



neOCampus DATA: user-guide

Dr Thiebolt François, thiebolt@irit.fr

This document is the DATA user guide for neOCampus end-users.

Modifications table

Date	Note
may.20	Initial release

Abstract

This guide will explain to you the various ways to gain access to the neOCampus DATA asset.

DRAFT

Table of contents

Abstract	1
Overview	3
neOCampus data lake first release	4
landing area	4
work zone	4
[LATER] golden area	5
Synthesis	6
neOCampus data sources	7
Data and API access	7
adding new data sources ?	7
[work zone] DATA access	8
'sensors' bucket structure	9
'data' measurement	9
[get data] UI & CSV	10
how to see 'outside' sensors ?	11
[get data] client library	12
[process data] Apache Zeppelin	13
[golden area] DATA access	14
[view data] Grafana	14
[landing area] DATA access	15
[get data] UI & CSV	15
[get data] client library	16
Annexe - A References	17

Overview

It's been a long time since neOCampus started to collect data from our sensors and actuators spread over our campus. The very first data was retrieved from U4 building, room 302 in 2017, june! It was our first attempt to collect these data in our mongoDB database.

Quite recently, we setup a new data storage architecture based on the **data lake** approach: the storage of our data is splitted across 3 zones as shown in fig.1:

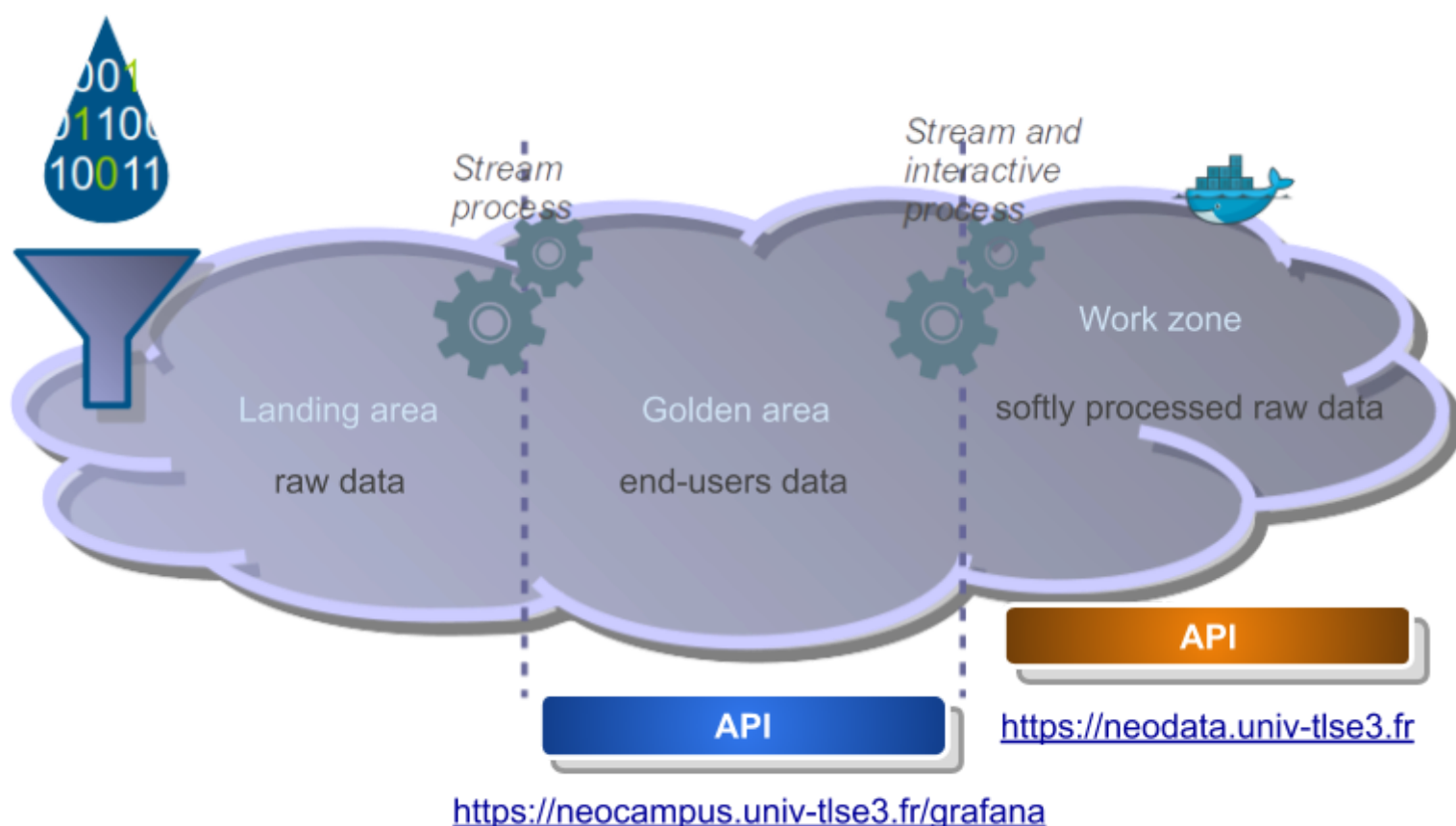


Fig. 1 - the neOCampus 'data lake' like.

Here is a short description of these 3 zones from our data lake:

Zone	Description
landing area	raw data coming from sensors; we just add a timestamp.
work zone	some soft processing has been applied to raw data producing 'more understandable' data (mostly data have been set with tags describing their location)
golden area	the real 'end-user' area with fully processed data showing a per location (i.e site, building, room) view featuring weekly, monthly and yearly synthesized data.

neOCampus data lake | first release

More than one year ago, we were facing issues accessing data from our mongoDB database: while the data was time indexed, it used to take, sometime, up to ten minutes to retrieve data for a specific room or building ... those days are gone :)

While talking with some of our colleagues and our internship applicant¹ a few months ago (feb.20), the idea of a multi-database storage area with stream processing in-between started to emerge. We tried to take into account the fact that different kinds of people will express different kinds of access to these data. Our end-users (university services, building managers etc) would like to have a synthetic view of these data while our PhD (kind of intermediate users) on their side would need to launch some specific processing on (almost) raw data.

While not really being a data lake as all of the storage and processing are self-contained in one server, we expressed a three tiered data storage architecture with stream processing in-between. One day, we'll probably switch to a real data lake, so it ought to be just a matter of switching from one local data source to a remote one ... at least we hope!

In the remaining part of this short overview, we'll give you some details about the different storage areas, we'll also present you the actual data sources, then next chapter will delve within the real data access procedures.

landing area

At this time of writing (may.20), we're using a **mongoDB** database as our **landing area**. These raw data (~100million entries ---may.20) are indexed according to the timestamp of the measure. Here are the ways to access these data:

- interactive GUI → <https://neocampus.univ-tlse3.fr/mongo>
- client library to mongoDB
(e.g Python3 pymongo <https://api.mongodb.com/python/current/python3.html>)

We must confess that the mongo-express UI is not really relevant for legacy data (i.e before apr.20 ---see fig. 2) access purposes. One way to circumvent this issue is to allow you to undertake your data processing in an offline manner: working directly on a mongoDB dump. If you're willing to go this way, please send an email to the neOCampus management team <https://neocampus.org>

work zone

Maybe the most interesting area for PhD and other people working on data to extract knowledge and other stuff. A few weeks ago we setup a new time series database: **InfluxDB 2.0**² This new release combines a data access API, (light) visualization tools, interactive query editor and many bindings for a broad range of client libraries. By the beginning of may, we imported all of the data available in mongoDB ⇒ you now have access to the broad range of data starting from 2017 :)

¹ Pr André Péninou, Pr Olivier Teste, Dr Jacques Thomazeau, M. Vincent Nam-Dang, Dr Hamdi BenHamou

² <https://www.influxdata.com/products/influxdb-overview/influxdb-2-0/>

There exists several ways to gain access to these data:

- interactive GUI → <https://neodata.univ-tlse3.fr>
- client library to influxDB 2.0
(e.g Python3 <https://github.com/influxdata/influxdb-client-python>)
- HTTP Restful API (better is to go through a client library),
- [later] Apache Zeppelin + Spark cluster,
- offline processing through InfluxDB dumps.



[LATER] golden area

This area will get dedicated to end-users that need a synthetic overview. You'll have a dashboard allowing you to select your data according to your needs (e.g location, building). Then more detailed graphs about rooms for example will feature: a week view (5mn resolution), a month view (1 day resolution with min, max, mean) and a year view (month resolution with min, max, mean). It will be something pretty close to what Domoticz offers.

Here are the following ways to access these data:

- Dashboard UI → <https://neocampus.univ-tlse3.fr/grafana>
- client library (may be a mix of InfluxDB and mongoDB)

Synthesis

Zone		
landing area (raw data)	before apr.20	[mongoDB] neOCampus DB: 'measure' collection 'failedData' collection
	after apr.20	[mongoDB] neocampus_datalake DB: 'sensors' collection
work zone (softly processed data)	from may.17	[influxDB] 'sensors' bucket (<i>legacy data were imported</i>) 'metrics' bucket (<i>MQTT metrics</i>)
golden area (synthesized data)		

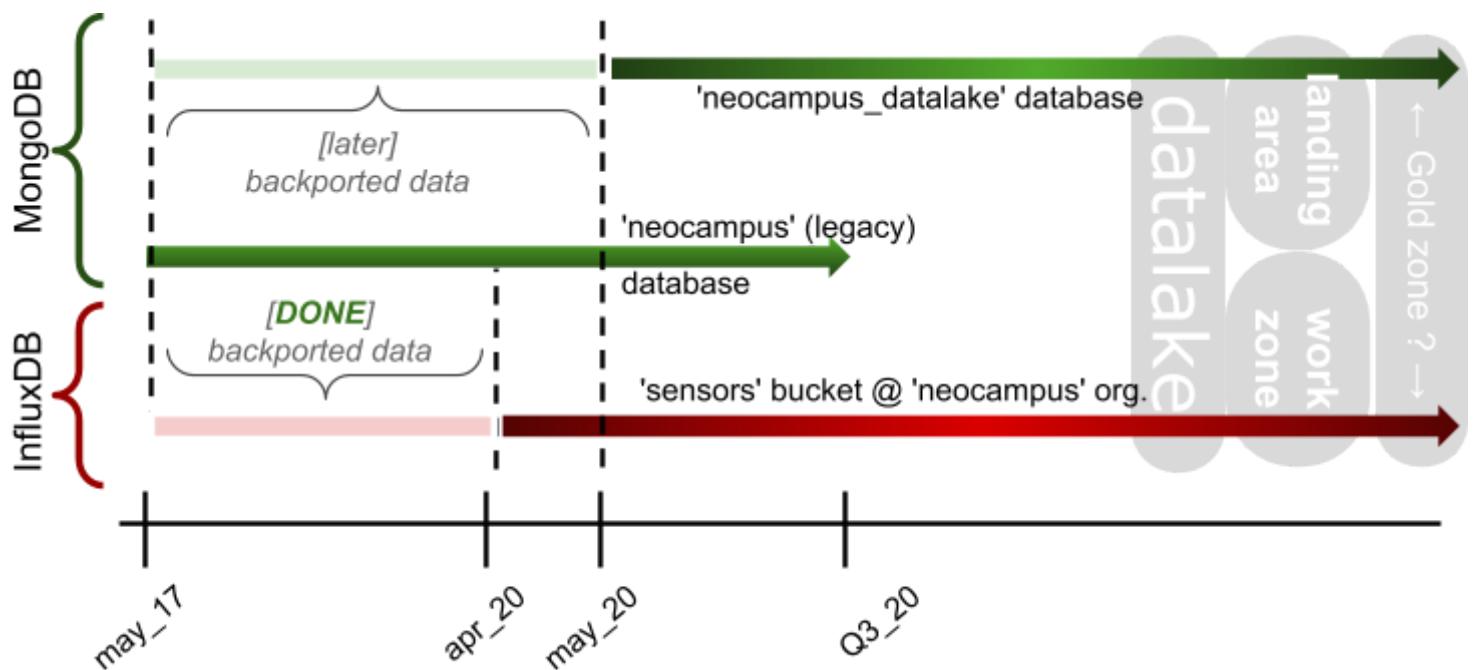


Fig. 2 - neOCampus data timeline.

As stated in fig. 2, a bit later we'll re-import all of the 'legacy' data to the neOCampus data lake 'landing area'.

neOCampus data sources

Below is a description of the various data sources brought to you.

Date	Description
	Apiary data
	LoRaWAN end-devices data
	SGE data
apr.20	UT3 / buildings BU science, IRIT2
apr.20	Toulouse metropole weather station
jun.17	UT3 / building U4

Note: greyed items means that data are not yet available.

Data and API access

Of course, all those access require credentials, either as login / password (e.g mongoDB) or tokens (e.g influxDB). Our students may just drop me an email. Pay attention to the fact that most of these API are only reachable from UT3 networks.

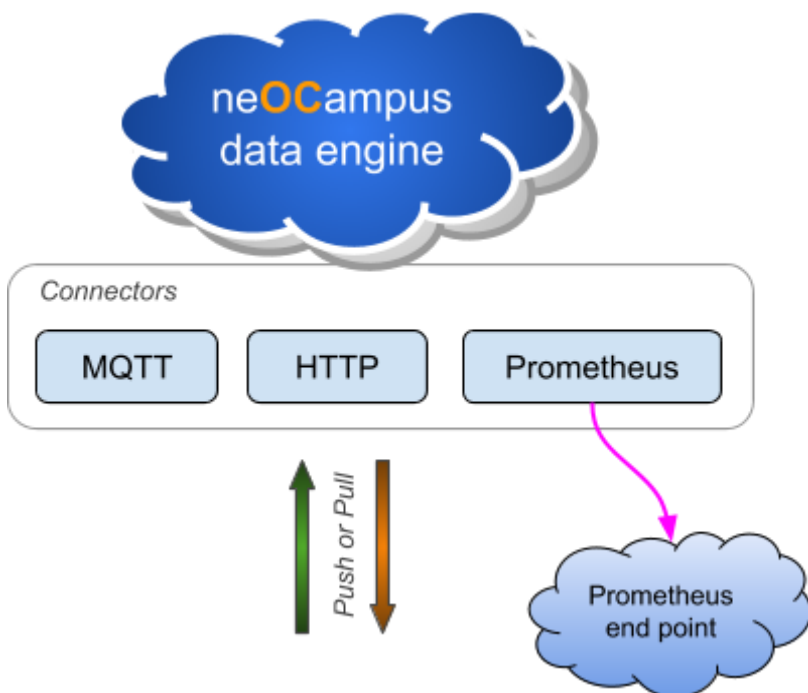
adding new data sources ?

We encourage you to propose new data sources that can be taken into account and that will benefit everyone.

We propose several connectors for both live and offline data:

- HTTP(s) in both PUSH or PULL modes,
- Prometheus protocol to collect data from end-points.

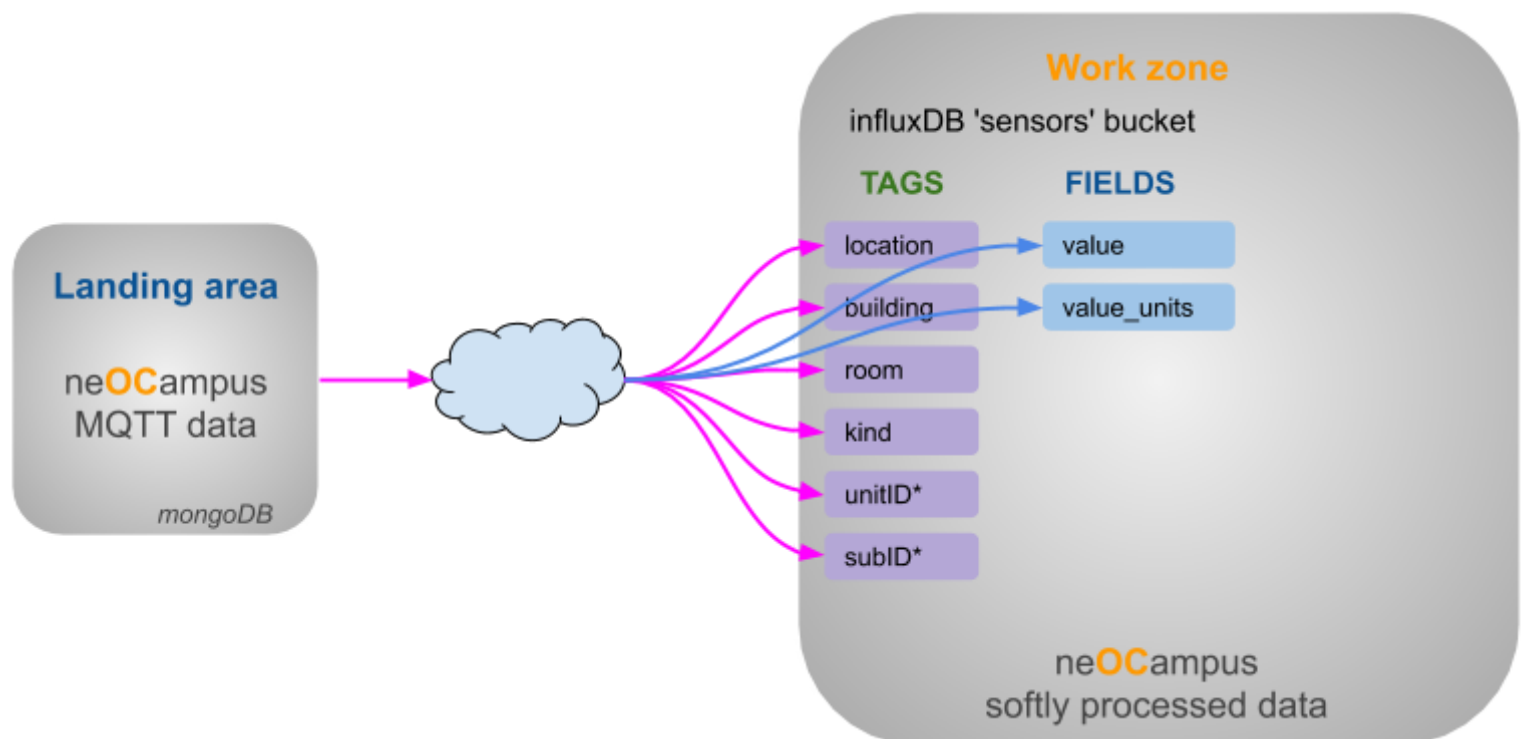
In case your app. does not behave as a prometheus end-point, we can add a prometheus gateway that will make your data easily reachable.



[work zone] DATA access

We're using influxDB to store our sensors' data. Before going on with data access, we need to specify their format.

As stated earlier, neOCampus MQTT data get stored in the 'landing area'. These data whose format is specified in [1] get slightly processed to finally arriving in the 'work zone'



'sensors' bucket structure

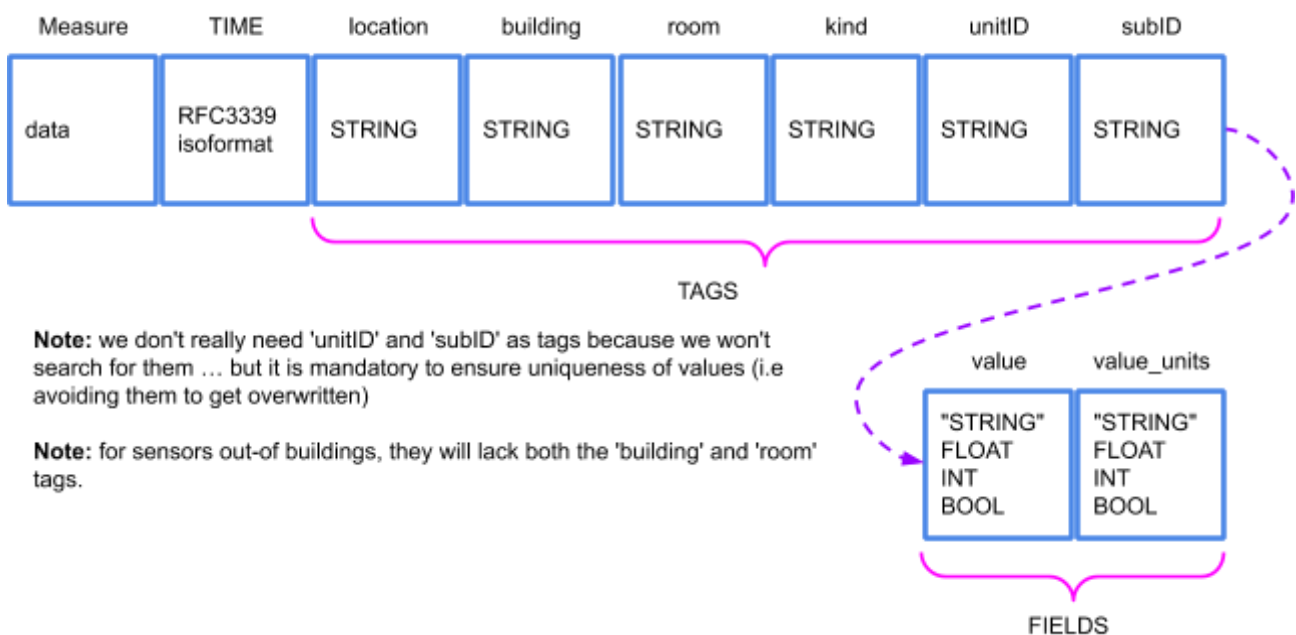
Hold all of the collected data from sensors / actuators.

This bucket only features one measurement:

- data

'data' measurement

TAGS	@data measures [sensors bucket]
location	e.g ut3, ensimag, carcassonne, ...
building	e.g u3, u4, ...
room	e.g 300, 301, 302, ...
kind	e.g temperature, wind, humidity, energy, ...
unitID	<i>as in payload</i>
subID	<i>as in payload</i>
FLAGS	@data measures [sensors bucket]
value	<i>as in payload</i>
value_unit	<i>as in payload</i>



Note: the various type (i.e kinds tag) are described in [1].

[get data] UI & CSV

InfluxDB natively proposes an UI which allows you to both see and export data in a CSV format.

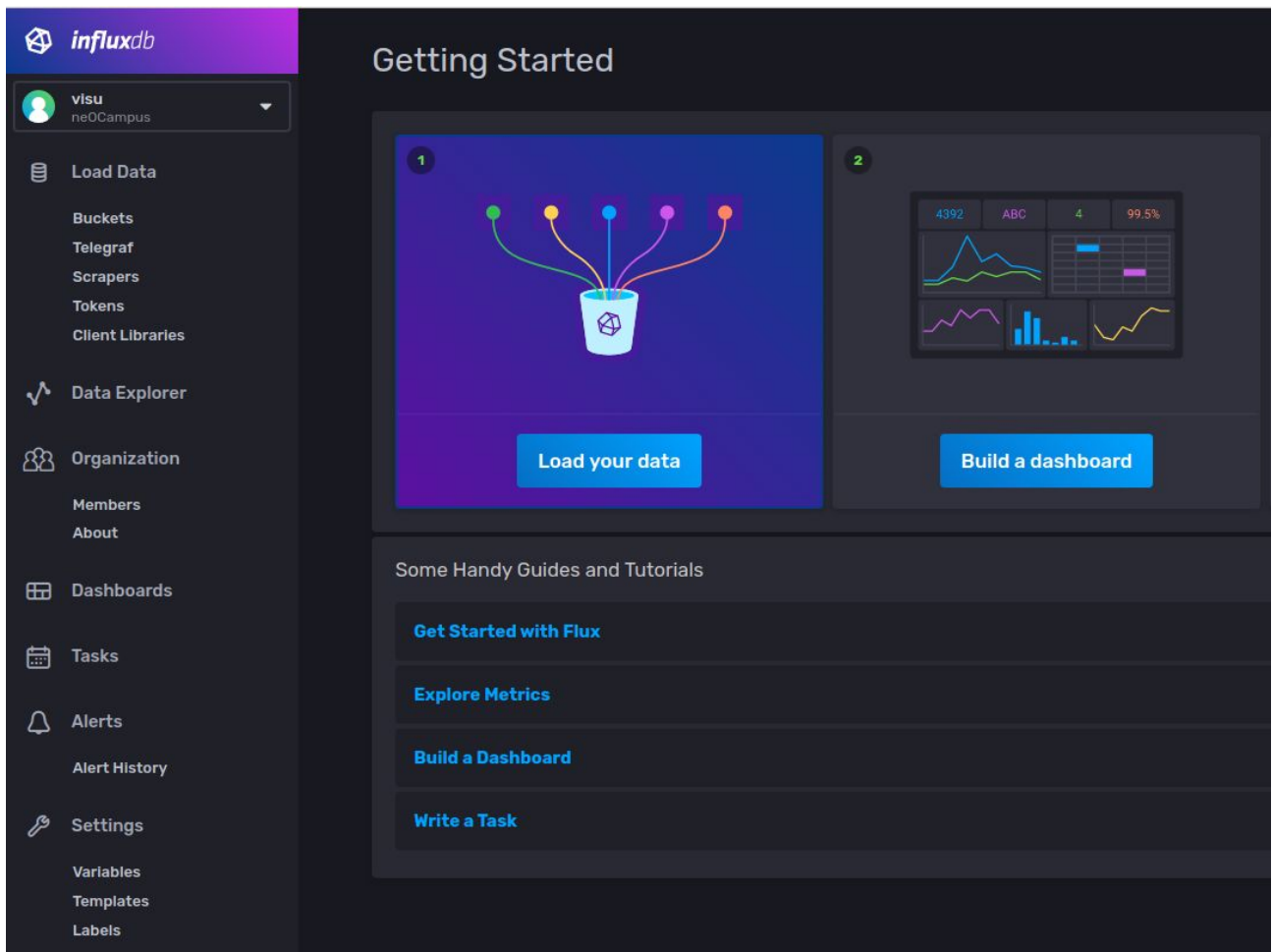


<https://neodata.univ-tlse3.fr>

login: visu

passwd: <ask us>

... you then get connected to influxDB

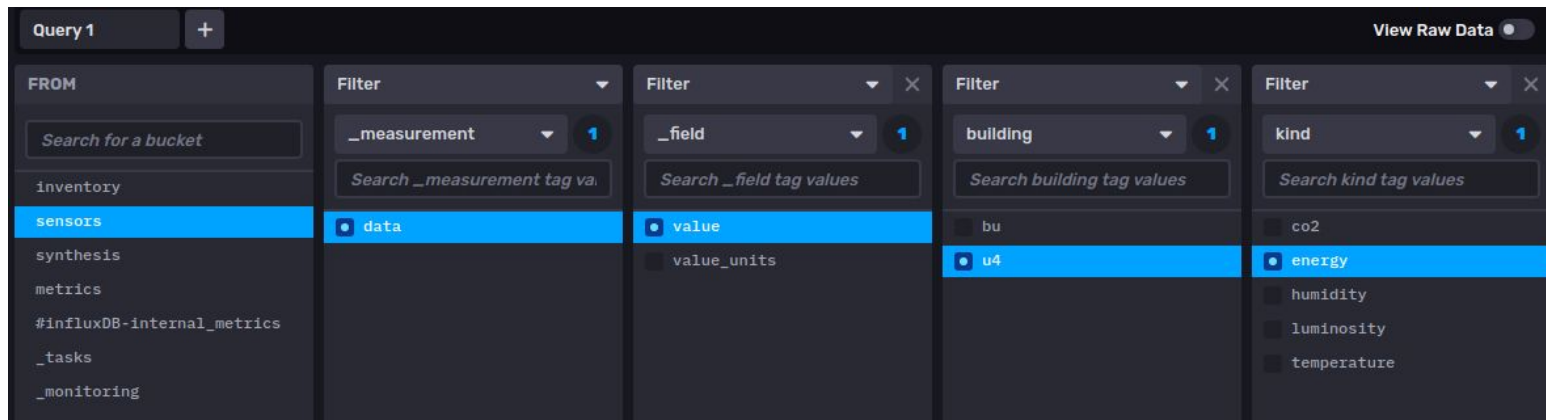


On the left side panel, select 'Data explorer'

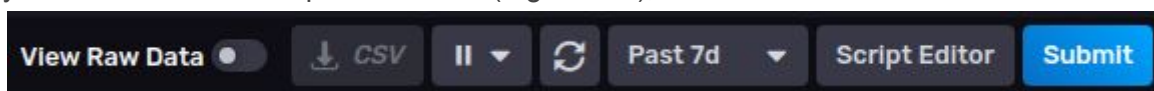
→ It will show you several data sources (buckets).

→ select 'sensors' bucket along with 'data' measurement (the time serie itself)

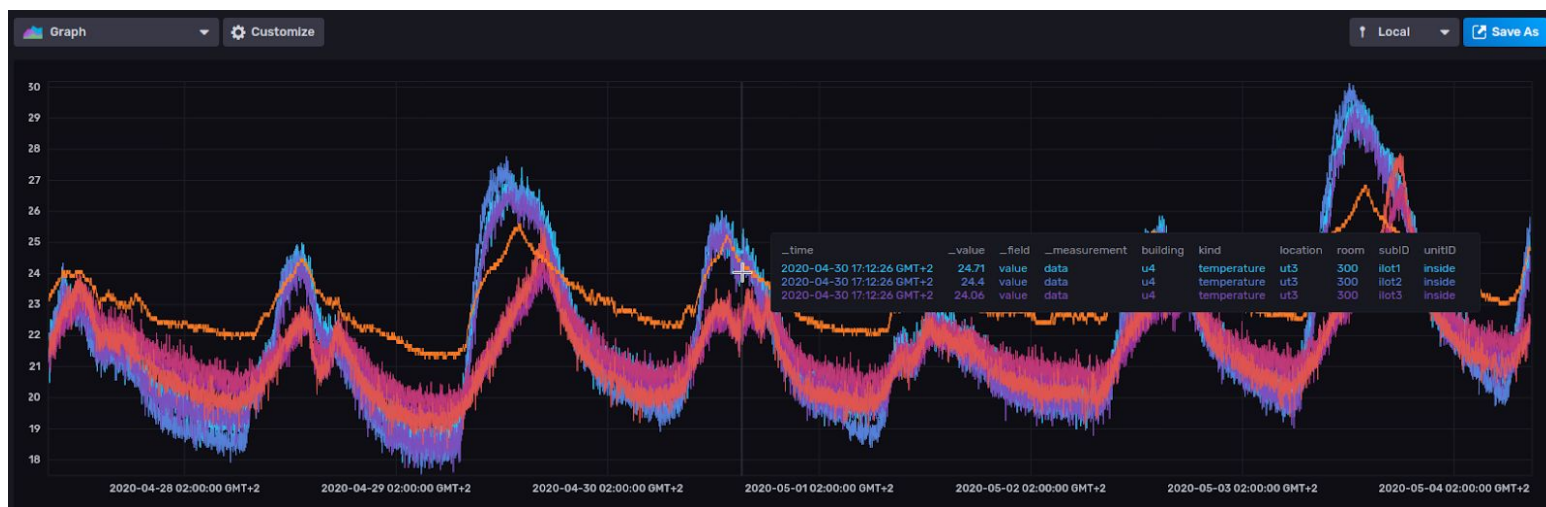
→ you can then start to explore the various buildings, rooms, kinds, ...



you can then choose a period of time (e.g last 7d), then click 'submit'



... tada !

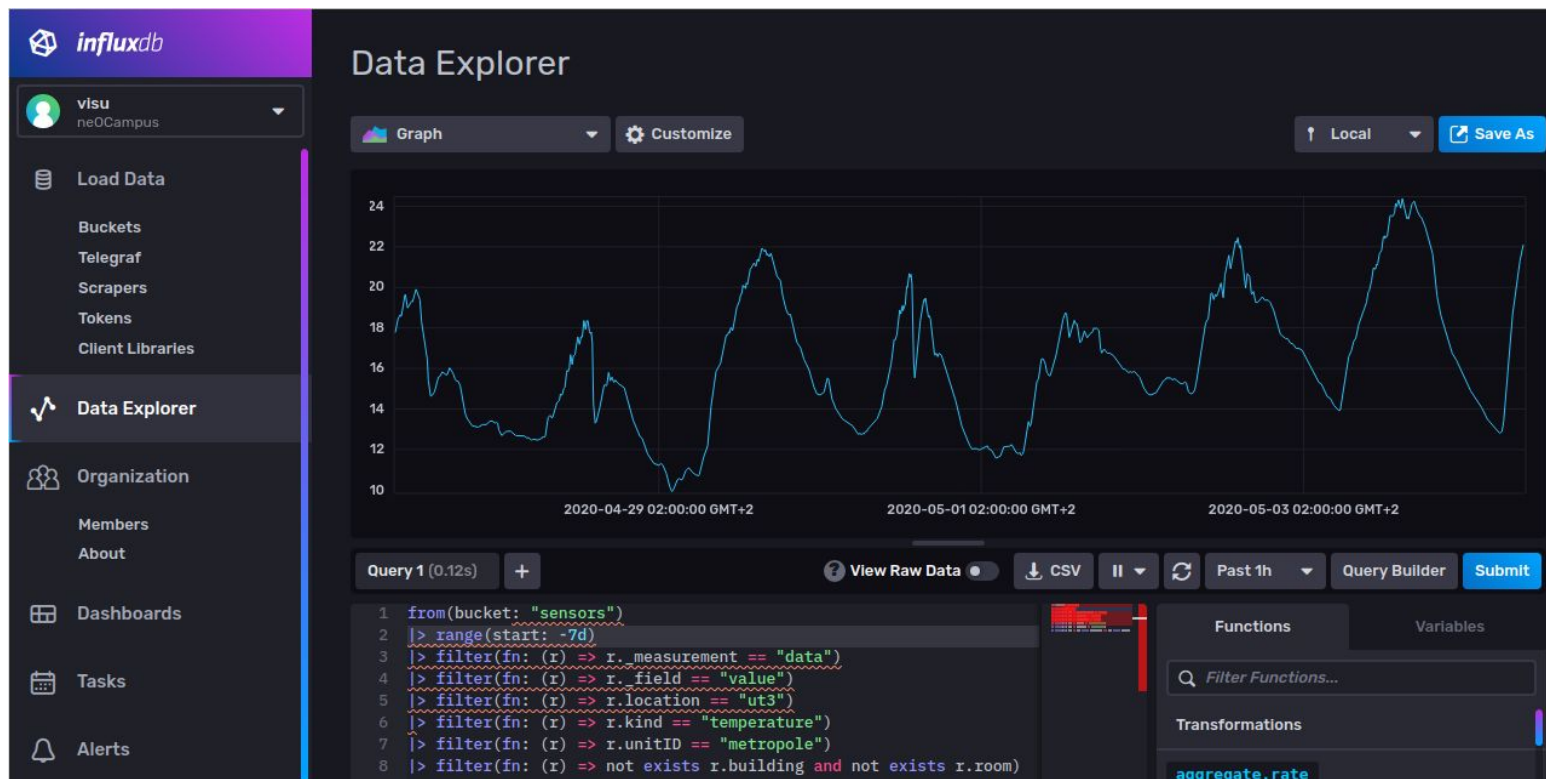


You can also choose to select 'RAW DATA' (tables) ... and export as CSV for further analysis.

how to see 'outside' sensors ?

As described in [1] along within the [dataCollector](#) document, outside sensors does not feature any building nor room tags → query editor

```
from(bucket: "sensors")
|> range(start: -24h)
|> filter(fn: (r) => r._measurement == "data")
|> filter(fn: (r) => r._field == "value")
|> filter(fn: (r) => r.location == "ut3")
|> filter(fn: (r) => r.kind == "temperature")
|> filter(fn: (r) => r.unitID == "metropole")
|> filter(fn: (r) => not exists r.building and not exists r.room)
```



Of course you can undertake more complex Flux queries, see languages details at <https://v2.docs.influxdata.com/v2.0/reference/flux/>

[get data] client library

Probably the simplest way to undertake your data analysis, there exists various client libraries to gain access to the influxDB database.



server: neodata.univ-tlse3.fr

port: 9999

read token: <ask us>

Our docker container 'dataCollector' features numerous apps written with the Python client <https://github.com/influxdata/influxdb-client-python>

```
pip3 install influxdb-client
```

[process data] Apache Zeppelin

TO BE CONTINUED

[golden area] DATA access

TO BE CONTINUED

[view data] Grafana

TO BE CONTINUED

[landing area] DATA access

TO BE CONTINUED

[get data] UI & CSV

mongoDB can be accessed through the mongo-express UI front-end.



server: <https://neocampus.univ-tlse3.fr/mongo/>

login: reader

passwd: <ask us>

Viewing Database: neocampus_datalake

Collections		Collection Name	+ Create collection
View	Export	[JSON]	Import
		sensors	Del

Database Stats	
Collections (incl. system.namespaces)	1
Data Size	47.2 MB
Storage Size	14.2 MB
Avg Obj Size #	183 Bytes
Objects #	257655
Extents #	0
Indexes #	1
Index Size	3.28 MB

[get data] client library

Probably the simplest way to undertake your data analysis, there exists various client libraries to gain access to the influxDB database.



server: neodata.univ-tlse3.fr
port: 27012
login / passwd: <ask us>

Our docker container 'dataCollector' features apps written with the Python client
<https://api.mongodb.com/python/current/tutorial.html>

```
pip3 install pymongo
```


Annexe - A | References

[1] neOCampus end-devices API and MQTT rules

This document describes the way end-devices are managed through the sensOCampus application. It also describes the various topics you may encounter in the 'landing area'. More important, all kinds of CLASS topics (i.e temperature, pressure, wind, energy ...) are fully specified in this document.

https://neocampus.univ-tlse3.fr/_media/sensocampus_end-devices_api.pdf