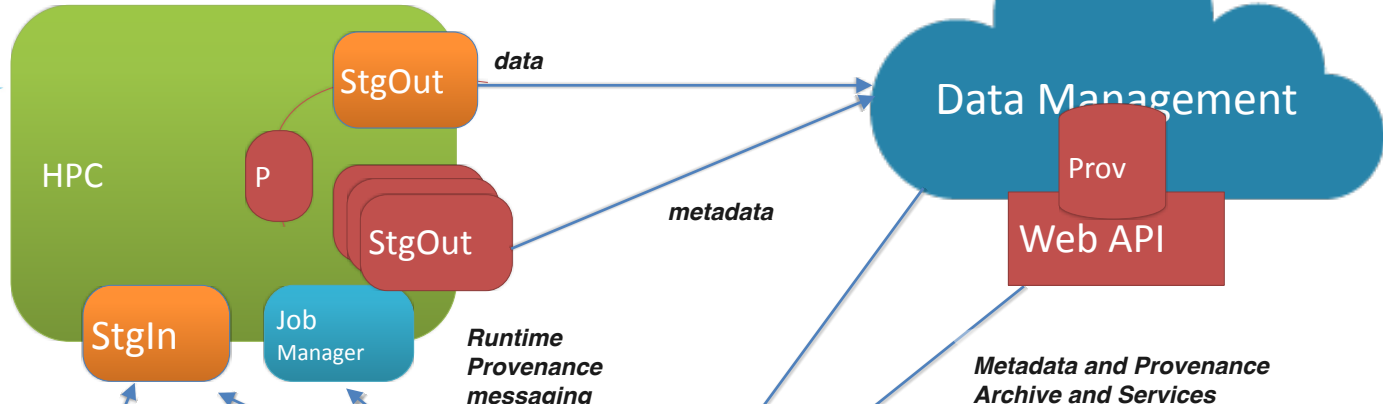
data

metadata

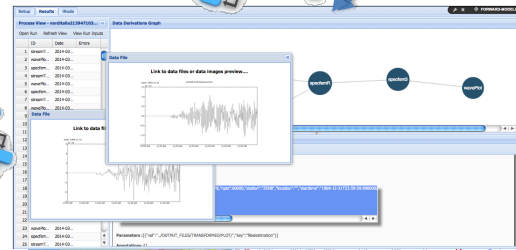
# VERCE Platform, Components Interaction

2 - HPC Simulation  
(model stage-in)

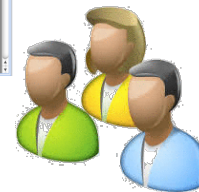
Results and provenance management



1 - Raw data acquisition  
3 - MISFIT



Interactive Validation and Visualisation throughout the process



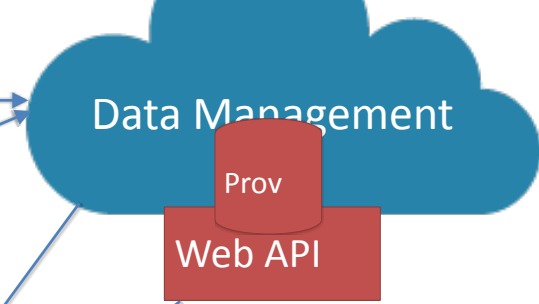
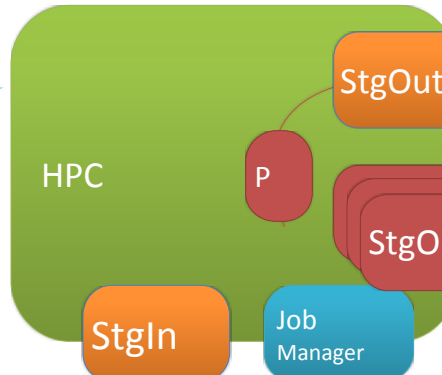
Science Gateway  
Community Applications



# VERCE Platform, Components Interaction

*2 - HPC Simulation  
(model stage-in)*

*Results and provenance management*



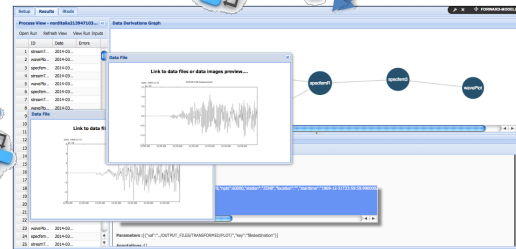
*1 - Raw data acquisition  
3 - MISFIT*



*GridFtp  
Globus  
UNICORE*

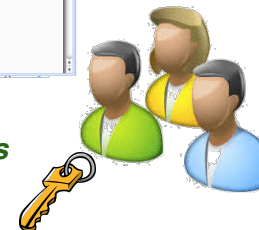
*Runtime  
Provenance  
messaging*

*Metadata and Provenance  
Archive and Services*



*Interactive Validation and Visualisation  
throughout the process*

*Science Gateway  
Community Applications*



VOMS X.509

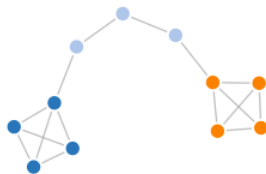


[dispel4py.org](http://dispel4py.org)

**Python library** used to describe **abstract workflows** for distributed data-intensive applications.

**Support for composition:** PEs may be defined by having their own internal workflows.

Abstract data-flows described in Dispel4Py can be automatically executed in **numerous parallel environments.**



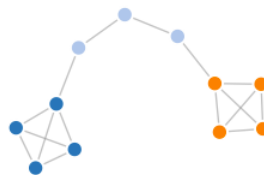


[dispel4py.org](http://dispel4py.org)

**Python library** used to describe **abstract workflows** for distributed data-intensive applications.

**Support for composition:** PEs may be defined by having their own internal workflows.

Abstract data-flows described in Dispel4Py can be automatically executed in **numerous parallel environments**.



MPI





[dispel4py.org](http://dispel4py.org)

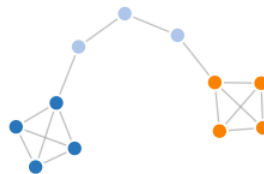
**Python library** used to describe **abstract workflows** for distributed data-intensive applications.

**Support for composition:** PEs may be defined by having their own internal workflows.

Abstract data-flows described in Dispel4Py can be automatically executed in **numerous parallel environments**.

**Storm**  
*Distributed and fault-tolerant realtime computation*

*Deployed on local  
Clouds*



MPI



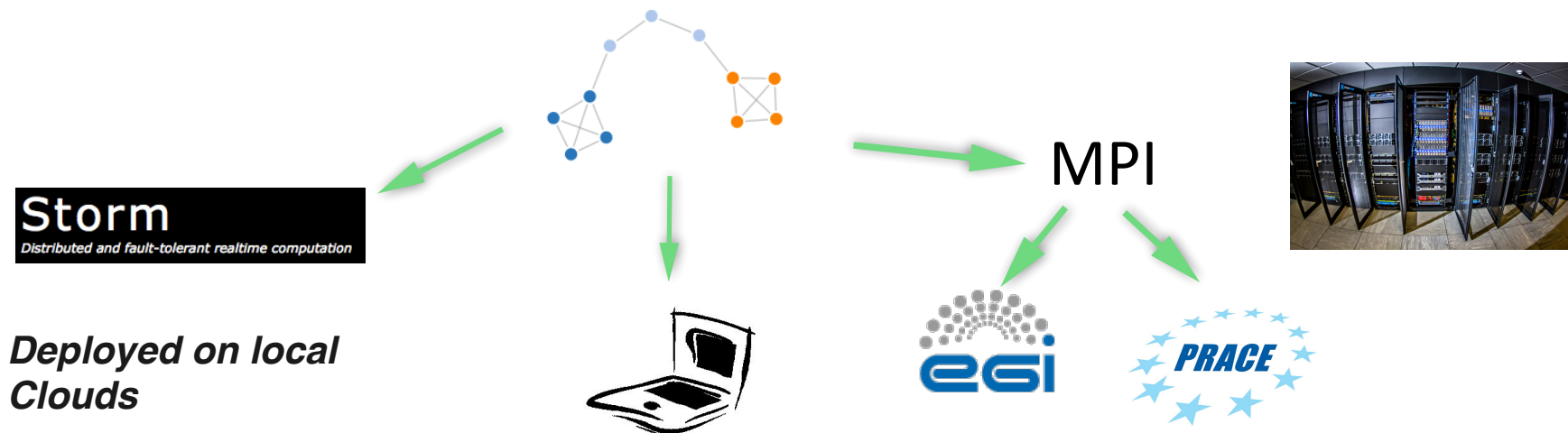


[dispel4py.org](http://dispel4py.org)

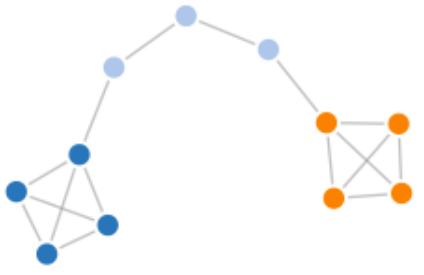
**Python library** used to describe **abstract workflows** for distributed data-intensive applications.

**Support for composition:** PEs may be defined by having their own internal workflows.

Abstract data-flows described in Dispel4Py can be automatically executed in **numerous parallel environments**.









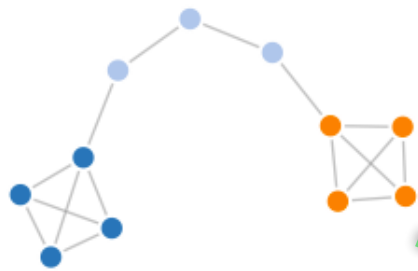
prov



WEB  
api

The screenshot shows a workflow management interface with the following components:

- Process View:** A table listing process steps with columns for ID, Date, and Error.
- Data Derivations Graph:** A flow diagram showing nodes for 'InputGen', 'decompose', 'speclemG', 'speclemF', and 'wavePlot' connected by arrows.
- Data File Preview:** A window titled 'Link to data files or data images preview....' displaying a waveform plot of '00000000 acceleration'.
- Log Output:** A scrollable area at the bottom showing command-line logs for various processes.



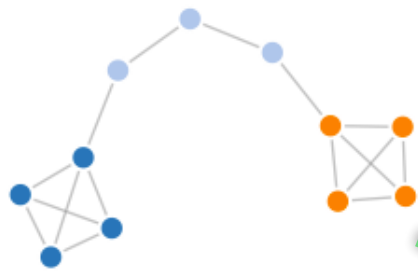
prov



WE  
B  
api

ID	Date	Error
8	2014-03-19 22:...	
9	2014-03-19 22:...	
10	2014-03-19 22:...	
11	2014-03-19 22:...	
12	2014-03-19 22:...	
13	2014-03-19 22:...	
14	2014-03-19 22:...	
15	2014-03-19 22:...	
16		
17		
18		
19		
20		
21		
22		
23		
24		
25		
26		
27		
28		
29		
30	2014-03-19 21:...	
31	2014-03-19 21:...	

- **Searches** over products metadata within and across runs
- **Data download** and preview
- **Diagnostic** / errors
- **Lineage**: Multi directional navigations across data dependencies
- **W3C PROV-DM** as reference model.



prov



WE  
B  
api

The screenshot shows a software interface for process management. On the left, a 'Process View' table lists runs with columns for ID, Date, and Error. The main area displays a 'Data Derivations Graph' with nodes representing processes: 'InputGen', 'decompose', 'specfemG', 'specfemF', and 'wavePlot'. A large double-headed arrow labeled 'source' and 'end' spans the graph. Below the graph, a 'Data File' window shows a waveform plot with the title 'Link to data files or data images preview...'. The interface also includes a 'Setup' tab and a 'FORWARD-MODELLING' window at the top right.

- **Searches** over products metadata within and across runs
- **Data download** and preview
- **Diagnostic** / errors
- **Lineage**: Multi directional navigations across data dependencies
- **W3C PROV-DM** as reference model.



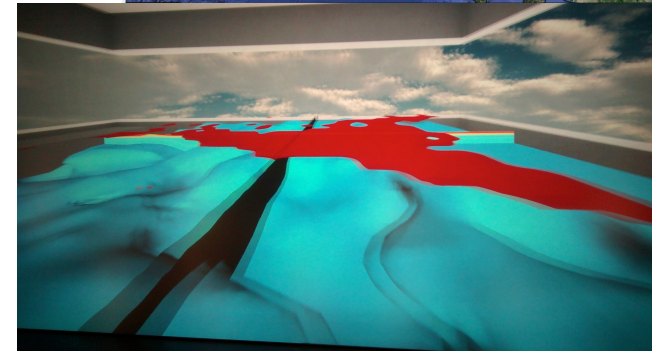
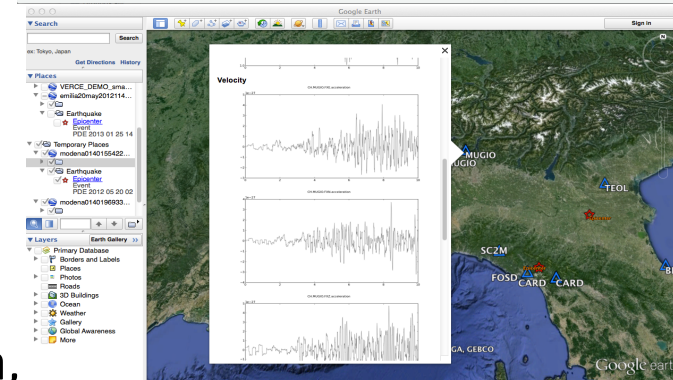
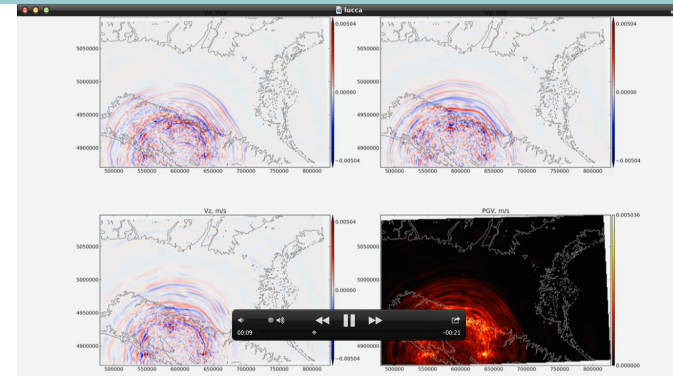
## Earthquake Simulation and Misfit Calculation Data overview:

**Synthetic Data:** seismograms, plots, 3D Geometry, Videos, KMZ packages, **meshes and models.**

(100 stations = 900 products and metadata )

**6-10 GB for a SPECIFEM3D simulation on 1000 cores**

**Raw Data:** on demand **access** and **staging** of **observational data** from **EIDA**: Earthquake Metadata, Sensors Metadata, waveform on regional scale.  
**(At the moment all via the FDSN WEB API)**





## Products

### Earthquake

**\_event\_id:** smi:local/ndk/C201302162116A/event  
m\_tp: -15000000000000000  
m\_tt: 5420000000000000  
m\_rr: -31300000000000000  
m\_rp: 29900000000000000  
m\_pp: 25900000000000000  
m\_rt: 87300000000000000  
description: SOUTHERN ITALY, C201302162116A  
depth\_in\_km: 21.5  
latitude: 41.7  
magnitude: 4.95  
longitude: 13.51  
origin\_time: 2013-02-16T21:16:13.300000

### Waveform

**MY\_TERM:VALUE**  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
network: IV  
longitude: 14.347  
latitude: 42.2273  
npts: 10000  
station: CAFR  
starttime: 1969-12-31T23:59:59.990000Z  
delta: 0.001  
calib: 1  
sampling\_rate: 1000  
endtime: 1970-01-01T00:00:09.989000Z  
type: velocity

## Processes

**name= SPECFEM3D**  
NPROC=1000  
STARTTIME=2015-03-06  
16:50:59.654342



## Products

### Earthquake

**\_event\_id:** smi:local/ndk/C201302162116A/event  
**m\_tp:** -15000000000000000  
**m\_tt:** 5420000000000000  
**m\_rr:** -31300000000000000  
**m\_rp:** 29900000000000000  
**m\_pp:** 25900000000000000  
**m\_rt:** 87300000000000000  
**description:** SOUTHERN ITALY, C201302162116A  
**depth\_in\_km:** 21.5  
**latitude:** 41.7  
**magnitude:** 4.95  
**longitude:** 13.51  
**origin\_time:** 2013-02-16T21:16:13.300000

### Waveform

**MY\_TERM:VALUE**  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**network:** IV  
**longitude:** 14.347  
**latitude:** 42.2273  
**npts:** 10000  
**station:** CAFR  
**starttime:** 1969-12-31T23:59:59.990000Z  
**delta:** 0.001  
**calib:** 1  
**sampling\_rate:** 1000  
**endtime:** 1970-01-01T00:00:09.989000Z  
**type:** velocity

## Processes

**name= SPECFEM3D**  
**NPROC=1000**  
**STARTTIME=2015-03-06**  
**16:50:59.654342**





## Products

### Earthquake

**\_event\_id:** smi:local/ndk/C201302162116A/event  
**m\_tp:** -15000000000000000  
**m\_tt:** 5420000000000000  
**m\_rr:** -31300000000000000  
**m\_rp:** 29900000000000000  
**m\_pp:** 25900000000000000  
**m\_rt:** 87300000000000000  
**description:** SOUTHERN ITALY, C201302162116A  
**depth\_in\_km:** 21.5  
**latitude:** 41.7  
**magnitude:** 4.95  
**longitude:** 13.51  
**origin\_time:** 2013-02-16T21:16:13.300000

### Waveform

**MY\_TERM:VALUE**  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**network:** IV  
**longitude:** 14.347  
**latitude:** 42.2273  
**npts:** 10000  
**station:** CAFR  
**starttime:** 1969-12-31T23:59:59.990000Z  
**delta:** 0.001  
**calib:** 1  
**sampling\_rate:** 1000  
**endtime:** 1970-01-01T00:00:09.989000Z  
**type:** velocity

## Processes

**name= SPECFEM3D**  
NPROC=1000  
STARTTIME=2015-03-06  
16:50:59.654342

**name= WavePlot**  
line-color=blue  
error=FileNotFound



## Products

### Earthquake

**\_event\_id:** smi:local/ndk/C201302162116A/event  
**m\_tp:** -15000000000000000  
**m\_tt:** 5420000000000000  
**m\_rr:** -31300000000000000  
**m\_rp:** 29900000000000000

**m\_pp:** 259000  
**m\_rt:** 873000  
**description:** S  
**depth\_in\_km:**  
**latitude:** 41.7  
**magnitude:** 4  
**longitude:** 13  
**origin\_time:** 2

### Images

**format:** image/png  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**network:** IV  
**longitude:** 14.347  
**latitude:** 42.2273  
**npts:** 10000  
**station:** CAFR  
**starttime:** 1969-12-31T23:59:59.990000Z  
**delta:** 0.001  
**calib:** 1  
**sampling\_rate:** 1000  
**endtime:** 1970-01-01T00:00:09.989000Z  
**type:** velocity

### Waveform

**MY\_TERM:** VALUE  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**network:** IV  
**longitude:** 14.347  
**latitude:** 42.2273  
**station:** CAFR  
**starttime:** 1969-12-31T23:59:59.990000Z  
**delta:** 0.001  
**sampling\_rate:** 1000  
**endtime:** 1970-01-01T00:00:09.989000Z  
**type:** velocity

## Processes

**name= SPECFEM3D**  
**NPROC=1000**  
**STARTTIME=2015-03-06**  
**16:50:59.654342**

**name= WavePlot**  
**line-color=blue**  
**error=FileNotFound**



## Products

### Earthquake

\_event\_id: smi:local/ndk/C201302162116A/event  
 m\_tp: -15000000000000000  
 m\_tt: 54200000000000000  
 m\_rr: -31300000000000000  
 m\_rp: 29900000000000000  
 m\_pp: 25900000000000000  
 m\_rt: 87300000000000000  
 description: S  
 depth\_in\_km: 10  
 latitude: 41.7  
 magnitude: 4  
 longitude: 13  
 origin\_time: 2

### Images

**format: image/png**  
**id: b9a3cbcc-060a-11e4-96c2-bcaec52d20d4**  
 network: IV  
 longitude: 14.347  
 latitude: 42.2273  
 npts: 10000  
 station: CAFR  
 starttime: 1969-12-31T23:59:59.990000Z  
 delta: 0.001  
 calib: 1  
 sampling\_rate: 1000  
 endtime: 1970-01-01T00:00:09.989000Z  
 type: velocity

### Waveform

**MY TERM-VALUE**  
**id: b9a3cbcc-060a-11e4-96c2-bcaec52d20d4**  
 Duration: 30s  
 Sampling-rate: 25fps  
 59.990000Z  
 1970-01-01T00:00:09.989000Z  
 velocity

### Video

**format: video/mpg**  
**id: b9a3cbcc-060a-11e4-96c2-bcaec52d20d4**  
 Duration: 30s  
 Sampling-rate: 25fps

## Processes

**name= SPECFEM3D**  
 NPROC=1000  
 STARTTIME=2015-03-06  
 16:50:59.654342

**name= WavePlot**  
 line-color=blue  
 error=FileNotFound



## Products

### Earthquake

**\_event\_id:** smi:local/ndk/C201302162116A/event  
**m\_tp:** -15000000000000000  
**m\_tt:** 54200000000000000  
**m\_rr:** -31300000000000000  
**m\_rp:** 29900000000000000  
**m\_pp:** 25900000000000000  
**m\_rt:** 87300000000000000  
**description:** S  
**depth\_in\_km:** 10  
**latitude:** 41.7  
**magnitude:** 4  
**longitude:** 13  
**origin\_time:** 2013-02-16T23:59:59.990000Z

### Images

**format:** image/png  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**network:** IV  
**longitude:** 14.347  
**latitude:** 42.2273  
**npts:** 10000  
**station:** CAFR  
**starttime:** 1969-12-31T23:59:59.990000Z  
**delta:** 0.001  
**calib:** 1  
**sampling\_rate:** 1000  
**endtime:** 1970-01-01T00:00:09.989000Z  
**type:** velocity

### Waveform

**format:** **MY\_TERM-VALUE**  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4

### Video

**format:** video/mpg  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**Duration:** 30s  
**Sampling-rate:** 25fps

### Mesh

**format:** binary  
**id:** b9a3cbcc-060a-11e4-96c2-bcaec52d20d4  
**maxlat:** 66.55  
**minlat:** 14.5  
**maxlon:** 33  
**minlong:** -12

## Processes

**name= SPECFEM3D**  
**NPROC=1000**  
**STARTTIME=2015-03-06**  
**16:50:59.654342**

**name= WavePlot**  
**line-color=blue**  
**error=FileNotFound**



W3C PROV-O offers a generic data model describing relations between data products and processes (**provenance**)



W3C PROV-O offers a generic data model describing relations between data products and processes (**provenance**)



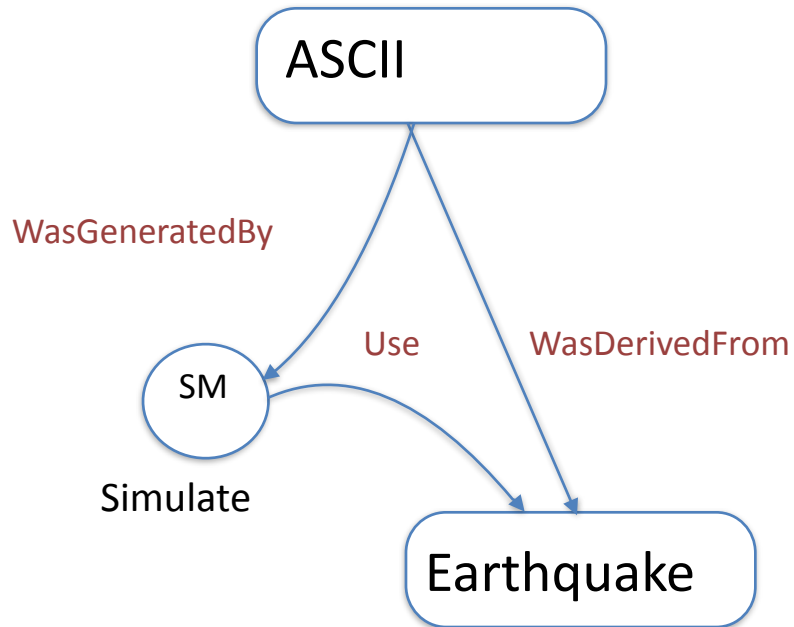
W3C PROV-O offers a generic data model describing relations between data products and processes (**provenance**)

Earthquake



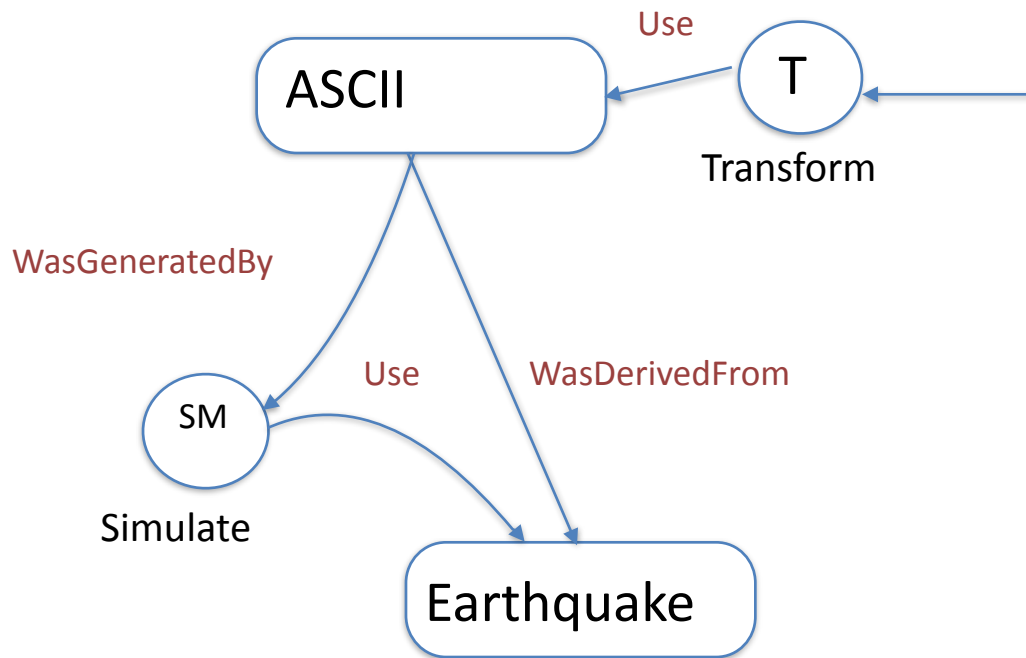


W3C PROV-O offers a generic data model describing relations between data products and processes (**provenance**)



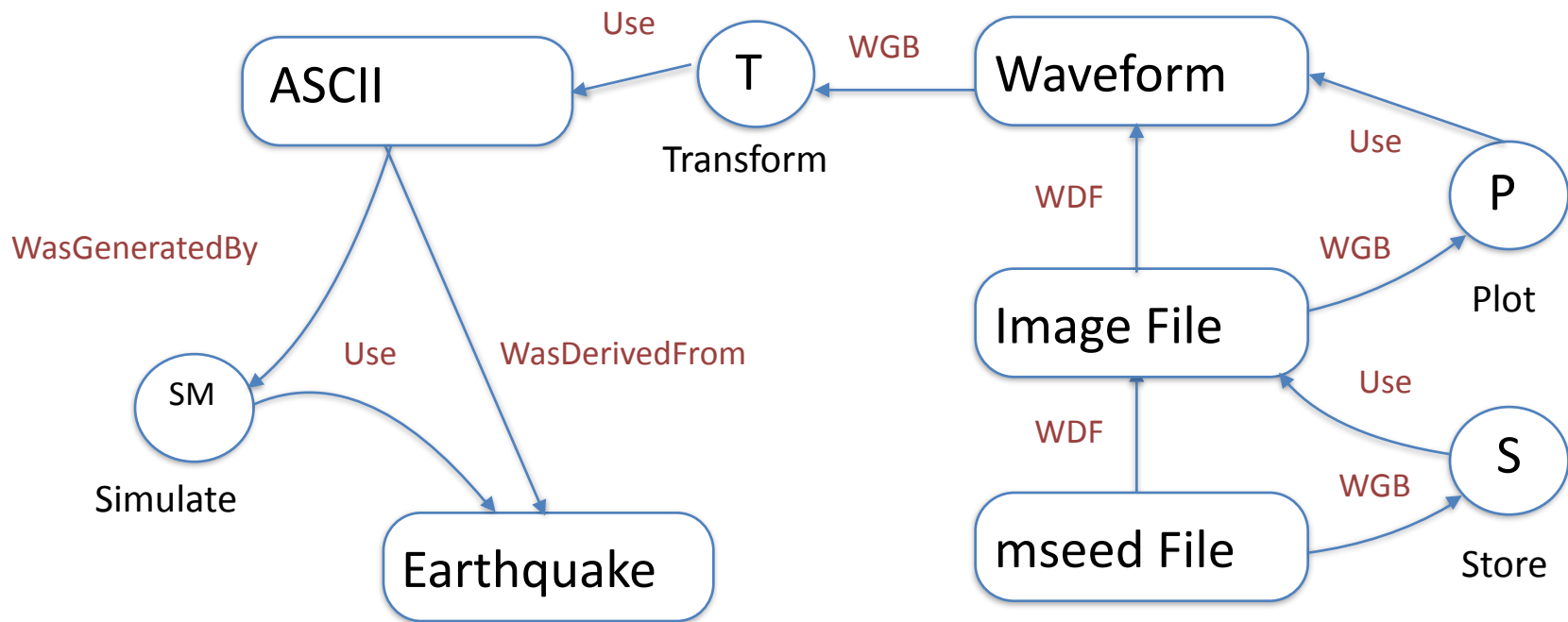


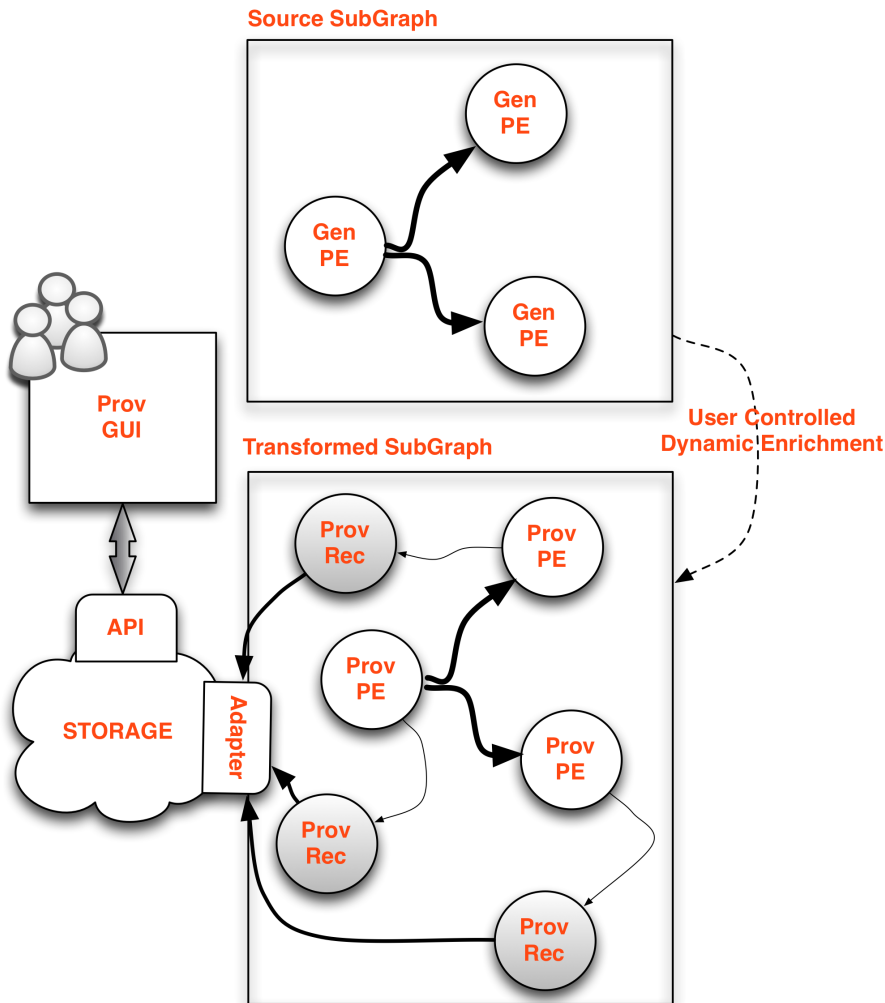
W3C PROV-O offers a generic data model describing relations between data products and processes (**provenance**)





W3C PROV-O offers a generic data model describing relations between data products and processes (provenance)





## Selective and extensible Provenance

The **GenericPEs** of a workflow subgraph are extended at runtime assuming the **ProvenancePE** type (*Dynamic Polymorphism*)

It gets connected to instances of **ProvenanceRecorderPE** for rapid analysis and transfer of the provenance traces to a dedicated catalogue.

The provenance data is then explored via a **GUI**, which connects to a web service interface.



## Selective and extensible Provenance

The *GenericPEs* of a workflow subgraph are extended at runtime assuming the *ProvenancePE* type (Dynamic Polymorphism)

### *ProvenancePE* extended functions:

- *Write(port, data, metadata, format, errors, control)*
- *ExtractItemMetadata(data)*

***port***: output port which is connect to the adjacent PEs.

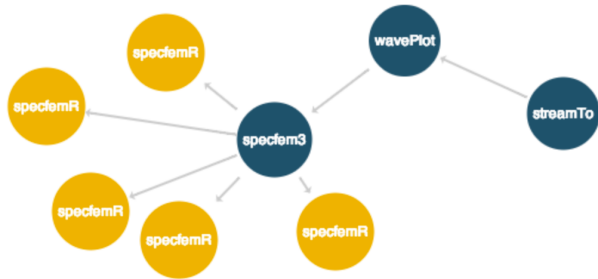
***data***: stream element to send to those adjacent PEs.

***metadata***: dictionary of metadata describing the data.

***format***: typically contains the mime-type of the data.

***errors***: erroneous situations and its description

***control***: are the control instructions like “*con:skip*” and “*con:immediateAccess*”



Process View - mordital0139524219...  
Data Derivations Graph

ID	Date	Error
8	2014-03-19 22:...	
9	2014-03-19 22:...	
10	2014-03-19 22:...	
11	2014-03-19 22:...	
12	2014-03-19 22:...	
13	2014-03-19 22:...	
14	2014-03-19 22:...	
15	2014-03-19 22:...	
16	2014-03-19 22:...	
17	2014-03-19 22:...	
18	2014-03-19 22:...	
19	2014-03-19 22:...	
20	2014-03-19 22:...	
21	2014-03-19 22:...	
22	2014-03-19 22:...	
23	2014-03-19 22:...	
24	2014-03-19 22:...	
25	2014-03-19 22:...	
26	2014-03-19 22:...	
27	2014-03-19 22:...	
28	2014-03-19 22:...	
29	2014-03-19 22:...	
30	2014-03-19 22:...	
31	2014-03-19 22:...	



Browse provenance  
and workflows output

VERCE  
portal.verce.eu Forward Modeling

forward-modelling

Name	Resource	Size	Date Modified
i19r01a36			May 4, 2014, 2:15 am
i19r01a35			May 4, 2014, 2:10 am
i19r01a03			May 4, 2014, 1:28 am
i19r01c17			May 4, 2014, 1:27 am
i19r01c16			May 4, 2014, 1:25 am
naslx			March 12, 2014, 11:14 am
work			March 10, 2014, 6:08 pm
pr49b			March 10, 2014, 6:08 pm
grid			March 10, 2014, 6:08 pm
hpc			March 10, 2014, 6:08 pm
modena01401549647956.mp4	uedinResc1	309.31 KB	May 31, 2014, 5:50 pm
modena01401471710308.mp4	uedinResc1	104.05 KB	May 31, 2014, 2:27 am

Access Data Products shipped from clusters  
to iRODS (e.g. **“con:immediateAccess”**)  
(Integrated Rule-Oriented Data System)



# Why ProvenanceRecorderPEs?

For a comprehensive provenance framework across DI-HPC models

1

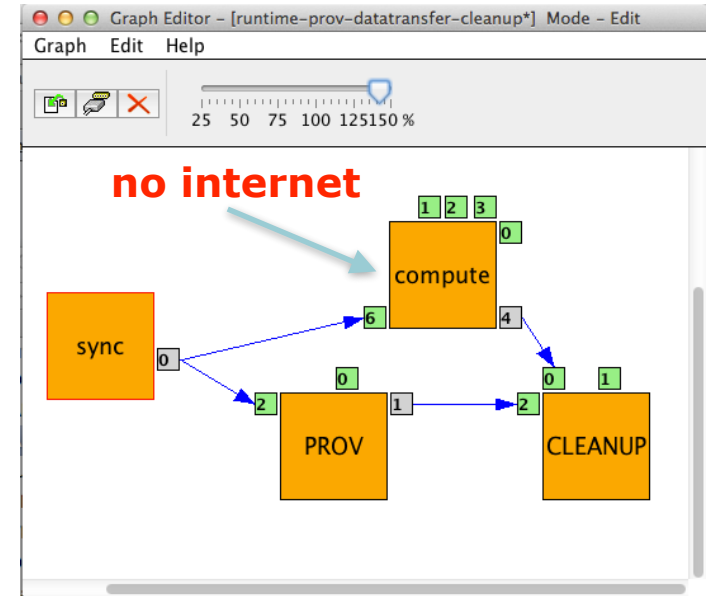
## Platform Control Workflow:

**sync:** metadata preparation and staging

**compute:** actual computation

**prov:** reads metadata, updates prov repository, intermediate data stage-out based on prov (**con:ImmediateAccess**)

**cleanup:** full data stageout and cleanups







# Why ProvenanceRecorderPEs?

For a comprehensive provenance framework across DI-HPC models

1

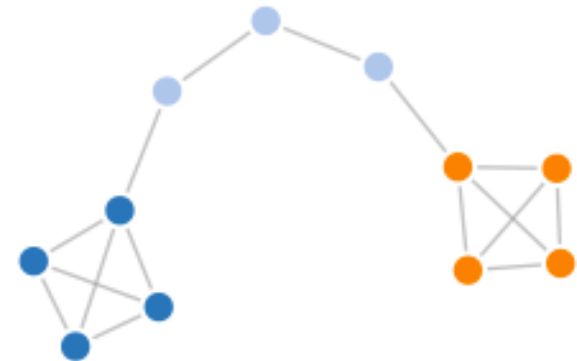
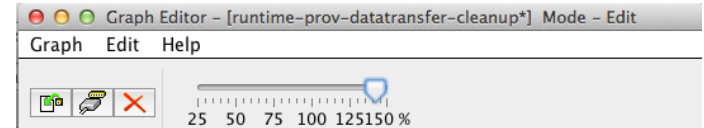
## Platform Control Workflow:

**sync:** metadata preparation and staging

**compute:** actual computation

**prov:** reads metadata, updates prov repository, intermediate data stage-out based on prov (**con:ImmediateAccess**)

**cleanup:** full data stageout and cleanups



**compute** receives and runs scientific workflows

2

## Science Case Workflow (According to community preferences)

Extraction of **user defined metadata** and **fine-grained Provenance**



## Use of *con:skip* for Selective provenance extraction

For a comprehensive provenance framework across DI-HPC models

**E.g. Large Streaming Cross-Correlations:** Data recorded by sensors are preprocessed in pipelines and cross correlated:  $n(n-1) / 2$  correlations  
*Full Provenance coverage would impact for 7-8% on our tests on ~350 sensors*

***con:skip*** enables provenance extraction and recording only for a subset of stations and for a specific set of metadata.

```
if stream.name not in list4prov:  
    self.write('output', stream, metadata=dic, control="con:skip")
```

**Note:** When using custom implementations of a ProvenancePE's function `ExtractItemMetadata(data)`, the metadata in *dic* will be added to the default.



# Use of Provenance for Diagnostics and Validation

For a comprehensive provenance framework across DI-HPC models

## Which PE in the pipeline is affecting the sampling rate?

The screenshot displays a provenance framework interface with the following components:

- Run activity monitor - concrete\_misfit\_preproc2**: A table listing simulation runs with columns for ID, Date, and Errors.
- Data Derivations Graph**: A directed graph showing data flow between processing elements (PE). Nodes include PE\_store, PE\_pre\_f, PE\_filto, and Rotation.
- Data Detail**: A window showing metadata for a specific run, including Date, Output Files, Output Metadata, Parameters, Annotations, and Errors.

The **Data Detail** window is open to the **Output Metadata** section, which contains the following information:

- calib: 1
- npts: 99
- endtime: 2013-02-16T21:16:12.189276Z

A red arrow points from the **npts: 99** value to the **PE\_pre\_f** node in the Data Derivations Graph.

### Metadata Dictionary



# Use of Provenance for Interactive Workflow Preparation

For a comprehensive provenance framework across DI-HPC models

## Simulations Runs

## Observed Data Pre-staging Runs

Simulation runs				raw-data download runs			
runId	Workflow name	Description	Date	runId	Workflow name	Description	Date
foggia01410428...	PRACE-LRZ-DIS...	foggia small	2014-09-11T09:...	test_download5	download	test	2015-05-28 13:...
foggia01410426...	PRACE-LRZ-DIS...	foggia con video	2014-09-11T09:...	test_download2	download	test	2015-05-22 12:...
modena014104...	PRACE-LRZ-DIS...	modena	2014-09-11T08:...	data_download_8	download	test	2015-05-21 13:...
abruzzo100001...	PRACE-LRZ-DIS...	d4py	2014-09-11T08:...	data_download_5	download	test	2015-05-21 13:...
abruzzo014104...	PRACE-LRZ-DIS...	d4py	2014-09-11T08:...	data_download...	download	test	2015-05-06 14:...
abruzzo014103...	PRACE-LRZ-DIS...	d4py	2014-09-10T12:...	data_download...	download	test	2015-05-06 14:...
modena014103...	PRACE-LRZ-DIS...	modena	2014-09-10T12:...	data_download...	download	test	2015-04-02 08:...

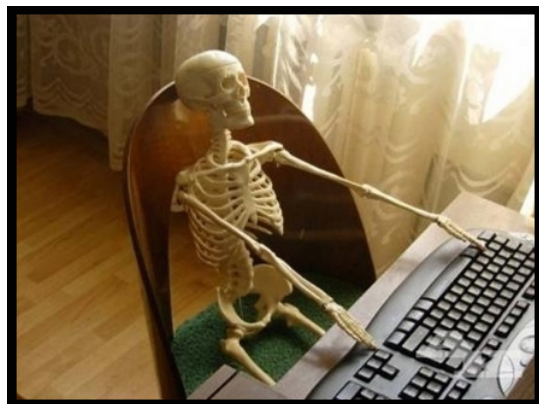
commons stations	
<input type="checkbox"/>	IV
<input checked="" type="checkbox"/>	<b>IV</b>
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input checked="" type="checkbox"/>	<b>IV</b>
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input type="checkbox"/>	IV
<input checked="" type="checkbox"/>	<b>IV</b>

**Sensors' data with common properties for MISFIT analysis (time-range, event)**

	TRIV
	<b>RNI2</b>
	T0204
	T0201
	SIRI
	GLCN
	SGA
	SOR
	<b>VULD</b>
	T1011
	VITU
	VULT
	<b>SGG</b>



## Runtime and Selective Provenance in Data Intensive Platforms



- Provides immediate feedback on the produced data with sensible and tuneable metadata.
- Useless waits for long and unfruitful runs can be reduce
- Fosters Dynamic Steering, Failure Recovery and Errors
- Detection, saving computing cycles, storage (\$\$) and energy!...



## References

P. Missier, S. Dey, K. Belhajjame, V. Cuevas-Vicenttin, and B. Ludascher. D-prov: extending the prov provenance model with workflow structure. In Proceedings of the 5th USENIX Workshop on the Theory and Practice of Provenance, TaPP '13, pages 9:1-9:7, Berkeley, CA, USA, 2013. USENIX Association.

M. Anand and S. Bowers. Exploring scientific workflow provenance using hybrid queries over nested data and lineage graphs. Scientific and Statistical Database Management, pages 237-254, 2009.

P. Missier and B. Ludascher. Linking multiple workflow provenance traces for interoperable collaborative science. Workflows in Support of Large-Scale Science (WORKS), 2010 5th Workshop on, 2010.

S. Madougou and S. Shahand. Characterizing workflow-based activity on a production e-infrastructure using provenance data. Future Generation Computer Systems, (Dci), 2013.

D. Crawl and I. Altintas. A provenance-based fault tolerance mechanism for scientific workflows. Provenance and Annotation of Data and Processes, pages 152-159, 2008.

P. Buneman, J. Cheney, and E.V. Kostylev. Hierarchical models of provenance. In Proceedings of the 4th USENIX conference on Theory and Practice of Provenance, TaPP'12, pages 10-10, Berkeley, CA, USA, 2012. USENIX Association.

E. Griffis, P. Martin, and J. Cheney. Semantics and provenance for processing element composition in dispel workflows. WORKS '13 Proceedings of the 8th Workshop on Workflows in Support of Large-Scale Science, 2013.

A. Spinuso, J. Cheney, and M. Atkinson. Provenance for seismological processing pipelines in a distributed streaming workflow. Proceedings of the Joint EDBT/ICDT 2013 Workshops, 2013.