

Data analytics workflows with Ophidia

Donatello Elia, Fabrizio Antonio, Alessandro D'Anca

Fondazione Centro Euro-Mediterraneo sui Cambiamenti Climatici (CMCC), Lecce, Italy



Ophidia

ENCCS/CMCC workshop:
*Training on HPDA for climate data with the
Ophidia framework*

11 November 2021

Session outline

Introduction to scientific workflows and motivations

Data analytics workflows in Ophidia

Ophidia workflows core concepts: JSON representation, workflow constructs, execution monitoring

Real-world examples of analytics workflow with the Ophidia framework

DEMO: Tutorial about workflow creation and execution with Ophidia

HANDS-ON: Data analytics workflows examples

Disclaimer: this material reflects only the authors' view, and the EU-Commission is not responsible for any use that may be made of the information it contains.



Large-scale climate analysis

Complexity of the analysis leads to the need for ***end-to-end workflow support***

- Typical approaches (mostly based on bash-like scripts) requires climate scientists to take care of implement and replicate workflow-like control logic
- Analyses can require the execution of *tens/hundreds of analytics operators*
 - *Efficient orchestration of the tasks is critical*
 - *Parallelism must be handled both at intra-task and inter-task level*
 - *Task failure should also be considered*

Workflows can represent a way to define ***portable*** and ***re-usable*** analyses
(targeting FAIR principles)



Session outline

Introduction to scientific workflows and motivations

Data analytics workflows in Ophidia

Ophidia workflows core concepts: JSON representation, workflow constructs, execution monitoring

Real-world examples of analytics workflow with the Ophidia framework

DEMO: DEMO: Tutorial about workflow creation and execution with Ophidia

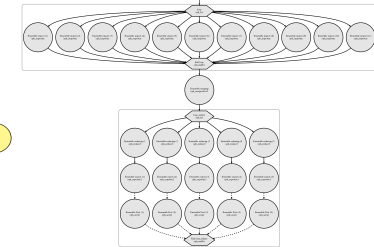
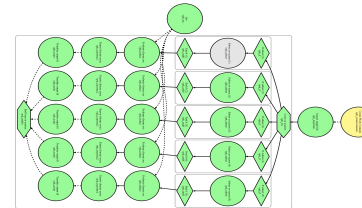
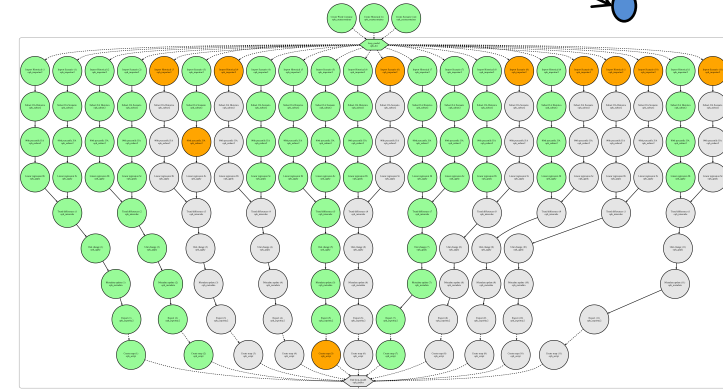
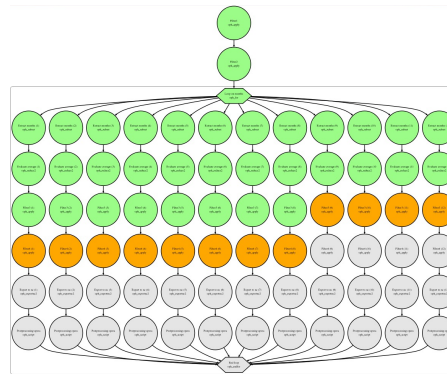
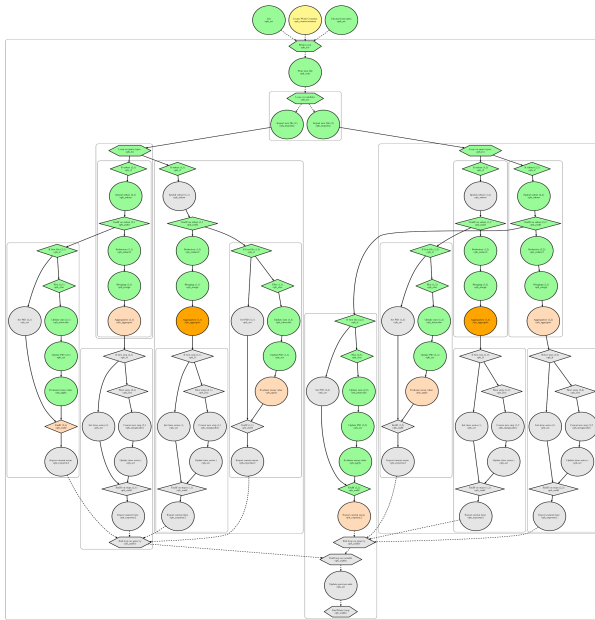
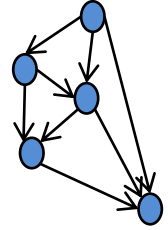
HANDS-ON: Data analytics workflows examples



Analytics workflows

Ophidia supports the execution of complex workflows of operators.

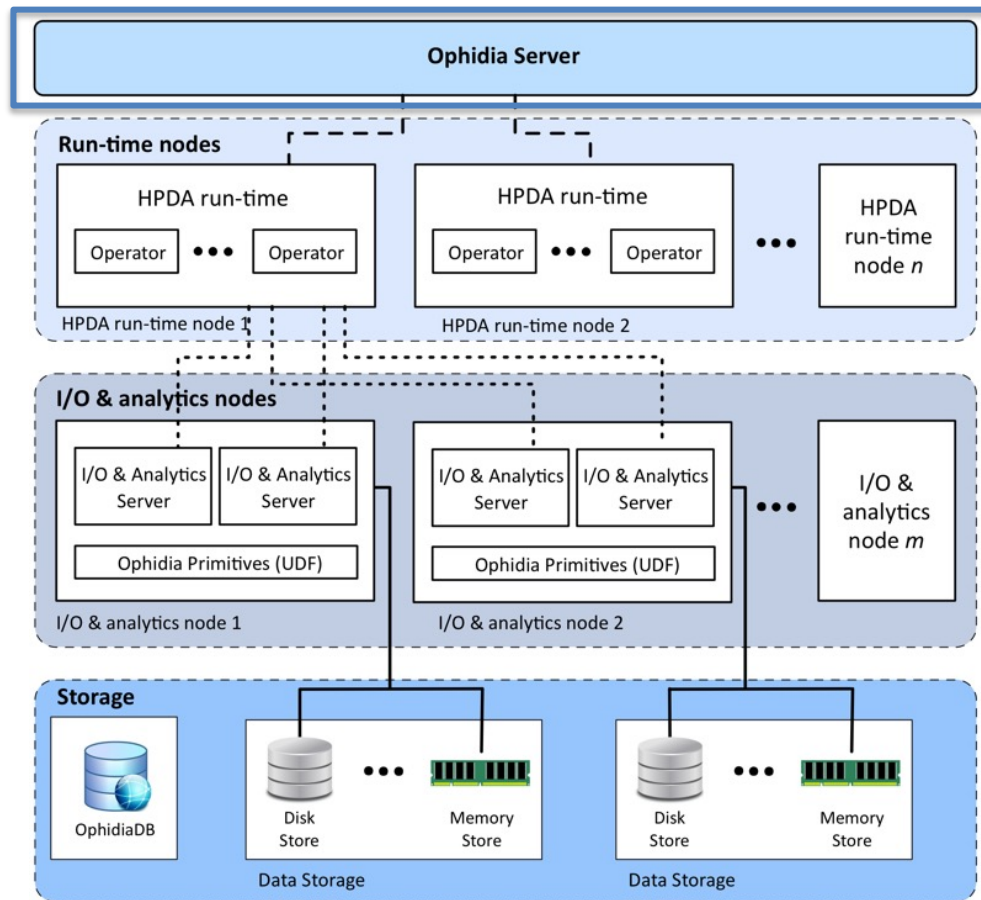
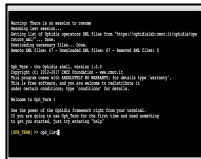
- Defines a **JSON representation** for the workflow DAG specification
- Supports different constructs: *dependencies; massive tasks; iterative (group of) tasks; parallel (group of) tasks; flow and error control*



C. Palazzo, A. Mariello, S. Fiore, A. D'Anca, D. Elia, D. N. Williams, G. Aloisio, "A Workflow-Enabled Big Data Analytics Software Stack for eScience", HPCS 2015, pp. 545-552



Ophidia architecture: front-end layer



The **Ophidia Server** is the **multi-interface** server front-end: OGC-WPS, WS-I

Manages user **authN/authZ**, **sessions** and enables server-side computation

Handles **single task** and **workflows** execution and monitors their execution on the server side

Remote interactions with:

- the Ophidia terminal CLI
- PyOphidia Python API
- WPS clients



Ophidia Terminal

The **Ophidia Terminal**, a CLI bash-like client for the Ophidia HPDA Framework:

- Executing *interactive* data analytics sessions;
- Submit *batch* data analytics tasks of *workflows*;
- Experiment and operators *debugging*;
- *File system exploration* and *environment management*.

```
[11..4495] >> oph_list level=2;
[Request]:
operator=oph_list;path=;level=2;sessionid=http://127.0.0.1/ophidia/sessions/1112
38695229505952271558621818154495/experiment;exec_mode=sync;cdd=/;

[JobID]:
http://127.0.0.1/ophidia/sessions/111238695229505952271558621818154495/experiment?2#45

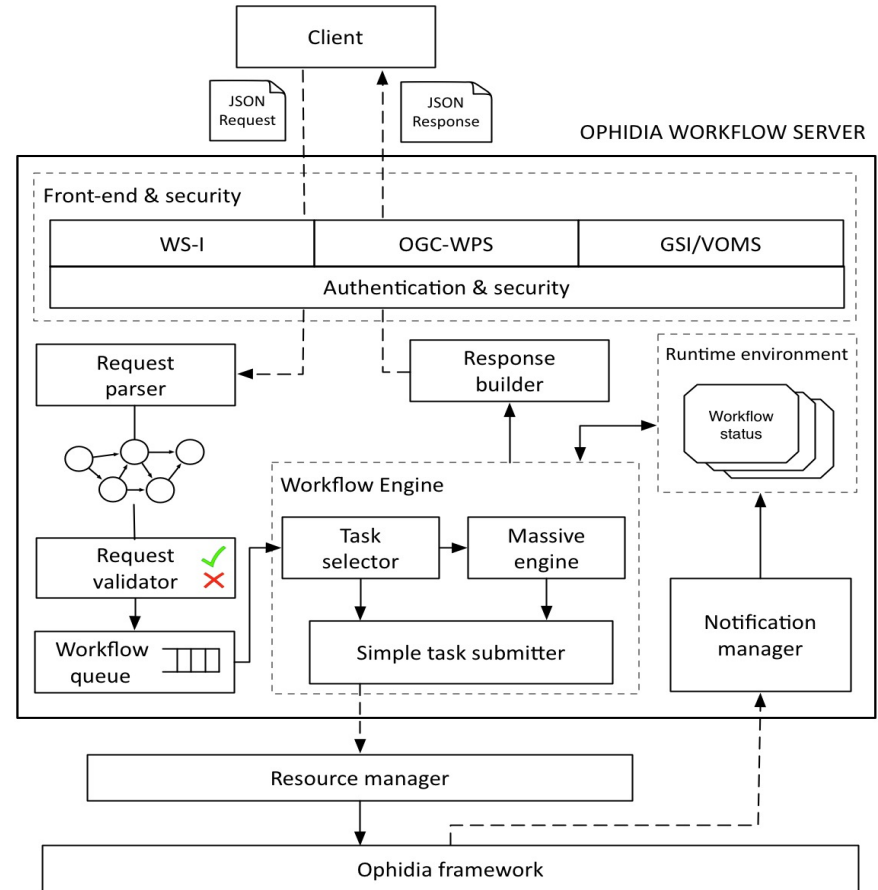
[Response]:
Ophidia Filesystem: /
-----
+==+=====+=====+=====+
| T | PATH                | DATACUBE PID                | DESCRIPTION |
+==+=====+=====+=====+
| f | testFolder/         |                              |             |
+==+=====+=====+=====+
| c | test                | http://127.0.0.1/ophidia/2917/374976 |             |
+==+=====+=====+=====+
```



The Ophidia Server

The **workflow management system (WMS)** is a core component of the Ophidia Server:

- *manages user request*
- *formats the commands for the analytics framework*
- *handles task dependencies and execution flow*
- *submits the tasks to the resource manager*
- *manages task status updates*
- *provides the proper response messages*



Session outline

Introduction to scientific workflows and motivations

Data analytics workflows in Ophidia

Ophidia workflows core concepts: JSON representation, workflow constructs, execution monitoring

Real-world examples of analytics workflow with the Ophidia framework

DEMO: DEMO: Tutorial about workflow creation and execution with Ophidia

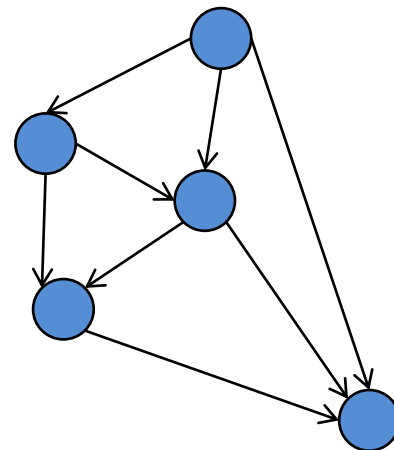
HANDS-ON: Data analytics workflows examples



Analytics Workflow Schema

Ophidia **workflows** schema:

- based on **JSON representation** for requests/responses
- defines application-level **semantic** and **syntactic rules**
- models scientific computations as **DAG**



Main supported abstractions:

- *Shared properties*
- *Flow/data dependencies*
- *Simple/massive tasks*
- *Iterative (group of) tasks*
- *Parallel (group of) tasks*
- *Flow and error control*
- *Interleaving and interactive tasks*



Behind the scene: workflow JSON representation

ophrpm@ophidiarpm:~/workflow

ophrpm@ophidiarpm:~/workflow

```
"tasks": [
  {
    "name": "Loop on tasmin and tasmax cubes",
    "operator": "oph_for",
    "arguments": [ "name=cube", "counter=1:2", "values=${1}|${2}", "parallel=yes" ]
  },
  {
    "name": "Compute operation over time",
    "operator": "oph_reduce2",
    "arguments": [
      "cube=@{cube}",
      "dim=time",
      "concept_level=M",
      "midnight=00",
      "operation=$3",
      "container=tmp"
    ],
    "dependencies": [
      { "task": "Loop on tasmin and tasmax cubes" }
    ]
  },
  {
    "name": "Conversion from Kelvin to Celsius degrees",
    "operator": "oph_apply",
    "arguments": [
      "query=oph_sum_scalar('oph_float','oph_float',measure,-273.15)"
    ],
    "dependencies": [{
      "task": "Compute operation over time",
      "type": "single"
    }]
  },
  {
    "name": "Loop for subset months",
    "operator": "oph_for",
    "arguments": [ "name=index", "counter=1:12", "values=Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec", "parallel=yes" ],
    "dependencies": [ { "task": "Conversion from Kelvin to Celsius degrees", "type": "single" } ]
  },
  {
    "name": "Subset on i-month",
    "operator": "oph_subset",
    "arguments": [
      "subset_dims=time",
      "subset_filter=&index:12:end"
    ],
    "dependencies": [
```

--More-- (65%)

Behind the scene: workflow JSON representation

ophrpm@ophidiarpm:~/workflow

ophrpm@ophidiarpm:~/workflow

```
"tasks": [
{
  "name": "Loop on tasmin and tasmax cubes",
  "operator": "oph_for",
  "arguments": [ "name=cube", "counter=1:2", "values=${1}|${2}", "parallel=yes" ]
},
```

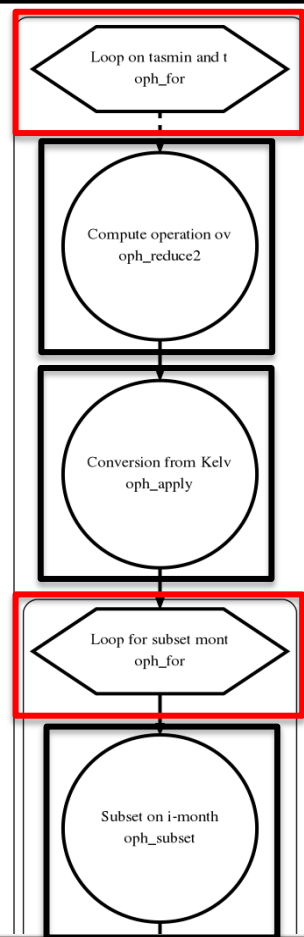
```
{
  "name": "Compute operation over time",
  "operator": "oph_reduce2",
  "arguments": [
    "cube=@{cube}",
    "dim=time",
    "concept_level=M",
    "midnight=00",
    "operation=$3",
    "container=tmp"
  ],
  "dependencies": [
    { "task": "Loop on tasmin and tasmax cubes" }
  ]
}
```

```
{
  "name": "Conversion from Kelvin to Celsius degrees",
  "operator": "oph_apply",
  "arguments": [
    "query=oph_sum_scalar('oph_float','oph_float',measure,-273.15)"
  ],
  "dependencies": [{
    "task": "Compute operation over time",
    "type": "single"
  }]
},
```

```
{
  "name": "Loop for subset months",
  "operator": "oph_for",
  "arguments": [ "name=index", "counter=1:12", "values=Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec", "parallel=y",
  "dependencies": [ { "task": "Conversion from Kelvin to Celsius degrees", "type": "single" } ]
},
```

```
{
  "name": "Subset on i-month",
  "operator": "oph_subset",
  "arguments": [
    "subset_dims=time",
    "subset_filter=&index:12:end"
  ],
  "dependencies": [
```

--More-- (65%)



Behind the scene: workflow JSON representation

ophrpm@ophidiarpm:~/workflow

ophrpm@ophidiarpm:~/workflow

```
{
  "name": "Loop on tasmin and tasmax cubes",
  "operator": "oph_for",
  "arguments": [ "name=cube", "counter=1:2", "values=${1}|${2}", "parallel=yes" ]
},
```

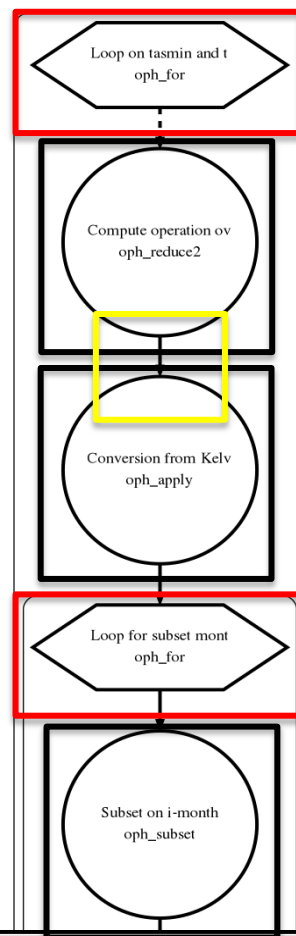
```
{
  "name": "Compute operation over time",
  "operator": "oph_reduce2",
  "arguments": [
    "cube=@{cube}",
    "dim=time",
    "concept_level=M",
    "midnight=00",
    "operation=$3",
    "container=tmp"
  ],
  "dependencies": [
    { "task": "Loop on tasmin and tasmax cubes" }
  ]
},
```

```
{
  "name": "Conversion from Kelvin to Celsius degrees",
  "operator": "oph_apply",
  "arguments": [
    "query=oph_sum_scalar('oph_float','oph_float',measure,-273.15)"
  ],
  "dependencies": [{
    "task": "Compute operation over time",
    "type": "single"
  }]
},
```

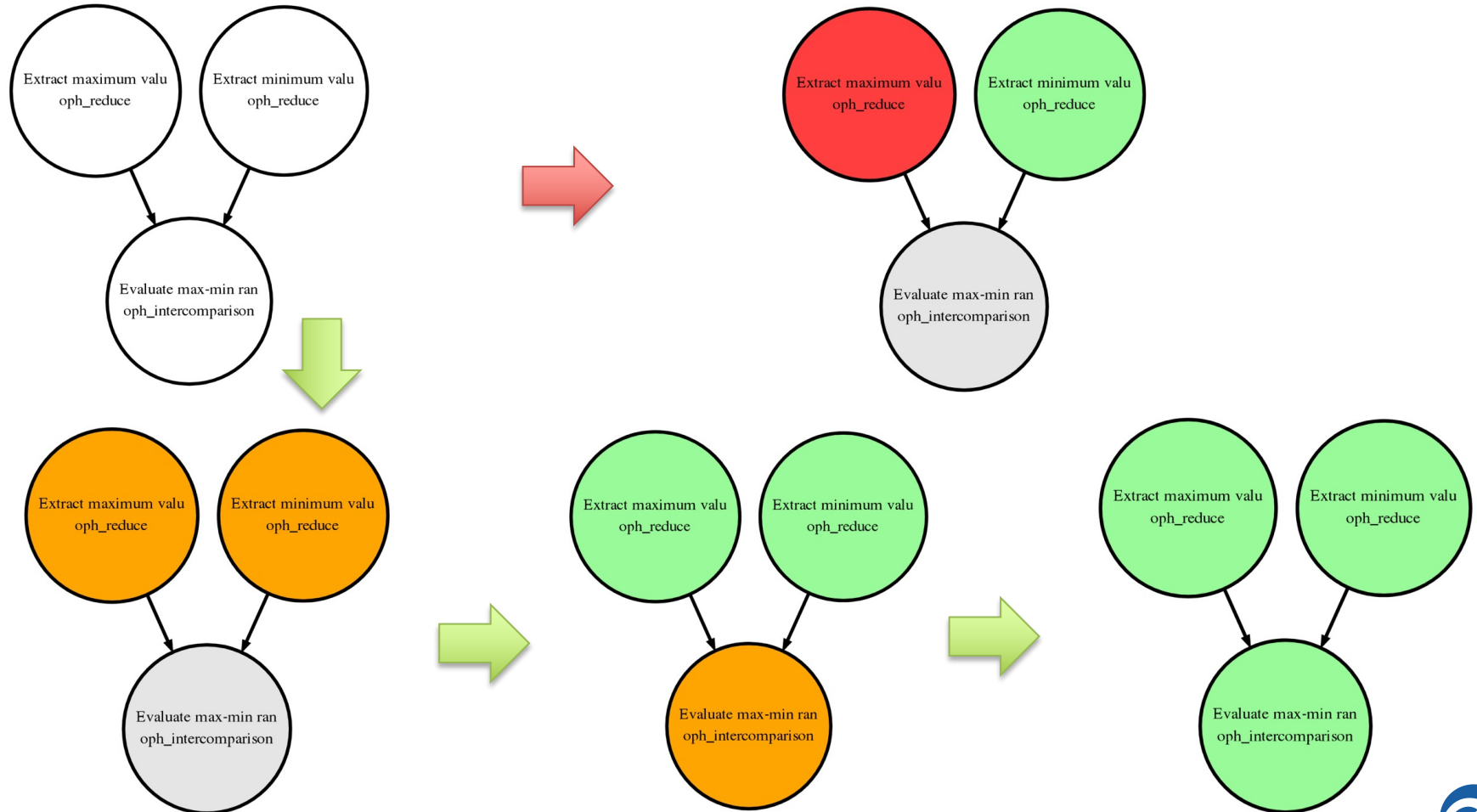
```
{
  "name": "Loop for subset months",
  "operator": "oph_for",
  "arguments": [ "name=index", "counter=1:12", "values=Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec", "parallel=yes" ],
  "dependencies": [ { "task": "Conversion from Kelvin to Celsius degrees", "type": "single" } ]
},
```

```
{
  "name": "Subset on i-month",
  "operator": "oph_subset",
  "arguments": [
    "subset_dims=time",
    "subset_filter=&index:12:end"
  ],
  "dependencies": [
```

--More-- (65%)



Workflow status monitoring



Workflow submission

```
ophrpm@ophidiarpm:~/devel/oph-client/res x oprpm@ophidiarpm:~/workflow
[37..6380] >>
[37..6380] >> ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max
[JobID]:
http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144

[37..6380] >> view 247
[247] ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max [http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144]

[Response]:
Workflow Status
-----
OPH_STATUS_COMPLETED

Workflow Progress
-----
+=====+
| NUMBER OF COMPLETED TASKS | TOTAL NUMBER OF TASKS |
+=====+
| 82 | 82 |
+=====+

Workflow Task List
-----
+=====+
| OPH JOB ID | SESSION CODE | WORKFL | MARKE | PARENT MA | TASK NAME | TYP | EXIT STATUS |
| | | OW ID | R ID | RKER ID | | E | |
+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3145 | 376699238311302232511449455166146380 | 247 | 3145 | 3144 | Loop on tasmin and tasmax cubes | SIM PLE | OPH_STATUS_COMPLETED |
+-----+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3146 | 376699238311302232511449455166146380 | 247 | 3146 | 3144 | Compute operation over time (1) | SIM PLE | OPH_STATUS_COMPLETED |
+-----+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3147 | 376699238311302232511449455166146380 | 247 | 3147 | 3144 | Compute operation over time (2) | SIM PLE | OPH_STATUS_COMPLETED |
+-----+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3148 | 376699238311302232511449455166146380 | 247 | 3148 | 3144 | Conversion from Kelvin to Celsius degrees (1) | SIM PLE | OPH_STATUS_COMPLETED |
+-----+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3149 | 376699238311302232511449455166146380 | 247 | 3149 | 3144 | Conversion from Kelvin to Celsius degrees (2) | SIM PLE | OPH_STATUS_COMPLETED |
+-----+
```

Workflow submission

```
ophrpm@ophidiarpm:~/devel/oph-client/res x oprpm@ophidiarpm:~/workflow
[37..6380] >> ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max
[37..6380] >> [JobID]: http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144
[37..6380] >> view 247
[247] ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max [http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144]

[Response]:
Workflow Status
-----
OPH_STATUS_COMPLETED

Workflow Progress
-----
=====+=====+
| NUMBER OF COMPLETED TASKS | TOTAL NUMBER OF TASKS |
| 82 | 82 |
=====+=====+

Workflow Task List
-----
=====+=====+=====+=====+=====+=====+=====+=====+
| OPH JOB ID | SESSION CODE | WORKFL | MARKE | PARENT MA | TASK NAME | TYP | EXIT STATUS |
| | | OW ID | R ID | RKER ID | | E | |
=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3145 | 376699238311302232511449455166146380 | 247 | 3145 | 3144 | Loop on tasmin and tasmax cubes | SIM | OPH_STATUS_COMPLETED |
| | | | | | | PLE | |
=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3146 | 376699238311302232511449455166146380 | 247 | 3146 | 3144 | Compute operation over time (1) | SIM | OPH_STATUS_COMPLETED |
| | | | | | | PLE | |
=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3147 | 376699238311302232511449455166146380 | 247 | 3147 | 3144 | Compute operation over time (2) | SIM | OPH_STATUS_COMPLETED |
| | | | | | | PLE | |
=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3148 | 376699238311302232511449455166146380 | 247 | 3148 | 3144 | Conversion from Kelvin to Celsius degrees (1) | SIM | OPH_STATUS_COMPLETED |
| | | | | | | PLE | |
=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3149 | 376699238311302232511449455166146380 | 247 | 3149 | 3144 | Conversion from Kelvin to Celsius degrees (2) | SIM | OPH_STATUS_COMPLETED |
| | | | | | | PLE | |
=====+=====+=====+=====+=====+=====+=====+=====+
```

Workflow submission

ophrpm@ophidiarpm:~/devel/oph-client/res

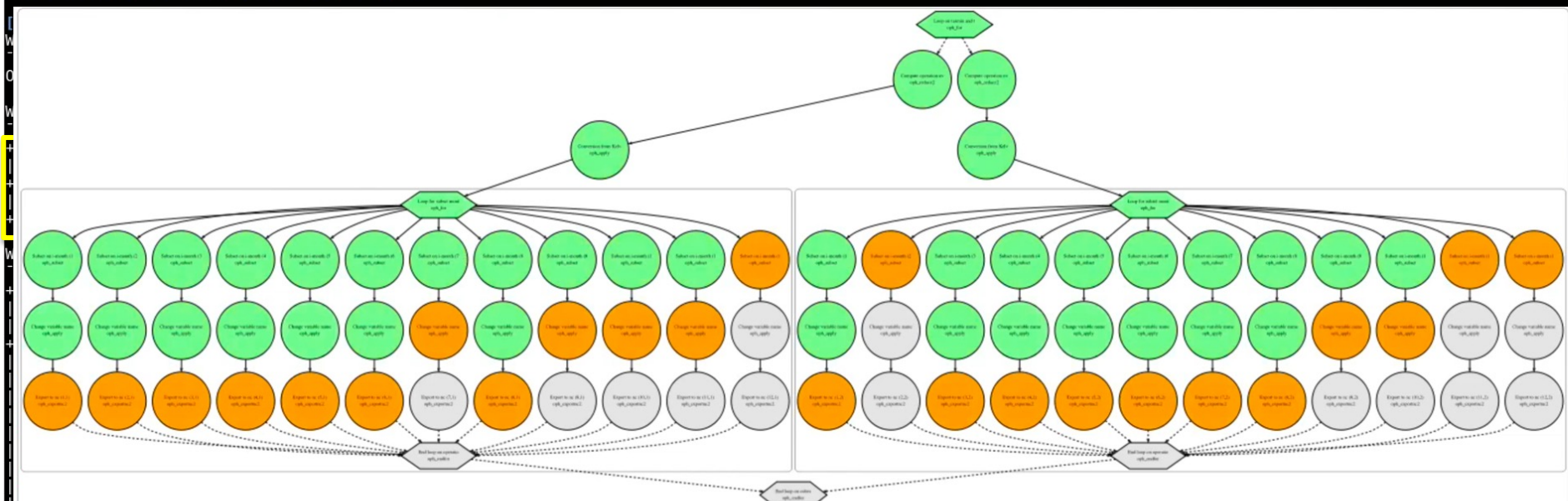
ophrpm@ophidiarpm:~/workflow

```
[37..6380] >> ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max
```

```
[JobID]: http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144
```

```
[37..6380] >> view 247
```

```
[247] ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max [http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144]
```



http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3148	376699238311302232511449455166146380	247	3148	3144	Conversion from Kelvin to Celsius degrees (1)	SIMPLE	OPH_STATUS_COMPLETED
http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3149	376699238311302232511449455166146380	247	3149	3144	Conversion from Kelvin to Celsius degrees (2)	SIMPLE	OPH_STATUS_COMPLETED

Analytics workflows constructs

Workflow Management

This group includes a number of flow control operators that could be used within an [Ophidia workflow](#) to implement complex data processing in batch mode. In particular, they implement several advanced features: [setting of run-time variables](#), [iterative and parallel interface](#), [selection interface](#), [interactive workflows](#), [interleaving workflows](#), etc.

NAME	DESCRIPTION
OPH_ELSE	Start the last sub-block of a selection block "if".
OPH_ELSEIF	Start a new sub-block of a selection block "if".
OPH_ENDFOR	Close a loop "for".
OPH_ENDIF	Close a selection block "if".
OPH_FOR	Implement a loop "for".
OPH_IF	Open a "if" selection block.
OPH_INPUT	It sends commands or data to an interactive task.
OPH_SET	Set a parameter in the workflow environment.
OPH_WAIT	Wait until an event occurs.



Iterative Interface

Allows to repeat the execution of a block of workflow **tasks** over different input data or over the result of the previous iteration.

Selection interface operators:

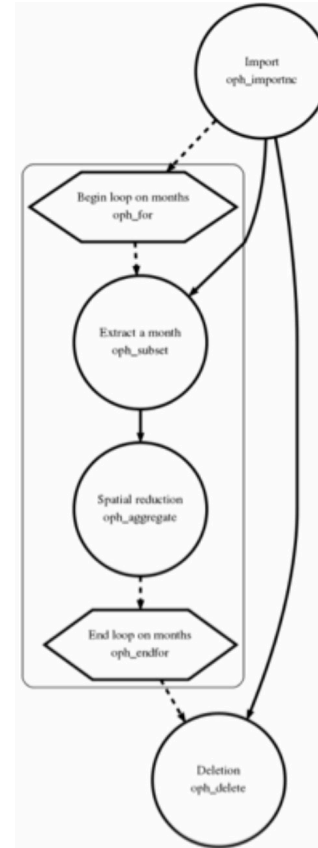
- *OPH_FOR*
- *OPH_ENDFOR*

The statement can be used in nested fashion

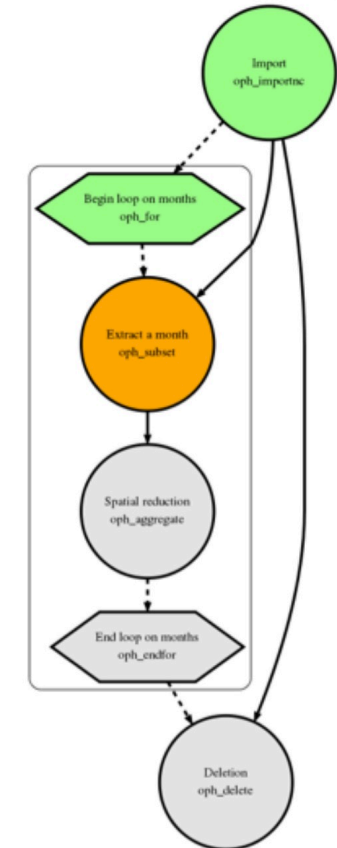
```
{
  "name": "Begin loop on months",
  "operator": "oph_for",
  "arguments":
  [
    "name=index",
    "counter=1:12",
    "values=Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec"
  ]
},
```

Workflow iterative interface documentation: http://ophidia.cmcc.it/documentation/users/workflow/workflow_for.html

AT DEFINITION TIME



AT RUNTIME



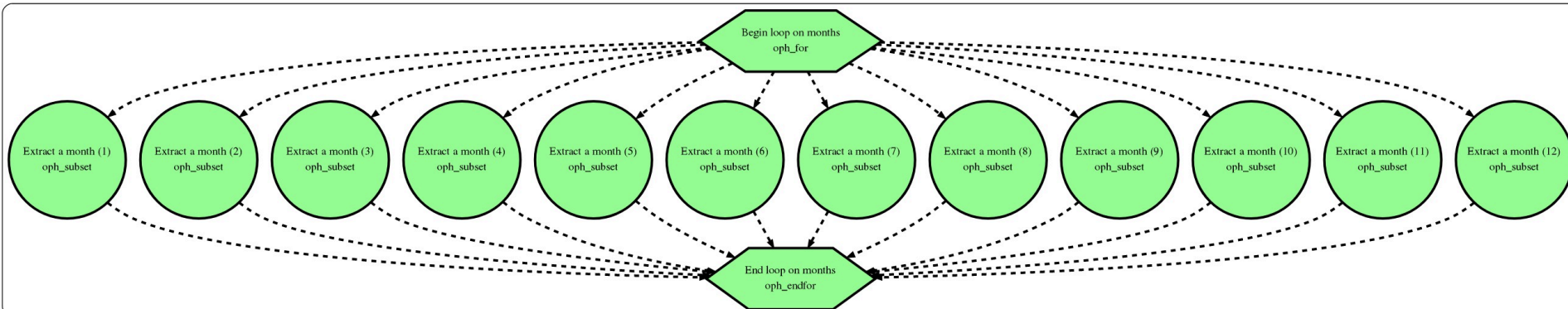
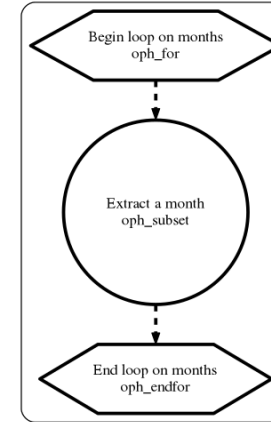
Parallel Interface

Extension of the OPH_FOR interface for parallel
(**concurrent**) execution of the loop iterations.

```
{  
  "name": "Begin loop on months",  
  "operator": "oph_for",  
  "arguments":  
  [  
    "parallel=yes",  
    "name=index",  
    "counter=1:12",  
    "values=Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec"  
  ]  
}
```

AT RUNTIME

AT DEFINITION TIME



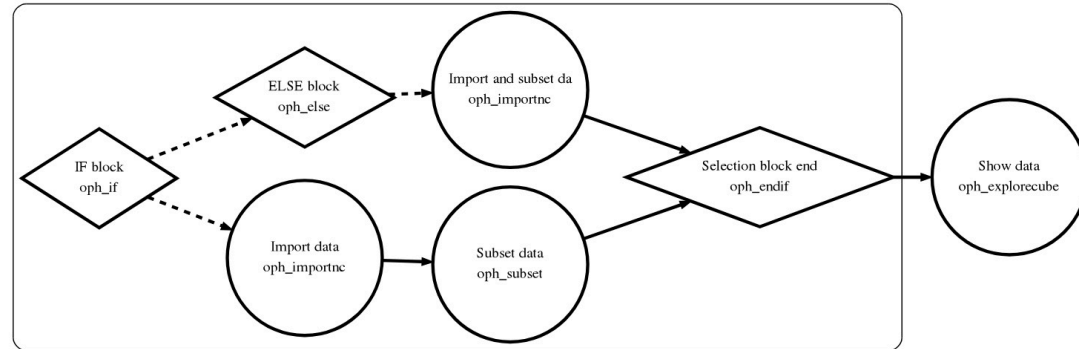
Selection Interface

Enables the workflow manager to **dynamically execute a block of tasks** based on boolean conditions evaluated at run-time.

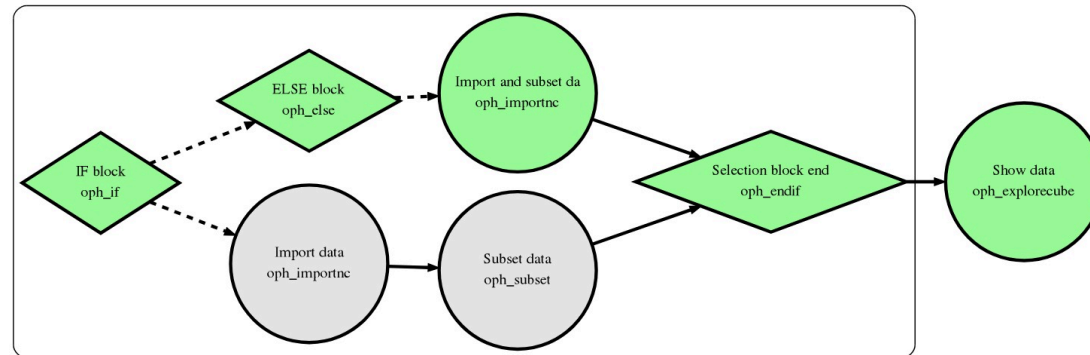
Selection interface operators:

- *OPH_IF*
- *OPH_ELSEIF*
- *OPH_ELSE*
- *OPH_ENDIF*

AT DEFINITION TIME



AT RUNTIME



Workflow error handling

In case of very large workflow executions **errors** in one of more **tasks** are likely.

Supported behaviours in case of task failure:

- *break*: the workflow is interrupted
- *skip*: the task is skipped and execution continues on the descendant tasks
- *continue*: the task and all depending task will be ignored, while other task will be executed
- *repeat N*: the task is re-executed N times

DEFINED AT GLOBAL WORKFLOW LEVEL

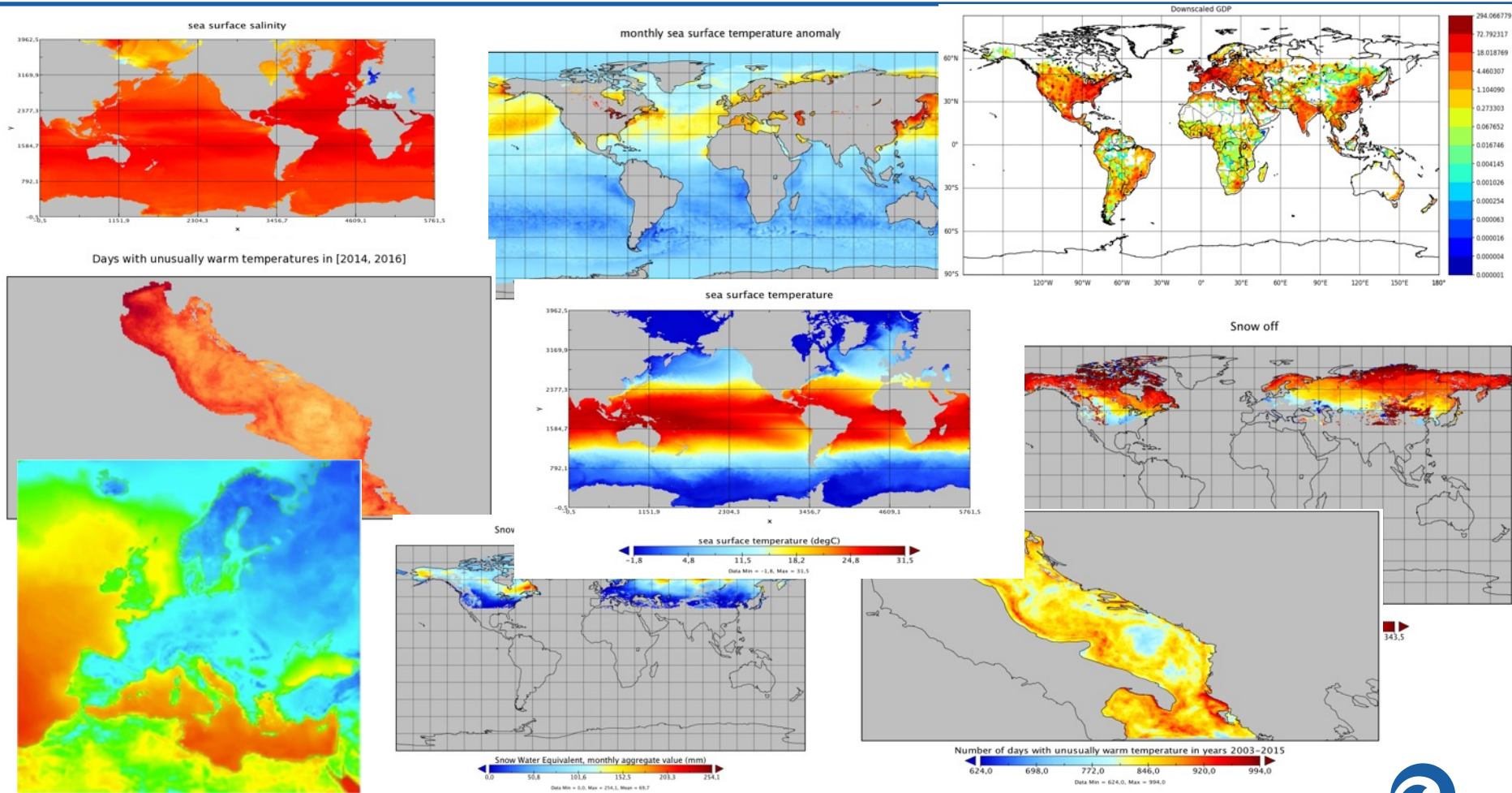
```
"name": "Example5",  
"author": "Foo",  
"abstract": "Simple workflow with automatic repetition",  
"exec_mode": "sync",  
"ncores": "1",  
"cube": "http://hostname/1/1",  
"on_error": "repeat 2",  
"tasks":
```

DEFINED AT TASK LEVEL (precedence)

```
{  
    "name": "Extract maximum value",  
    "operator": "oph_reduce",  
    "arguments": [ "operation=max" ],  
    "on_error": "repeat 5"  
},
```



Efficient support for advanced analytics experiments



Session outline

Introduction to scientific workflows and motivations

Data analytics workflows in Ophidia

Ophidia workflows core concepts: JSON representation, workflow constructs, execution monitoring

Real-world examples of analytics workflow with the Ophidia framework

DEMO: DEMO: Tutorial about workflow creation and execution with Ophidia

HANDS-ON: Data analytics workflows examples

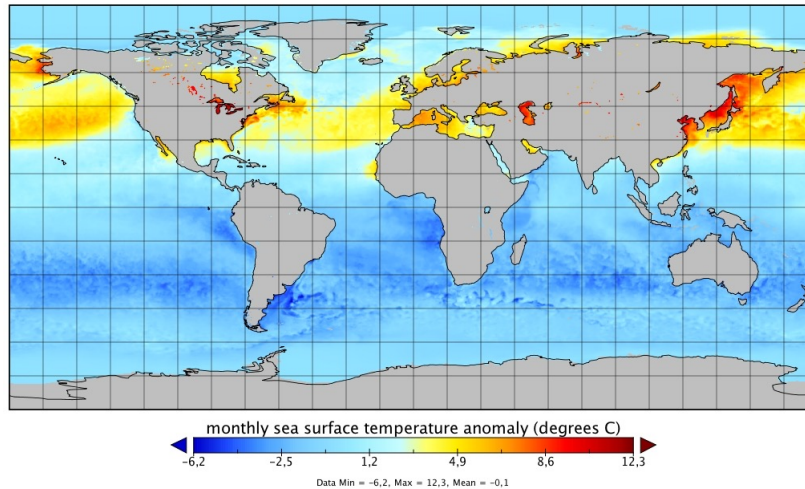


Workflow example I: climate indicators processing

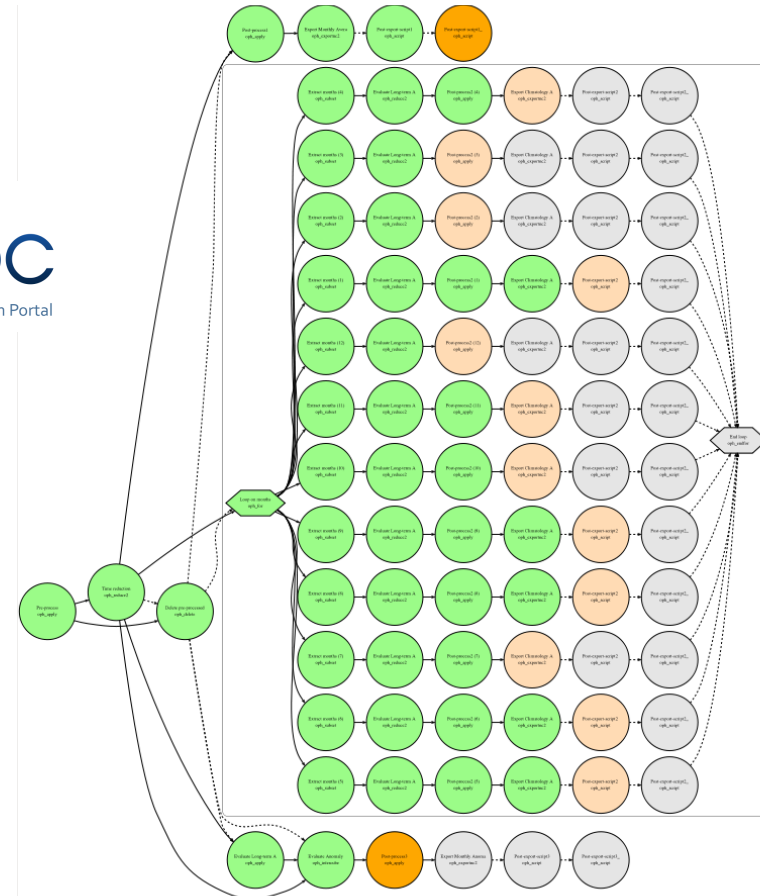
SST (monthly) mean, anomaly, climatological mean

- Dataset time range: 1991-2010
- **7062** nc files
- **350GB** of input data
- **87 tasks** performed
- **12x51MB + 2x12GB** of output files

monthly sea surface temperature anomaly



Clipc
Climate Information Portal



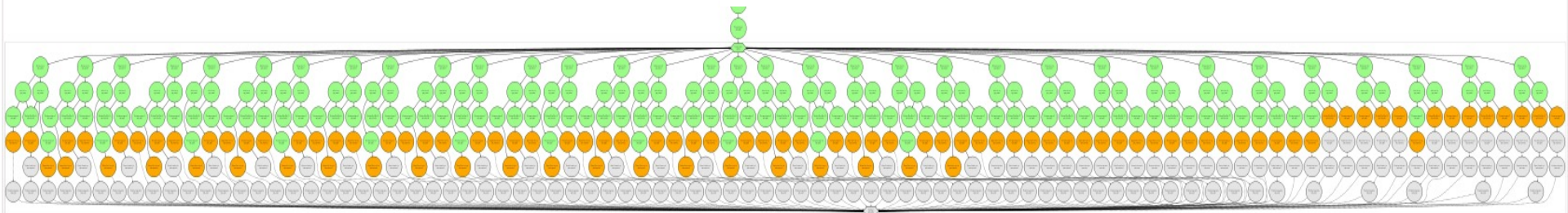
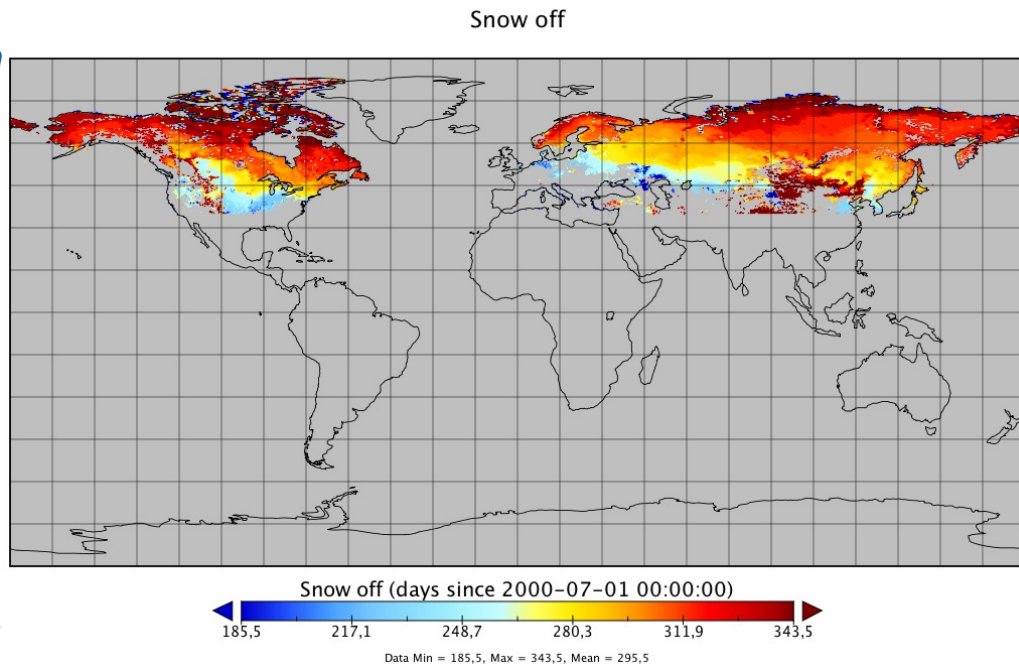
A. D'Anca, et al., "On the Use of In-memory Analytics Workflows to Compute eScience Indicators from Large Climate Datasets," 2017 17th IEEE/ACM Int. Symposium on Cluster, Cloud and Grid Computing (CCGRID), pp. 1035-1043



Workflow example II: climate indicators processing

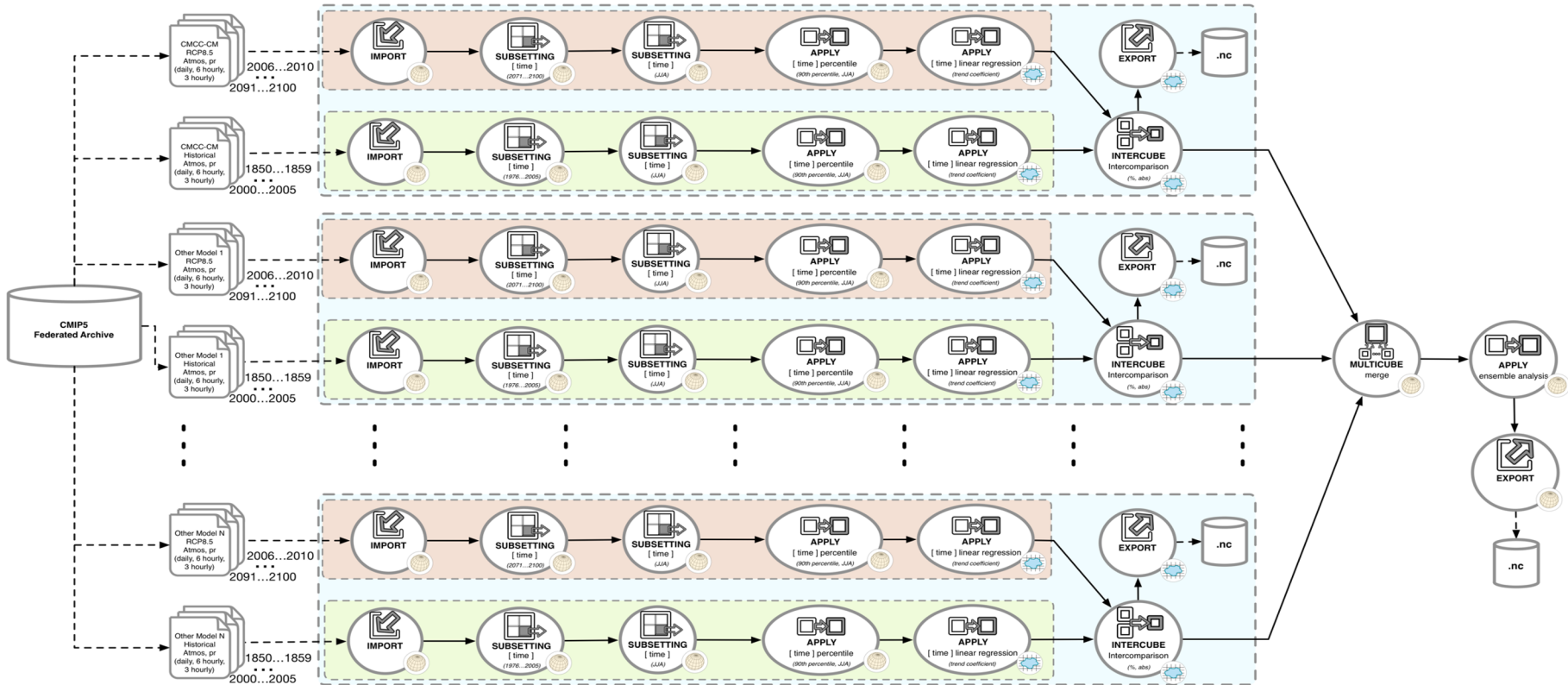
Snow on/off – Length of snow season (single workflow for 3 indicators)

- Dataset time range: 1979-2012
- **6341** nc files
- **50 GB** of input data
- **599 tasks** performed
- **99 NetCDF output files** (6MB each)
- **21 tasks** in the exp. description



Workflow example III: Multi-model experiment design

Precipitation Trend Analysis use case implemented as an Ophidia workflow



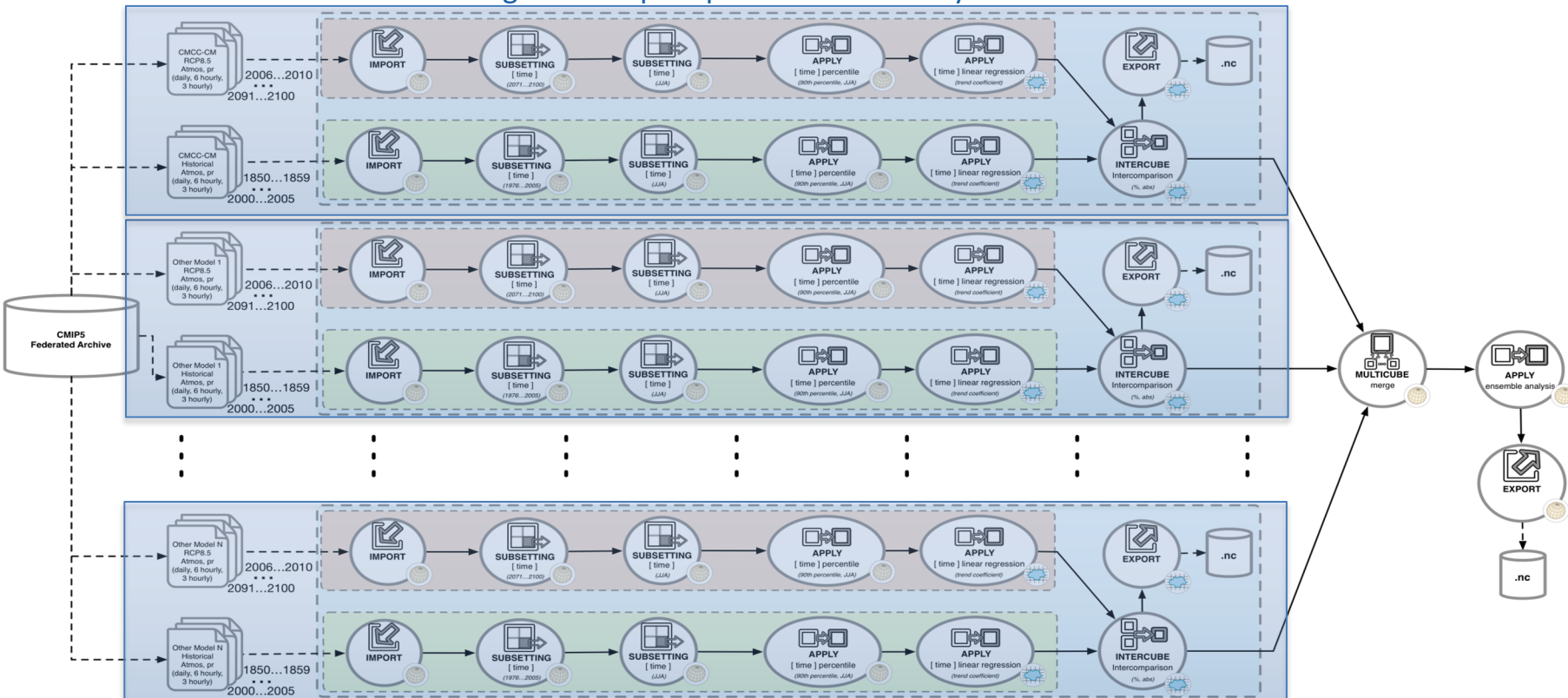
S. Fiore, et al., "Distributed and cloud-based multi-model analytics experiments on large volumes of climate change data in the earth system grid federation eco-system". In *Big Data (Big Data)*, 2016 IEEE Int. Conference on. IEEE, 2016. pp. 2911-2918



Workflow example III: Multi-model experiment design

Precipitation Trend Analysis use case implemented as an Ophidia workflow

Single model precipitation trend analysis



S. Fiore, et al., "Distributed and cloud-based multi-model analytics experiments on large volumes of climate change data in the earth system grid federation eco-system". In *Big Data (Big Data)*, 2016 IEEE Int. Conference on. IEEE, 2016. pp. 2911-2918

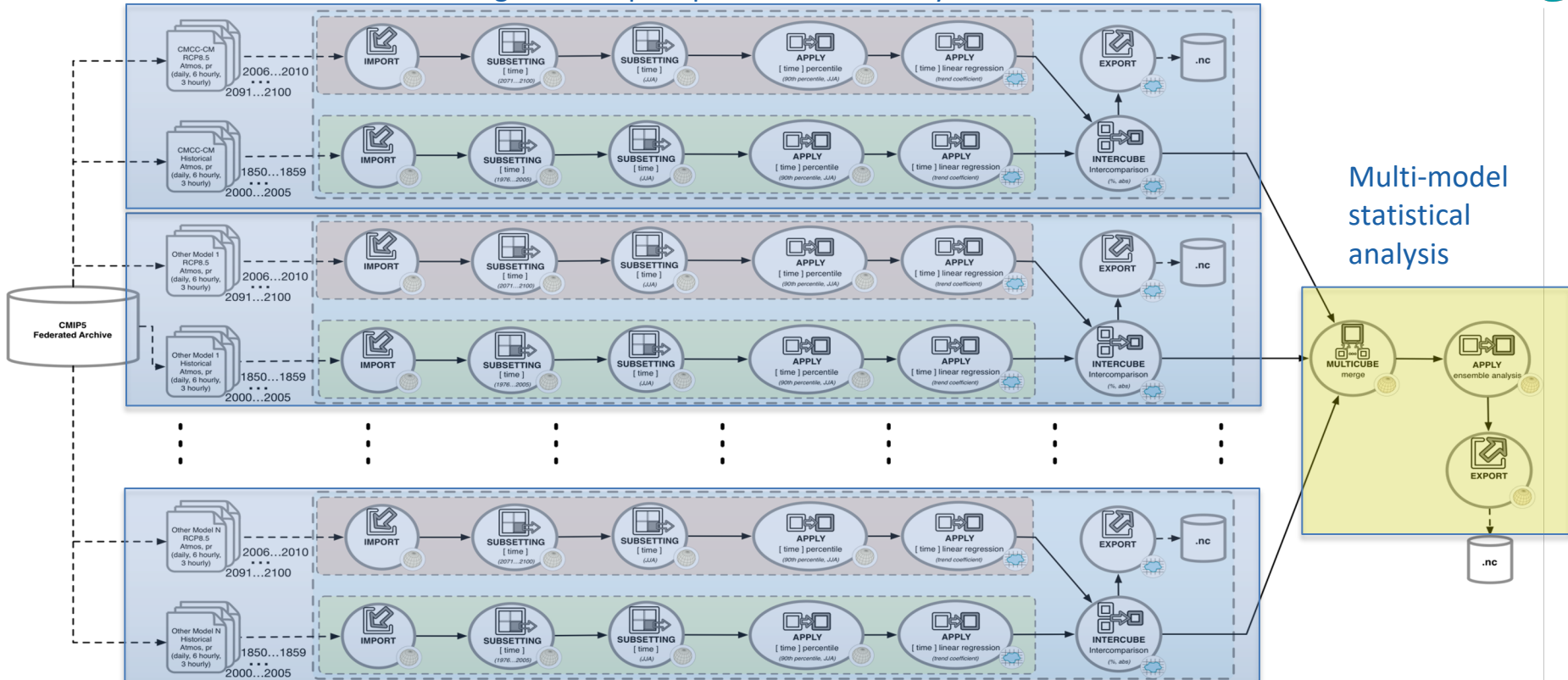


Workflow example III: Multi-model experiment design

Precipitation Trend Analysis use case implemented as an Ophidia workflow

Single model precipitation trend analysis

Multi-model statistical analysis

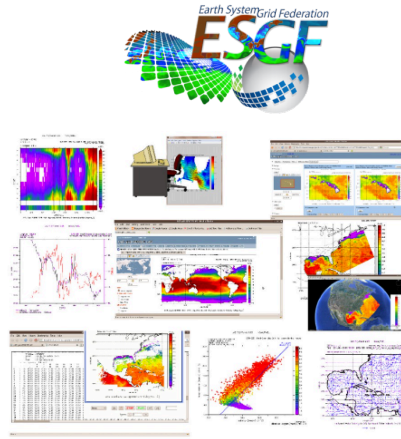


S. Fiore, et al., "Distributed and cloud-based multi-model analytics experiments on large volumes of climate change data in the earth system grid federation eco-system". In *Big Data (Big Data)*, 2016 IEEE Int. Conference on. IEEE, 2016. pp. 2911-2918



Multi-model experiment input data

ESGF¹ is a coordinated multiagency, international collaboration of institutions that continually develop, deploy, and maintain software needed to facilitate and empower the study of climate.

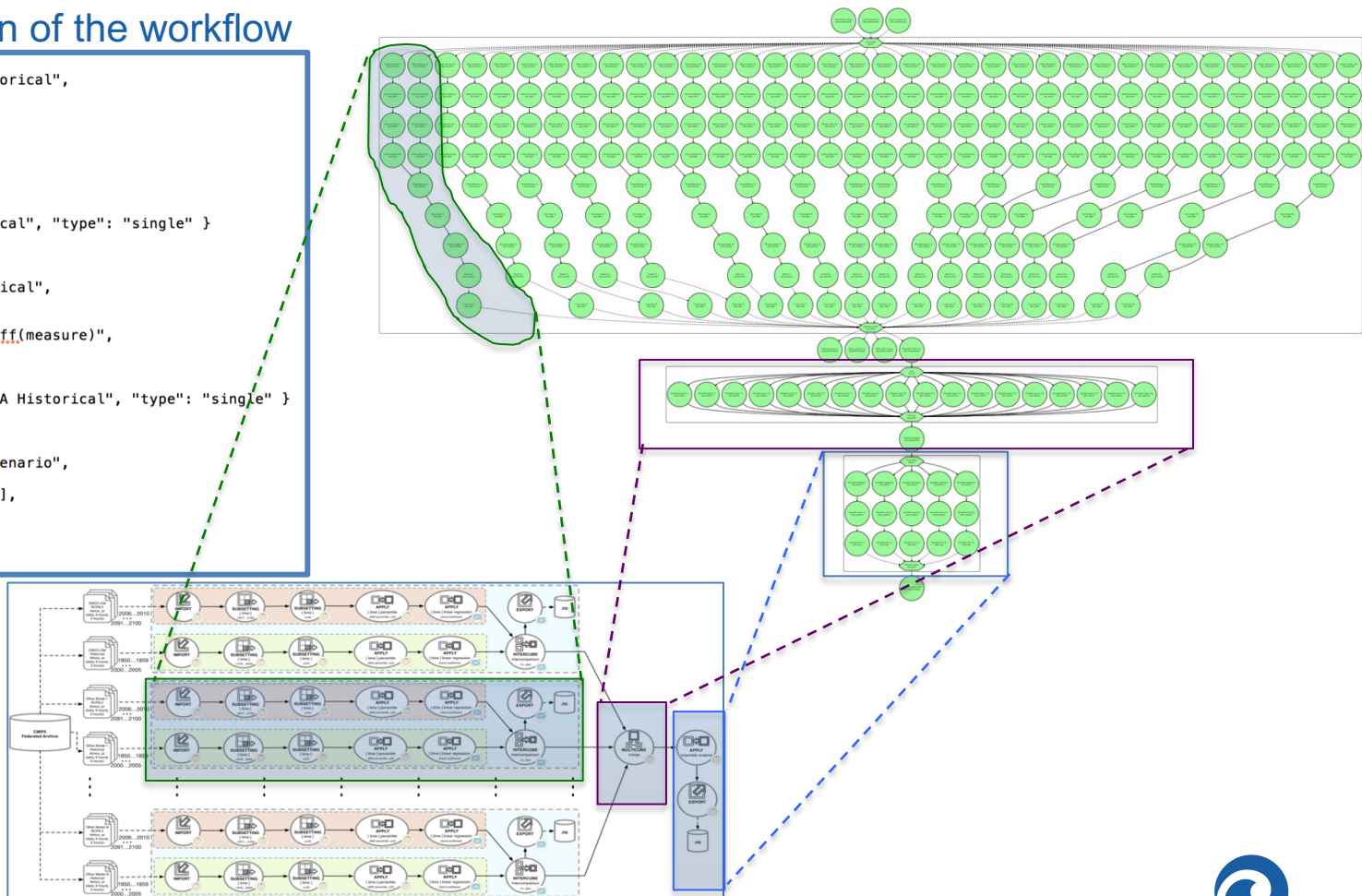


Model acronym	Model expansion	Institute
CCSM4	Community Climate System Model, v4	National Center for Atmospheric Research (NCAR)
CMCC-CESM	CMCC - Community Earth System Model	Euro-Mediterranean Center on Climate Change (CMCC)
CMCC-CMS	CMCC - Coupled Modeling System	Euro-Mediterranean Center on Climate Change (CMCC)
CMCC-CM	CMCC - Climate Model	Euro-Mediterranean Center on Climate Change (CMCC)
CNRM-CM5	CNRM - Coupled Global Climate Model, v5	Centre National de Recherches Météorologiques (CNRM)/Centre Européen de Recherche et de Formation Avancée en Calcul Scientifique (CERFACS)
CSIRO Mk3.6.0	CSIRO Mark, v3.6.0	Commonwealth Scientific and Industrial Research Organisation (CSIRO) in collaboration with Queensland Climate-Change Centre of Excellence (QCCCE)
CanESM2	Second Generation Canadian Earth System Model	Canadian Centre for Climate Modelling and Analysis (CCCMA)
GFDL-CM3	GFDL Climate Model, v3	National Oceanic and Atmospheric Administration (NOAA)/Geophysical Fluid Dynamics Laboratory (GFDL)
GFDL-ESM2G	GFDL Earth System Model with Generalized Ocean Layer Dynamics (GOLD) component	National Oceanic and Atmospheric Administration (NOAA)/Geophysical Fluid Dynamics Laboratory (GFDL)
GFDL-ESM2M	GFDL Earth System Model with Modular Ocean Model 4 (MOM4) component	National Oceanic and Atmospheric Administration (NOAA)/Geophysical Fluid Dynamics Laboratory (GFDL)
HadGEM2-CC	Hadley Centre Global Environment Model, v2 (Carbon Cycle)	Met Office (UKMO) Hadley Centre (HC)
HadGEM2-ES	Hadley Centre Global Environment Model, v2 (Earth System)	Met Office (UKMO) Hadley Centre (HC)
INM-CM4.0	INM Coupled Model, v4.0	Institute of Numerical Mathematics (INM)
IPSL-CM5A-MR	IPSL Coupled Model, version 5, coupled with NEMO, mid resolution	L'Institut Pierre-Simon Laplace (IPSL)
MIROC5	Model for Interdisciplinary Research on Climate, v5	Atmosphere and Ocean Research Institute (The University of Tokyo), National Institute for Environmental Studies, and Japan Agency for Marine-Earth Science and Technology
MPI-ESM-MR	MPI Earth System Model, medium resolution	Max Planck Institute for Meteorology (MPI-M)
MRI-CGCM3	MRI Coupled Atmosphere - Ocean General Circulation Model, v3	Meteorological Research Institute (MRI)
NorESM1-M	Norwegian Earth System Model, v1 (intermediate resolution)	Norwegian Climate Centre (NCC)

Multi-model experiment implementation & execution

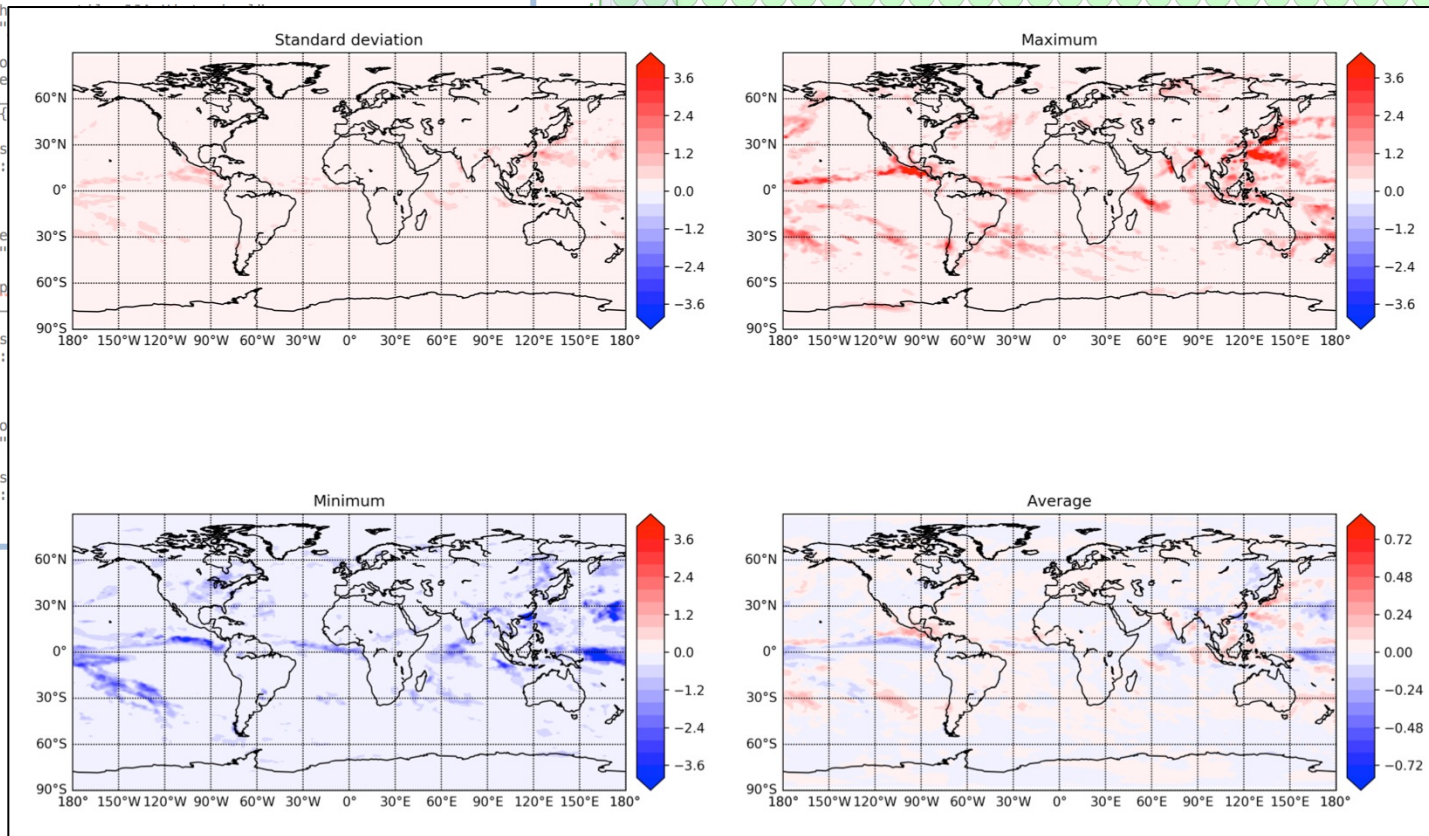
JSON implementation of the workflow

```
{
  {
    "name": "90th percentile JJA Historical",
    "operator": "oph_reduce2",
    "arguments": [
      "operation=quantile",
      "dim=time",
      "concept_level=y",
      "order=${5}"
    ],
    "dependencies": [
      { "task": "Subset JJA Historical", "type": "single" }
    ]
  },
  {
    "name": "Linear regression Historical",
    "operator": "oph_apply",
    "arguments": [
      "query=oph_gsl_fit_linear_coeff(measure)",
      "measure_type=auto"
    ],
    "dependencies": [
      { "task": "90th percentile JJA Historical", "type": "single" }
    ]
  },
  {
    "name": "Import Type Selection Scenario",
    "operator": "oph_if",
    "arguments": [ "condition=${10}" ],
    "dependencies": [
      { "task": "loop_model" }
    ]
  }
},
}
```



JSON implementation of the workflow

```
{
  "name": "get",
  "operator": ":",
  "arguments": {
    "operation": "dim=times",
    "concept": "order=${",
  },
  "dependencies": {
    "task": "line"
  }
},
{
  "name": "Line",
  "operator": ":",
  "arguments": {
    "query=op",
    "measure_": "order=${",
  },
  "dependencies": {
    "task": "get"
  }
},
{
  "name": "Impo",
  "operator": ":",
  "arguments": {
    "dependencies": {
      "task": "line"
    }
  },
}
```



Two approaches for the implementation

The screenshot shows a Jupyter Notebook with two code cells. The first cell, labeled 'In [35]:', defines two functions: `merge_results` and `final_reduce`. `merge_results` takes a list of single model cubes and a client, joins them into a single cube, and returns its PID. `final_reduce` takes a cube PID, an operation, and a client, performs a reduce2 operation with 1 core and 4 threads, and exports the result. The second cell, labeled 'In [34]:', shows a loop over a list of models. For each model, it computes trend analysis on historical and RCP datasets, compares them, and appends the results to a list. Finally, it merges the results and computes statistics (avg, max, min, var, std) for the merged data. The notebook interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help), a toolbar with icons for saving, running, and other actions, and a status bar indicating 'Notebook saved', 'Trusted', and 'Python 2'.

```
In [35]: def merge_results(single_model_cubes, my_client):
cube.Cube.client = my_client
cubesList = '|'.join(single_model_cubes)
merged_cube = cube.Cube.mergecubes2(cubes=cubesList, dim='new_dim', description='ensemble_r
return merged_cube.pid

def final_reduce(in_cube_pid, operation, my_client):
cube.Cube.client = my_client
cube_input = cube.Cube(pid = in_cube_pid)
cube_output = cube_input.reduce2(operation=operation, dim='new_dim', ncores=1, nthreads=4)
session_code = cube.Cube.client.session.split('/')
cube_output.exportnc2(force='yes', output_name=operation, output_path='/INDIGO/precip_trend
return cube_output

In [34]: %%time
single_model_cubes = []
#For each model
for l in list of models:
#Compute trend analysis on historical dataset
hist_cube_pid = delayed(historical_scenario_function)(l, 'historical', myClient)
#Compute trend analysis on RCP dataset
scenario_cube_pid = delayed(historical_scenario_function)(l, 'scenario', myClient)
#Compare trend from historical and RCP
model_cube = delayed(single_model_calculation)(l, hist_cube_pid, scenario_cube_pid, myClient)
single_model_cubes.append(model_cube)

merged_cube_pid = delayed(merge_results)(single_model_cubes, myClient)

stats = ['avg', 'max', 'min', 'var', 'std']
stat_cubes = []
for s in stats:
stat_cube = delayed(final_reduce)(merged_cube_pid, s, myClient)
stat_cubes.append(stat_cube)

final_result = compute(*stat_cubes)
```

Single
model
analysis

Multi-model
statistical
analysis

The screenshot shows a JSON configuration file for a workflow. It defines three tasks: '90th percentile JJA Historical', 'Linear regression Historical', and 'Import Type Selection Scenario'. Each task specifies its name, operator, arguments, and dependencies. The '90th percentile JJA Historical' task depends on 'Subset JJA Historical'. The 'Linear regression Historical' task depends on '90th percentile JJA Historical'. The 'Import Type Selection Scenario' task depends on 'loop_model'.

```
{
  "name": "90th percentile JJA Historical",
  "operator": "oph_reduce2",
  "arguments": [
    "operation=quantile",
    "dim=time",
    "concept_level=y",
    "order=${5}"
  ],
  "dependencies": [
    { "task": "Subset JJA Historical", "type": "single" }
  ]
},
{
  "name": "Linear regression Historical",
  "operator": "oph_apply",
  "arguments": [
    "query=oph_gsl_fit_linear_coeff(measure)",
    "measure_type=auto"
  ],
  "dependencies": [
    { "task": "90th percentile JJA Historical", "type": "single" }
  ]
},
{
  "name": "Import Type Selection Scenario",
  "operator": "oph_if",
  "arguments": [
    "condition=${10}"
  ],
  "dependencies": [
    { "task": "loop_model" }
  ]
},
{
  "name": "loop_model",
  "operator": "oph_loop",
  "arguments": [
    "loop_model"
  ],
  "dependencies": [
    { "task": "loop_model" }
  ]
}
```

	Approach	Mode	Library	Code	ExecTime
Workflow	SS - SI*	Batch	Ophida WF	JSON	~170s (1.35x)
Notebook	SS - MI*	Interactive	PyOphidia	Python	~230s

* SS: Server Side; SI: Single Interaction, MI: Multiple Interactions



What have we learned so far?

Complex climate data analysis requires workflow support

*The **Ophidia HPDA framework** provides **workflow management features**:*

- *Target large-scale analysis and parallel execution of tasks*
- *Support for different constructs and workflow resiliency*
- *Integrated job orchestration, management and monitoring features*

Real case studies can be effectively modeled as (complex) workflows composed of hundreds of tasks

Next: Demo and hands-on of Ophidia workflows



References and further readings

- Asch, M., et al. (2018). *Big data and extreme-scale computing: Pathways to convergence-toward a shaping strategy for a future software and data ecosystem for scientific inquiry*. *Int. J. High Perform. Comput. Appl.*, 32(4), 435-479.
- E. Deelman, et al. (2018) 'The future of scientific workflows', *The International Journal of High Performance Computing Applications*, 32(1), pp. 159–175.
- S. Fiore, et al. (2013). *Ophidia: Toward Big Data Analytics for eScience*. *ICCS 2013*, volume 18 of *Procedia Computer Science*, pp. 2376-2385.
- S. Fiore, et al. (2014). "Ophidia: A Full Software Stack for Scientific Data Analytics", *proc. of the 2014 Int. Conference on High Performance Computing & Simulation (HPCS 2014)*, pp. 343-350.
- S. Fiore, D. Elia, C. Palazzo, F. Antonio, A. D'Anca, I. Foster and G. Aloisio (2019), "Towards High Performance Data Analytics for Climate Change", *ISC High Performance 2019. Lecture Notes in Computer Science*, vol. 11887, pp. 240-257.
- D. Elia, S. Fiore and G. Aloisio, "Towards HPC and Big Data Analytics Convergence: Design and Experimental Evaluation of a HPDA Framework for eScience at Scale," in *IEEE Access*, vol. 9, pp. 73307-73326, 2021
- D. Elia, et al. (2016). "An in-memory based framework for scientific data analytics". In *Proc. of the ACM Int. Conference on Computing Frontiers (CF '16)*, pp. 424-429.
- C. Palazzo, et al. (2015), "A Workflow-Enabled Big Data Analytics Software Stack for eScience", *HPCS 2015*, pp. 545-552
- A. D'Anca, et al. (2017), "On the Use of In-memory Analytics Workflows to Compute eScience Indicators from Large Climate Datasets," *2017 17th IEEE/ACM Int. Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, pp. 1035-1043.
- S. Fiore, et al. (2016). "Distributed and cloud-based multi-model analytics experiments on large volumes of climate change data in the earth system grid federation eco-system". In *Big Data (Big Data)*, 2016 *IEEE Int. Conference on*. IEEE. pp. 2911-2918.



Acknowledgements

ESiWACE2 has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 823988



eFlows4HPC this project has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 955558. The JU receives support from the European Union's Horizon 2020 research and innovation programme and Spain, Germany, France, Italy, Poland, Switzerland, Norway



IS-ENES3 has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 824084



Thank you!

Questions?

More about Ophidia?

Ophidia website: <http://ophidia.cmcc.it>

GitHub repo: <https://github.com/OphidiaBigData>

Contact: *ophidia-info [at] cmcc.it*

Twitter channel: <https://twitter.com/OphidiaBigData>

