



EOSC-hub

Report about the evaluation of OneData by the SeaDataNet community

Authors:	Jordan Maduro, Dick Schaap (MARIS) Gergely Sipos, Björn Backeberg, Andrea Manzi (EGI Foundation)
Date:	11/Mar/2020
Dissemination Level:	PUBLIC
Document Link:	https://documents.egi.eu/document/3583

Abstract

In this report we provide a summary of the findings of the OneData evaluation activity that was carried out by MARIS, member of the Marine Competence Centre of the EOSC-hub project during 2019. The evaluation was focussed on different data movement, access and caching use cases that are also described in the report.



COPYRIGHT NOTICE

This work by Parties of the EOSC-hub Consortium is licensed under a Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>). The EOSC-hub project is co-funded by the European Union Horizon 2020 programme under grant number 777536.

Contents

1	Introduction	4
1.1	Challenges	4
2	Case 1: Cloud migration for legacy applications.....	5
2.1	Discussion.....	7
2.2	Summary	9
3	Case 2: Reducing redundant data transfer	11
4	Case 3: Virtual space for user data storage and delivery.....	12
5	Case 4: Interface for distributed search using metadata queries	13
6	Learnings/experience.....	16
6.1	Unexplained behavior.....	17

1 Introduction

MARIS is technical coordinator of the SeaDataNet infrastructure and as such currently involved in the EU SeaDataCloud project. MARIS also represents SeaDataNet in the EOSC-hub Marine Competence Center (MCC). In this position MARIS regularly encounters software development challenges related to the distributed nature of the data and the large number of partners involved in the SeaDataNet network. A large portion of the challenges is related to the various aspects of data storage management. Currently, MARIS has multiple applications and systems created internally and by third parties, which make specific platform assumptions. For example, files locally available, synchronous file transfers, proprietary log formats, and log locations. These assumptions are difficult to overcome when these applications do not provide an interface to modify their behavior. A way to overcome these limitations is to develop specialized middleware or wrappers. However, developing these applications are time-consuming and often not reusable.

Onedata offers many advanced features related to data storage management. These features seem promising in solving many of our recurring challenges. By means of an amended EOSC-hub MCC workplan, MARIS aims to explore the opportunities available by leveraging the Onedata platform.

In this report, we examine the implementation of the use-cases proposed in the EOSC-hub MCC workplan.

1.1 Challenges

Each proposed use case illustrates an abstract challenge inspired by concrete cases we currently have or are faced with in multiple projects. The average number of files requested and processed is 500. Each file has a size of tens of KB but occasionally some larger files that average 500 MB are processed. A request is usually for ease of transfer and is usually between 50 to 100 MB. We list the Onedata concepts associated with each case.

2 Case 1: Cloud migration for legacy applications

To simplify the processes of migrating mostly desktop applications to the cloud we can use Onedata as a file storage abstraction layer providing seamless access to files distributed in the cloud environments. For this case, we used Octopus as the main application. We set up the testbed using Onedata version 19.02. We installed the Oneproviders using the onedatify installation script. The simple workflow consists of a number of steps.

Simple workflow

- Without Onedata
 - Download zip
 - Unpack zip
 - Run Octopus
 - Upload log
 - Run conversion
- With Onedata
 - Run Octopus
 - Move log
 - Run conversion

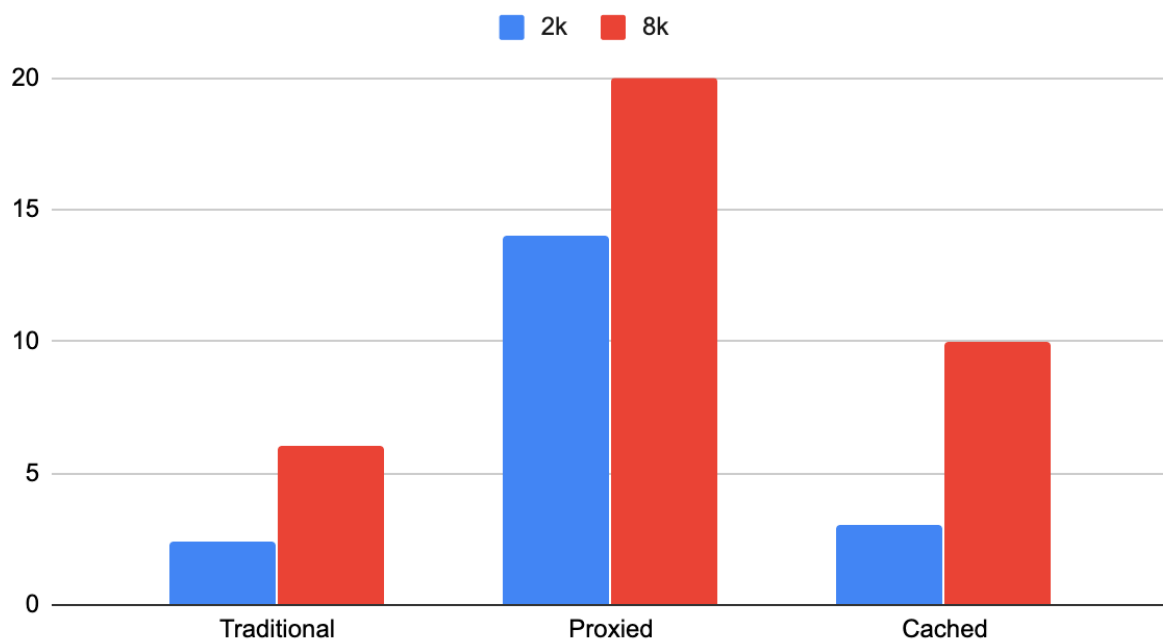
Setup

- **1 OneProvider with datasets**
- **1 OneProvider at computing site**

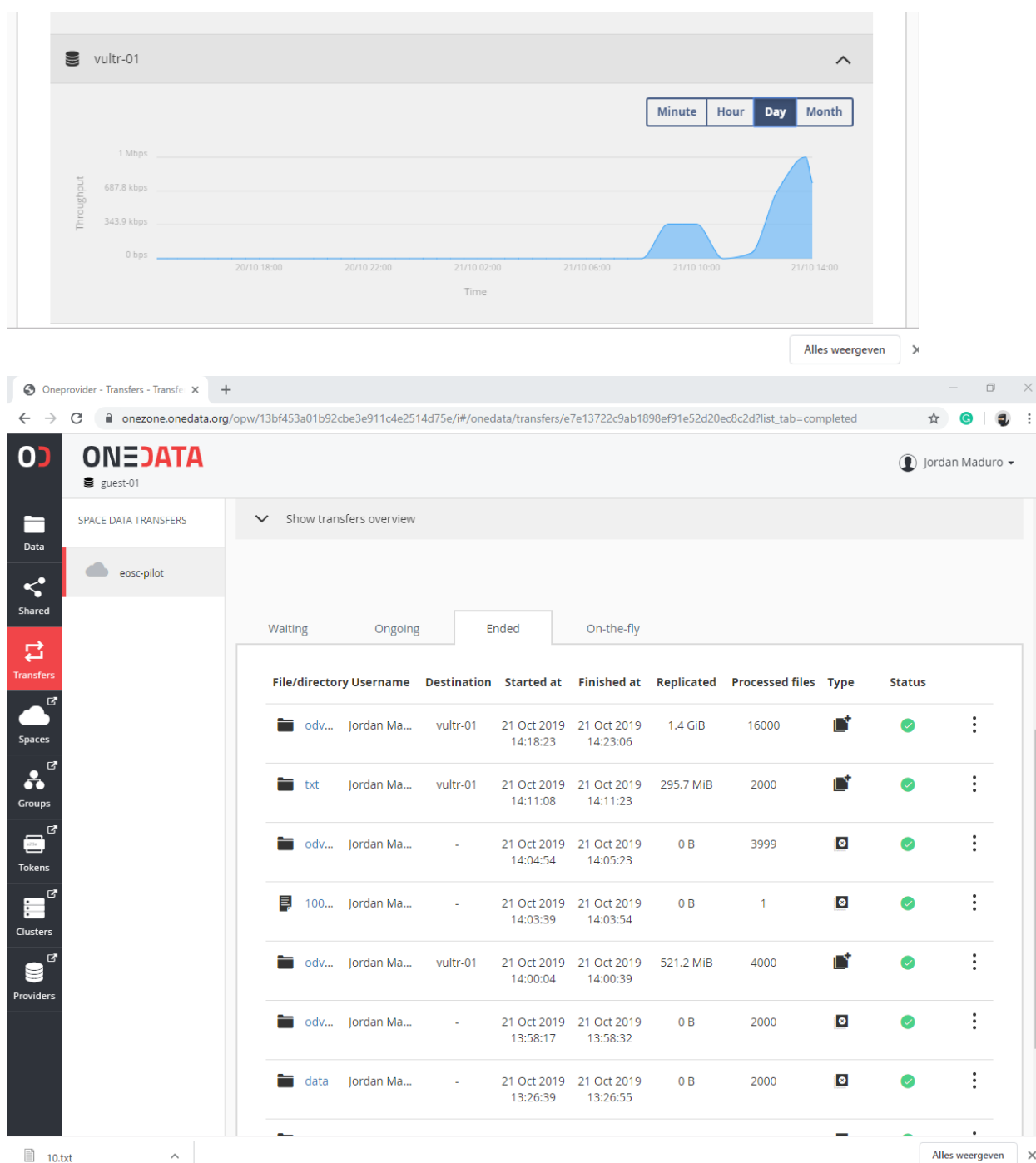
The Oneproviders are hosted in two different countries. The compute instance is hosted in the Netherlands and the storage instance is hosted in Poland. The experiment consists of multiple tests using different sample sizes of files. The workflow is semantically the same. A number of files undergo quality control by first using Octopus in batch mode to check the files. Then the files are converted to NetCDF.

In this case, we ran the workflow multiple times under different conditions. We benchmark the workflows using the time command. The benchmark compares the performance of proxied (uncached remote files), Oneprovider cached files and traditional. In the chart below, the average runtime can be seen in minutes.

2k and 8k files



We experimented with the replication API to fine-tune the optimized setup. By first transferring the files in bulk to the compute instance you reduce the transfer times significantly.



2.1 Discussion

Using Onedata for this use case greatly reduces complexity. Once setup onedata is very transparent and does not require any special modifications to existing programs. However, there are drawbacks. Using the unoptimized setup can be extremely slow depending on the available bandwidth. Every file requires a network call that takes considerably longer time than an I/O. This can be mitigated if we use the replication feature of Onedata. By replicating the dataset before running the workflow

you only add about 15 seconds overhead instead of the 10+ minutes using the proxied network calls. The replication of data before running the workflow saves a considerable amount of time. However, running through a Oneclient still has an overhead. In the case of 'odv2k' a little less than 1 minute.

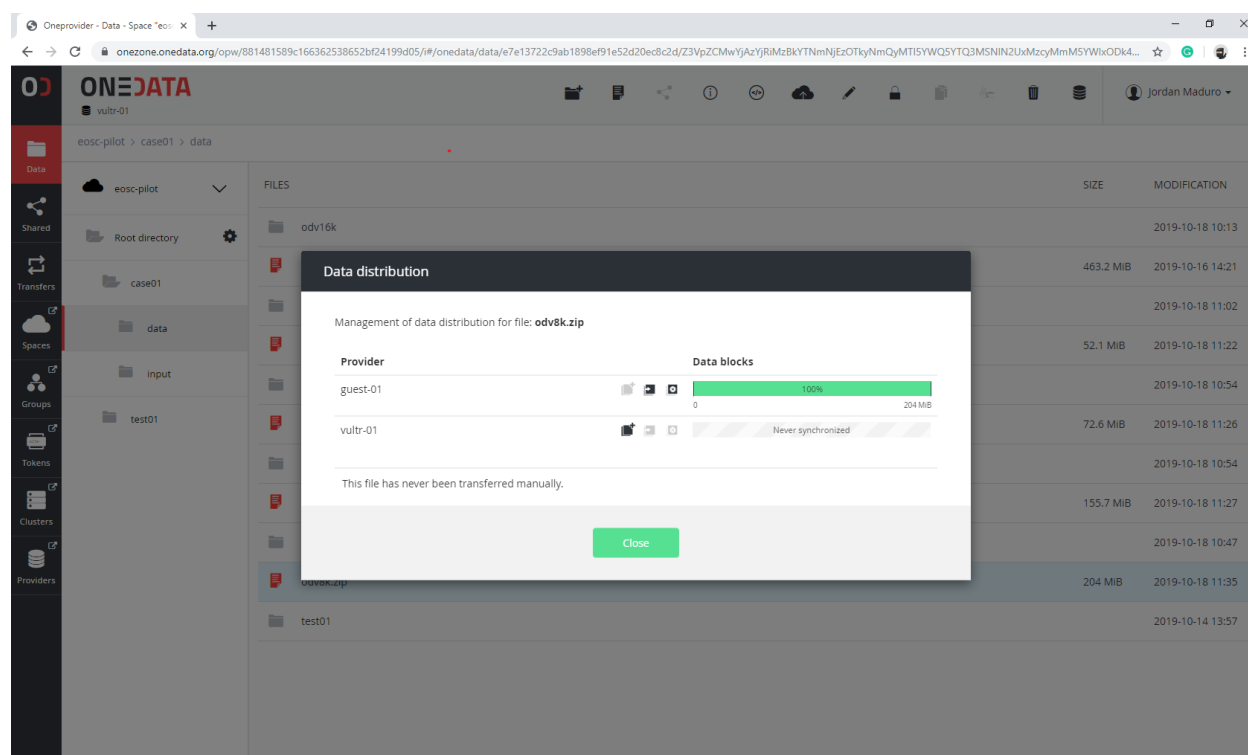


Figure 1 Data distribution before replication

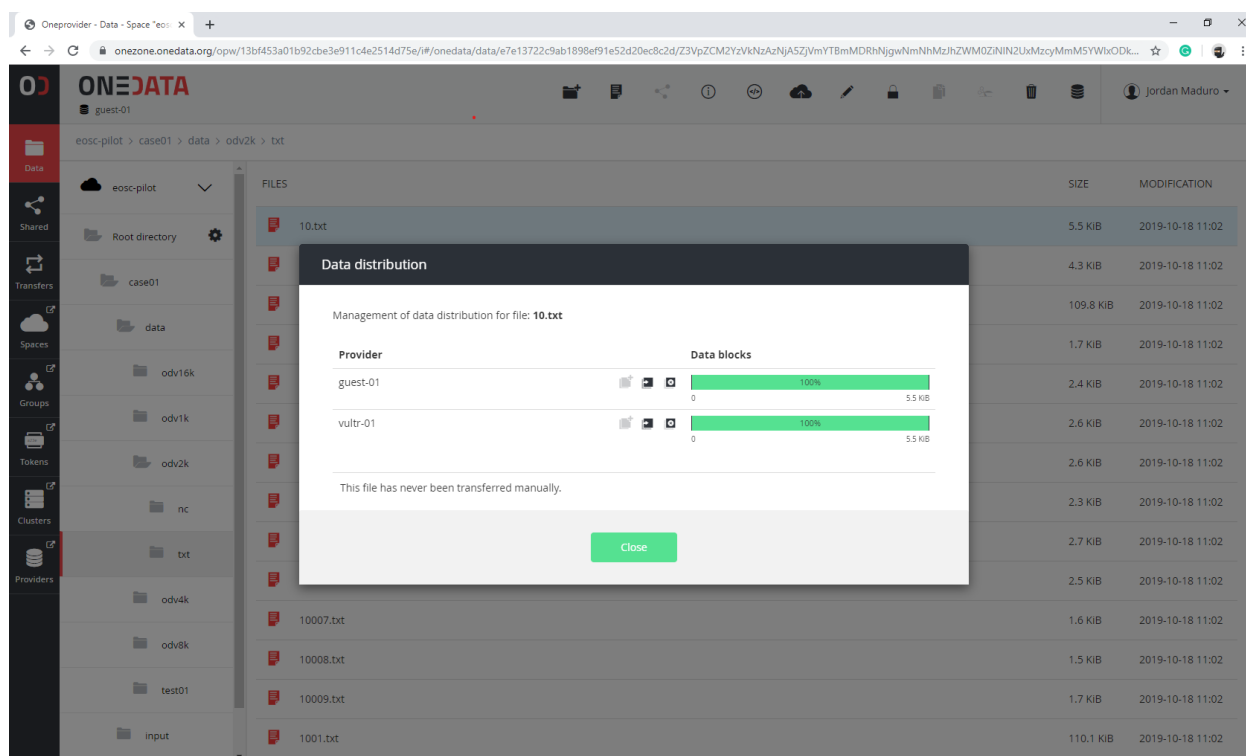


Figure 2 Data distribution after replication

2.2 Summary

Use case 1 is a great candidate for using Oneddata because of the reduced complexity for the user. While Oneddata is very sophisticated, it does not place any more cognitive load on the user. Thinking in terms of “removable storage” is enough to understand the basic concept of Oneddata. We implemented the use case 1 by removing steps in the traditional workflow. The only considerations are if the overhead incurred by Oneddata is acceptable and if the optimized setup can be used.



3 Case 2: Reducing redundant data transfer

Operations on large datasets are not always successful. In the case of quality control, certain datasets may be rejected and have to be revised before being submitted again.

This means that some complete datasets are transferred multiple times before being accepted. Using the direct access provided by Onedata to these files we can process only the required amount of data.

The implementation of this use case did not lead to a working solution. The functionality of viewing a zip file works. However, there is some unexpected behavior that made this use case hard to implement. To illustrate, the zip file is visible using the ``ls`` command. Which shows the files within the space, in this instance, located on a remote server. However, it is not found using the `unzip -l` command which lists the contents of the zip file. The command only works once you cache the file on the local Oneprovider. This defeats the point of using the Oneprovider to reduce redundant data transfer.

Another method that works is extracting the data on the remote Oneprovider. Then you can access each file individually. While it would not require files to be transferred multiple times it still transfers the uncompressed file. In the case of 500mb file zip, you get 3.4GB uncompressed files.

A possible solution is to first transfer the zip file to the compute site. After processing, any update operations should be performed on the remote site and OneProvider should sync the changed blocks. This approach was not attempted.

4 Case 3: Virtual space for user data storage and delivery

After a user searches for specific data he sends a request for a subset of datasets, a process is started to collect the datasets from many partners. This collection is an asynchronous process. The process can take weeks to collect all the requested files from the partners. This process is dependent on the resources available at the partners. We intend to use the features of Space and privileges management provided by Onedata to streamline these processes. We would provide the end-users access to his requested files through a shared space. Ideally, such a space can be used to make his requested files available for further processing in a cloud environment.

The Onedata platform offers the functionality to implement this use case. However, it was not possible to implement this use case due to a lack of resources, timing, and technical issues. The use cases were implemented using the latest version of Onedata. This version was only available at the Onedata.org Onezone. The partners were connected to the EGI Onezone which ran an older incompatible version of Onezone. This meant that the testbed created for the other use cases could not be used for this use case. Furthermore, technical issues arose during the migration of the EGI Onezone to the latest version of Onedata. The Onedata onedatify installation script was broken and therefore no new instances could be installed.

5 Case 4: Interface for distributed search using metadata queries

In order to find specific files in a distributed environment, we use proprietary search indexes. These indexes are inflexible and only allow querying of predetermined fields. This increases the time required to process and index all available datasets. With the current search interface, we process changes daily. However, the use of datasets within workflows would benefit from up to date information on the available datasets. To extend the discovery capabilities of a cloud application we can leverage the advanced metadata querying functionalities of the Onedata platform.

The implementation of the case started with the definition of views. The views are created using a MapReduce function to index specific fields in a space. In this case, we expect to have an array of subjects. Each subject is an array with 3 entries that represent the subject, object and units metadata fields in the ODV files. The function was written in a file and posted to the API using Postman. Once posted, the view would be available at the following URL.

```

Z: > eoschub > JS subjects.js > <function>
1  function(id, type, meta, ctx) {
2      if(type === "custom_metadata"){
3          if(meta['onedata_json']['subjects']) {
4              return {
5                  list: meta['onedata_json']['subjects'].map(function(r) {
6                      return [r, id]
7                  })
8              };
9          }
10     }
11 }
  
```

Figure 5 MapReduce function written in JavaScript

We then created a script to index files within a given folder within a space. This file is written in JavaScript and runs on Node v12. The script loads all files and then extracts the metadata lines, transforms and parses them to JSON. The JSON contains the 3 fields subject, object, and units. The JSON is then posted to the Oneprovider using the API.

After running the script we can query the API to see all the files that contain one or more combinations of the parameters. We first use Postman to query the view.

```

https://{HOST}/api/v3/oneprovider/spaces/e7e13722c9ab1898ef91e52d20ec8c2d/indexes/subjects/query?keys=[
  "SDN:LOCAL:AMON_7",

```

```

"SDN:P01::AMONAAZX",
"SDN:P06::UPOX"
],[
"SDN:LOCAL:ASLV_52",
"SDN:P01::ASLVZZ01",
"SDN:P06::ULAA"
]]

```

Query URL

Which results in the following output. A list of JSON objects which contain a “key”, “value”, and “id” property. The most important property is the “value” property because this is the file id.

The file id can then be used to retrieve the file or the complete metadata.

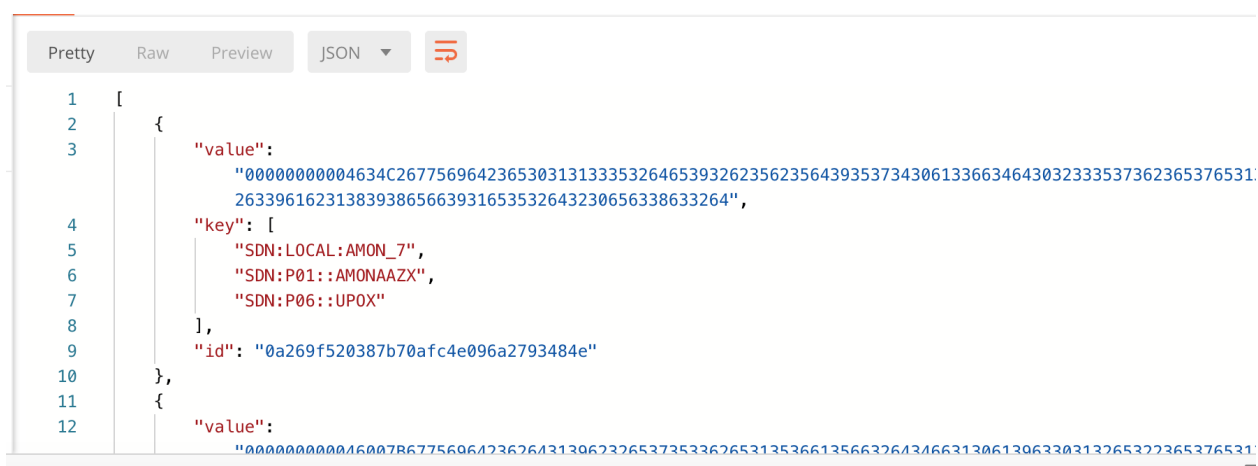


Figure 6 Pretty formatted search results

```

Pretty  Raw  Preview
[{"value":"00000000004634C2677569642365303131333532646539326235623564393537343061336634643032333537362365376531",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"0a269f520387b70afc4e096a2793484e"},
{"value":"000000000046007B677569642362643139623265373533626531353661356632643466313061396330313265322365376531",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"16a839aebcea9fb2d2d1c9a6cf1683b6"},
{"value":"0000000000466FB867756964233763373665353564376265666263363386430396338646163333835353561392365376531",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"25d586f91ef09a5130d7fd5ec7931cf9"},
{"value":"000000000046D7A9677569642338393766646431626333343936303062373139313833386133316266663464342365376531",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"2ceb54cfed366f00d44098ba2b975ece"},
{"value":"000000000046A1F6677569642334626435626131666266356662363637396330373739623335666165376563652365376531",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"496083746c5481f4b72c8a3ad6b880d6"},
{"value":"000000000046B14B677569642334303866343536326430343365376365613839356564333963313030373239362365376531",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"519a89ddce8deb7aafda9d186d9c4267"},
{"value":"00000000004657116775696423663566646533663061663738636234306437363932373133306438383937643523653765313",
["SDN:LOCAL:AMON_7","SDN:P01::AMONAAZX","SDN:P06::UPOX"],"id":"6385c1f4db7d49c9de15128d6e9c5cae"},

```

Figure 7 Raw search results

This use case depends largely on how the metadata is structured and added to each file. The primary method of annotation would be during the import of a file. Unlike the traditional method of having a central database with the metadata. The data providers can annotate their own files with metadata which can then be queried using the views.

6 Learnings/experience

The Onedata platform allows you to build workflows that eliminate explicit data transfer actions. Using it in the exact same way as the traditional workflows do not automatically provide a performance boost. It actually incurs a slight overhead. It does, however, reduce the complexity and the points of failure within a workflow. This helps in creating, maintaining, and understanding workflows.

Implicit cross space transfers can be very slow. Having your application request a file that is located somewhere else causes that file to first be transferred to your local Oneprovider. This implicit transfer is a great mechanism that reduces the complexity of a program. However, when your application requests thousands of files this can become a bottleneck. In this case, using the transfer API can reduce the amount of time spent transferring files by doing it in bulk.

Onedata is not geared towards handling files on a low level. It is intended to be an easy-access storage. The data should already be processed and organized in the desired way before uploading it to onedata. Otherwise, it can be a slow and tedious process. One way to bypass this limitation is by turning on the file system sync. With sync turned on and having direct access to the files' initial location will allow performing fs operations at native speed.

The Onedata platform is well documented and offers multiple tools to install and manage your Onedata ecosystem. Once installed the systems work seamlessly and have a lot of functionality out of the box. The only issue is the complexity of the installation. If the onedatify script is broken manual installation would be a challenge for even the most experienced users.

Optimized and Unoptimized setup

The unoptimized setup uses the Oneclient to connect to the space. This means that even though the files are available to the system they first have to be transferred before processing can begin. Subsequent reads suffer the same latency issue. This setup is only recommended for a very low number of files or prototypes.

The optimized setup is the same as the unoptimized version however the one difference is that there is Oneprovider installed on the system and connected to the space. The Oneprovider caches the files locally allowing subsequent reads to be at near-native speeds.

Vultr VM (Compute Instance)

CPU: 6 vCore

RAM: 16384 MB

Storage: 320 GB SSD

6.1 Unexplained behavior

The server suddenly showed an increase in CPU usage. This change happened without any traffic or operation and it was holding steady at about 80%. Using `top` we could see that the process `beam.smp` is the culprit. This behavior was monitored for two weeks and resulted in over 100% CPU average usage (162% to be precise). This irregularity was reported but not investigated further during the implementation of the cases.

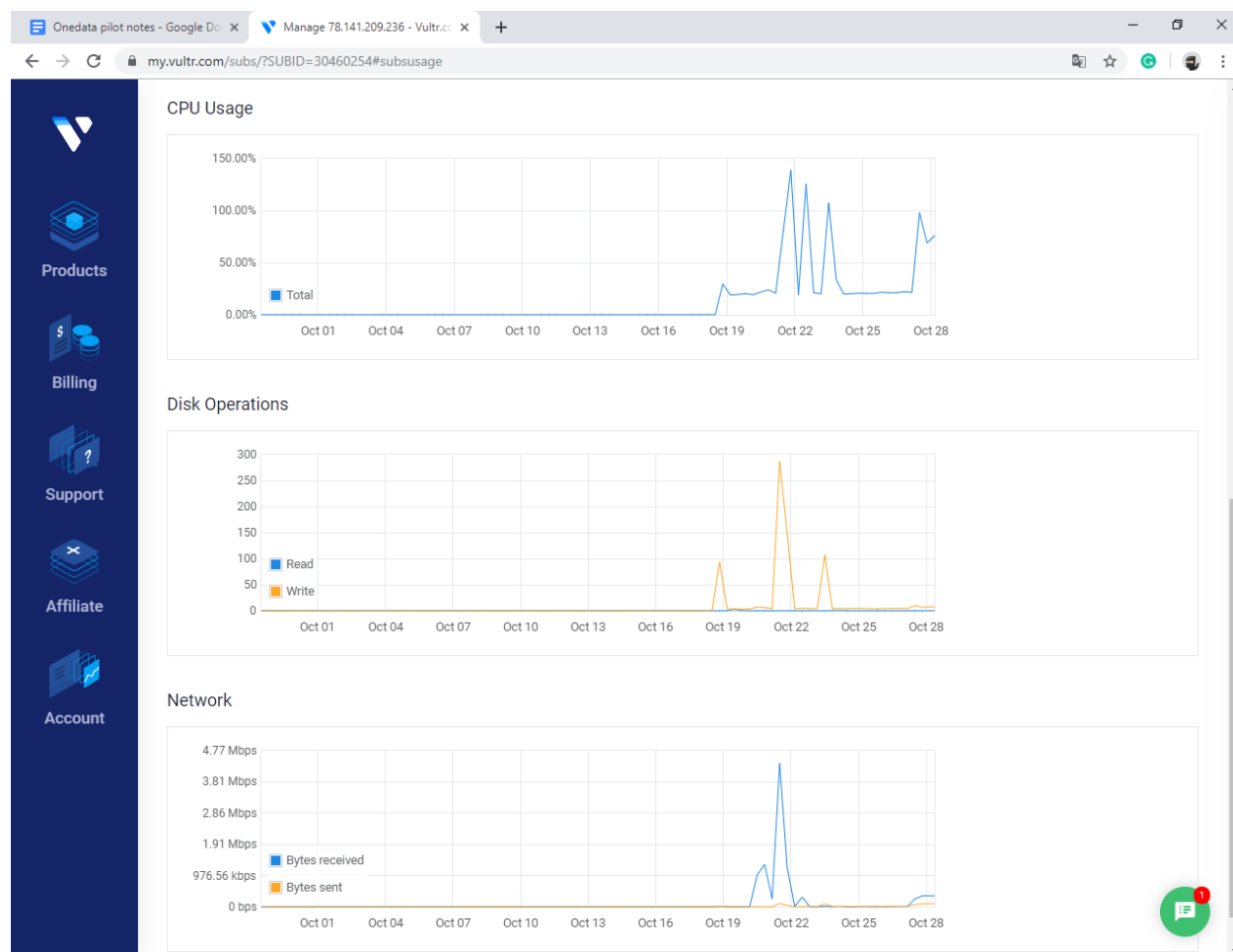


Figure 8 Unexplained CPU usage spikes

```

root@vultr:~# top
top - 09:29:25 up 9 days, 16:37, 2 users, load average: 0.68, 0.81, 0.84
Tasks: 183 total, 1 running, 124 sleeping, 1 stopped, 0 zombie
%Cpu(s): 8.5 us, 1.0 sy, 0.0 ni, 90.4 id, 0.0 wa, 0.0 hi, 0.1 si, 0.0 st
KiB Mem : 16425152 total, 641044 free, 4565896 used, 11218212 buff/cache
KiB Swap: 0 total, 0 free, 0 used, 11551828 avail Mem

```

PID	USER	PR	NI	VIRT	RES	SHR	S	%CPU	%MEM	TIME+	COMMAND
9103	root	20	0	8358604	2.704g	39004	S	54.2	17.3	1189:46	beam.smp
7393	999	20	0	1648668	261100	5868	S	8.6	1.6	1042:02	beam.smp
8249	152	20	0	3318624	38244	5104	S	1.7	0.2	32:06.30	beam.smp
7444	999	20	0	1614496	85648	19528	S	1.3	0.5	150:59.86	beam.smp
7537	999	20	0	788868	425128	13708	S	0.7	2.6	167:22.89	memcached
854	root	20	0	72296	6188	5440	S	0.3	0.0	0:20.90	sshd
3589	root	20	0	1498760	46204	24632	S	0.3	0.3	91:15.95	containerd
6878	root	20	0	3385080	76644	8640	S	0.3	0.5	11:15.48	beam.smp
7536	999	20	0	5848	2504	1896	S	0.3	0.0	23:43.35	goport
7541	999	20	0	763356	13132	10192	S	0.3	0.1	16:32.51	goxdcr
7638	999	20	0	5592	2456	1896	S	0.3	0.0	24:22.28	goport
29695	root	20	0	44520	4100	3340	R	0.3	0.0	0:01.66	top
1	root	20	0	225444	9276	6788	S	0.0	0.1	0:17.12	systemd
2	root	20	0	0	0	0	S	0.0	0.0	0:00.71	kthreadd
4	root	0	-20	0	0	0	I	0.0	0.0	0:00.00	kworker/0:0H
6	root	0	-20	0	0	0	I	0.0	0.0	0:00.00	mm_percpu_wq
7	root	20	0	0	0	0	S	0.0	0.0	0:07.67	ksoftirqd/0
8	root	20	0	0	0	0	I	0.0	0.0	7:54.22	rcu_sched
9	root	20	0	0	0	0	I	0.0	0.0	0:00.00	rcu_bh
10	root	rt	0	0	0	0	S	0.0	0.0	0:01.33	migration/0
11	root	rt	0	0	0	0	S	0.0	0.0	0:01.53	watchdog/0

Figure 9 Unexplained CPU usage by beam.smp