



#### OVERVIEW

TAO (stands for Tool Augmentation by user enhancements and Orchestration) is a lightweight, generic integration and distributed orchestration framework. It allows to reuse (i.e. integrate) commonly used toolboxes (such as, but not limited to, SNAP, Orfeo Toolbox, GDAL, PolSARPro, etc.). This framework allows for processing composition and distribution in such a way that end users could define by themselves processing workflows and easily integrate additional processing modules (by processing module it is understood either a standalone executable or a script).

In terms of use, the TAO platform provides a mean for orchestration of heterogeneous processing components and libraries in order to process scientific data. This is achieved in following steps:

- Preparation of resources (including processing components) and data input,
- Definition of a workflow as a processing chain,
- Execution of workflows,
- Retrieval / visualization of the results.

To have a simple view of the TAO platform, the platform model is split among four main macrocomponents. Such a macro-component is a logical collection of components with related functions. It has no direct relationship to the software implementation.









The TAO platform was designed and built to be a multi-user platform. Therefore, one of its basic features relates to the user management.

User can authenticate to the platform by several configurable mechanisms (local user base, LDAP or even single sign-on). User accounts can be linked to groups with appropriate system rights.

Since many years, EO Data and EO Data processing as well,

140   Dashbeard	× +	– n ×
€⇒ເພ	0 192.168.61.101/tao web u/sap.html#my/worldowstedit=1jpage=0	□ ··· ♥☆ II\ Œ Ξ
ΤΑΟ		🕫 📽 🖉 💩 Hax Mustermann 😽
Max Mustermann • Ontre	Users management Full list of users.	🚯 Home > Usos
Search Q	All users Active users Disabled Pending	٥
😆 Dashboard	+	≡ ⊞
My WORKSPINCE C	User list 23 from 23	Order by:
SHARED WORKSPACE	Anton Fan (anton.pan@c.s.ro)	Onese
% RESOURCES ✓	registration date: Wednesday 26, 2017 - 14:04:21 Jast Joan date: Wednesday 20, 2017 - 14:04:21	
Datasources	number of workflows: xxxx	
Processing components	phone: 0764666666666	
Users	Queens	👁 Vera 🥒 Edit 🔢 Disable 👔 Delete
EDCUMENTATION		
7/0 Documentation	Corina Vasile (cv@c.s.ro)	Onver
O OLA	registration date: Thursday 20, 2017 - 16:09:043	
User quota: unknown 🔽	last login date: Thursday 20, 2017 - 16:09:48 number of workflows: xxxxx email address: cx/pr-s.ro	

use increasingly disk space and computing power. Even if technical infrastructures evolve likewise, computing and storage resources have a not negligible provisioning and usage cost. Therefore, we have to limit resource usage in order to offer a fair operation of the platform.

The platform, through its Authorization module, defines **quota management** by user. The quota consists in the pool of resources that may be allocated for a single user, in terms of: number of CPUs, storage space and, possibly, RAM memory amount for the execution of user workflows.

A dedicated activity monitoring service captures the required information and collects it in a database. This database with tracked processing activity and disk space consumption is used then to automatically restrict the access to resources according to the user quota.

# DATA ACCESS

In general, data access refers to activities related to storing, retrieving or acting on data housed in a repository.

For TAO, the Data Access feature deals with the provisioning of data (either EO products or ancillary data) that exhibits certain criteria for workflows execution.

Two categories of data sources can be distinguished: local and remote.

#### LOCAL DATA ACCESS

A local **data product** is nothing but a descriptor for one or more files (for example, a GeoTIFF data product can be just a single raster file, while a Sentinel-2 data product is a structured collection of raster and metadata files). Consequently, a local data source is responsible with providing local data to the TAO platform to the requesting processing component(s).

Even if the local data source (comprising in the local database and file system) is not an external interface, it is important to layout its organization before describing the external interfaces of the system.



As it can be seen in the figure beside, the local data source consists of a product database, in which metadata associated with concrete data products are stored, and a local (in the platform context) file system, in which the data products are physically stored.

The product database stores basic metadata about the EO products, such as acquisition date, geographical footprint, product type, etc.

This metadata allows the users to query the local data source for products satisfying certain criteria, and is created when either products are initially imported or new products are downloaded from remote data sources.

The file system is organized such that an easier

distinction can be made between public products (products that are visible/usable by all users) and private products (products visible/usable only by a specific user). This separation further allows the implementation of user quota management.

The visibility of the products is implemented at a logical level (i.e. not by physical operating system rights) in the database.

The file system structure depicted above is visible to (i.e. shared with) all the processing nodes in the TAO cluster. This is necessary in order to allow a uniform way of accessing products by processing components from remote nodes.

#### EXTERNAL DATA ACCESS

External data sources represent external repositories that are not under the direct control of the TAO platform (an example of such an external data source is the Copernicus Scientific Data Hub). They may implement different access interfaces and authentication/authorization schemes. This is transparent for the user of the TAO platform.

Eventually, all needed data would become local, taking also into account the defined user quotas.

Contrary to the differences that may exist between repositories access and the data formats that may be provided, several common features may be derived:

- All data can be expressed in binary form;
- There are common operations (actions) that can be performed on any data source:
  - Authenticate/authorize the user connection to the data source;
  - Query the data source for data satisfying several criteria;
  - Retrieve the data in a binary form that can be interpreted by TAO.









The framework is capable to abstract the data format and to expose them in a unitary way. In addition, it makes the user (or a requesting component) unaware of the original location of data (i.e. a remote repository). This is accomplished by an interface abstraction with querying capabilities.



This approach is illustrated in the following figure:

The two types of data sources (local and remote) share the same interface, the location of the data being transparent for the TAO API. It is the implementation of the data source that takes care of properly connecting to the repository and querying and retrieving the data products. The platform can thus seamlessly access data from local repositories, such as the Sentinel-1 and Sentinel-2 repositories found on the DIAS platforms (CreoDIAS, Mundi and Onda), but also from remote repositories, such as the USGS Landsat repository or Alaska Satellite Facility.

The data source interface exposes operations (i.e. methods) allowing to:

- Connect to the remote repository;
- Authenticate to the remote repository;
- Create a query to be executed against the repository.

Given the diversity of EO products (and product providers), the parameters of a query may be bound to a specific repository. Nevertheless, there is a small subset of parameters that are supported by different providers, namely:

- The product name or identifier;
- The acquisition date (and time);
- The product footprint.





esa



(€) ⇒ œ @	3 # 192.168.61.101/tao web ulibaa/htmi#sharet/datascuree					10 合	₩.CD Ξ
T A O							👵 Hechaterran 🗠 ⊄
Hackletomann s beles	Datasources Filt list of dataseaseas	Datasource		×			🚓 Harrie II. Datasa ras
panda, Q	Al dataseuras - Baiblici dataseuras	Sentinel1					0
Eastbard	τQ	Scientific Data Hub					
B BYROEKSANCE K	ROPTION REATORER X		Tan natio 🤤				
SHORED INDRESINCE	Blatassurors 9 to 16 of 17	Terrs Terrs	Tres	Palast votes			Oraba by: (912)(913)-
montows	SENTINEL-3	priorisationWode	Swing		8		8
components todasources		to organize	Palipatita				
Summer e	•	beginPosition	Cate		-	•	
DOCUMENTS) DI		endPecifian	Cate				
# No occumentation	SENTINEL1	strationalitate	string		8		89
0 QM		platformkarre	song	Sentinel 1			
User çuntat unikniwi 🛛 🔤		product/spc	song	SLC		_	
		Notice					
		Paramatas willing pa definition stage.	gulated with specific valu	is only in the workflow			
				Ctte			
	Copyright > 2013 . This largest entities by lists indexected and $\mathbf{O}$	wheshvilian (TEO)				Page times	ANNJO 2018 OD 18728/48/06/5412

Currently, several data source plugins are available in the TAO platform, namely:

Plugin for	Supported Sensors	Remarks
Scientific Data Hub	Sentinel-1, Sentinel-2, Sentinel-3	Supports also local archives
Amazon Web Services	Sentinel-2, Landsat-8	Supports also local archives
PEPS	Sentinel-1, Sentinel-2	
USGS	Landsat-7, Landsat-8	Supports also local archives
FedEO	ALOS, ALOS-2, CryoSAT, MetOP, DEIMOS-1, ERS-1, ERS-2, ENVISAT, FORMOSAT-2, GeoEYE-1, GOES, IKONOS, IRS-1D, IRS-P5, IRS-P6, JASON, KANOPUS-V1, KOMPSAT- 2, METEOR-3M, OCEANSAT, OCEANSAT-2, Pleiades, Proba-V, QuickBird, RadarSAT-1, RapidEYE, SPOT 1-5, SeaSAT, TerraSAR-X, WorldView-1, WorldView-2,	Collections limited to ESA archive. Not all products may be retrieved
Alaska Satellite Facility	Sentinel-1, ALOS	
CreoDIAS	Sentinel-1, Sentinel-2, Landsat-8	Supports also local archives
Mundi DIAS	Sentinel-1, Sentinel-2, Landsat-8	Supports also local archives









All data sources can be intuitively queried by means of the TAO web interface, which dynamically adapts to the parameters of the respective data source. The query results can be then selectively downloaded, and used as input for further processing.

# PROCESSING COMPONENTS AND MODULES

A **processing component** (on short called component) represents a standalone **application** (or module) defined by the following parameters:

- Input description: type of data that the component accept as input source (e.g. image, raster maps, vector maps, sensors, etc.)
- Processing operation with execution parameters: the operation that the component executes with the list of accepted parameters.
- Output description: type of data provided as processing operation result. Can consist in one or more files.

A component is viewed as a generic execution resource and should permit, through its parameters list, the possibility to define *input location* (the place where it expects to retrieve input data for consuming) and *output location* (the place where the processing results are persisted).



Each processing component is described in terms of

its expected input, processing module (see the below paragraphs for details about what a processing module is) and output.

	0					ROMA	S.	La force de l'innovation	BROCKMANN CONSULT	(Cees
TAO   Dashboard X	+									– ø ×
-) → ଙ ŵ	I 22.168.61.101/tao-web-ui/sap LABEL LIKE: CALIBRATION ×	.html#my/components							··· 🛡 🏠	
My WORKSPACE Y	Processing components 1 to	E dia manda					×	_		Order by:
	Processing components - 1 to	Edit processing con	nponent				^			TAGE
	OPTICALCALIBRATI	Please check all the tabs	before saving a new component.							<b>%</b>
	Version: 6.6.0	General desc. Co	nfiguration Parameters Sys. vars.	Sources & Targets						
	Authors: OTB Team Copyright: (C) CNES Apache License	Parameter list:								
	Perform optical calibration TOA/TO	Parameter ID	Description	Parameter Label	Parameter Type	Data Type	Default value	pplication also allows providin		
		SourceBands	sourceBands	PsourceBands	Regular	String	default value		≠ Edit (8 Di	lete
	SARCALIBRATION SARCalibration	🌶 🗃 auxFile	auxFile	PauxFile	Regular	String ~	default value			<b>#</b>
	Version: 6.6.0 Authors: OTB Team Copyright: (C) CNES Apache License Perform radiometric calibration of	🖋 🖹 externalAuxF	externalÄuxFile	PexternalAuxFile	Regular 🗸	String ~	default value	i input.		
		🖋 🖹 outputimage	outputimageInComplex	Poutputimagein	Regular 🕑	Boole 🗸	false		🖌 Edit 🔒 De	iete
	CALIBRATION snap-calibration		outputimageScaleInDb	PoutputimageSc	Regular 🗸	Boole ~	default value		(20AT) (20AT)	(TAG6)-
	Version: 1.0 Authors: SNAP Team Copyright: (C) SNAP Team Calibration operator	🖋 🔒 createGamm	createGammaBand	PcreateGammai	Regular 🗸	Boole ~	false			
		🖋 🗃 createBetaBa	createBetaBand	PcreateBetaBan	Regular 🗸	Boole ~	false		<b>∕</b> Edit 🗐 G	dete
		🖋 🗃 selectedPola	selectedPolarisations	PselectedPolaria	Regular 💟	String ~	[default value]			
	Copyright © 2018 - Tool Augmentation by		outputSigmaBand	PoutputSigmaB:	Regular 🗸	Boole 🗸	false			8708:45:40.415Z

There are three kinds of components that could be used in a workflow:

- **System components**: these are the components that ship with the system and can be used by everyone. Currently, TAO ships with all the components of Orfeo Toolbox 6.4 and of SNAP 6.0.0.
- User components: a user of the platform also has the possibility to upload his own processing algorithms as script in Python or R. These user components can be used only by the owner (private mode) or by everyone (public mode - in which case it will be part of a contributor list). Users who upload components into the system are responsible for the management of their components.
- **Contributor components list**: all user components loaded into the system and declared for public use will be a part of this list.

The person who upload a component into the system is responsible for providing also all the necessary dependencies. The component is saved into a repository either as an archive or as an installation kit. The repository can be a tree folder into the file system or a database.

The processing modules are heterogeneous regarding their implementation language or operating system. The TAO framework does not target to integrate at once the processing modules from several operating systems (i.e. multiple modules deployed on different operating systems at once). Nevertheless, TAO aims to be as OS-independent as possible and to allow the integration of modules written in different programming languages (such as C/C++, Java or Python).

In an ideal scenario, a processing module may be an independent executable (i.e. not having any additional runtime dependencies). However, this is seldom the case and different modules have different runtime dependencies. Therefore, it is critical that each module has its own appropriate dependencies when the module would execute. Moreover, dependencies of one module may break the execution of another module if they are executed on the same machine, in the same memory space.

attillus





) esa

With these constraints in mind, the framework is capable of providing to each module its appropriate runtime dependencies, without disrupting the functioning of other modules, by using Docker execution containers. Each container hosts a processing module together with its runtime dependencies. Each container is then isolated from other containers that may execute on the same machine.

Out of the box, TAO ships with predefined Docker images for OTB 6.4.0, SNAP 6.0.0, Python (with most common python scientific libraries) and R. Of course, additional Docker images can be registered if available.

# WORKFLOWS

By **workflow** it is understood a sequence of processing operations performed on a given input (EO data and/or ancillary data), having at least one output. A workflow is described by a formal flow diagramming technique, with directed flows between components.



01 Completion

The processing operations in a workflow are handled by processing components.

The Workflow Management feature of TAO copes with all operations necessary to define such workflows and parameters of their components.

After defining / describing processing components, these can be glued together to form a workflow. A workflow definition

then comprises of the set of constituent modules and the rules that link them.

For a better user experience, a workflow is graphically / visually created by dragging and dropping processing components (boxes) on a drawing canvas, and then connecting them with directed (arrowed) lines.



Then, for each component, the parameters of the component can be modified.

A workflow can be *created* from scratch, can be *edited* or *deleted*, and/or can be created (*cloned*) from an existing workflow.





# COMPUTING RESOURCES AND EXECUTION



The TAO platform can scale up from one to as many nodes as made available. A node represents a computer resource (a machine) with a defined amount of processing power (number of processors, amount of RAM and disk space), a container for execution of processing components. The platform provides a dedicated administrative interface to manage the processing nodes (even via remote shells).

Currently, the framework deals with already existing nodes (Linux-based with default installation – i.e. no particular configuration) that can be registered with it. There is work in progress of integrating the OpenStack Nova API (used by several cloud providers) to allow the dynamic creation and registration of computing machines/resources). Once done, this will allow the dynamic expansion of resources for achieving computation scalability as needed, without any operator intervention.

The actual execution of workflows is handled by dedicated software, usually named Cluster Resource Managers. In order to decouple the core of the TAO framework from the actual Cluster Resource Manager software, the latter component has to provide a DRMAA-compliant interface.

The DRMAA (Distributed Resource Management Application API) provides a standardized access to the DRM systems for execution resources. It is focused on job submission, job control, reservation management, and retrieval of jobs and machine monitoring information.

Currently, there are several DRMAA implementations, most of them implemented in the C programming language and for Linux operating systems.



The bridging plugin approach taken by the TAO framework allows to easily swap such implementations, making the framework loosely coupled from the actual DRM system used.

The platform provides two such DRMAA-compliant plugins, namely for SLURM (Simple LinUx Resource Manager) and Torque (the open-source version of PBS/Torque).

#### USER WORKSPACES

The results of the workflow executions (and not only) can be found in the user workspace. This is a system view that allows a user to:





ΤΑΟ						🖓 🜔 Cara Cosmin 🛛 🕫
Cara Cosmin • Altern	Shared files Pull list of public files.					B Home > NyTiles
Search Q	✓ Se public			88-90-output_s	nap-ndviop.tif	×
n Dashboard	Standard John Markey North 10125 Note 100 - 100 Note 1			Name	Value	^
S My WORKSTRICE (	> 05, rop.into > 05, HTML			acquisitionDate visibility	2018-10-08 14:58:29:548 PUBLIC	
SHARED WORKSPACE <	> , GRANULE			cm	EPSG/WGS 84 / UTM zone 30N	
Data sources				sensor Type	UNKNOWN	
Components	MTD. MSILICXII	ami	43.87 KB	width	10980	
WerkRows	B INSPIRANT	ami	18.43 KB	geometry	POLYDON ([-0.5337109 43.328257, -0.5727045 42.3024536, 0.8181237 43.2890818, 0.533710	42.3383721, 0.7578584 (9.43.326257))
Common files		.5370	52.07 KB	id .	88.90 output_snap-ndvlop	
% RESOURCES <	<ul> <li>         AL_MOBILE_CONTROLTIONSSI_MODE_NON_INTERNITIONSSI_MAPE     </li> <li>         SIA_MOBILE_CONTROLTIONSSI_MODE_NON_INTERNITIONSSI_MAPE     </li> <li>         SIA_MOBILE_CONTROLTIONSSI_MODE_NON_INTERNITIONSSI_MAPE     </li> <li>         SIA_MOBILE_CONTROLTIONSSI_MODE_NON_INTERNITIONSSI_MAPE     </li> </ul>			pixelType	FLOAT32	
too weather				formatType	RASTER	
	🗸 💺 admin			productType	GeoTIFF	<ul> <li>Toggle buildkivlew</li> </ul>
TAO Documentation	🛩 🛼 83-30-output_snap-ndviep			height	10080	
O How To	<ul> <li>Bis00-output_snap-advlop.tf</li> <li>See files</li> </ul>	BL.	919.81 MB			
User quotec 100000 🕓						
ung occ Office, I follows				Actions Download	Described size: 202.81 MB	u ▲ Download He
ana Tanti masta (1972 julium bimitt	Copyright © 2018 Teol Augmentation by User Enhancements and Orchestration (TAO)				Page ti	mostamp: 2018-20-00710:30:37.074Z

• View the details about a local data product (accessible for both administrator and user roles, with the difference that the user can only view details about his local data products as well as public data products),

esa

• View all local data products (accessible for both administrator and user roles, with the difference that the user can only view his local data products as well as public data products),

- Upload additional files that can be further used in workflows (such as model files, shape files, etc.),
- Perform a check-style selection from the existing products that can be used as input for his/her workflows

# EXTERNAL PLATFORM ACCESS

A closed platform is of limited use and impales its adoption by the community. Therefore, the TAO platform allows for the interoperation with external processing systems, and also for a standardized access by external clients. To accomplish this, it relies on several OGC WPS – compliant components, each having a specific role:

- WPS processing component: this is the TAO component that can invoke external processors via a WPS interface;
- WPS server: this is the TAO component that allows external clients to invoke TAO workflows. The other TAO processing components are not exposed via WPS interfaces. Instead, only public workflows are exposed as WPS capabilities of this WPS interface.







esa



The platform comes with a minimal, but sufficient monitoring dashboard, which allows, at any time, to see what is executing in the system and where it is executing the system resource usage (CPU, memory, storage space).

Besides platform information, users and administrators are visually notified, in real-time,

ΤΑΟ			Ko <sup>1</sup> 😒 🜔 Caro Coortin
ra Cosnin 1949	Dashboard Overview of TAG schildy		🗱 linne - Dub
	Notifications History	commit page 1 + +	🕍 Master node monitor
	Job [205] for workflow [Local DB, SMAP NEW, SNAP RM, SNAP SAPLand OTE Becample (or group)] completed in 76a	Construction and the	CPU usage 🔲 16%
	DONE	Contractor Contractor	Memory Ph (NU)/12885, and load, megalyte) Storage R76 (NUN4/DEX, and (Nuts, ggdyte)
	RUNNING	A DESCRIPTION OF THE OWNER OWNER OF THE OWNER	
	DONE	Contractor and the second second	
	DONE	CONTRACTOR OF STREET, STRE	
	DONE	CONTRACTOR OF STREET,	
	RUNNING	OCCUPATION AND AND AND AND AND AND AND AND AND AN	
ientation.	RUNNING		••••
	ICHNNINE	Constant Constant Constant	やけらとう りょう りょう きょう りょう りょう すい
	IXINI	Construction Constant	
10G8 💽		sametpage 1 + +	Running jobs orretarge1 + +
		1.1.1	Job name: Local DB, SNAP NDVI and OTB Resample, 09/10/2018 15:35:28
	Executions History	concrepage1 + * D	Workfrow: Local DB, SNAP NDVI and OTB Resample
	Inh name: Test workflow 9		ador: vincont, status: RUNNING CONSISTENCESSES - Coxe Task sumary:
	Workflow: Local DB, SNAP NDVI and OTB Resample		✓ snap-ndvlop (2018-10-00700:00-2018-10-00716:35:30.662)
	user: admin, status: FALED Case-Co-contractor see 0 Onais-as-contractor-sec.or		<ul> <li>Right Transformillesample [2018-10-00708/02/02/2018-10-00718/30150.039]</li> <li>Sentimel2-Local Database [2018-10-00118/15/04.541-2018-10-00118/35/85.005], on host sen2agri-</li> </ul>
	Task sumary:		hung
	SentineD-Local Batabase [2010-10-09107:53:38.216-2018-10-09107:53:38.238], on host sen2agri-prod		i and in the second sec
	@ RigidTransformResample [#/a (\0]		Job name: Test workflow 9
			man where the second se
	JOD name: Local DB, SNAP NDVI and OTB Resample, 09/10/2018 15:34:50		Task surrory:
	Workdow Local DB, SMAP NOW and OTB REsumption		@ smap-mbylep (2010-10-0010000200-e/s)

via the web interface, about events occurring in the system.

Currently, the monitoring dashboard presents the following indicators:

Name	Purpose	Туре
Status of an execution node	Indicates the status of a particular execution node	Possible statuses: online, offline, not configured.
Workflow status	Indicates the status of the execution of a particular workflow	Possible statuses: not run, running, paused, completed, error.
Task status	Indicates the status of the execution of a particular workflow task	Possible statuses: not run, running, completed, error.
Quota status	Indicates the status of the user quota	Percentage and value of used disk space.