

---

**vantage6**

**A. van Gestel, B. van Beusekom, D. Smits, F. Martin, J. van Soest,**

**Jun 20, 2023**



# CONTENTS

<b>1</b>	<b>What is vantage6?</b>	<b>1</b>
<b>2</b>	<b>Overview of this documentation</b>	<b>3</b>
<b>3</b>	<b>Vantage6 resources</b>	<b>5</b>
<b>4</b>	<b>Index</b>	<b>7</b>
4.1	Concepts . . . . .	7
4.1.1	Architecture . . . . .	7
4.1.2	Components . . . . .	8
4.1.2.1	End to end encryption . . . . .	9
4.2	Installation . . . . .	10
4.2.1	Requirements . . . . .	10
4.2.1.1	Node & Server . . . . .	10
4.2.2	Client . . . . .	11
4.2.2.1	User interface . . . . .	11
4.2.2.2	Python client library . . . . .	12
4.2.2.3	R client library . . . . .	12
4.2.2.4	Server API . . . . .	12
4.2.3	Node . . . . .	12
4.2.4	Server . . . . .	12
4.2.4.1	Local Installation . . . . .	12
4.2.4.2	Cloud Service Provider . . . . .	13
4.2.4.3	Optional components . . . . .	13
4.2.4.4	Deployment . . . . .	18
4.3	How to use . . . . .	20
4.3.1	Client . . . . .	20
4.3.1.1	Introduction . . . . .	20
4.3.1.2	User Interface . . . . .	20
4.3.1.3	Python client . . . . .	22
4.3.1.4	R Client . . . . .	28
4.3.1.5	Server API . . . . .	29
4.3.2	Node . . . . .	29
4.3.2.1	Introduction . . . . .	29
4.3.2.2	Configure . . . . .	30
4.3.2.3	Security . . . . .	34
4.3.2.4	Logging . . . . .	35
4.3.3	Server . . . . .	35
4.3.3.1	Introduction . . . . .	35
4.3.3.2	Configure . . . . .	36

	4.3.3.3	Logging	42
	4.3.3.4	Shell	42
4.4		Algorithm Development	48
	4.4.1	Concepts	48
	4.4.1.1	Input & output	49
	4.4.1.2	Wrappers	50
	4.4.1.3	Mock client	53
	4.4.1.4	Child containers	53
	4.4.1.5	Networking	53
	4.4.1.6	Cross language	54
	4.4.1.7	Package & distribute	55
	4.4.2	Classic Tutorial	57
	4.4.2.1	Mathematical decomposition	57
	4.4.2.2	Federated implementation	57
	4.4.2.3	Vantage6 integration	59
	4.4.2.4	Cross-language serialization	65
4.5		Developer community	66
	4.5.1	Contribute	66
	4.5.1.1	Support questions	66
	4.5.1.2	Reporting issues	66
	4.5.1.3	Security vulnerabilities	67
	4.5.1.4	Community Planning	67
	4.5.1.5	Submitting patches	68
	4.5.2	Documentation	70
	4.5.2.1	How this documentation is created	70
	4.5.2.2	API Documenation with OAS3+	71
	4.5.3	Release	71
	4.5.3.1	Version format	71
	4.5.3.2	Create a release	72
	4.5.3.3	The release pipeline	72
	4.5.3.4	Distribute release	73
	4.5.3.5	User Interface release	73
4.6		Glossary	74
4.7		Release notes	77
	4.7.1	3.6.1	77
	4.7.2	3.5.2	78
	4.7.3	3.5.1	78
	4.7.4	3.5.0	78
	4.7.5	3.4.2	79
	4.7.6	3.4.0 & 3.4.1	79
	4.7.7	3.3.7	80
	4.7.8	3.3.6	80
	4.7.9	3.3.5	80
	4.7.10	3.3.3	80
	4.7.11	3.3.2	81
	4.7.12	3.3.1	81
	4.7.13	3.3.0	81
	4.7.14	3.2.0	82
	4.7.15	3.1.0	83
	4.7.16	3.0.0	83
	4.7.17	2.3.0 - 2.3.4	84
	4.7.18	2.2.0	84
	4.7.19	2.1.2 & 2.1.3	85
	4.7.20	2.1.1	85

4.7.21	2.1.0 . . . . .	85
4.7.22	2.0.0.post1 . . . . .	85
4.7.23	2.0.0 . . . . .	86
4.7.24	1.2.3 . . . . .	86
4.7.25	1.2.2 . . . . .	86
4.7.26	1.2.1 . . . . .	86
4.7.27	1.2.0 . . . . .	87
4.7.28	1.1.0 . . . . .	87
4.7.29	1.0.0 . . . . .	88
4.8	Partners . . . . .	89



## WHAT IS VANTAGE6?

Vantage6 stands for **privacy preserving federated learning infrastructure** for **secure insight exchange**.

The project is inspired by the [Personal Health Train](#) (PHT) concept. In this analogy vantage6 is the *tracks* and *stations*. Compatible algorithms are the *trains*, and computation tasks are the *journey*. Vantage6 is completely open source under the [Apache License](#).

What vantage6 does:

- delivering algorithms to data stations and collecting their results
- managing users, organizations, collaborations, computation tasks and their results
- providing control (security) at the data-stations to their owners

What vantage6 does *not* (yet) do:

- formatting the data at the data station
- aligning data across the data stations (for the vertical partitioned use case)

The vantage6 infrastructure is designed with three fundamental functional aspects of federated learning.

1. **Autonomy**. All involved parties should remain independent and autonomous.
2. **Heterogeneity**. Parties should be allowed to have differences in hardware and operating systems.
3. **Flexibility**. Related to the latter, a federated learning infrastructure should not limit the use of relevant data.





## OVERVIEW OF THIS DOCUMENTATION

This documentation space consists of the following main sections:

- **Introduction** → *You are here now*
- *Installation* → *How to install vantage6 servers, nodes and clients*
- *How to use* → *How to use vantage6 servers, nodes and clients*
- */technical-documentation/index* (Under construction) → *Implementation details of the vantage6 platform*
- *Developer community* → *How to collaborate on the development of the vantage6 infrastructure*
- *Algorithm Development* → *Develop algorithms that are compatible with vantage6*
- *Glossary* → *A dictionary of common terms used in these docs*
- *Release notes* → *Log of what has been released and when*



## VANTAGE6 RESOURCES

This is a - non-exhaustive - list of vantage6 resources.

### Documentation

- [docs.vantage6.ai](https://docs.vantage6.ai) → *This documentation.*
- [vantage6.ai](https://vantage6.ai) → *vantage6 project website*
- [Academic papers](#) → *Technical insights into vantage6*

### Source code

- [vantage6](#) → *Contains all components (and the python-client).*
- [Planning](#) → *Contains all features, bugfixes and feature requests we are working on. To submit one yourself, you can create a [new issue](#).*

### Community

- [Discord](#) → *Chat with the vantage6 community*
  - [Community meetings](#) → *Bi-monthly developer community meeting*
-



## 4.1 Concepts

### 4.1.1 Architecture

In vantage6, a **client** can pose a question to the **server**, which is then delivered as an **algorithm** to the **node** (Fig. 4.1). When the algorithm completes, the node sends the results back to the client via the server. An algorithm may be enabled to communicate directly with twin algorithms running on other nodes.

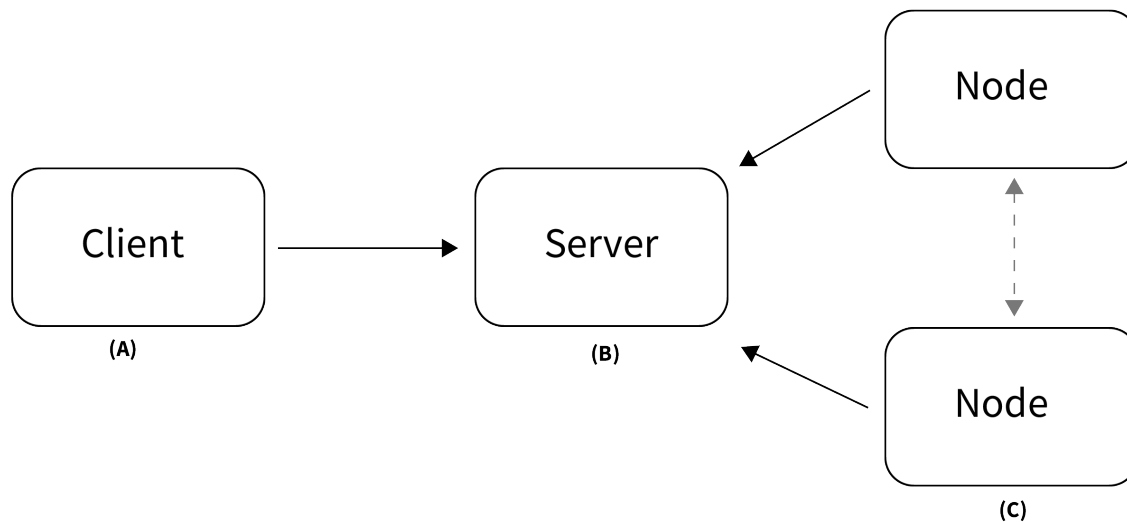


Fig. 4.1: Vantage6 has a client-server architecture. (A) The client is used by the researcher to create computation requests. It is also used to manage users, organizations and collaborations. (B) The server contains users, organizations, collaborations, tasks and their results. (C) The nodes have access to data and handle computation requests from the server.

Conceptually, vantage6 consists of the following parts:

- A (central) **server** that coordinates communication with clients and nodes. The server is in charge of processing tasks as well as handling administrative functions such as authentication and authorization.
- One or more **node(s)** that have access to data and execute algorithms
- **Users** (i.e. researchers or other applications) that request computations from the nodes via the client

- **Organizations** that are interested in collaborating. Each user belongs to one of these organizations.
- A **Docker registry** that functions as database of algorithms

On a technical level, vantage6 may be seen as a container orchestration tool for privacy preserving analyses. It deploys a network of containerized applications that together ensure insights can be exchanged without sharing record-level data.

### 4.1.2 Components

There are several entities in vantage6, such as users, organizations, tasks, etc. The following statements should help you understand their relationships.

- A **collaboration** is a collection of one or more **organizations**.
- For each collaboration, each participating organization needs a **node** to compute tasks.
- Each organization can have **users** who can perform certain actions.
- The permissions of the user are defined by the assigned **rules**.
- It is possible to collect multiple rules into a **role**, which can also be assigned to a user.
- Users can create **tasks** for one or more organizations within a collaboration.
- A task should produce a **result** for each organization involved in the task.

The following schema is a *simplified* version of the database:

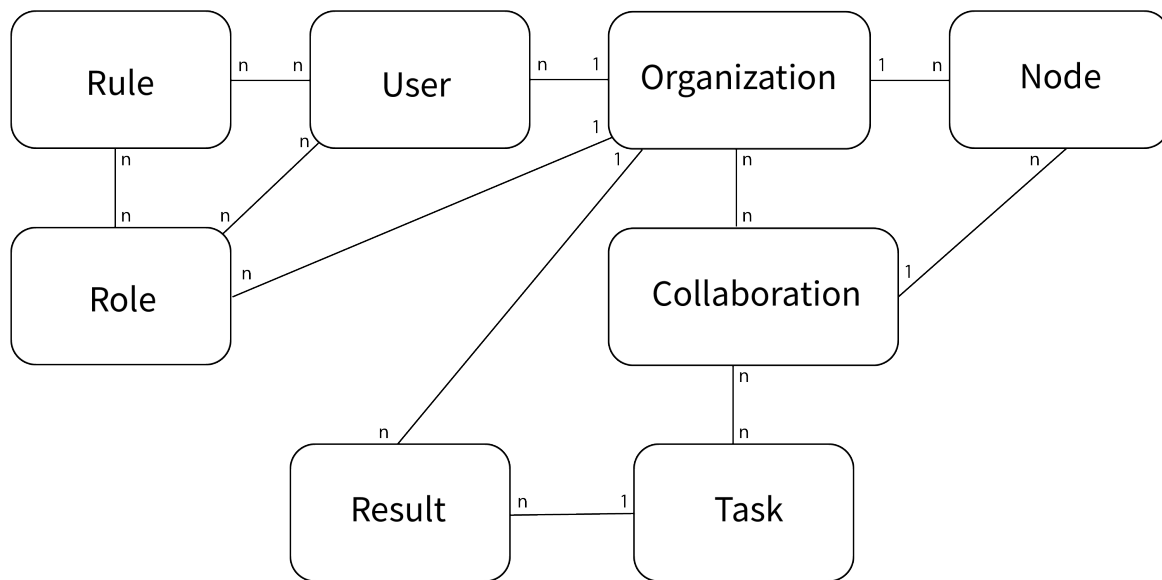


Fig. 4.2: Simplified database model

#### 4.1.2.1 End to end encryption

Encryption in vantage6 is handled at organization level. Whether encryption is used or not, is set at collaboration level. All the nodes in the collaboration need to agree on this setting. You can enable or disable encryption in the node configuration file, see the example in [Configuration file structure](#).

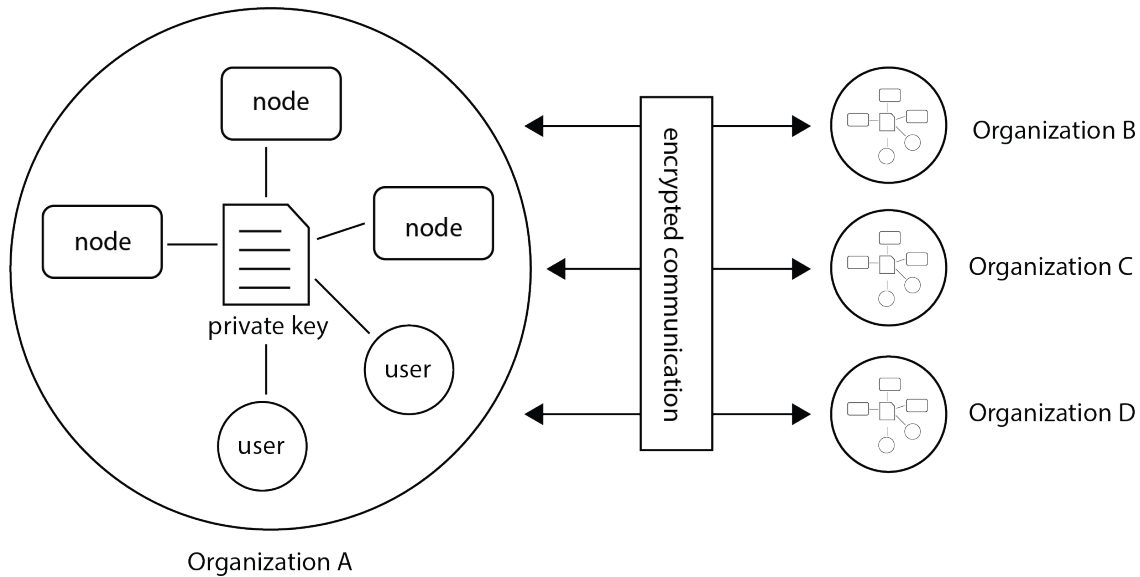


Fig. 4.3: Encryption takes place between organizations therefore all nodes and users from the a single organization should use the same private key.

The encryption module encrypts data so that the server is unable to read communication between users and nodes. The only messages that go from one organization to another through the server are computation requests and their results. Only the algorithm input and output are encrypted. Other metadata (e.g. time started, finished, etc), can be read by the server.

The encryption module uses RSA keys. The public key is uploaded to the vantage6-server. Tasks and other users can use this public key (this is automatically handled by the python-client and R-client) to send messages to the other parties.

---

**Note:** The RSA key is used to create a shared secret which is used for encryption and decryption of the payload.

---

When the node starts, it checks that the public key stored at the server is derived from the local private key. If this is not the case, the node will replace the public key at the server.

**Warning:** If an organization has multiple nodes and/or users, they must use the same private key.

In case you want to generate a new private key, you can use the command `vnode create-private-key`. If a key already exists at the local system, the existing key is reused (unless you use the `--force` flag). This way, it is easy to configure multiple nodes to use the same key.

It is also possible to generate the key yourself and upload it by using the endpoint `https://SERVER[/api_path]/organization/<ID>`.

## 4.2 Installation

The vantage6 framework consists of several components that should be installed. Which component(s) you need depends on your use case. For example, if you only need to communicate to an existing vantage6 server, you don't have to install a server.

We will first detail the requirements for installing a node and/or server, and then explain the installation process per component.

### 4.2.1 Requirements

#### 4.2.1.1 Node & Server

The (minimal) requirements of the node and server are similar. Note that these are recommendations: it may also work on other hardware, operating systems, versions of Python etc. (but they are not tested as much).

##### Hardware

- x86 CPU architecture + virtualization enabled
- 1 GB memory
- 50GB+ storage
- Stable and fast (1 Mbps+ internet connection)
- Public IP address

##### Software

- Operating system
  - Ubuntu 18.04+ or
  - MacOS Big Sur+ or
  - Windows 10+
- *Python*
- *Docker*

**Warning:** The hardware requirements of the node also depend on the algorithms that the node will run. For example, you need much less compute power for a descriptive statistical algorithm than for a machine learning model.

### Python

Installation of any of the vantage6 packages requires Python 3.7. For installation instructions, see [python.org](https://python.org), [anaconda.com](https://anaconda.com) or use the package manager native to your OS and/or distribution.

---

**Note:** We recommend you install vantage6 in a new, clean Python (Conda) environment.

Other version of Python  $\geq 3.6$  will most likely also work, but may give issues with installing dependencies. For now, we test vantage6 on version 3.7, so that is a safe choice.

---



## Docker

Docker facilitates encapsulation of applications and their dependencies in packages that can be easily distributed to diverse systems. Algorithms are stored in Docker images which nodes can download and execute. Besides the algorithms, both the node and server are also running from a docker container themselves.

Please refer to [this page](#) on how to install Docker. To verify that Docker is installed and running you can run the hello-world example from Docker.

```
docker run hello-world
```

**Warning:** Note that for **Linux**, some [post-installation steps](#) may be required. Vantage6 needs to be able to run docker without sudo, and these steps ensure just that.

---

### Note:

- Always make sure that Docker is running while using vantage6!
  - We recommend to always use the latest version of Docker.
- 

## 4.2.2 Client

We provide four ways in which you can interact with the server to manage your vantage6 resources: the *User interface* (UI), the *Python client*, the *R client*, and the server API. Below are installation instructions for each of them.

For most use cases, we recommend to use the UI (for anything except creating tasks - this is coming soon) and/or the Python Client. The latter covers the server functionality completely, but is more convenient for most users than sending HTTP requests directly to the API.

**Warning:** Depending on your algorithm it *may* be required to use a specific language to post a task and retrieve the results. This could happen when the output of an algorithm contains a language specific datatype and or serialization.

### 4.2.2.1 User interface

The UI is available as a website, so you don't have to install anything! Just go to the webpage and login with your user account. If you are using the Petronas server, simply go to <https://portal.petronas.vantage6.ai>.

If you are a server admin and want to set up a user interface, see *User Interface*.

#### 4.2.2.2 Python client library

Before you install the Python client, we check the version of the server you are going to interact with first. If you are using an existing server, check `https://<server_url>/version` (e.g. `https://petronas.vantage6.ai/version` or `http://localhost:5000/api/version`) to find its version.

Then you can install the `vantage6-client` with:

```
pip install vantage6==<version>
```

where you add the version you want to install. You may also leave out the version to install the most recent version.

#### 4.2.2.3 R client library

The R client currently only supports creating tasks and retrieving their results. It can not (yet) be used to manage resources, such as creating and deleting users and organizations.

You can install the R client by running:

```
devtools::install_github('IKNL/vtg', subdir='src')
```

#### 4.2.2.4 Server API

The API can be called via HTTP requests from a programming language of your choice. Hence, what you need to install, depends on you!

You can explore how to use the server API on `https://<serverdomain>/apidocs` (e.g. `https://petronas.vantage6.ai/apidocs` for our Petronas server).

### 4.2.3 Node

To install the **vantage6-node** make sure you have met the [requirements](#). Then install the latest version:

```
pip install vantage6
```

This will install the CLI in order to configure, start and stop the node. The node software itself will be downloaded when you start the node for the first time. See [here](#) for more details on how to do that.

### 4.2.4 Server

#### 4.2.4.1 Local Installation

This installs the vantage6 server at a VM or your local machine. First, make sure you have met the [Requirements](#). Then install the latest version:

```
pip install vantage6
```

This command will install the vantage6 command line interface (CLI), from which you can create new servers (see [Use Server](#)).

#### 4.2.4.2 Cloud Service Provider

To use vantage6 at a cloud service provider, you can use the Docker image we provide. Check the [Deployment](#) section for deployment examples.

---

**Note:** We recommend to provide the latest version. Should you have reasons to deploy an older VERSION, use the image `harbor2.vantage6.ai/infrastructure/server:<VERSION>`.

If you deploy an older version, it is also recommended that the nodes match that version.

---

#### 4.2.4.3 Optional components

There are several optional components that you can set up apart from the vantage6 server itself:

##### *User Interface*

A web application that will allow your users to interact more easily with your vantage6 server.

##### *EduVPN*

If you want to enable algorithm containers that are running on different nodes, to directly communicate with one another, you require a VPN server. Refer to on how to install the VPN server.

##### *rabbitmq*

If you have a server with a high workload whose performance you want to improve, you may want to set up a RabbitMQ service which enables horizontal scaling of the Vantage6 server.

##### *Docker registry*

A docker registry can be used to store algorithms but it is also possible to use [Docker hub](#) for this.

### User Interface

The User Interface (UI) is a web application that will make it easier for your users to interact with the server. It allows you to manage all your resources (such as creating collaborations, editing users, or viewing tasks), except for creating new tasks. We aim to incorporate this functionality in the near future.

If you plan on creating your own server and want to use interact with it via the UI, follow the instructions on the [UI Github page](#). We also provide a Docker image that runs the UI.

The UI is not compatible with older versions (<3.3) of vantage6.

### EduVPN

EduVPN is an optional server component that enables the use of algorithms that require node-to-node communication.

[EduVPN](#) provides an API for the OpenVPN server, which is required for automated certificate retrieval by the nodes. Like vantage6, it is an open source platform.

The following documentation shows you how to install EduVPN:

- [Debian](#)
- [Centos](#)
- [Fedora](#)

After the installation is done, you need to configure the server to:

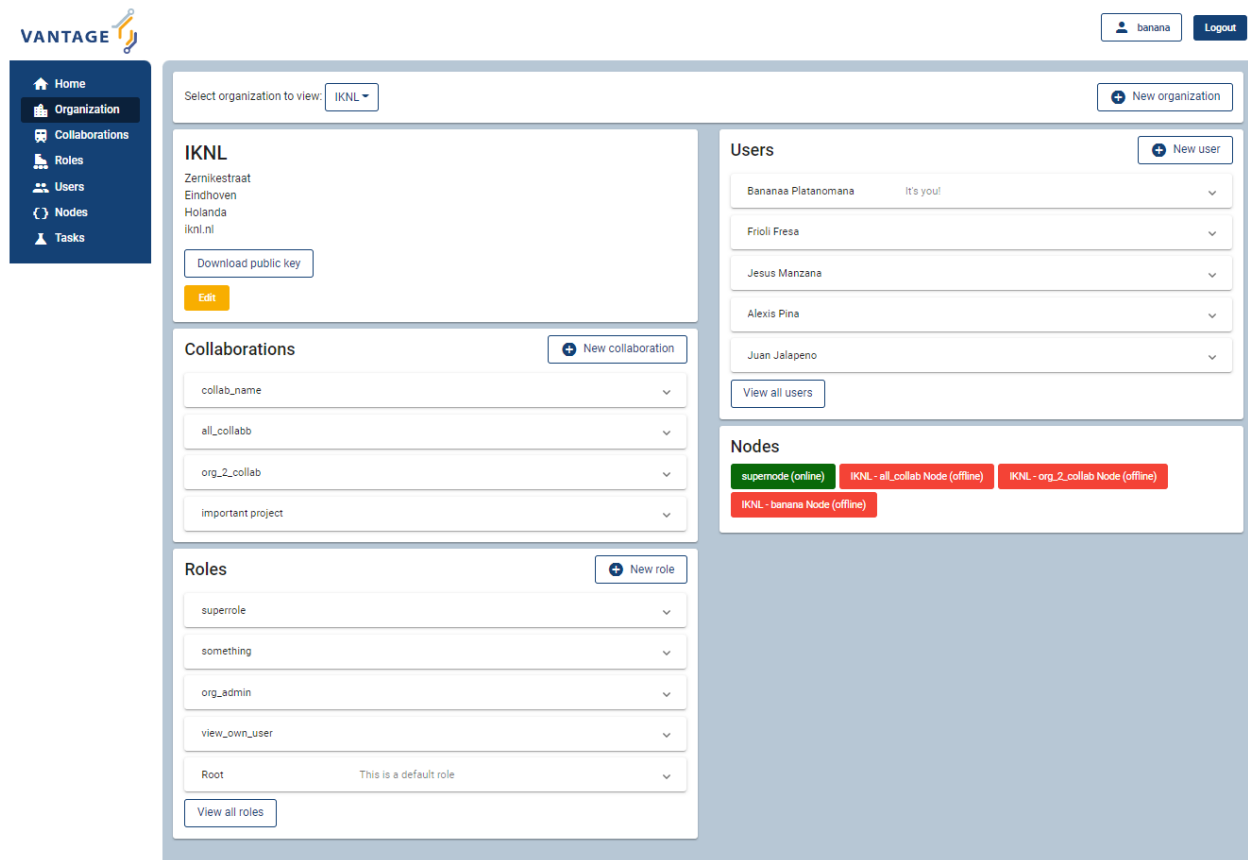


Fig. 4.4: Screenshot of the vantage6 UI

1. Enable client-to-client communication. This can be achieved in the configuration file by the `clientToClient` setting (see [here](#)).
2. Do not block LAN communication (set `blockLan` to `false`). This allows your docker subnetworks to continue to communicate, which is required for vantage6 to function normally.
3. Enable [port sharing](#) (Optional). This may be useful if the nodes are behind a strict firewall. Port sharing allows nodes to connect to the VPN server only using outgoing `tcp/443`. Be aware that [TCP meltdown](#) can occur when using the TCP protocol for VPN.
4. Create an application account.

**Warning:** EduVPN enables listening to multiple protocols (UDP/TCP) and ports at the same time. Be aware that all nodes need to be connected using the same protocol and port in order to communicate with each other.

**Warning:** The EduVPN server should usually be available to the public internet to allow all nodes to find it. Therefore, it should be properly secured, for example by closing all public ports (except `http/https`).

Additionally, you may want to explicitly allow *only* VPN traffic between nodes, and not between a node and the VPN server. You can achieve that by updating the firewall rules on your machine.

On Debian machines, these rules can be found in `/etc/iptables/rules.v4` and `/etc/iptables/rules.v6`, on CentOS, Red Hat Enterprise Linux and Fedora they can be found in `/etc/sysconfig/iptables` and `/etc/sysconfig/ip6tables`. You will have to do the following:

```
# In the firewall rules, below INPUT in the #SSH section, add this line
# to drop all VPN traffic with the VPN server as final destination:
-I INPUT -i tun+ -j DROP

# We only want to allow nodes to reach other nodes, and not other
# network interfaces available in the VPN.
# To achieve, replace the following rules:
-A FORWARD -i tun+ ! -o tun+ -j ACCEPT
-A FORWARD ! -i tun+ -o tun+ -j ACCEPT
# with:
-A FORWARD -i tun+ -o tun+ -j ACCEPT
-A FORWARD -i tun+ -j DROP
```

### Example configuration

The following configuration makes a server listens to `TCP/443` only. Make sure you set `clientToClient` to `true` and `blockLan` to `false`. The `range` needs to be supplied to the node configuration files. Also note that the server configured below uses [port-sharing](#).

```
// /etc/vpn-server-api/config.php
<?php

return [
    // List of VPN profiles
    'vpnProfiles' => [
        'internet' => [
            // The number of this profile, every profile per instance has a
            // unique number
```

(continues on next page)

(continued from previous page)

```

// REQUIRED
'profileNumber' => 1,

// The name of the profile as shown in the user and admin portals
// REQUIRED
'displayName' => 'vantage6 :: vpn service',

// The IPv4 range of the network that will be assigned to clients
// REQUIRED
'range' => '10.76.0.0/16',

// The IPv6 range of the network that will be assigned to clients
// REQUIRED
'range6' => 'fd58:63db:3245:d20d::/64',

// The hostname the VPN client(s) will connect to
// REQUIRED
'hostName' => 'eduvpn.vantage6.ai',

// The address the OpenVPN processes will listen on
// DEFAULT = ':::'
'listen' => ':::',

// The IP address to use for connecting to OpenVPN processes
// DEFAULT = '127.0.0.1'
'managementIp' => '127.0.0.1',

// Whether or not to route all traffic from the client over the VPN
// DEFAULT = false
'defaultGateway' => true,

// Block access to local LAN when VPN is active
// DEFAULT = false
'blockLan' => false,

// IPv4 and IPv6 routes to push to the client, only used when
// defaultGateway is false
// DEFAULT = []
'routes' => [],

// IPv4 and IPv6 address of DNS server(s) to push to the client
// DEFAULT = []
// Quad9 (https://www.quad9.net)
'dns' => ['9.9.9.9', '2620:fe::fe'],

// Whether or not to allow client-to-client traffic
// DEFAULT = false
'clientToClient' => true,

// Whether or not to enable OpenVPN logging
// DEFAULT = false
'enableLog' => false,

```

(continues on next page)

(continued from previous page)

```

// Whether or not to enable ACLs for controlling who can connect
// DEFAULT = false
'enableAcl' => false,

// The list of permissions to allow access, requires enableAcl to
// be true
// DEFAULT = []
'aclPermissionList' => [],

// The protocols and ports the OpenVPN processes should use, MUST
// be either 1, 2, 4, 8 or 16 proto/port combinations
// DEFAULT = ['udp/1194', 'tcp/1194']
'vpnProtoPorts' => [
    'tcp/1195',
],

// List the protocols and ports exposed to the VPN clients. Useful
// for OpenVPN port sharing. When empty (or missing), uses list
// from vpnProtoPorts
// DEFAULT = []
'exposedVpnProtoPorts' => [
    'tcp/443',
],

// Hide the profile from the user portal, i.e. do not allow the
// user to choose it
// DEFAULT = false
'hideProfile' => false,

// Protect to TLS control channel with PSK
// DEFAULT = tls-crypt
'tlsProtection' => 'tls-crypt',
// 'tlsProtection' => false,
],
],

// API consumers & credentials
'apiConsumers' => [
    'vpn-user-portal' => '***',
    'vpn-server-node' => '***',
],
];

```

The following configuration snippet can be used to add an API user. The username and the `client_secret` have to be added to the `vantage6-server` configuration file.

```

...
'Api' => [
    'consumerList' => [
        'vantage6-user' => [
            'redirect_uri_list' => [

```

(continues on next page)

(continued from previous page)

```
        'http://localhost',
    ],
    'display_name' => 'vantage6',
    'require_approval' => false,
    'client_secret' => '***'
]
]
...
```

## RabbitMQ

RabbitMQ is an optional component that enables the server to handle more requests at the same time. This is important if a server has a high workload.

There are several options to host your own RabbitMQ server. You can run [RabbitMQ in Docker](#) or host [RabbitMQ on Azure](#). When you have set up your RabbitMQ service, you can connect the server to it by adding the following to the server configuration:

```
rabbitmq_uri: amqp://<username>:<password@<hostname>:5672/<vhost>
```

Be sure to create the user and vhost that you specify exist! Otherwise, you can add them via the [RabbitMQ management console](#).

## Docker registry

A Docker registry or repository provides storage and versioning for Docker images. Installing a private Docker registry is useful if you want don't want to share your algorithms.

## Docker Hub

Docker itself provides a registry as a turn-key solution on Docker Hub. Instructions for setting it up can be found here: [https://hub.docker.com/\\_/registry](https://hub.docker.com/_/registry).

## Harbor

[Harbor](#) is another option for running a registry. Harbor provides access control, a user interface and automated scanning on vulnerabilities.

### 4.2.4.4 Deployment

The vantage6 server is a Flask application, together with [python-socketio](#) for websocket support. The server runs as a standalone process (listening on its own ip address/port).

There are many deployment options. We simply provide a few examples.

- [NGINX](#)
- [Azure app service](#)
- ...



**Note:** From version 3.2+ it is possible to horizontally scale the server (This upgrade is also made available to version 2.3.4)

Documentation on how to deploy it will be shared here soon. Reach out to us on Discord for now.

---

## NGINX

A basic setup is shown below. Note that SSL is not configured in this example.

```
server {  
  
    # Public port  
    listen 80;  
    server_name _;  
  
    # vantage6-server. In the case you use a sub-path here, make sure  
    # to forward also it to the proxy_pass  
    location /subpath {  
        include proxy_params;  
  
        # internal ip and port  
        proxy_pass http://127.0.0.1:5000/subpath;  
    }  
  
    # Allow the websocket traffic  
    location /socket.io {  
        include proxy_params;  
        proxy_http_version 1.1;  
        proxy_buffering off;  
        proxy_set_header Upgrade $http_upgrade;  
        proxy_set_header Connection "Upgrade";  
        proxy_pass http://127.0.0.1:5000/socket.io;  
    }  
}
```

**Note:** When you *Configure* the server, make sure to include the `/subpath` that has been set in the NGINX configuration into the `api_path` setting (e.g. `api_path: /subpath/api`)

---

## Azure app service

**Note:** We still have to document this. Reach out to us on Discord for now.

---

## 4.3 How to use

In this section of the documentation, we aim to explain to you how to use vantage6. The first part, *Client*, will inform you how to interact with a vantage6 server through the several types of clients that are available. The *Node* section will explain how to install a node for your organization and finally, in *Server* you will be informed how to setup your own vantage6 server.

### 4.3.1 Client

#### 4.3.1.1 Introduction

We provide four ways in which you can interact with the server to manage your vantage6 resources:

- *User Interface*
- *Python client*
- *R Client*
- *Server API*

The UI and the clients make it much easier to interact with the server than directly interacting with the server API through HTTP requests, especially as data is serialized and encrypted automatically. For most use cases, we recommend to use the UI and/or the Python client.

---

**Note:** The R client is only suitable for creating tasks and retrieve their results. With the Python client it is possible to use the entire API.

---

Note that whenever you interact with the server, you are limited by your permissions. For instance, if you try to create another user but do not have permission to do so, you will receive an error message. All permissions are described by rules, which can be aggregated in roles. Contact your server administrator if you find your permissions are inappropriate.

---

**Note:** There are predefined roles such as ‘Researcher’ and ‘Organization Admin’ that are automatically created by the server. These can be assigned to any new user by the administrator that is creating the user.

---

#### 4.3.1.2 User Interface

The User Interface (UI) is a web application that aims to make it easy to interact with the server. At present, it provides all functionality except for creating tasks. We aim to incorporate this functionality in the near future.

Using the UI should be relatively straightforward. There are buttons that should help you e.g. create a collaboration or delete a user. If anything is unclear, please contact us via [Discord](#).

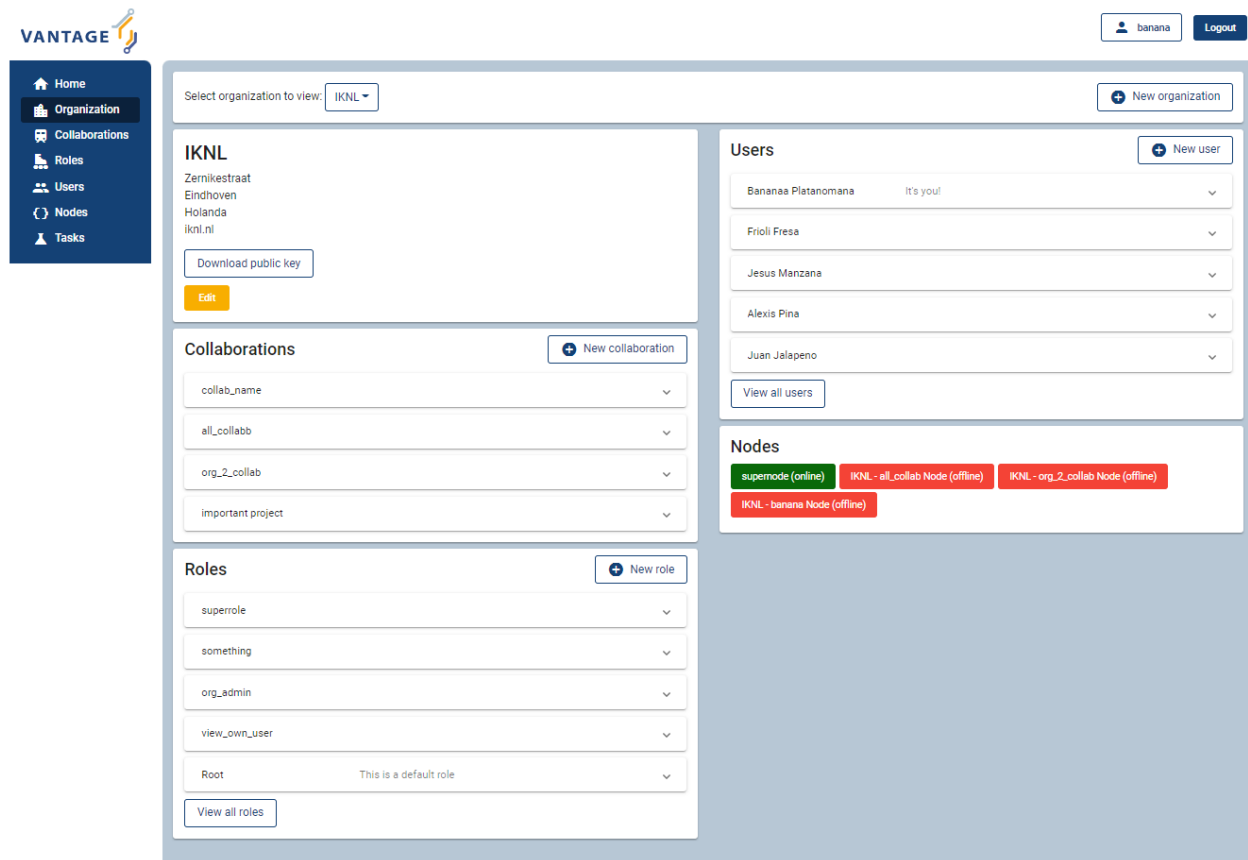


Fig. 4.5: Screenshot of the vantage6 UI

### 4.3.1.3 Python client

It is assumed you installed the *Client*. The Python client aims to completely cover the vantage6-server communication possibilities. It can create computation tasks and collect their results, manage organizations, collaborations, users, etc. The server hosts an API which the client uses for this purpose.

The methods in the library are all documented in their docstring, you can view them using `help(...)` , e.g. `help(client.user.create)` will show you the parameters needed to create a new user:

```
help(client.task.create)
#Create a new task
#
#   Parameters
#   -----
#   collaboration : int
#       Id of the collaboration to which this task belongs
#   organizations : list
#       Organization ids (within the collaboration) which need
#       to execute this task
#   name : str
#       Human readable name
#   image : str
#       Docker image name which contains the algorithm
#   description : str
#       Human readable description
#   input : dict
#       Algorithm input
#   data_format : str, optional
#       IO data format used, by default LEGACY
#   database: str, optional
#       Name of the database to use. This should match the key
#       in the node configuration files. If not specified the
#       default database will be tried.
#
#   Returns
#   -----
#   dict
#       Containing the task information
```

In *Authentication* and sections after that, there are more examples on how to use the Python client.

The following groups (related to the components) of methods are available, most of them have a `list()`, `create()`, `delete()` and `get()` method attached.

- `client.user`
- `client.organization`
- `client.rule`
- `client.role`
- `client.collaboration`
- `client.task`
- `client.result`
- `client.util`

- `client.node`

## Authentication

This section and the following sections introduce some minimal examples for administrative tasks that you can perform with our *Python client*. We start by authenticating.

To authenticate, we create a config file to store our login information. We do this so we do not have to define the `server_url`, `server_port` and so on every time we want to use the client. Moreover, it enables us to separate the sensitive information (login details, organization key) that you do not want to make publicly available, from other parts of the code you might write later (e.g. on submitting particular tasks) that you might want to share publicly.

```
# config.py

server_url = "https://MY VANTAGE6 SERVER" # e.g. https://petronas.vantage6.ai or
                                           # http://localhost for a local dev server
server_port = 443 # This is specified when you first created the server
server_api = "" # This is specified when you first created the server

username = "MY USERNAME"
password = "MY PASSWORD"

organization_key = "FILEPATH TO MY PRIVATE KEY" # This can be empty if you do not want
↳ to set up encryption
```

Note that the `organization_key` should be a filepath that points to the private key that was generated when the organization to which your login belongs was first created (see *Creating an organization*).

Then, we connect to the vantage 6 server by initializing a Client object, and authenticating

```
from vantage6.client import Client

# Note: we assume here the config.py you just created is in the current directory.
# If it is not, then you need to make sure it can be found on your PYTHONPATH
import config

# Initialize the client object, and run the authentication
client = Client(config.server_url, config.server_port, config.server_api,
               verbose=True)
client.authenticate(config.username, config.password)

# Optional: setup the encryption, if you have an organization_key
client.setup_encryption(config.organization_key)
```

**Note:** Above, we have added `verbose=True` as additional argument when creating the `Client(...)` object. This will print much more information that can be used to debug the issue.

## Creating an organization

After you have authenticated, you can start generating resources. The following also assumes that you have a login on the Vantage6 server that has the permissions to create a new organization. Regular end-users typically do not have these permissions (typically only administrators do); they may skip this part.

The first (optional, but recommended) step is to create an RSA keypair. A keypair, consisting of a private and a public key, can be used to encrypt data transfers. Users from the organization you are about to create will only be able to use encryption if such a keypair has been set up and if they have access to the private key.

```
from vantage6.common import warning, error, info, debug, bytes_to_base64s
from vantage6.client.encryption import RSACryptor
from pathlib import Path

# Generated a new private key
# Note that the file below doesn't exist yet: you will create it
private_key_filepath = r'/path/to/private/key'
private_key = RSACryptor.create_new_rsa_key(Path(private_key_filepath))

# Generate the public key based on the private one
public_key_bytes = RSACryptor.create_public_key_bytes(private_key)
public_key = bytes_to_base64s(public_key_bytes)
```

Now, we can create an organization

```
client.organization.create(
    name = 'The_Shire',
    address1 = '501 Buckland Road',
    address2 = 'Matamata',
    zipcode = '3472',
    country = 'New Zealand',
    domain = 'the_shire.org',
    public_key = public_key # use None if you haven't set up encryption
)
```

Users can now be created for this organization. Any users that are created and who have access to the private key we generated above can now use encryption by running

```
client.setup_encryption('/path/to/private/key')

# or, if you don't use encryption
client.setup_encryption(None)
```

after they authenticate.

## Creating a collaboration

Here, we assume that you have a Python session with an authenticated Client object, as created in [Authentication](#). We also assume that you have a login on the Vantage6 server that has the permissions to create a new collaboration (regular end-users typically do not have these permissions, this is typically only for administrators).

A collaboration is an association of multiple organizations that want to run analyses together. First, you will need to find the organization id's of the organizations you want to be part of the collaboration.

```
client.organization.list(fields=['id', 'name'])
```

Once you know the id's of the organizations you want in the collaboration (e.g. 1 and 2), you can create the collaboration:

```
collaboration_name = "fictional_collab"
organization_ids = [1,2] # the id's of the respective organizations
client.collaboration.create(name = collaboration_name,
                           organizations = organization_ids,
                           encrypted = True)
```

Note that a collaboration can require participating organizations to use encryption, by passing the `encrypted = True` argument (as we did above) when creating the collaboration. It is recommended to do so, but requires that a keypair was created when [Creating an organization](#) and that each user of that organization has access to the private key so that they can run the `client.setup_encryption(...)` command after [Authentication](#).

## Registering a node

Here, we again assume that you have a Python session with an authenticated Client object, as created in [Authentication](#), and that you have a login that has the permissions to create a new node (regular end-users typically do not have these permissions, this is typically only for administrators).

A node is associated with both a collaboration and an organization (see components). You will need to find the collaboration and organization id's for the node you want to register:

```
client.organization.list(fields=['id', 'name'])
client.collaboration.list(fields=['id', 'name'])
```

Then, we register a node with the desired organization and collaboration. In this example, we create a node for the organization with id 1 and collaboration with id 1.

```
# A node is associated with both a collaboration and an organization
organization_id = 1
collaboration_id = 1
api_key = client.node.create(collaboration = collaboration_id, organization = organization_id)
print(f"Registered a node for collaboration with id {collaboration_id}, organization with id {organization_id}. The API key that was generated for this node was {api_key}")
```

Remember to save the `api_key` that is returned here, since you will need it when you [Configure](#) the node.

## Creating a task

### Preliminaries

Here we assume that

- you have a Python session with an authenticated Client object, as created in [Authentication](#).
- you already have the algorithm you want to run available as a container in a docker registry (see [here](#) for more details on developing your own algorithm)
- the nodes are configured to look at the right database

In this manual, we'll use the averaging algorithm from `harbor2.vantage6.ai/demo/average`, so the second requirement is met. This container assumes a comma-separated (\*.csv) file as input, and will compute the average over one of the named columns. We'll assume the nodes in your collaboration have been configured to look at a comma-separated database, i.e. their config contains something like

```
databases:
  default: /path/to/my/example.csv
  my_other_database: /path/to/my/example2.csv
```

so that the third requirement is also met. As an end-user running the algorithm, you'll need to align with the node owner about which database name is used for the database you are interested in. For more details, see how to [Configure](#) your node.

### Determining which collaboration / organizations to create a task for

First, you'll want to determine which collaboration to submit this task to, and which organizations from this collaboration you want to be involved in the analysis

```
>>> client.collaboration.list(fields=['id', 'name', 'organizations'])
[
  {'id': 1, 'name': 'example_collab1',
   'organizations': [
     {'id': 2, 'link': '/api/organization/2', 'methods': ['GET', 'PATCH']},
     {'id': 3, 'link': '/api/organization/3', 'methods': ['GET', 'PATCH']},
     {'id': 4, 'link': '/api/organization/4', 'methods': ['GET', 'PATCH']}
   ]
}]
```

In this example, we see that the collaboration called `example_collab1` has three organizations associated with it, of which the organization id's are 2, 3 and 4. To figure out the names of these organizations, we run:

```
>>> client.organization.list(fields=['id', 'name'])
[{'id': 1, 'name': 'root'}, {'id': 2, 'name': 'example_org1'},
 {'id': 3, 'name': 'example_org2'}, {'id': 4, 'name': 'example_org3'}]
```

i.e. this collaboration consists of the organizations `example_org1` (with id 2), `example_org2` (with id 3) and `example_org3` (with id 4).

### Creating a task that runs the master algorithm

Now, we have two options: create a task that will run the master algorithm (which runs on one node and may spawns subtasks on other nodes), or create a task that will (only) run the RPC methods (which are run on each node). Typically, the RPC methods only run the node local analysis (e.g. compute the averages per node), whereas the master algorithms performs aggregation of those results as well (e.g. starts the node local analyses and then also computes the overall average). First, let us create a task that runs the master algorithm of the `harbor2.vantage6.ai/demo/average` container



```
input_ = {'method': 'master',
          'kwargs': {'column_name': 'age'},
          'master': True}

average_task = client.task.create(collaboration=1,
                                  organizations=[2,3],
                                  name="an-awesome-task",
                                  image="harbor2.vantage6.ai/demo/average",
                                  description='',
                                  input=input_,
                                  data_format='json')
```

Note that the kwargs we specified in the `input_` are specific to this algorithm: this algorithm expects an argument `column_name` to be defined, and will compute the average over the column with that name. Furthermore, note that here we created a task for collaboration with id 1 (i.e. our `example_collab1`) and the organizations with id 2 and 3 (i.e. `example_org1` and `example_org2`). I.e. the algorithm need not necessarily be run on *all* the organizations involved in the collaboration. Finally, note that `client.task.create()` has an optional argument called `database`. Suppose that we would have wanted to run this analysis on the database called `my_other_database` instead of the default database, we could have specified an additional `database = 'my_other_database'` argument. Check `help(client.task.create)` for more information.

### Creating a task that runs the RPC algorithm

You might be interested to know output of the RPC algorithm (in this example: the averages for the ‘age’ column for each node). In that case, you can run only the RPC algorithm, omitting the aggregation that the master algorithm will normally do:

```
input_ = {'method': 'average_partial',
          'kwargs': {'column_name': 'age'},
          'master': False}

average_task = client.task.create(collaboration=1,
                                  organizations=[2,3],
                                  name="an-awesome-task",
                                  image="harbor2.vantage6.ai/demo/average",
                                  description='',
                                  input=input_,
                                  data_format='json')
```

### Inspecting the results

Of course, it will take a little while to run your algorithm. You can use the following code snippet to run a loop that checks the server every 3 seconds to see if the task has been completed:

```
print("Waiting for results")
task_id = average_task['id']
task_info = client.task.get(task_id)
while not task_info.get("complete"):
    task_info = client.task.get(task_id, include_results=True)
    print("Waiting for results")
    time.sleep(3)

print("Results are ready!")
```

When the results are in, you can get the `result_id` from the task object:

```
result_id = task_info['id']
```

and then retrieve the results

```
result_info = client.result.list(task=result_id)
```

The number of results may be different depending on what you run, but for the master algorithm in this example, we can retrieve it using:

```
>>> result_info['data'][0]['result']
{'average': 53.25}
```

while for the RPC algorithm, dispatched to two nodes, we can retrieve it using

```
>>> result_info['data'][0]['result']
{'sum': 253, 'count': 4}
>>> result_info['data'][1]['result']
{'sum': 173, 'count': 4}
```

#### 4.3.1.4 R Client

It is assumed you installed the *R client library*. The R client can create tasks and retrieve their results. If you want to do more administrative tasks, either use the API directly or use the *Python client*.

Initialization of the R client can be done by:

```
setup.client <- function() {
  # Username/password should be provided by the administrator of
  # the server.
  username <- "username@example.com"
  password <- "password"

  host <- 'https://petronas.vantage6.ai:443'
  api_path <- ''

  # Create the client & authenticate
  client <- vtg::Client$new(host, api_path=api_path)
  client$authenticate(username, password)

  return(client)
}

# Create a client
client <- setup.client()
```

Then this client can be used for the different algorithms. Refer to the README in the repository on how to call the algorithm. Usually this includes installing some additional client-side packages for the specific algorithm you are using.

**Warning:** The R client is subject to change. We aim to make it more similar to the Python client.

## Example

This example shows how to run the vantage6 implementation of a federated Cox Proportional Hazard regression model. First you need to install the client side of the algorithm by:

```
devtools::install_github('iknl/vtg.coxph', subdir="src")
```

This is the code to run the coxph:

```
print( client$getCollaborations() )

# Should output something like this:
#   id   name
# 1  1 ZEPPELIN
# 2  2 PIPELINE

# Select a collaboration
client$setCollaborationId(1)

# Define explanatory variables, time column and censor column
expl_vars <- c("Age", "Race2", "Race3", "Mar2", "Mar3", "Mar4", "Mar5", "Mar9",
              "Hist8520", "hist8522", "hist8480", "hist8501", "hist8201",
              "hist8211", "grade", "ts", "nne", "nnp", "er2", "er4")
time_col <- "Time"
censor_col <- "Censor"

# vtg.coxph contains the function `dcoxph`.
result <- vtg.coxph::dcoxph(client, expl_vars, time_col, censor_col)
```

### 4.3.1.5 Server API

The server API is documented in the path [https://SERVER\[/api\\_path\]/apidocs](https://SERVER[/api_path]/apidocs). For Petronas, the API docs can thus be found at <https://petronas.vantage6.ai/apidocs>.

This page will show you which API endpoints exist and how you can use them. All endpoints communicate via HTTP requests, so you can communicate with them using any platform or programming language that supports HTTP requests.

## 4.3.2 Node

### 4.3.2.1 Introduction

It is assumed you have successfully installed `vantage6-node`. To verify this you can run the command `vnode --help`. If that prints a list of commands, the installation is completed. Also, make sure that Docker is running.

---

**Note:** An organization runs a node for each of the collaborations it participates in

---

## Quick start

To create a new node, run the command below. A menu will be started that allows you to set up a node configuration file. For more details, check out the [Configure](#) section.

```
vnode new
```

To run a node, execute the command below. The `--attach` flag will cause log output to be printed to the console.

```
vnode start --name <your_node> --attach
```

Finally, a node can be stopped again with:

```
vnode stop --name <your_node>
```

## Available commands

Below is a list of all commands you can run for your node(s). To see all available options per command use the `--help` flag, i.e. `vnode start --help`.

Command	Description
<code>vnode new</code>	Create a new node configuration file
<code>vnode start</code>	Start a node
<code>vnode stop</code>	Stop one or all nodes
<code>vnode files</code>	List the files of a node
<code>vnode attach</code>	Print the node logs to the console
<code>vnode list</code>	List all available nodes
<code>vnode create-private-key</code>	Create and upload a new public key for your organization

See the following sections on how to configure and maintain a vantage6-node instance:

- [Configure](#)
- [Security](#)
- [Logging](#)

### 4.3.2.2 Configure

The vantage6-node requires a configuration file to run. This is a `yaml` file with a specific format. To create an initial configuration file, start the configuration wizard via: `vnode new`. You can also create and/or edit this file manually.

The directory where the configuration file is stored depends on your operating system (OS). It is possible to store the configuration file at **system** or at **user** level. By default, node configuration files are stored at **user** level. The default directories per OS are as follows:

Operating System	System-folder	User-folder
Windows	<code>C:\ProgramData \vantage\node</code>	<code>C:\Users\&lt;user&gt; \AppData\Local\vantage\node</code>
MacOS	<code>/Library/Application Support/vantage6/node</code>	<code>/Users/&lt;user&gt;/Library/Application Support/vantage6/node</code>
Linux	<code>/etc/vantage6/node</code>	<code>/home/&lt;user&gt; /.config/vantage6/node</code>

**Warning:** The command `vnode` looks in certain directories by default. It is possible to use any directory and specify the location with the `--config` flag. However, note that using a different directory requires you to specify the `--config` flag every time!

## Configuration file structure

Each node instance (configuration) can have multiple environments. You can specify these under the key `environments` which allows four types: `dev`, `test`, `acc` and `prod`. If you do not want to specify any environment, you should only specify the key `application` (not within `environments`).

```
application:

# API key used to authenticate at the server.
api_key: ***

# URL of the vantage6 server
server_url: https://petronas.vantage6.ai

# port the server listens to
port: 443

# API path prefix that the server uses. Usually '/api' or an empty string
api_path: ''

# subnet of the VPN server
vpn_subnet: 10.76.0.0/16

# add additional environment variables to the algorithm containers.
# this could be usefull for passwords or other things that algorithms
# need to know about the node it is running on
# OPTIONAL
algorithm_env:

    # in this example the environment variable 'player' has
    # the value 'Alice' inside the algorithm container
    player: Alice

# specify custom Docker images to use for starting the different
# components.
# OPTIONAL
images:
    node: harbor2.vantage6.ai/infrastructure/node:petronas
    alpine: harbor2.vantage6.ai/infrastructure/alpine
    vpn_client: harbor2.vantage6.ai/infrastructure/vpn_client
    network_config: harbor2.vantage6.ai/infrastructure/vpn_network

# path or endpoint to the local data source. The client can request a
# certain database to be used if it is specified here. They are
# specified as label:local_path pairs.
databases:
    default: D:\data\datafile.csv
```

(continues on next page)

```
# end-to-end encryption settings
encryption:

    # whenever encryption is enabled or not. This should be the same
    # as the `encrypted` setting of the collaboration to which this
    # node belongs.
    enabled: false

    # location to the private key file
    private_key: /path/to/private_key.pem

# To control which algorithms are allowed at the node you can set
# the allowed_images key. This is expected to be a valid regular
# expression
allowed_images:
    - ^harbor.vantage6.ai/[a-zA-Z]+/[a-zA-Z]+

# credentials used to login to private Docker registries
docker_registries:
    - registry: docker-registry.org
      username: docker-registry-user
      password: docker-registry-password

# Create SSH Tunnel to connect algorithms to external data sources. The
# `hostname` and `tunnel:bind:port` can be used by the algorithm
# container to connect to the external data source. This is the address
# you need to use in the `databases` section of the configuration file!
ssh-tunnels:

    # Hostname to be used within the internal network. I.e. this is the
    # hostname that the algorithm uses to connect to the data source. Make
    # sure this is unique and the same as what you specified in the
    # `databases` section of the configuration file.
    - hostname: my-data-source

    # SSH configuration of the remote machine
    ssh:

        # Hostname or ip of the remote machine, in case it is the docker
        # host you can use `host.docker.internal` for Windows and MacOS.
        # In the case of Linux you can use `172.17.0.1` (the ip of the
        # docker bridge on the host)
        host: host.docker.internal
        port: 22

        # fingerprint of the remote machine. This is used to verify the
        # authenticity of the remote machine.
        fingerprint: "ssh-rsa ..."

        # Username and private key to use for authentication on the remote
        # machine
```

(continues on next page)

(continued from previous page)

```

identity:
    username: username
    key: /path/to/private_key.pem

# Once the SSH connection is established, a tunnel is created to
# forward traffic from the local machine to the remote machine.
tunnel:

    # The port and ip on the tunnel container. The ip is always
    # 0.0.0.0 as we want the algorithm container to be able to
    # connect.
    bind:
        ip: 0.0.0.0
        port: 8000

    # The port and ip on the remote machine. If the data source runs
    # on this machine, the ip most likely is 127.0.0.1.
    dest:
        ip: 127.0.0.1
        port: 8000

# Settings for the logger
logging:
    # Controls the logging output level. Could be one of the following
    # levels: CRITICAL, ERROR, WARNING, INFO, DEBUG, NOTSET
    level:          DEBUG

    # Filename of the log-file, used by RotatingFileHandler
    file:           my_node.log

    # whenever the output needs to be shown in the console
    use_console:    True

    # The number of log files that are kept, used by RotatingFileHandler
    backup_count:   5

    # Size kb of a single log file, used by RotatingFileHandler
    max_size:       1024

    # format: input for logging.Formatter,
    format:         "%(asctime)s - %(name)-14s - %(levelname)-8s - %(message)s"
    datefmt:        "%Y-%m-%d %H:%M:%S"

# directory where local task files (input/output) are stored
task_dir: C:\Users\<your-user>\AppData\Local\vantage6\node\tno1

```

**Note:** We use [DTAP](#) for [key environments](#). In short:

- dev: Development environment. It is ok to break things here
- test: Testing environment. Here, you can verify that everything works as expected. This environment should resemble the target environment where the final solution will be deployed as much as possible.

- **acc:** Acceptance environment. If the tests were successful, you can try this environment, where the final user will test his/her analysis to verify if everything meets his/her expectations.
  - **prod:** Production environment. The version of the proposed solution where the final analyses are executed.
- 

## Configure using the Wizard

The most straightforward way of creating a new server configuration is using the command `vnode new` which allows you to configure the most basic settings.

By default, the configuration is stored at user level, which makes this configuration available only for your user. In case you want to use a system directory you can add the `--system` flag when invoking the `vnode new` command.

## Update configuration

To update a configuration you need to modify the created `yaml` file. To see where this file is located, you can use the command `vnode files`. Do not forget to specify the flag `--system` in case of a system-wide configuration or the `--user` flag in case of a user-level configuration.

## Local test setup

Check the section on [Local test setup](#) of the server if you want to run both the node and server on the same machine.

### 4.3.2.3 Security

As a data owner it is important that you take the necessary steps to protect your data. Vantage6 allows algorithms to run on your data and share the results with other parties. It is important that you review the algorithms before allowing them to run on your data.

Once you approved the algorithm, it is important that you can verify that the approved algorithm is the algorithm that runs on your data. There are two important steps to be taken to accomplish this:

- Set the (optional) `allowed_images` option in the node-configuration file. You can specify a regex expression here. For example
  1. `^harbor2.vantage6.ai/[a-zA-Z]+/[a-zA-Z]+`: allows all images from the vantage6 registry
  2. `^harbor2.vantage6.ai/algorithms/glm`: only allows this specific image, but all builds of this image
  3. `^harbor2.vantage6.ai/algorithms/glm@sha256:82becede498899ec668628e7cb0ad87b6e1c371cb8a1e597d83a47fac21d6af3`: allows only this specific build from the GLM to run on your data
- Enable `DOCKER_CONTENT_TRUST` to verify the origin of the image. For more details see the [documentation from Docker](#).

**Warning:** By enabling `DOCKER_CONTENT_TRUST` you might not be able to use certain algorithms. You can check this by verifying that the images you want to be used are signed.

In case you are using our Docker repository you need to use `harbor2.vantage6.ai` as `harbor.vantage6.ai` does not have a notary.



#### 4.3.2.4 Logging

Logging is enabled by default. To configure the logger, look at the logging section in the example configuration file in *Configuration file structure*.

Useful commands:

1. `vnode files`: shows you where the log file is stored
2. `vnode attach`: shows live logs of a running server in your current console. This can also be achieved when starting the node with `vnode start --attach`

### 4.3.3 Server

#### 4.3.3.1 Introduction

It is assumed that you successfully installed the vantage6 *Server*. To verify this, you can run the command `vserver --help`. If that prints a list of commands, your installation is successful. Also, make sure that Docker is running.

#### Quick start

To create a new server, run the command below. A menu will be started that allows you to set up a server configuration file.

```
vserver new
```

For more details, check out the *Configure* section.

To run a server, execute the command below. The `--attach` flag will copy log output to the console.

```
vserver start --name <your_server> --attach
```

**Warning:** When the server is run for the first time, the following user is created:

- username: root
- password: root

It is recommended to change this password immediately.

Finally, a server can be stopped again with:

```
vserver stop --name <your_server>
```

## Available commands

The following commands are available in your environment. To see all the options that are available per command use the `--help` flag, e.g. `vserver start --help`.

Command	Description
<code>vserver new</code>	Create a new server configuration file
<code>vserver start</code>	Start a server
<code>vserver stop</code>	Stop a server
<code>vserver files</code>	List the files that a server is using
<code>vserver attach</code>	Show a server's logs in the current terminal
<code>vserver list</code>	List the available server instances
<code>vserver shell</code>	Run a server instance python shell
<code>vserver import</code>	Import server entities as a batch
<code>vserver version</code>	Shows the versions of all the components of the running server

The following sections explain how to use these commands to configure and maintain a vantage6-server instance:

- *Configure*
- *Batch import*
- *Deployment*
- *Logging*
- *Shell*

### 4.3.3.2 Configure

The vantage6-server requires a configuration file to run. This is a `yaml` file with specific contents. You can create and edit this file manually. To create an initial configuration file you can also use the configuration wizard: `vserver new`.

The directory where to store the configuration file depends on your operating system (OS). It is possible to store the configuration file at **system** or at **user** level. By default, server configuration files are stored at **system** level. The default directories per OS are as follows:

OS	System	User
Windows	<code>C:\ProgramData \vantage6\server</code>	<code>C:\Users\&lt;user&gt; \AppData\Local\vantage6\server\</code>
Macos	<code>/Library/Application Support/vantage6/server/</code>	<code>/Users/&lt;user&gt;/Library/Application Support/vantage6/server/</code>
Ubuntu	<code>/etc/xdg/vantage6/ server/</code>	<code>~/.config/vantage6/server/</code>

**Warning:** The command `vserver` looks in certain directories by default. It is possible to use any directory and specify the location with the `--config` flag. However, note that using a different directory requires you to specify the `--config` flag every time!

## Configuration file structure

Each server instance (configuration) can have multiple environments. You can specify these under the key environments which allows four types: dev ,test, acc and prod . If you do not want to specify any environment, you should only specify the key application (not within environments) .

```

application:
...
environments:
  test:

    # Human readable description of the server instance. This is to help
    # your peers to identify the server
    description: Test

    # Should be prod, acc, test or dev. In case the type is set to test
    # the JWT-tokens expiration is set to 1 day (default is 6 hours). The
    # other types can be used in future releases of vantage6
    type: test

    # IP adress to which the server binds. In case you specify 0.0.0.0
    # the server listens on all interfaces
    ip: 0.0.0.0

    # Port to which the server binds
    port: 5000

    # API path prefix. (i.e. https://yourdomain.org/api_path/<endpoint>). In the
    # case you use a referse proxy and use a subpath, make sure to include it
    # here also.
    api_path: /api

    # The URI to the server database. This should be a valid SQLAlchemy URI,
    # e.g. for an Sqlite database: sqlite:///database-name.sqlite,
    # or Postgres: postgresql://username:password@172.17.0.1/database).
    uri: sqlite:///test.sqlite

    # This should be set to false in production as this allows to completely
    # wipe the database in a single command. Useful to set to true when
    # testing/developing.
    allow_drop_all: True

    # The secret key used to generate JWT authorization tokens. This should
    # be kept secret as others are able to generate access tokens if they
    # know this secret. This parameter is optional. In case it is not
    # provided in the configuration it is generated each time the server
    # starts. Thereby invalidating all previous distributed keys.
    # OPTIONAL
    jwt_secret_key: super-secret-key! # recommended but optional

    # Settings for the logger
    logging:

```

(continues on next page)

(continued from previous page)

```

# Controls the logging output level. Could be one of the following
# levels: CRITICAL, ERROR, WARNING, INFO, DEBUG, NOTSET
level:          DEBUG

# Filename of the log-file, used by RotatingFileHandler
file:          test.log

# Whether the output is shown in the console or not
use_console:   True

# The number of log files that are kept, used by RotatingFileHandler
backup_count: 5

# Size in kB of a single log file, used by RotatingFileHandler
max_size:     1024

# format: input for logging.Formatter,
format:       "%(asctime)s - %(name)-14s - %(levelname)-8s - %(message)s"
datefmt:      "%Y-%m-%d %H:%M:%S"

# Configure a smtp mail server for the server to use for administrative
# purposes. e.g. allowing users to reset their password.
# OPTIONAL
smtp:
  port: 587
  server: smtp.yourmailserver.com
  username: your-username
  password: super-secret-password

# Set an email address you want to direct your users to for support
# (defaults to the address you set above in the SMTP server or otherwise
# to support@vantage6.ai)
support_email: your-support@email.com

# set how long reset token provided via email are valid (default 1 hour)
email_token_validity_minutes: 60

# If algorithm containers need direct communication between each other
# the server also requires a VPN server. (!) This must be a EduVPN
# instance as vantage6 makes use of their API (!)
# OPTIONAL
vpn_server:
  # the URL of your VPN server
  url: https://your-vpn-server.ext

  # OATH2 settings, make sure these are the same as in the
  # configuration file of your EduVPN instance
  redirect_url: http://localhost
  client_id: your_VPN_client_user_name
  client_secret: your_VPN_client_user_password

  # Username and password to access the EduVPN portal

```

(continues on next page)

(continued from previous page)

```
portal_username: your_eduvpn_portal_user_name
portal_userpass: your_eduvpn_portal_user_password
```

```
prod:
  ...
```

**Note:** We use [DTAP](#) for key environments. In short:

- **dev** Development environment. It is ok to break things here
- **test** Testing environment. Here, you can verify that everything works as expected. This environment should resemble the target environment where the final solution will be deployed as much as possible.
- **acc** Acceptance environment. If the tests were successful, you can try this environment, where the final user will test his/her analysis to verify if everything meets his/her expectations.
- **prod** Production environment. The version of the proposed solution where the final analyses are executed.

## Configuration wizard

The most straightforward way of creating a new server configuration is using the command `vserver new` which allows you to configure most settings. The [Configure](#) section details what each setting represents.

By default, the configuration is stored at system level, which makes this configuration available for *all* users. In case you want to use a user directory you can add the `--user` flag when invoking the `vserver new` command.

## Update configuration

To update a configuration you need to modify the created `yaml` file. To see where this file is located you can use the command `vserver files`. Do not forget to specify the flag `--system` in case of a system-wide configuration or the flag `--user` in case of a user-level configuration.

## Local test setup

If the nodes and the server run at the same machine, you have to make sure that the node can reach the server.

### Windows and MacOS

Setting the server IP to `0.0.0.0` makes the server reachable at your localhost (this is also the case when the dockerized version is used). In order for the node to reach this server, set the `server_url` setting to `host.docker.internal`.

#### warning

On the **M1** mac the local server might not be reachable from your nodes as `host.docker.internal` does not seem to refer to the host machine. Reach out to us on [Discourse](#) for a solution if you need this!

### Linux

You should bind the server to `0.0.0.0`. In the node configuration files, you can then use `http://172.17.0.1`, assuming you use the default docker network settings.

## Batch import

**Warning:** All users that are imported using `vserver import` receive the superuser role. We are looking into ways to also be able to import roles. For more background info refer to this [issue](#).

You can easily create a set of test users, organizations and collaborations by using a batch import. To do this, use the `vserver import /path/to/file.yaml` command. An example yaml file is provided here:

### organizations:

```
- name:      IKNL
  domain:    iknl.nl
  address1:  Godebaldkwartier 419
  address2:
  zipcode:   3511DT
  country:   Netherlands
  users:
    - username: admin
      firstname: admin
      lastname: robot
      password: password
    - username: frank@iknl.nl
      firstname: Frank
      lastname: Martin
      password: password
    - username: melle@iknl.nl
      firstname: Melle
      lastname: Sieswerda
      password: password
  public_key: LS0tLS1CRUdJTiBQVUJMSUMgS0VZLS0tLS0KTU1JQ01qQU5CZ2txaGtpRz13MEJB
  BUUVGQUFPQ0FnoEFNSU1DQ2dLQ0FnRUF2eU4wW

- name:      Small Organization
  domain:    small-organization.example
  address1:  Big Ambitions Drive 4
  address2:
  zipcode:   1234AB
  country:   Nowhereland
  users:
    - username: admin@small-organization.example
      firstname: admin
      lastname: robot
      password: password
    - username: stan
      firstname: Stan
      lastname: the man
      password: password
  public_key: LS0tLS1CRUdJTiBQVUJMSUMgS0VZLS0tLS0KTU1JQ01qQU5CZ2txaGtpRz13MEJB
  BUUVGQUFPQ0FnoEFNSU1DQ2dLQ0FnRUF2eU4wW

- name:      Big Organization
```

(continues on next page)

(continued from previous page)

```

domain:      big-organization.example
address1:    Offshore Accounting Drive 19
address2:
zipcode:     54331
country:     Nowhereland
users:
  - username: admin@big-organization.example
    firstname: admin
    lastname: robot
    password: password
public_key: ↵
↵LS0tLS1CRUdJTiBQVUJMSUMgS0VZLS0tLS0KTU1JQ0lqQU5CZ2txaGtpRzl3MEJBUUVGQUFPQ0FnoEFNSU1DQ2dLQ0FnRUF2eU4wW
collaborations:

- name: ZEPPELIN
  participants:
    - name: IKNL
      api_key: 123e4567-e89b-12d3-a456-426614174001
    - name: Small Organization
      api_key: 123e4567-e89b-12d3-a456-426614174002
    - name: Big Organization
      api_key: 123e4567-e89b-12d3-a456-426614174003
  tasks: ["hello-world"]
  encrypted: false

- name: PIPELINE
  participants:
    - name: IKNL
      api_key: 123e4567-e89b-12d3-a456-426614174004
    - name: Big Organization
      api_key: 123e4567-e89b-12d3-a456-426614174005
  tasks: ["hello-world"]
  encrypted: false

- name: SLIPPERS
  participants:
    - name: Small Organization
      api_key: 123e4567-e89b-12d3-a456-426614174006
    - name: Big Organization
      api_key: 123e4567-e89b-12d3-a456-426614174007
  tasks: ["hello-world", "hello-world"]
  encrypted: false

```

#### 4.3.3.3 Logging

Logging is enabled by default. To configure the logger, look at the `logging` section in the example configuration in *Configuration file structure*.

Useful commands:

1. `vserver files`: shows you where the log file is stored
2. `vserver attach`: show live logs of a running server in your current console. This can also be achieved when starting the server with `vserver start --attach`

#### 4.3.3.4 Shell

**Warning:** The preferred method of managing entities is using the API, instead of the shell interface. This because the API will perform validation of the input, whereas in the shell all inputs are accepted.

Through the shell it is possible to manage all server entities. To start the shell, use `vserver shell [options]`.

In the next sections the different database models that are available are explained. You can retrieve any record and edit any property of it. Every `db.` object has a `help()` method which prints some info on what data is stored in it (e.g. `db.Organization.help()`).

---

**Note:** Don't forget to call `.save()` once you are done editing an object.

---

### Organizations

---

**Note:** Organizations have a public key that is used for end-to-end encryption. This key is automatically created and/or uploaded by the node the first time it runs.

---

To store an organization you can use the `db.Organization` model:

```
# create new organization
organization = db.Organization(
    name="IKNL",
    domain="iknl.nl",
    address1="Zernikestraat 29",
    address2="Eindhoven",
    zipcode="5612HZ",
    country="Netherlands"
)

# store organization in the database
organization.save()
```

Retrieving organizations from the database:



```
# get all organizations in the database
organizations = db.Organization.get()

# get organization by its unique id
organization = db.Organization.get(1)

# get organization by its name
organization = db.Organization.get_by_name("IKNL")
```

A lot of entities (e.g. users) at the server are connected to an organization. E.g. you can see which (computation) tasks are issued by the organization or see which collaborations it is participating in.

```
# retrieve organization from which we want to know more
organization = db.Organization.get_by_name("IKNL")

# get all collaborations in which the organization participates
collaborations = organization.collaborations

# get all users from the organization
users = organization.users

# get all created tasks (from all users)
tasks = organization.created_tasks

# get the results of all these tasks
results = organization.results

# get all nodes of this organization (for each collaboration
# an organization participates in, it needs a node)
nodes = organization.nodes
```

## Roles and Rules

A user can have multiple roles and rules assigned to them. These are used to determine if the user has permission to view, edit, create or delete certain resources using the API. A role is a collection of rules.

```
# display all available rules
db.Rule.get()

# display rule 1
db.Rule.get(1)

# display all available roles
db.Role.get()

# display role 3
db.Role.get(3)

# show all rules that belong to role 3
db.Role.get(3).rules

# retrieve a certain rule from the DB
```

(continues on next page)

(continued from previous page)

```
rule = db.Rule.get_by_("node", Scope, Operation)

# create a new role
role = db.Role(name="role-name", rules=[rule])
role.save()

# or assign the rule directly to the user
user = db.User.get_by_username("some-existing-username")
user.rules.append(rule)
user.save()
```

## Users

Users belong to an organization. So if you have not created any *Organizations* yet, then you should do that first. To create a user you can use the `db.User` model:

```
# first obtain the organization to which the new user belongs
org = db.Organization.get_by_name("IKNL")

# obtain role 3 to assign to the new user
role_3 = db.Role.get(3)

# create the new users, see section Roles and Rules on how to
# deal with permissions
new_user = db.User(
    username="root",
    password="super-secret",
    firstname="John",
    lastname="Doe",
    roles=[role_3],
    rules=[],
    organization=org
)

# store the user in the database
new_user.save()
```

You can retrieve users in the following ways:

```
# get all users
db.User.get()

# get user 1
db.User.get(1)

# get user by username
db.User.get_by_username("root")

# get all users from the organization IKNL
db.Organization.get_by_name("IKNL").users
```

To modify a user, simply adjust the properties and save the object.

```

user = db.User.get_by_username("some-existing-username")

# update the firstname
user.firstname = "Brandnew"

# update the password; it is automatically hashed.
user.password = "something-new"

# store the updated user in the database
user.save()

```

## Collaborations

A collaboration consists of one or more organizations. To create a collaboration you need at least one but preferably multiple *Organizations* in your database. To create a collaboration you can use the `db.Collaboration` model:

```

# create a second organization to collaborate with
other_organization = db.Organization(
    name="IKNL",
    domain="iknl.nl",
    address1="Zernikestraat 29",
    address2="Eindhoven",
    zipcode="5612HZ",
    country="Netherlands"
)
other_organization.save()

# get organization we have created earlier
iknl = db.Organization.get_by_name("IKNL")

# create the collaboration
collaboration = db.Collaboration(
    name="collaboration-name",
    encrypted=False,
    organizations=[iknl, other_organization]
)

# store the collaboration in the database
collaboration.save()

```

Tasks, nodes and organizations are directly related to collaborations. We can obtain these by:

```

# obtain a collaboration which we like to inspect
collaboration = db.Collaboration.get(1)

# get all nodes
collaboration.nodes

# get all tasks issued for this collaboration
collaboration.tasks

```

(continues on next page)

(continued from previous page)

```
# get all organizations
collaboration.organizations
```

**Warning:** Setting the encryption to False at the server does not mean that the nodes will send encrypted results. This is only the case if the nodes also agree on this setting. If they don't, you will receive an error message.

## Nodes

Before nodes can login, they need to exist in the server's database. A new node can be created as follows:

```
# we'll use a uuid as the API-key, but you can use anything as
# API key
from uuid import uuid1

# nodes always belong to an organization *and* a collaboration,
# this combination needs to be unique!
iknl = db.Organization.get_by_name("IKNL")
collab = iknl.collaborations[0]

# generate and save
api_key = str(uuid1())
print(api_key)

node = db.Node(
    name = f"IKNL Node - Collaboration {collab.name}",
    organization = iknl,
    collaboration = collab,
    api_key = api_key
)

# save the new node to the database
node.save()
```

**Note:** API keys are hashed before stored in the database. Therefore you need to save the API key immediately. If you lose it, you can reset the API key later via the shell, API, client or UI.

## Tasks and Results

**Warning:** Tasks(/results) created from the shell are not picked up by nodes that are already running. The signal to notify them of a new task cannot be emitted this way. We therefore recommend sending tasks via the Python client.

A task is intended for one or more organizations. For each organization the task is intended for, a corresponding (initially empty) result should be created. Each task can have multiple results, for example a result from each organization.

```

# obtain organization from which this task is posted
iknl = db.Organization.get_by_name("IKNL")

# obtain collaboration for which we want to create a task
collaboration = db.Collaboration.get(1)

# obtain the next run_id. Tasks sharing the same run_id
# can share the temporary volumes at the nodes. Usually this
# run_id is assigned through the API (as the user is not allowed
# to do so). All tasks from a master-container share the
# same run_id
run_id = db.Task.next_run_id()

task = db.Task(
    name="some-name",
    description="some human readable description",
    image="docker-registry.org/image-name",
    collaboration=collaboration,
    run_id=run_id,
    database="default",
    initiator=iknl,
)
task.save()

# input the algorithm container (docker-registry.org/image-name)
# expects
input_ = {
}

import datetime

# now create a result model for each organization within the
# collaboration. This could also be a subset
for org in collaboration.organizations:
    res = db.Result(
        input=input_,
        organization=org,
        task=task,
        assigned_at=datetime.datetime.now()
    )
    res.save()

```

Tasks can have a child/parent relationship. Note that the `run_id` is for parent and child tasks the same.

```

# get a task to which we want to create some
# child tasks
parent_task = db.Task.get(1)

child_task = db.Task(
    name="some-name",
    description="some human readable description",
    image="docker-registry.org/image-name",
    collaboration=collaboration,

```

(continues on next page)

(continued from previous page)

```
run_id=parent_task.run_id,
database="default",
initiator=iknl,
parent=parent_task
)
child_task.save()
```

---

**Note:** Tasks that share a `run_id` have access to the same temporary folder at the node. This allows for multi-stage algorithms.

---

Obtaining results:

```
# obtain all Results
db.Result.get()

# obtain only completed results
[result for result in db.Result.get() if result.complete]

# obtain result by its unique id
db.Result.get(1)
```

## 4.4 Algorithm Development

This section helps you to develop MPC and FL algorithms that are compatible with vantage6. You are **not** going to find a list of algorithms here or help on how to use them. In the components the basic concepts and interface between node and algorithm is explained. Then in the *Classic Tutorial* a FL algorithm is build from scratch.

This section is to be extended with more examples in the future.

### 4.4.1 Concepts

Algorithms are executed at the (vantage6-)node. The node receives a computation task from the vantage6-server. The node will then retrieve the algorithm, execute it and return the results to the server.

Algorithms are shared using [Docker images](#) which are stored in a [Docker image registry](#) which is accessible to the nodes. In the following sections we explain the fundamentals of algorithm containers.

1. *Input & output* Interface between the node and algorithm container
2. *Wrappers* Library to simplify and standardized the node-algorithm input and output
3. *Child containers* Creating subtasks from an algorithm container
4. *Networking* Communicate with other algorithm containers and the vantage6-server
5. *Cross language* Cross language data serialization
6. *Package & distribute* Packaging and shipping algorithms

#### 4.4.1.1 Input & output

The algorithm runs in an isolated environment within the data station (node). As it is important to limit the connectivity and accessibility for obvious security reasons. In order for the algorithm to do its work, it is provided with several resources.

---

**Note:** This section describes the current process. Keep in mind that this is subjected to be changed. For more information, please see this [Github](#)

---

#### Environment variables

The algorithms have access to several environment variables, see [Section 4.4.1.1](#). These can be used to locate certain files or to add local configuration settings into the container.

Table 4.1: Environment variables

Variable	Description
INPUT_FILE	path to the input file. The input file contains the user defined input for the algorithms.
TOKEN_FILE	Path to the token file. The token file contains a JWT token which can be used to access the vantage6-server. This way the algorithm container is able to post new tasks and retrieve results.
TEMPORARY_FOLDER	Path to the temporary folder. This folder can be used to store intermediate results. These intermediate results are shared between all containers that have the same run_id. Algorithm containers which are created from an algorithm container themselves share the same run_id.
HOST	Contains the URL to the vantage6-server.
PORT	Contains the port to which the vantage6-server listens. Is used in combination with HOST and API_PATH.
API_PATH	Contains the api base path from the vantage6-server.
[*]_DATABASE_URI	Contains the URI of the local database. The * is replaced by the key specified in the node configuration file.

---

**Note:** Additional environment variables can be specified in the node configuration file using the `algorithm_env` key. These additional variables are forwarded to all algorithm containers.

---

#### File mounts

The algorithm container has access to several file mounts.

##### *Input*

The input file contains the user defined input. The user specifies this when a task is created.

##### *Output*

The algorithm should write its output to this file. When the docker container exits the contents of this file will be send back to the vantage6-server.

##### *Token*

The token file contains a JWT token which can be used by the algorithm to communicate with the central server.

The token can only be used to create a new task with the same image, and is only valid while the task has not yet been completed.

#### Temporary directory

The temporary directory can be used by an algorithm container to share files with other algorithm containers that:

- run on the same node
- have the same `run_id`

Algorithm containers that origin from another container (a.k.a master container or parent container) share the same `run_id`. i.o. if a user creates a task a new `run_id` is assigned.

The paths to these files and directories are stored in the environment variables, which we will explain now.

#### 4.4.1.2 Wrappers

The algorithm wrapper simplifies and standardizes the interaction between algorithm and node. The [client libraries](#) and the algorithm wrapper are tied together and use the same standards. The algorithm wrapper:

- reads the environment variables and file mounts and supplies these to your algorithm.
- provides an [entrypoint](#) for the docker container
- allows to write a single algorithm for multiple types of data sources

The wrapper is language specific and currently we support Python and R. Extending this concept to other languages is not so complex.

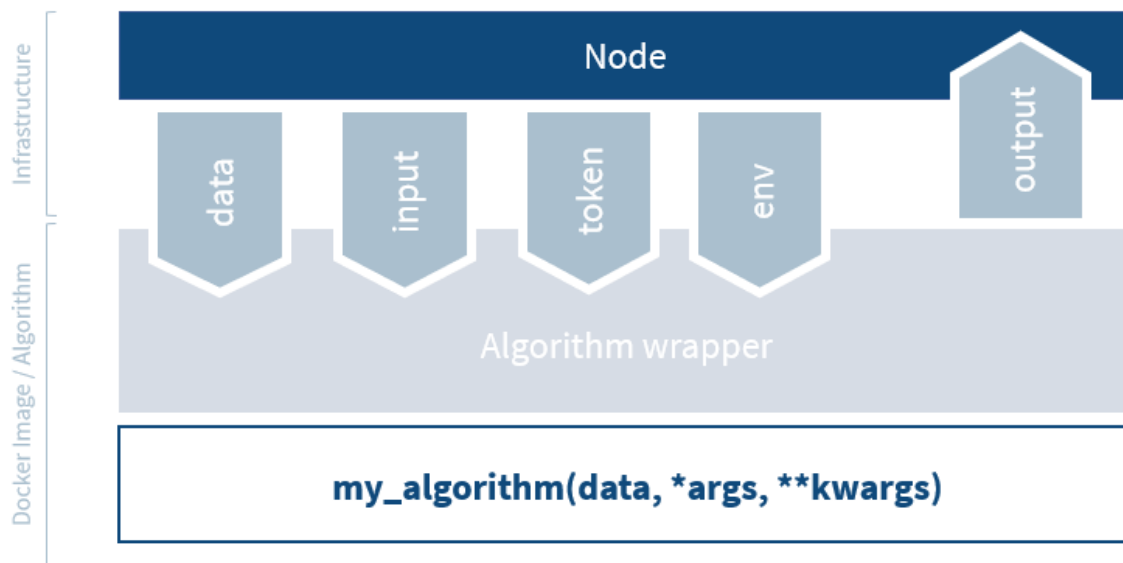


Fig. 4.6: The algorithm wrapper handles algorithm input and output.



## Federated function

The signature of your function has to contain `data` as the first argument. The method name should have a `RPC_` prefix. Everything that is returned by the function will be written to the output file.

*Python:*

```
def RPC_my_algorithm(data, *args, **kwargs):
    pass
```

*R:*

```
RPC_my_algorithm <- function(data, ...) {
}
```

## Central function

It is quite common to have a central part of your federated analysis which orchestrates the algorithm and combines the partial results. A common pattern for a central function would be:

1. Request partial models from all participants
2. Obtain the partial models
3. Combine the partial models to a global model
4. (optional) Repeat step 1-3 until the model converges

It is possible to run the central part of the analysis on your own machine, but it is also possible to let vantage6 handle the central part. There are several advantages to letting vantage6 handle this:

- You don't have to keep your machine running during the analysis
- You don't need to use the same programming language as the algorithm in case a language specific serialization is used in the algorithm

---

**Note:** Central functions also run at a node and *not* at the server.

---

In contrast to the federated functions, central functions are not prefixed. The first argument needs to be `client` and the second argument needs to be `data`. The `data` argument contains the local data and the `client` argument provides an interface to the vantage6-server.

**Warning:** The argument `data` is not present in the R wrapper. This is a consistency issue which will be solved in a future release.

```
def main(client, data, *args, **kwargs):
    # Run a federated function. Note that we omit the
    # RPC_ prefix. This prefix is added automatically
    # by the infrastructure
    task = client.create_new_task(
        {
            "method": "my_algorithm",
            "args": [],
```

(continues on next page)

(continued from previous page)

```

        "kwargs": {}
    },
    organization_ids=[...]
)

# wait for the federated part to complete
# and return
results = wait_and_collect(task)

return results

```

```

main <- function(client, ...) {
  # Run a federated function. Note that we omit the
  # RPC_prefix. This prefix is added automatically
  # by the infrastructure
  result <- client$call("my_algorithm", ...)

  # Optionally do something with the results

  # return the results
  return(result)
}

```

## Different wrappers

The docker wrappers read the local data source and supplies this to your functions in your algorithm. Currently CSV and SPARQL for Python and a CSV wrapper for R is supported. Since the wrapper handles the reading of the data, you need to rebuild your algorithm with a different wrapper to make it compatible with a different type of data source. You do this by updating the CMD directive in the dockerfile.

*CSV wrapper (Python)*

```

...
CMD python -c "from vantage6.tools.docker_wrapper import docker_wrapper; docker_wrapper('
↳ ${PKG_NAME}')"

```

*CSV wrapper (R)*

```

...
CMD Rscript -e "vtg::docker.wrapper('${PKG_NAME}')"

```

*SPARQL wrapper (Python)*

```

...
CMD python -c "from vantage6.tools.docker_wrapper import sparql_wrapper; sparql_wrapper('
↳ ${PKG_NAME}')"

```

## Data serialization

TODO

### 4.4.1.3 Mock client

TODO

### 4.4.1.4 Child containers

When a user creates a task, one or more nodes spawn an algorithm container. These algorithm containers can create new tasks themselves.

Every algorithm is supplied with a JWT token (see *Input & output*). This token can be used to communicate with the vantage6-server. In case you use a algorithm wrapper, you simply can use the supplied `Client` in case you use a *Central function*.

A child container can be a parent container itself. There is no limit to the amount of task layers that can be created. It is common to have only a single parent container which handles many child containers.

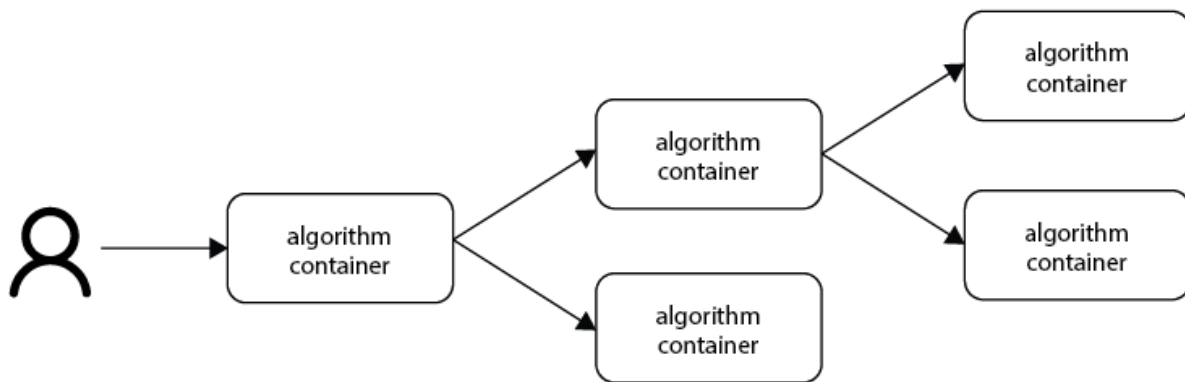


Fig. 4.7: Each container can spawn new containers in the network. Each container is provided with a unique token which they can use to communicate to the vantage6-server.

The token to which the containers have access supplies limited permissions to the container. For example, the token can be used to create additional tasks, but only in the same collaboration, and using the same image.

### 4.4.1.5 Networking

The algorithm container is deployed in an isolated network to reduce their exposure. Hence, the algorithm it cannot reach the internet. There are two exceptions:

1. When the VPN feature is enabled on the server all algorithm containers are able to reach each other using an ip and port over VPN.
2. The central server is reachable through a local proxy service. In the algorithm you can use the `HOST`, `POST` and `API_PATH` to find the address of the server.

**Note:** We are working on a whitelisting feature which allows a node to configure addresses that the algorithm container is able to reach.

---

## VPN connection

Algorithm containers can expose one or more ports. These ports can then be used by other algorithm containers to exchange data. The infrastructure uses the Dockerfile from which the algorithm has been build to determine to which ports are used by the algorithm. This is done by using the EXPOSE and LABEL directives.

For example when an algorithm uses two ports, one port for communication `com` and one port for data exchange `data`. The following block should be added to you algorithm Dockerfile:

```
# port 8888 is used by the algorithm for communication purposes
EXPOSE 8888
LABEL p8888 = "com"

# port 8889 is used by the algorithm for data-exchange
EXPOSE 8889
LABEL p8889 = "data"
```

Port 8888 and 8889 are the internal ports to which the algorithm container listens. When another container want to communicate with this container it can retrieve the IP and external port from the central server by using the `result_id` and the label of the port you want to use (`com` or `data` in this case)

### 4.4.1.6 Cross language

Because algorithms are exchanged through Docker images they can be written in any language. This is an advantage as developers can use their preferred language for the problem they need to solve.

**Warning:** The `wrappers` are only available for R and Python, so when you use different language you need to handle the IO yourself. Consult the *Input & Output* section on what the node supplies to your algorithm container.

When data is exchanged between the user and the algorithm they both need to be able to read the data. When the algorithm uses a language specific serialization (e.g. a `pickle` in the case of Python or `RData` in the case of R) the user needs to use the same language to read the results. A better solution would be to use a type of serialization that is not specific to a language. For our wrappers we use JSON for this purpose.

**Note:** Communication between algorithm containers can use language specific serialization as long as the different parts of the algorithm use the same language.

---

#### 4.4.1.7 Package & distribute

Once the algorithm is completed it needs to be packaged and made available for retrieval by the nodes. The algorithm is packaged in a Docker image. A Docker image is created from a Dockerfile, which acts as blue-print. Once the Docker image is created it needs to be uploaded to a registry so that nodes can retrieve it.

#### Dockerfile

A minimal Dockerfile should include a base image, injecting your algorithm and execution command of your algorithm. Here are several examples:

```
# python3 image as base
FROM python:3

# copy your algorithm in the container
COPY . /app

# maybe your algorithm is installable.
RUN pip install /app

# execute your application
CMD python /app/app.py
```

When using the Python Wrappers, the Dockerfile needs to follow a certain format. You should only change the PKG\_NAME value to the Python package name of your algorithm.

```
# python vantage6 algorithm base image
FROM harbor.vantage6.ai/algorithms/algorithm-base

# this should reflect the python package name
ARG PKG_NAME="v6-summary-py"

# install federated algorithm
COPY . /app
RUN pip install /app

ENV PKG_NAME=${PKG_NAME}

# Tell docker to execute `docker_wrapper()` when the image is run.
CMD python -c "from vantage6.tools.docker_wrapper import docker_wrapper; docker_wrapper('
↳ ${PKG_NAME}')
```

**Note:** When using the python wrapper your algorithm file needs to be installable. See [here](#) for more information on how to create a python package.

When using the R Wrappers, the Dockerfile needs to follow a certain format. You should only change the PKG\_NAME value to the R package name of your algorithm.

```
# The Dockerfile tells Docker how to construct the image with your algorithm.
# Once pushed to a repository, images can be downloaded and executed by the
# network hubs.
```

(continues on next page)

(continued from previous page)

```
FROM harbor2.vantage6.ai/base/custom-r-base

# this should reflect the R package name
ARG PKG_NAME='vtg.package'

LABEL maintainer="Main Tainer <m.tainer@vantage6.ai>"

# Install federated glm package
COPY . /usr/local/R/${PKG_NAME}/

WORKDIR /usr/local/R/${PKG_NAME}
RUN Rscript -e 'library(devtools)' -e 'install_deps(".")'
RUN R CMD INSTALL --no-multiarch --with-keep.source .

# Tell docker to execute `docker.wrapper()` when the image is run.
ENV PKG_NAME=${PKG_NAME}
CMD Rscript -e "vtg::docker.wrapper('$PKG_NAME')"
```

---

**Note:** Additional Docker directives are needed when using direct communication between different algorithm containers, see *Networking*.

---

## Build & upload

If you are in the folder containing the Dockerfile, you can build the project as follows:

```
docker build -t repo/image:tag .
```

The `-t` indicated the name of your image. This name is also used as reference where the image is located on the internet. If you use Docker hub to store your images, you only specify your username as `repo` followed by your image name and tag: `USERNAME/IMAGE_NAME:IMAGE_TAG`. When using a private registry `repo` should contain the URL of the registry also: e.g. `harbor2.vantage6.ai/PROJECT/IMAGE_NAME:TAG`.

Then you can push you image:

```
docker push repo/image:tag
```

Now that is has been uploaded it is available for nodes to retrieve when they need it.

## Signed images

It is possible to use the Docker the framework to create signed images. When using signed images, the node can verify the author of the algorithm image adding an additional protection layer.

Dockerfile

- Build project
- CMD
- Expose

## 4.4.2 Classic Tutorial

In this section the basic steps for creating an algorithm for horizontally partitioned data are explained.

**Note:** The final code of this tutorial is published on [Github](#). The algorithm is also published in our Docker registry: [harbor2.vantage6.ai/demo/average](https://harbor2.vantage6.ai/demo/average)

It is assumed that it is mathematically possible to create a federated version of the algorithm you want to use. In the following sections we create a federated algorithm to compute the average of a distributed dataset. An overview of the steps that we are going through:

1. Mathematically decompose the model
2. Federated implementation and local testing
3. Vantage6 algorithm wrapper
4. Dockerize and push to a registry

This tutorial shows you how to create a **federated mean** algorithm.

### 4.4.2.1 Mathematical decomposition

The mean of  $Q = [q_1, q_2 \dots q_n]$  is computed as:

$$Q_{mean} = \frac{1}{n} \sum_{i=1}^n q_i = \frac{q_1 + q_2 + \dots + q_n}{n}$$

When dataset  $Q$  is **horizontally partitioned** in dataset  $A$  and  $B$ :

$$\begin{aligned} A &= [a_1, a_2 \dots a_j] = [q_1, q_2 \dots q_j] \\ B &= [b_1, b_2 \dots b_k] = [q_{j+1}, q_{j+2} \dots q_n] \end{aligned}$$

We would like to compute  $Q_{mean}$  from dataset  $A$  and  $B$ . This could be computed as:

$$Q_{mean} = \frac{(a_1 + a_2 + \dots + a_j) + (b_1 + b_2 + \dots + b_k)}{j + k} = \frac{\sum A + \sum B}{j + k}$$

Both the number of samples in each dataset and the total sum of each dataset is needed. Then we can compute the global average of dataset  $A$  and  $B$ .

**Note:** We cannot simply compute the average on each node and combine them, as this would be mathematically incorrect. This would only work if dataset **A** and **B** contain the exact same number of samples.

### 4.4.2.2 Federated implementation

**Warning:** In this example we use python, however you are free to use any language. The only requirements are: 1) It has to be able to create HTTP-requests, and 2) has to be able to read and write to files.

However, if you use a different language you are not able to use our wrapper. Reach out to us on [Discord](#) to discuss how this works.

A federated algorithm consist of two parts:

1. A federated part of the algorithm which is responsible for creating the partial results. In our case this would be computing (1) the sum of the observations, and (2) the number of observations.
2. A central part of the algorithm which is responsible for combining the partial results from the nodes. In the case of the federated mean that would be dividing the total sum of the observations by the total number of observations.

---

**Note:** The central part of the algorithm can either be run on the machine of the researcher himself or in a master container which runs on a node. The latter is the preferred method.

In case the researcher runs this part, he/she needs to have a proper setup to do so (i.e. Python 3.5+ and the necessary dependencies). This can be useful when developing new algorithms.

---

## Federated part

The node that runs this part contains a CSV-file with one column (specified by the argument *column\_name*) which we want to use to compute the global mean. We assume that this column has no *NaN* values.

```
import pandas

def federated_part(path, column_name="numbers"):
    """Compute the sum and number of observations of a column"""

    # extract the column numbers from the CSV
    numbers = pandas.read_csv(path)[column_name]

    # compute the sum, and count number of rows
    local_sum = numbers.sum()
    local_count = len(numbers)

    # return the values as a dict
    return {
        "sum": local_sum,
        "count": local_count
    }
```

## Central part

The central algorithm receives the sums and counts from all sites and combines these to a global mean. This could be from one or more sites.

```
def central_part(node_outputs):
    """Combine the partial results to a global average"""
    global_sum = 0
    global_count = 0
    for output in node_outputs:
        global_sum += output["sum"]
        global_count += output["count"]

    return {"average": global_sum / global_count}
```



## Local testing

To test, simply create two datasets **A** and **B**, both having a numerical column **numbers**. Then run the following:

```
outputs = [
    federated_part("path/to/dataset/A"),
    federated_part("path/to/dataset/B")
]
Q_average = central_part(outputs)["average"]
print(f"global average = {Q_average}.")
```

### 4.4.2.3 Vantage6 integration

**Note:** A good starting point would be to use the boilerplate code from our [Github](#). This section outlines the steps needed to get to this boilerplate but also provides some background information.

**Note:** In this example we use a **csv**-file. It is also possible to use other types of data sources. This tutorial makes use of our algorithm wrapper which is currently only available for **csv**, **SPARQL** and **Parquet** files.

Other wrappers like **SQL**, **OMOP**, etc. are under consideration. Let us now if you want to use one of these or other datasources.

Now that we have a federated implementation of our algorithm we need to make it compatible with the vantage6 infrastructure. The infrastructure handles the communication with the server and provides data access to the algorithm.

The algorithm consumes a file containing the input. This contains both the method name to be triggered as well as the arguments provided to the method. The algorithm also has access to a CSV file (in the future this could also be a database) on which the algorithm can run. When the algorithm is finished, it writes back the output to a different file.

The central part of the algorithm has to be able to create (sub)tasks. These subtasks are responsible for executing the federated part of the algorithm. The central part of the algorithm can either be executed on one of the nodes in the vantage6 network or on the machine of a researcher. In this example we only show the case in which one of the nodes executes the central part of the algorithm. The node provides the algorithm with a JWT token so that the central part of the algorithm has access to the server to post these subtasks.

## Algorithm Structure

The algorithm needs to be structured as a Python [package](#). This way the algorithm can be installed within the Docker image. The minimal file-structure would be:

```
project_folder
├── Dockerfile
├── setup.py
├── algorithm_pkg
│   └── __init__.py
```

We also recommend adding a `README.md`, `LICENSE` and `requirements.txt` to the *project\_folder*.

## setup.py

Contains the setup method to create a package from your algorithm code. Here you specify some details about your package and the dependencies it requires.

```
from os import path
from codecs import open
from setuptools import setup, find_packages

# we're using a README.md, if you do not have this in your folder, simply
# replace this with a string.
here = path.abspath(path.dirname(__file__))
with open(path.join(here, 'README.md'), encoding='utf-8') as f:
    long_description = f.read()

# Here you specify the meta-data of your package. The `name` argument is
# needed in some other steps.
setup(
    name='v6-average-py',
    version="1.0.0",
    description='vantage6 average',
    long_description=long_description,
    long_description_content_type='text/markdown',
    url='https://github.com/IKNL/v6-average-py',
    packages=find_packages(),
    python_requires='>=3.6',
    install_requires=[
        'vantage6-client',
        # list your dependencies here:
        # pandas, ...
    ]
)
```

---

**Note:** The setup.py above is sufficient in most cases. However if you want to do more advanced stuff (like adding static data, or a CLI) you can use the [extra arguments](#) from setup.

---

## Dockerfile

The Dockerfile contains the recipe for building the Docker image. Typically you only have to change the argument PKG\_NAME to the name of you package. This name should be the same as as the name you specified in the setup.py. In our case that would be v6-average-py.

```
# This specifies our base image. This base image contains some commonly used
# dependancies and an install from all vantage6 packages. You can specify a
# different image here (e.g. python:3). In that case it is important that
# `vantage6-client` is a dependancy of you project as this contains the wrapper
# we are using in this example.
FROM harbor.vantage6.ai/algorithms/algorithm-base

# Change this to the package name of your project. This needs to be the same
```

(continues on next page)

(continued from previous page)

```
# as what you specified for the name in the `setup.py`.
ARG PKG_NAME="v6-average-py"

# This will install your algorithm into this image.
COPY . /app
RUN pip install /app

# This will run your algorithm when the Docker container is started. The
# wrapper takes care of the IO handling (communication between node and
# algorithm). You dont need to change anything here.
ENV PKG_NAME=${PKG_NAME}
CMD python -c "from vantage6.tools.docker_wrapper import docker_wrapper; docker_wrapper('
↳ ${PKG_NAME}')"

```

### \_\_init\_\_.py

This contains the code for your algorithm. It is possible to split this into multiple files, however the methods that should be available to the researcher should be in this file. You can do that by simply importing them into this file (e.g. `from .average import my_nested_method`)

We can distinguish two types of methods that a user can trigger:

name	description	pre- fix	arguments
master	Central part of the algorithm. Receives a <code>client</code> as argument which provides an interface to the central server. This way the master can create tasks and collect their results.		( <code>client</code> , <code>data</code> , <code>*args</code> , <code>**kwargs</code> )
Remote proce- dure call	Consumes the data at the node to compute the partial.	<i>RPC</i>	( <code>data</code> , <code>*args</code> , <code>**kwargs</code> )

**Warning:** Everything that is returned by the `return` statement is sent back to the central vantage6-server. This should never contain any privacy-sensitive information.

**Warning:** The client the master method receives is a `ContainerClient` which is different than the client you use as a user.

For our average algorithm the implementation will look as follows:

```
import time

from vantage6.tools.util import info

def master(client, data, column_name):
    """Combine partials to global model

```

(continues on next page)

(continued from previous page)

First we collect the parties that participate in the collaboration. Then we send a task to all the parties to compute their partial (the row count and the column sum). Then we wait for the results to be ready. Finally when the results are ready, we combine them to a global average.

Note that the master method also receives the (local) data of the node. In most usecases this data argument is not used.

The client, provided in the first argument, gives an interface to the central server. This is needed to create tasks (for the partial results) and collect their results later on. Note that this client is a different client than the client you use as a user.

```
"""
```

```
# Info messages can help you when an algorithm crashes. These info
# messages are stored in a log file which is send to the server when
# either a task finished or crashes.
```

```
info('Collecting participating organizations')
```

```
# Collect all organization that participate in this collaboration.
# These organizations will receive the task to compute the partial.
organizations = client.get_organizations_in_my_collaboration()
ids = [organization.get("id") for organization in organizations]
```

```
# Request all participating parties to compute their partial. This
# will create a new task at the central server for them to pick up.
# We've used a kwarg but is is also possible to use `args`. Although
# we prefer kwargs as it is clearer.
```

```
info('Requesting partial computation')
```

```
task = client.create_new_task(
    input_={
        'method': 'average_partial',
        'kwargs': {
            'column_name': column_name
        }
    },
    organization_ids=ids
)
```

```
# Now we need to wait untill all organizations(/nodes) finished
# their partial. We do this by polling the server for results. It is
# also possible to subscribe to a websocket channel to get status
# updates.
```

```
info("Waiting for results")
```

```
task_id = task.get("id")
```

```
task = client.get_task(task_id)
```

```
while not task.get("complete"):
```

```
    task = client.get_task(task_id)
```

```
    info("Waiting for results")
```

```
    time.sleep(1)
```

(continues on next page)

(continued from previous page)

```

# Once we now the partials are complete, we can collect them.
info("Obtaining results")
results = client.get_results(task_id=task.get("id"))

# Now we can combine the partials to a global average.
global_sum = 0
global_count = 0
for result in results:
    global_sum += result["sum"]
    global_count += result["count"]

return {"average": global_sum / global_count}

def RPC_average_partial(data, column_name):
    """Compute the average partial

    The data argument contains a pandas-dataframe containing the local
    data from the node.
    """

    # extract the column_name from the dataframe.
    info(f'Extracting column {column_name}')
    numbers = data[column_name]

    # compute the sum, and count number of rows
    info('Computing partials')
    local_sum = numbers.sum()
    local_count = len(numbers)

    # return the values as a dict
    return {
        "sum": local_sum,
        "count": local_count
    }

```

## Local testing

Now that we have a vantage6 implementation of the algorithm it is time to test it. Before we run it in a vantage6 setup we can test it locally by using the ClientMockProtocol which simulates the communication with the central server.

Before we can locally test it we need to (editable) install the algorithm package so that the Mock client can use it. Simply go to the root directory of your algorithm package (with the setup.py file) and run the following:

```
pip install -e .
```

Then create a script to test the algorithm:

```

from vantage6.tools.mock_client import ClientMockProtocol

# Initialize the mock server. The datasets simulate the local datasets from
# the node. In this case we have two parties having two different datasets:

```

(continues on next page)

(continued from previous page)

```

# a.csv and b.csv. The module name needs to be the name of your algorithm
# package. This is the name you specified in `setup.py`, in our case that
# would be v6-average-py.
client = ClientMockProtocol(
    datasets=["local/a.csv", "local/b.csv"],
    module="v6-average-py"
)

# to inspect which organization are in your mock client, you can run the
# following
organizations = client.get_organizations_in_my_collaboration()
org_ids = ids = [organization["id"] for organization in organizations]

# we can either test a RPC method or the master method (which will trigger the
# RPC methods also). Lets start by triggering an RPC method and see if that
# works. Note that we do *not* specify the RPC_ prefix for the method! In this
# example we assume that both a.csv and b.csv contain a numerical column `age`.
average_partial_task = client.create_new_task(
    input_={
        'method': 'average_partial',
        'kwargs': {
            'column_name': 'age'
        }
    },
    organization_ids=org_ids
)

# You can directly obtain the result (we dont have to wait for nodes to
# complete the tasks)
results = client.get_results(average_partial_task.get("id"))
print(results)

# To trigger the master method you also need to supply the `master`-flag
# to the input. Also note that we only supply the task to a single organization
# as we only want to execute the central part of the algorithm once. The master
# task takes care of the distribution to the other parties.
average_task = client.create_new_task(
    input_={
        'master': 1,
        'method': 'master',
        'kwargs': {
            'column_name': 'age'
        }
    },
    organization_ids=[org_ids[0]]
)
results = client.get_results(average_task.get("id"))
print(results)

```

## Building and Distributing

Now that we have a fully tested algorithm for the vantage6 infrastructure. We need to package it so that it can be distributed to the data-stations/nodes. Algorithms are delivered in Docker images. So that's where we need the Dockerfile for. To build an image from our algorithm (make sure you have docker installed and it's running) you can run the following command from the root directory of your algorithm project.

```
docker build -t harbor2.vantage6.ai/demo/average .
```

The option `-t` specifies the (unique) identifier used by the researcher to use this algorithm. Usually this includes the registry address (harbor2.vantage6.ai) and the project name (demo).

---

**Note:** In case you are using docker hub as registry, you do not have to specify the registry or project as these are set by default to the Docker hub and your docker hub username.

---

```
docker push harbor2.vantage6.ai/demo/average
```

---

**Note:** Reach out to us on [Discord](#) if you want to use our registries (harbor.vantage6.ai and harbor2.vantage6.ai).

---

### 4.4.2.4 Cross-language serialization

It is possible that a vantage6 algorithm is developed in one programming language, but you would like to run the task from another language. For these use-cases, the Python algorithm wrapper and client support cross-language serialization. By default, input to the algorithms and output back to the client are serialized using pickle. However, it is possible to define a different serialization format.

Input and output serialization can be specified as follows:

```
client.post_task(  
    name='mytask',  
    image='harbor2.vantage6.ai/testing/v6-test-py',  
    collaboration_id=COLLABORATION_ID,  
    organization_ids=ORGANIZATION_IDS,  
    data_format='json', # Specify input format to the algorithm  
    input_={  
        'method': 'column_names',  
        'kwargs': {'data_format': 'json'}, # Specify output format  
    }  
)
```

## 4.5 Developer community

As an open-source platform, we welcome anyone who would like to contribute to the vantage6 code and/or documentation. The following sections are meant to clarify our processes in development, documentation and releasing.

### 4.5.1 Contribute

#### 4.5.1.1 Support questions

If you have questions, you can use

- [Github discussions](#)
- Ask us on [Discord](#)

We prefer that you ask questions via these routes rather than creating Github issues. The issue tracker is intended to address bugs, feature requests, and code changes.

#### 4.5.1.2 Reporting issues

Issues can be posted at our [Github issue page](#), or, if you have issues that are specific to the user interface, please post them to the [UI issue page](#).

We distinguish between the following types of issues:

- Bug report: you encountered broken code
- Feature request: you want something to be added
- Change request: there is a something you would like to be different but it is not considered a new feature nor is something broken
- Security vulnerabilities: you found a security issue

Each issue type has its own template. Using these templates makes it easier for us to manage them.

**Warning:** Security vulnerabilities should not be reported in the Github issue tracker as they should not be publicly visible. To see how we deal with security vulnerabilities read our [policy](#).

See the [Security vulnerabilities](#) section when you want to release a security patch yourself.

We distribute the open issues in sprints and hotfixes. You can check out these boards here:

- [Sprints](#)
- [Hotfixes](#)

When a high impact bug is reported, we will put it on the hotfix board and create a patch release as soon as possible.

The sprint board tracks which issues we plan to fix in which upcoming release. Low-impact bugs, new features and changes will be scheduled into a sprint periodically. We automatically assign the label 'new' to all newly reported issues to track which issues should still be scheduled.

If you would like to fix an existing bug or create a new feature, check [Submitting patches](#) for more details on e.g. how to set up a local development environment and how the release process works. We prefer that you let us know you what are working on so we prevent duplicate work.



### 4.5.1.3 Security vulnerabilities

If you are a member of the Vantage6 Github organization, you can create an security advisory in the [Security](#) tab. See [Table 4.2](#) on what to fill in.

If you are not a member, please reach out directly to Frank Martin and/or Bart van Beusekom, or any other project member. They can then create a security advisory for you.

Table 4.2: Advisory details

Name	Details
Ecosystem	Set to <code>pip</code>
Package name	Set to <code>vantage6</code>
Affected versions	Specify the versions (or set of verions) that are affected
Patched version	Version where the issue is addressed, you can fill this in later when the patch is released.
Severity	Determine severity score using <a href="#">this</a> tool. Then use table <a href="#">Table 4.3</a> to determine the level from this score.
Common weakness enumerator (CWE)	Find the CWE (or multiple) on <a href="#">this</a> website.

Table 4.3: Severity

Score	Level
0.1-3.9	Low
4.0-6.9	Medium
7.0-8.9	High
9.0-10.0	Critical

Once the advisory has been created it is possible to create a private fork from there (Look for the button `Start a temporary private fork`). This private fork should be used to solve the issue.

From the same page you should request a CVE number so we can alert dependent software projects. Github will review the request. We are not sure what this entails, but so far they approved all advisories.

### 4.5.1.4 Community Planning

We host bi-monthly community meetings intended for aligning development efforts. Anyone is welcome to join although they are mainly intended for infrastructure and algorithm developers. There is an opportunity to present what your team is working on and find collaboration partners.

Community meetings schedule:

Table 4.4: Community meetings

Date	Time
17 November 2022	10:00 - 12:00
19 January 2023	10:00 - 12:00
16 March 2023	10:00 - 12:00

Reach out on [Discord](#) if you want to join the community meeting.

#### 4.5.1.5 Submitting patches

If there is not an open issue for what you want to submit, please open one for discussion before submitting the PR. We encourage you to reach out to us on [Discord](#), so that we can work together to ensure your contribution is added to the repository.

The workflow below is specific to the [vantage6 infrastructure repository](#). However, the concepts for our other repositories are the same. Then, modify the links below and ignore steps that may be irrelevant to that particular repository.

#### Setup your environment

- Make sure you have a Github account
- Install and configure git
- (Optional) install and configure Miniconda
- Clone the main repository locally:

```
git clone https://github.com/vantage6/vantage6
cd vantage6
```

- Add your fork as a remote to push your work to. Replace {username} with your username.

```
git remote add fork https://github.com/{username}/vantage6
```

- Create a virtual environment to work in. For miniconda:

```
conda create -n vantage6 python=3.7
conda activate vantage6
```

It is also possible to use `virtualenv` if you do not have a conda installation.

- Update pip and setuptools

```
python -m pip install --upgrade pip setuptools
```

- Install vantage6 as development environment with the `-e` flag.

```
pip install -e .
```

#### Coding

First, create a branch you can work on. Make sure you branch of the latest main branch:

```
git fetch origin
git checkout -b your-branch-name origin/main
```

Then you can create your bugfix, change or feature. Make sure to commit frequently. Preferably include tests that cover your changes.

Finally, push your commits to your fork on Github and create a pull request.

```
git push --set-upstream fork your-branch-name
```

Please apply the [PEP8](#) standards to your code.

## Local test setup

To test your code changes, it may be useful to create a local test setup. There are several ways of doing this.

1. Use the command `vserver-local` and `vnode-local`. This runs the application in your current activated Python environment.
2. Use the command `vserver` and `vnode` in combination with the options `--mount-src` and optionally `--image`.
  - The `--mount-src` option will run your current code in the docker image. The provided path should point towards the root folder of the [vantage6 repository](#).
  - The `--image` can be used to point towards a custom build infrastructure image. Note that when your code update includes dependency upgrades you need to build a custom infrastructure image as the 'old' image does not contain these and the `--mount-src` option will only overwrite the source and not re-install dependencies.

---

**Note:** If you are using Docker Desktop (which is usually the case if you are on Windows or MacOS) and want to setup a test environment, you should use `http://host.docker.internal` for the server address in the node configuration file. You should not use `http://localhost` in that case as that points to the localhost within the docker container instead of the system-wide localhost.

---

## Unit tests & coverage

You can execute unit tests using the `test` command in the Makefile:

```
make test
```

If you want to execute a specific unit test (e.g. the one you just created or one that is failing), you can use a command like:

```
python -m unittest tests_folder.test_filename.TestClassName.test_name
```

This command assumes you are in the directory above `tests_folder`. If you are inside the `tests_folder`, then you should remove that part.

## Pull Request

Please consider first which branch you want to merge your contribution into. **Patches** are usually directly merged into `main`, but **features** are usually merged into a development branch (e.g. `dev3` for version 3) before being merged into the `main` branch.

Before the PR is merged, it should pass the following requirements:

- At least one approved review of a code owner
- All [unit tests](#) should complete
- [CodeQL](#) (vulnerability scanning) should pass
- [Codacy](#) - Code quality checks - should be OK
- [Coveralls](#) - Code coverage analysis - should not decrease

## Documentation

Depending on the changes you made, you may need to add a little (or a lot) of documentation. For more information on how and where to edit the documentation, see the section [Documentation](#).

Consider which documentation you need to update:

- **User documentation.** Update it if your change led to a different experience for the end-user
- **Technical documentation.** Update it if you added new functionality. Check if your function docstrings have also been added (see last bullet below).
- **OAS (Open API Specification).** If you changed input/output for any of the API endpoints, make sure to add it to the docstrings. See [API Documentation with OAS3+](#) for more details.
- **Function docstrings** These should always be documented using the [numpy format](#). Such docstrings can then be used to automatically generate parts of the technical documentation space.

### 4.5.2 Documentation

The vantage6 framework is documented on this website. Additionally, there is [API Documentation with OAS3+](#). This documentation is shipped directly with the server instance. All of these documentation pages are described in more detail below.

#### 4.5.2.1 How this documentation is created

The source of the documentation you are currently reading is located [here](#), in the docs folder of the *vantage6* repository itself.

To build the documentation locally, there are two options. To build a static version, you can do `make html` when you are in the docs directory. If you want to automatically refresh the documentation whenever you make a change, you can use `sphinx-autobuild`. Assuming you are in the main directory of the repository, run the following commands:

```
pip install -r docs/requirements.txt
sphinx-autobuild docs docs/_build/html --watch .
```

Of course, you only have to install the requirements if you had not done so before.

Then you can access the documentation on `http://127.0.0.1:8000`. The `--watch` option makes sure that if you make changes to either the documentation text or the docstrings, the documentation pages will also be reloaded.

This documentation is automatically built and published on a commit (on certain branches, including `main`). Both Frank and Bart have access to the vantage6 project when logged into readthedocs. Here they can manage which branches are to be synced, manage the webhook used to trigger a build, and some other -less important- settings.

The files in this documentation use the `rst` format, to see the syntax view [this cheatsheet](#).

### 4.5.2.2 API Documentation with OAS3+

The API documentation is hosted at the server at the `/apidocs` endpoint. This documentation is generated from the docstrings using `Flasgger`. The source of this documentation can be found in the docstrings of the API functions.

If you are unfamiliar with OAS3+, note that it was formerly known as Swagger.

An example of such a docstring:

```
"""Summary of the endpoint
---
description: >-
    Short description on what the endpoint does, and which users have
    access or which permissions are required.

parameters:
    - in: path
      name: id
      schema:
          type: integer
          description: some identifier
          required: true

responses:
    200:
        description: Ok
    401:
        description: Unauthorized or missing permission

security:
    - bearerAuth: []

tags: ["Group"]
"""
```

### 4.5.3 Release

This page is intended to provide information about our release process. First, we discuss the version formatting, after which we discuss the actual creation and distribution of a release.

#### 4.5.3.1 Version format

Semantic versioning is used: `Major.Minor.Patch.Pre[N].Post<n>`.

**Major** is used for releasing breaking changes. For example, when the database model has changed, a new major version should be issued.

**Minor** is used for releasing new features, enhancements and other changes that are compatible with all other components. An example is the release of a new endpoint.

**Patch** is used for bugfixes and other minor changes

**Pre[N]** is used for alpha (a), beta (b) and release candidates (rc) releases and the build number is appended (e.g. `2.0.1b1` indicates the first beta-build of version `2.0.1`)

**Post[N]** is used for a rebuild where no code changes have been made, but where, for example, a dependency has been updated and a rebuild is required.

**Warning:** Post releases are only used by versioning the Docker images. Code changes should never be released with a `.post[N]` version.

#### 4.5.3.2 Create a release

To create a new release, one should go through the following steps:

- Check out the correct branch of the [vantage6](#) repository and pull the latest version:

```
git checkout main
git pull
```

*Make sure the branch is up-to-date.* **Patches** are usually directly merged into main, but for **minor** or **major** releases you usually need to execute a pull request from a development branch.

- Create a tag for the release. See [Version format](#) for more details on version names:

```
git tag version/x.y.z
```

- Push the tag to the remote. This will trigger the release pipeline on Github:

```
git push origin version/x.y.z
```

---

**Note:** The release process is protected and can only be executed by certain people. Reach out if you have any questions regarding this.

---

#### 4.5.3.3 The release pipeline

The release pipeline executes the following steps:

1. It checks if the tag contains a valid version specification. If it does not, the process it stopped.
2. Update the version in the repository code to the version specified in the tag and commit this back to the main branch.
3. Install the dependencies and build the Python package.
4. Upload the package to PyPi.
5. Build and push the Docker image to [harbor2.vantage6.ai](#).
6. Post a message in Discord to alert the community of the new release. This is not done if the version is a pre-release (e.g. `version/x.y.0rc1`).

---

**Note:** If you specify a tag with a version that already exists, the build pipeline will fail as the upload to PyPi is rejected.

---

The release pipeline uses a number of environment variables to, for instance, authenticate to PyPi and Discord. These variables are listed and explained in the table below.

Table 4.5: Environment variables

Secret	Description
COMMIT_PAT	Github Personal Access Token with commit privileges. This is linked to an individual user with admin right as the commit on the main needs to bypass the protections. There is unfortunately not -yet- a good solution for this.
ADD_TO_PROJECT_PAT	Github Personal Access Token with project management privileges. This token is used to add new issues to project boards.
COVERALLS_TOKEN	Token from coveralls to post the test coverage stats.
DOCKER_TOKEN	Token used together DOCKER_USERNAME to upload the container images to our <a href="https://harbor2.vantage6.ai">https://harbor2.vantage6.ai</a> .
DOCKER_USERNAME	See DOCKER_TOKEN.
PYPI_TOKEN	Token used to upload the Python packages to PyPi.
DISCORD_RELEASE_TOKEN	Token to post a message to the Discord community when a new release is published.

#### 4.5.3.4 Distribute release

Nodes and servers that are already running will automatically be upgraded to the latest version of their major release when they are restarted. This happens by pulling the newly released docker image. Note that the major release is never automatically updated: for example, a node running version 2.1.0 will update to 2.1.1 or 2.2.0, but never to 3.0.0. Depending on the version of Vantage6 that is being used, there is a reserved Docker image tag for distributing the upgrades. These are the following:

Tag	Description
petronas	3.x.x release
harukas	2.x.x release
troltungu	1.x.x release

Docker images can be pulled manually with e.g.

```
docker pull harbor2.vantage6.ai/infrastructure/server:petronas
docker pull harbor2.vantage6.ai/infrastructure/node:3.1.0
```

#### 4.5.3.5 User Interface release

The release process for the user interface (UI) is very similar to the release of the infrastructure detailed above. The same versioning format is used, and when you push a version tag, the automated release process is triggered.

We have semi-synchronized the version of the UI with that of the infrastructure. That is, we try to release major and minor versions at the same time. For example, if we are currently at version 3.5 and release version 3.6, we release it both for the infrastructure and for the UI. However, there may be different patch versions for both: the latest version for the infrastructure may then be 3.6.2 while the UI may still be at 3.6.

The release pipeline for the UI executes the following steps:

1. Version tag is verified (same as infrastructure).
2. Version is updated in the code (same as infrastructure).
3. Application is built.
4. Docker images are built and released to harbor2.
5. Application is pushed to our UI deployment slot (an Azure app service).

## 4.6 Glossary

The following is a list of definitions used in vantage6.

### A

- **Autonomy:** the ability of a party to be in charge of the control and management of its own data.

### C

- **Collaboration:** an agreement between two or more parties to participate in a study (i.e., to answer a research question).

### D

- **Distributed learning:** see *Federated Learning*
- **Docker:** a platform that uses operating system virtualization to deliver software in packages called containers. It is worth noting that although they are often confused, [Docker containers are not virtual machines](#).
- **Data Station:** Virtual Machine containing the vantage6-node application and a database.

### F

- **FAIR data:** data that are Findable, Accessible, Interoperable, and Reusable. For more information, see [the original paper](#).
- **Federated learning:** an approach for analyzing data that are spread across different parties. Its main idea is that parties run computations on their local data, yielding either aggregated parameters or encrypted values. These are then shared to generate a global (statistical) model. In other words, instead of bringing the data to the algorithms, federated learning brings the algorithms to the data. This way, patient-sensitive information is not disclosed. Federated learning is some times known as *distributed learning*. However, we try to avoid this term, since it can be confused with distributed computing, where different computers share their processing power to solve very complex calculations.

### H

- **Heterogeneity:** the condition in which in a federated learning scenario, parties are allowed to have differences in hardware and software (i.e., operating systems).
- **Horizontally-partitioned data:** data spread across different parties where the latter have the same features of different instances (i.e., patients). See also vertically-partitioned data.

### N

- **Node:** vantage6 node application that runs at a **Data Station** which has access to the local data.

### M

- **Multi-party computation:** an approach to perform analyses across different parties by performing operations on encrypted data.

### P

- **Party:** an entity that takes part in one (or more) collaborations
- **Python:** a high-level general purpose programming language. It aims to help programmers write clear, logical code. vantage6 is [written in Python](#).

### S

- **Secure multi-party computation:** see *Multi-party computation*



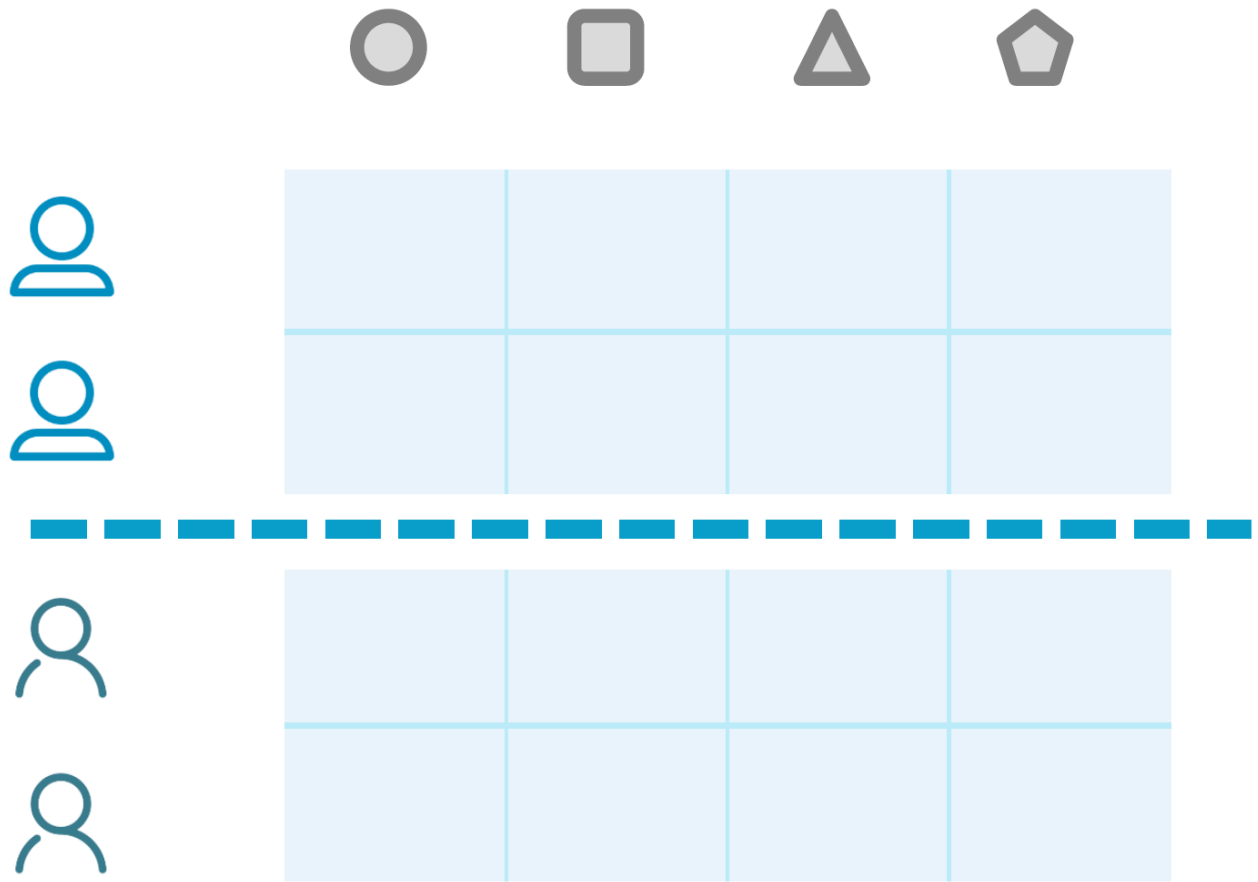


Fig. 4.8: Horizontally-partitioned data

- **Server:** Public access point of the vantage6 infrastructure. Contains at least the **vantage6-server** application but can also host the optional components: Docker registry, VPN server and RabbitMQ. In this documentation space we try to be explicit when we talk about `_server_` and `_vantage6-server_`, however you might encounter `_server_` where `_vantage6-server_` should have been.

## V

- **vantage6:** priVAcy preserviNg federaTed leArninG infrastruCTurE for Secure Insight eXchange. In short, vantage6 is an infrastructure for executing federated learning analyses. However, it can also be used as a FAIR data station and as a model repository.
- **Vertically-partitioned data:** data spread across different parties where the latter have different features of the same instances (i.e., patients). See also horizontally-partitioned data.

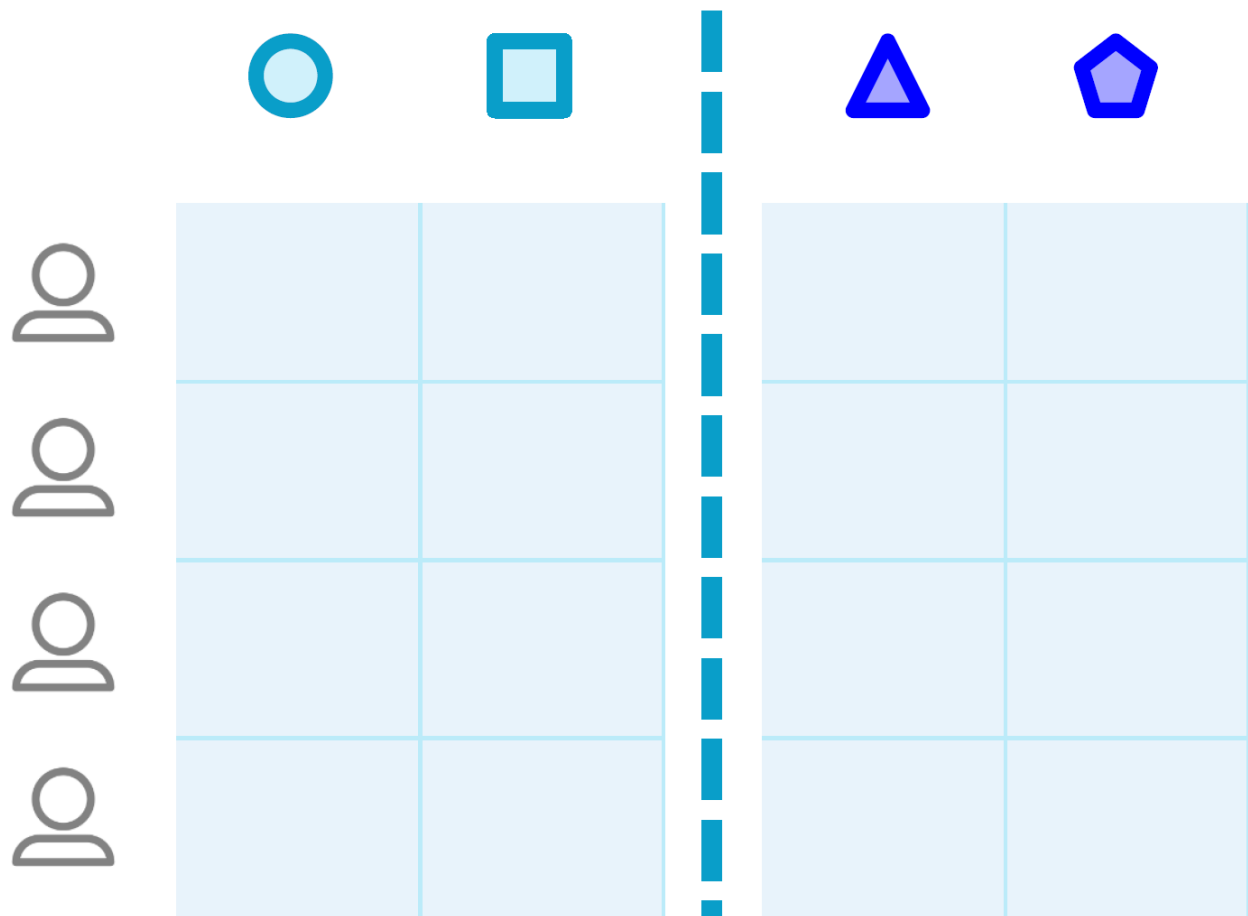


Fig. 4.9: Vertically partitioned data

## 4.7 Release notes

### 4.7.1 3.6.1

12 january 2023

- **Feature**

- Algorithm containers can be killed from the client. This can be done for a specific task or it possible to kill all tasks running at a specific node. (PR#417, Issue#167)
- Added a status field for an algorithm, that tracks if an algorithm has yet to start, is started, has finished, or has failed. In the latter case, it also indicates how/when the algorithm failed (PR#417).
- The UI has been connected to the socket, and gives messages about node and task status changes (UI PR#84, UI Issue #73). There are also new permissions for socket events on the server to authorize users to see events from their (or all) collaborations (PR#417).
- It is now possible to create tasks in the UI (UI version >3.6.0). Note that all tasks are then JSON serialized and you will not be able to run tasks in an encrypted collaboration (as that would require uploading a private key to a browser) (PR#90).

**Warning:** If you want to run the UI Docker image, note that from this version onwards, you have to define the SERVER\_URL and API\_PATH environment variables (compared to just a API\_URL before).

- There is a new multi-database wrapper that will forward a dictionary of all node databases and their paths to the algorithm. This allows you to use multiple databases in a single algorithm easily. (PR#424, Issue #398).
- New rules are now assigned automatically to the default root role. This ensures that rules that are added in a new version are assigned to system administrators, instead of them having to change the database (PR#456, Issue #442).
- There is a new command `vnode set-api-key` that facilitates putting your API key into the node configuration file (PR#428, Issue #259).
- Logging in the Python client has been improved: instead of all or nothing, log level is now settable to one of debug, info, warn, error, critical (PR#453, Issue #340).
- When there is an error in the VPN server configuration, the user receives clearer error messages than before (PR#444, Issue #278).

- **Change**

- The node status (online/offline) is now checked periodically over the socket connection via a ping/pong construction. This is an improvement over the older version where a node's status was changed only when it connected or disconnected (PR#450, Issue #40).

**Warning:** If a server upgrades to 3.6.1, the nodes should also be upgraded. Otherwise, the node status will be incorrect and the logs will show errors periodically with each attempted ping/pong.

- It is no longer possible for any user to change the username of another user, as this would be confusing for that user when logging in (PR#433, Issue #396).
- The server has shorter log messages when someone calls a non-existing route. The resulting 404 exception is no longer logged (PR#452, Issue #393).
- Removed old, unused scripts to start a node (PR#464).

- **Bugfix**

- Node was unable to pull images from Docker Hub; this has been corrected. ([PR#432](#), [Issue#422](#)).
- File-based database extensions were always converted to .csv when they were mounted to a node. Now, files keep their original file extensions ([PR#426](#), [Issue #397](#)).
- When a node configuration defined a wrong VPN subnet, the VPN connection didn't work but this was not detected until VPN was used. Now, the user is alerted immediately and VPN is turned off ([PR#444](#)).
- If a user tries to write a node or server config file to a non-existing directory, they are now getting a clear error message instead of an incorrect one ([PR#455](#), [Issue #1](#)).
- There was a circular import in the infrastructure code, which has now been resolved ([PR#451](#), [Issue #53](#)).
- In PATCH /user, it was not possible to set some fields (e.g. `firstname`) to an empty string if there was a value before. ([PR#439](#), [Issue #334](#)).

---

**Note:** Release 3.6.0 was skipped due to an issue in the release process.

---

## 4.7.2 3.5.2

30 november 2022

- **Bugfix**

- Fix for automatic addition of column. This failed in some SQL dialects because reserved keywords (i.e. 'user' for PostgreSQL) were not escaped ([PR#415](#))
- Correct installation order for uWSGI in node and server docker file ([PR#414](#))

## 4.7.3 3.5.1

30 november 2022

- **Bugfix**

- Backwards compatibility for which organization initiated a task between v3.0-3.4 and v3.5 ([PR#412](#))
- Fixed VPN client container. Entry script was not executable in Github pipelines ([PR#413](#))

## 4.7.4 3.5.0

30 november 2022

**Warning:** When upgrading to 3.5.0, you might need to add the `otp_secret` column to the `user` table manually in the database. This may be avoided by upgrading to 3.5.2.

- **Feature**

- Multi-factor authentication via TOTP has been added. Admins can enforce that all users enable MFA ([PR#376](#), [Issue#355](#)).
- You can now request all tasks assigned by a given user ([PR#326](#), [Issue#43](#)).

- The server support email is now settable in the configuration file, used to be fixed at `support@vantage6.ai` (PR#330, Issue#319).
- When pickles are used, more task info is shown in the node logs (PR#366, Issue#171).
- **Change**
- The `harbor2.vantag6.ai/infrastructure/algorithm-base:[TAG]` is tagged with the `vantage6-client` version that is already in the image (PR#389, Issue#233).
- The infrastructure base image has been updated to improve build time (PR#406, Issue#250).

### 4.7.5 3.4.2

3 november 2022

- **Bugfix**
- Fixed a bug in the local proxy server which made algorithm containers crash in case the `client.create_new_task` method was used (PR#382).
- Fixed a bug where the node crashed when a non existing image was sent in a task (PR#375).

### 4.7.6 3.4.0 & 3.4.1

25 oktober 2022

- **Feature**
- Add columns to the SQL database on startup (PR#365, ISSUE#364). This simplifies the upgrading proces when a new column is added in the new release, as you do no longer need to manually add columns. When downgrading the columns will **not** be deleted.
- Docker wrapper for Parquet files (PR#361, ISSUE#337). Parquet provides a way to store tabular data with the datatypes included which is an advantage over CSV.
- When the node starts, or when the client is verbose initialized a banner to cite the vantage6 project is added (PR#359, ISSUE#356).
- In the client a waiting for results method is added (PR#325, ISSUE#8). Which allows you to automatically poll for results by using `client.wait_for_results(...)`, for more info see `help(client.wait_for_results)`.
- Added Github releases (PR#358, ISSUE#357).
- Added option to filter GET /role by user id in the Python client (PR#328, ISSUE#213). E.g.: `client.role.list(user=...)`.
- In release process, build and release images for both ARM and x86 architecture.
- **Change**
- Unused code removed from the Makefile (PR#324, ISSUE#284).
- Pandas version is frozen to version 1.3.5 (PR#363, ISSUE#266).
- **Bugfix**
- Improve checks for non-existing resources in unittests (PR#320, ISSUE#265). Flask did not support negative ints, so the tests passed due to another 404 response.
- `client.node.list` does no longer filter by offline nodes (PR#321, ISSUE#279).

---

**Note:** 3.4.1 is a rebuild from 3.4.0 in which the all dependencies are fixed, as the build led to a broken server image.

---

### 4.7.7 3.3.7

- **Bugfix**
- The function `client.util.change_my_password()` was updated ([Issue #333](#))

### 4.7.8 3.3.6

- **Bugfix**
- Temporary fix for a bug that prevents the master container from creating tasks in an encrypted collaboration. This temporary fix disables the parallel encryption module in the local proxy. This functionality will be restored in a future release.

---

**Note:** This version is also the first version where the User Interface is available in the right version. From this point onwards, the user interface changes will also be part of the release notes.

---

### 4.7.9 3.3.5

- **Feature**
- The release pipeline has been expanded to automatically push new Docker images of node/server to the harbor2 service.
- **Bugfix**
- The VPN IP address for a node was not saved by the server using the PATCH `/node` endpoint, while this functionality is required to use the VPN

---

**Note:** Note that 3.3.4 was only released on PyPi and that version is identical to 3.3.5. That version was otherwise skipped due to a temporary mistake in the release pipeline.

---

### 4.7.10 3.3.3

- **Bugfix**
- Token refresh was broken for both users and nodes. ([Issue#306](#), [PR#307](#))
- Local proxy encryption was broken. This prevented algorithms from creating sub tasks when encryption was enabled. ([Issue#305](#), [PR#308](#))

### 4.7.11 3.3.2

- **Bugfix**
- `vpn_client_image` and `network_config_image` are settable through the node configuration file. ([PR#301](#), [Issue#294](#))
- The option `--all` from `vnode stop` did not stop the node gracefully. This has been fixed. It is possible to force the nodes to quit by using the `--force` flag. ([PR#300](#), [Issue#298](#))
- Nodes using a slow internet connection (high ping) had issues with connecting to the websocket channel. ([PR#299](#), [Issue#297](#))

### 4.7.12 3.3.1

- **Bugfix**
- Fixed faulty error status codes from the `/collaboration` endpoint ([PR#287](#)).
- *Default* roles are always returned from the `/role` endpoint. This fixes the error when a user was assigned a *default* role but could not reach anything (as it could not view its own role) ([PR#286](#)).
- Performance upgrade in the `/organization` endpoint. This caused long delays when retrieving organization information when the organization has many tasks ([PR#288](#)).
- Organization admins are no longer allowed to create and delete nodes as these should be managed at collaboration level. Therefore, the collaboration admin rules have been extended to include create and delete nodes rules ([PR#289](#)).
- Fixed some issues that made 3.3.0 incompatible with 3.3.1 ([Issue#285](#)).

### 4.7.13 3.3.0

- **Feature**
- Login requirements have been updated. Passwords are now required to have sufficient complexity (8+ characters, and at least 1 uppercase, 1 lowercase, 1 digit, 1 special character). Also, after 5 failed login attempts, a user account is blocked for 15 minutes (these defaults can be changed in a server config file).
- Added endpoint `/password/change` to allow users to change their password using their current password as authentication. It is no longer possible to change passwords via `client.user.update()` or via a `PATCH /user/{id}` request.
- Added the default roles ‘viewer’, ‘researcher’, ‘organization admin’ and ‘collaboration admin’ to newly created servers. These roles may be assigned to users of any organization, and should help users with proper permission assignment.
- Added option to filter get all roles for a specific user id in the `GET /role` endpoint.
- RabbitMQ has support for multiple servers when using `vserver start`. It already had support for multiple servers when deploying via a Docker compose file.
- When exiting server logs or node logs with `Ctrl+C`, there is now an additional message alerting the user that the server/node is still running in the background and how they may stop them.
- **Change**
- Node proxy server has been updated
- Updated PyJWT and related dependencies for improved JWT security.

- When nodes are trying to use a wrong API key to authenticate, they now receive a clear message in the node logs and the node exits immediately.
- When using `vserver import`, API keys must now be provided for the nodes you create.
- Moved all swagger API docs from YAML files into the code. Also, corrected errors in them.
- API keys are created with UUID4 instead of UUID1. This prevents that UUIDs created milliseconds apart are not too similar.
- Rules for users to edit tasks were never used and have therefore been deleted.
- **Bugfix**
- In the Python client, `client.organization.list()` now shows pagination metadata by default, which is consistent all other `list()` statements.
- When not providing an API key in `vnode new`, there used to be an unclear error message. Now, we allow specifying an API key later and provide a clearer error message for any other keys with inadequate values.
- It is now possible to provide a name when creating a name, both via the Python client as via the server.
- A GET `/role` request crashed if parameter `organization_id` was defined but not `include_root`. This has been resolved.
- Users received an ‘unexpected error’ when performing a GET `/collaboration?organization_id=<id>` request and they didn’t have global collaboration view permission. This was fixed.
- GET `/role/<id>` didn’t give an error if a role didn’t exist. Now it does.

#### 4.7.14 3.2.0

- **Feature**
- Horizontal scaling for the vantage6-server instance by adding support for RabbitMQ.
- It is now possible to connect other docker containers to the private algorithm network. This enables you to attach services to the algorithm network using the `docker_services` setting.
- Many additional select and filter options on API endpoints, see swagger docs endpoint (`/apidocs`). The new options have also been added to the Python client.
- Users are now always able to view their own data
- Usernames can be changed through the API
- **Bugfix**
- (Confusing) SQL errors are no longer returned from the API.
- Clearer error message when an organization has multiple nodes for a single collaboration.
- Node no longer tries to connect to the VPN if it has no `vpn_subnet` setting in its configuration file.
- Fix the VPN configuration file renewal
- Superusers are no longer able to post tasks to collaborations its organization does not participate in. Note that superusers were never able to view the results of such tasks.
- It is no longer possible to post tasks to organization which do not have a registered node attach to the collaboration.
- The `vnode create-private-key` command no longer crashes if the `ssh` directory does not exist.
- The client no longer logs the password



- The version of the `alpine` docker image (that is used to set up algorithm runs with VPN) was fixed. This prevents that many versions of this image are downloaded by the node.
- Improved reading of username and password from docker registry, which can be capitalized differently depending on the docker version.
- Fix error with multiple-database feature, where default is now used if specific database is not found

#### 4.7.15 3.1.0

- **Feature**
  - Algorithm-to-algorithm communication can now take place over multiple ports, which the algorithm developer can specify in the Dockerfile. Labels can be assigned to each port, facilitating communication over multiple channels.
  - Multi-database support for nodes. It is now also possible to assign multiple data sources to a single node in Petronas; this was already available in Harukas 2.2.0. The user can request a specific data source by supplying the `database` argument when creating a task.
  - The CLI commands `vserver new` and `vnode new` have been extended to facilitate configuration of the VPN server.
  - Filter options for the client have been extended.
  - Roles can no longer be used across organizations (except for roles in the default organization)
  - Added `vnode remove` command to uninstall a node. The command removes the resources attached to a node installation (configuration files, log files, docker volumes etc).
  - Added option to specify configuration file path when running `vnode create-private-key`.
- **Bugfix**
  - Fixed swagger docs
  - Improved error message if docker is not running when a node is started
  - Improved error message for `vserver version` and `vnode version` if no servers or nodes are running
  - Patching user failed if users had zero roles - this has been fixed.
  - Creating roles was not possible for a user who had permission to create roles only for their own organization - this has been corrected.

#### 4.7.16 3.0.0

- **Feature**
  - Direct algorithm-to-algorithm communication has been added. Via a VPN connection, algorithms can exchange information with one another.
  - Pagination is added. Metadata is provided in the headers by default. It is also possible to include them in the output body by supplying an additional parameter `include=metadata`. Parameters `page` and `per_page` can be used to paginate. The following endpoints are enabled:
    - `GET /result`
    - `GET /collaboration`
    - `GET /collaboration/{id}/organization`
    - `GET /collaboration/{id}/node`

- GET /collaboration/{id}/task
- GET /organization
- GET /role
- GET /role/{id}/rule
- GET /rule
- GET /task
- GET /task/{id}/result
- GET /node

- API keys are encrypted in the database
- Users cannot shrink their own permissions by accident
- Give node permission to update public key
- Dependency updates
- **Bugfix**
- Fixed database connection issues
- Don't allow users to be assigned to non-existing organizations by root
- Fix node status when node is stopped and immediately started up
- Check if node names are allowed docker names

#### 4.7.17 2.3.0 - 2.3.4

- **Feature**
- Allows for horizontal scaling of the server instance by adding support for RabbitMQ. Note that this has not been released for version 3(!)
- **Bugfix**
- Performance improvements on the /organization endpoint

#### 4.7.18 2.2.0

- **Feature**
- Multi-database support for nodes. It is now possible to assign multiple data sources to a single node. The user can request a specific data source by supplying the *database* argument when creating a task.
- The mailserver now supports TLS and SSL options
- **Bugfix**
- Nodes are now disconnected more gracefully. This fixes the issue that nodes appear offline while they are in fact online
- Fixed a bug that prevented deleting a node from the collaboration
- A role is now allowed to have zero rules
- Some http error messages have improved
- Organization fields can now be set to an empty string

#### 4.7.19 2.1.2 & 2.1.3

- **Bugfix**
- Changes to the way the application interacts with the database. Solves the issue of unexpected disconnects from the DB and thereby freezing the application.

#### 4.7.20 2.1.1

- **Bugfix**
- Updating the country field in an organization works again\
- The `client.result.list(...)` broke when it was not able to deserialize one of the in- or outputs.

#### 4.7.21 2.1.0

- **Feature**
- Custom algorithm environment variables can be set using the `algorithm_env` key in the configuration file. [See this Github issue.](#)
- Support for non-file-based databases on the node. [See this Github issue.](#)
- Added flag `--attach` to the `vserver start` and `vnode start` command. This directly attaches the log to the console.
- Auto updating the node and server instance is now limited to the major version. [See this Github issue.](#)
  - e.g. if you've installed the Trolltunga version of the CLI you will always get the Trolltunga version of the node and server.
  - Infrastructure images are now tagged using their version major. (e.g. `trolltunga` or `harukas` )
  - It is still possible to use intermediate versions by specifying the `--image` option when starting the node or server. (e.g. `vserver start --image harbor.vantage6.ai/infrastructure/server:2.0.0.post1` )
- **Bugfix**
- Fixed issue where node crashed if the database did not exist on startup. [See this Github issue.](#)

#### 4.7.22 2.0.0.post1

- **Bugfix**
- Fixed a bug that prevented the usage of secured registry algorithms

### 4.7.23 2.0.0

- **Feature**

- Role/rule based access control
  - Roles consist of a bundle of rules. Rules provide access to certain API endpoints at the server.
  - By default 3 roles are created: 1) Container, 2) Node, 3) Root. The root role is assigned to the root user on the first run. The root user can assign rules and roles from there.
- Major update on the *python*-client. The client also contains management tools for the server (i.e. to creating users, organizations and managing permissions. The client can be imported from `from vantage6.client import Client` .
- You can use the argument `verbose` on the client to output status messages. This is useful for example when working with Jupyter notebooks.
- Added CLI `vserver version` , `vnode version` , `vserver-local version` and `vnode-local version` commands to report the version of the node or server they are running
- The logging contains more information about the current setup, and refers to this documentation and our Discourse channel

- **Bugfix**

- Issue with the DB connection. Session management is updated. Error still occurs from time to time but can be reset by using the endpoint `/health/fix` . This will be patched in a newer version.

### 4.7.24 1.2.3

- **Feature**

- The node is now compatible with the Harbor v2.0 API

### 4.7.25 1.2.2

- **Bug fixes**

- Fixed a bug that ignored the `--system` flag from `vnode start`
- Logging output muted when the `--config` option is used in `vnode start`
- Fixed config folder mounting point when the option `--config` option is used in `vnode start`

### 4.7.26 1.2.1

- **Bug fixes**

- starting the server for the first time resulted in a crash as the root user was not supplied with an email address.
- Algorithm containers could still access the internet through their host. This has been patched.

### 4.7.27 1.2.0

- **Features**

- Cross language serialization. Enabling algorithm developers to write algorithms that are not language dependent.
- Reset password is added to the API. For this purpose two endpoints have been added: `/recover/lostandrecover/reset` . The server config file needs to be extended to be connected to a mail-server in order to make this work.
- User table in the database is extended to contain an email address which is mandatory.

- **Bug fixes**

- Collaboration name needs to be unique
- API consistency and bug fixes:
  - GET `/organization` was missing domain key
  - PATCH `/organization` could not patch domain
  - GET `/collaboration/{id}/node` has been made consistent with `/node`
  - GET `/collaboration/{id}/organization` has been made consistent with `/organization`
  - PATCH `/user root-user` was not able to update users
  - DELETE `/user root-user` was not able to delete users
  - GET `/task` null values are now consistent: `[]` is replaced by `null`
  - POST, PATCH, DELETE `/node root-user` was not able to perform these actions
  - GET `/node/{id}/task` output is made consistent with the

- **other**

- questionnaire dependency is updated to 1.5.2
- vantage6-toolkit repository has been merged with the vantage6-client as they were very tight coupled.

### 4.7.28 1.1.0

- **Features**

- new command `vnode clean` to clean up temporary docker volumes that are no longer used
- Version of the individual packages are printed in the console on startup
- Custom task and log directories can be set in the configuration file
- Improved **CLI** messages
- Docker images are only pulled if the remote version is newer. This applies both to the node/server image and the algorithm images
- Client class names have been simplified (`UserClientProtocol` -> `Client`)

- **Bug fixes**

- Removed defective websocket watchdog. There still might be disconnection issues from time to time.

### 4.7.29 1.0.0

- **Updated Command Line Interface (CLI)**

- The commands `vnode list`, `vnode start` and the new command `vnode attach` are aimed to work with multiple nodes at a single machine.
- System and user-directories can be used to store configurations by using the `--user/--system` options. The node stores them by default at user level, and the server at system level.
- Current status (online/offline) of the nodes can be seen using `vnode list`, which also reports which environments are available per configuration.
- Developer container has been added which can inject the container with the source. `vnode start --develop [source]`. Note that this Docker image needs to be build in advance from the `development.Dockerfile` and tag `devcon`.
- `vnode config_file` has been replaced by `vnode files` which not only outputs the config file location but also the database and log file location.

- **New database model**

- Improved relations between models, and with that, an update of the Python API.
- Input for the tasks is now stored in the result table. This was required as the input is encrypted individually for each organization (end-to-end encryption (E2EE) between organizations).
- The `Organization` model has been extended with the `public_key` (String) field. This field contains the public key from each organization, which is used by the E2EE module.
- The `Collaboration` model has been extended with the `encrypted` (Boolean) field which keeps track if all messages (tasks, results) need to be E2EE for this specific collaboration.
- The `Task` keeps track of the initiator (organization) of the organization. This is required to encrypt the results for the initiator.

- **End to end encryption**

- All messages between all organizations are by default be encrypted.
- Each node requires the private key of the organization as it needs to be able to decrypt incoming messages. The private key should be specified in the configuration file using the `private_key` label.
- In case no private key is specified, the node generates a new key and uploads the public key to the server.
- If a node starts (using `vnode start`), it always checks if the `public_key` on the server matches the private key the node is currently using.
- In case your organization has multiple nodes running they should all point to the same private key.
- Users have to encrypt the input and decrypt the output, which can be simplified by using our client `vantage6.client.Client` for Python or `vtg::Client` for R.
- Algorithms are not concerned about encryption as this is handled at node level.

- **Algorithm container isolation**

- Containers have no longer an internet connection, but are connected to a private docker network.
- Master containers can access the central server through a local proxy server which is both connected to the private docker network as the outside world. This proxy server also takes care of the encryption of the messages from the algorithms for the intended receiving organization.
- In case a single machine hosts multiple nodes, each node is attached to its own private Docker network.

- **Temporary Volumes**

- Each algorithm mounts temporary volume, which is linked to the node and the `run_id` of the task
- The mounting target is specified in an environment variable `TEMPORARY_FOLDER`. The algorithm can write anything to this directory.
- These volumes need to be cleaned manually. (`docker rm VOLUME_NAME`)
- Successive algorithms only have access to the volume if they share the same `run_id`. Each time a **user** creates a task, a new `run_id` is issued. If you need to share information between containers, you need to do this through a master container. If a master container creates a task, all slave tasks will obtain the same `run_id`.
- **RESTful API**
- **All RESTful API output is HATEOS formatted.**  
([wiki](#))
- **Local Proxy Server**
- Algorithm containers no longer receive an internet connection. They can only communicate with the central server through a local proxy service.
- It handles encryption for certain endpoints (i.e. `/task`, the input or `/result` the results)
- **Dockerized the Node**
- All node code is run from a Docker container. Build versions can be found at our Docker repository: `harbor.distributedlearning.ai/infrastructure/node`. Specific version can be pulled using tags.
- For each running node, a Docker volume is created in which the data, input and output is stored. The name of the Docker volume is: `vantage-NODE_NAME-vol`. This volume is shared with all incoming algorithm containers.
- Each node is attached to the public network and a private network: `vantage-NODE_NAME-net`.

## 4.8 Partners

Our community is open to everyone. The following people and organizations made a significant contribution to the design and implementation of vantage6.



- Anja van Gestel
- Bart van Beusekom
- Frank Martin
- Hasan Alradhi
- Melle Sieswerda
- Gijs Geleijnse



- Djura Smits
- Lourens Veen



- Johan van Soest

**Would you like to contribute?** Check out [how to contribute!](#) Find and chat with us via the [Discord](#) chat!