



Royal Netherlands  
Meteorological Institute  
*Ministry of Infrastructure  
and Water Management*

# KNMI.nl Dataloader

Quinten Cabo



# Finding a Dataset



## Source of bad datasets?

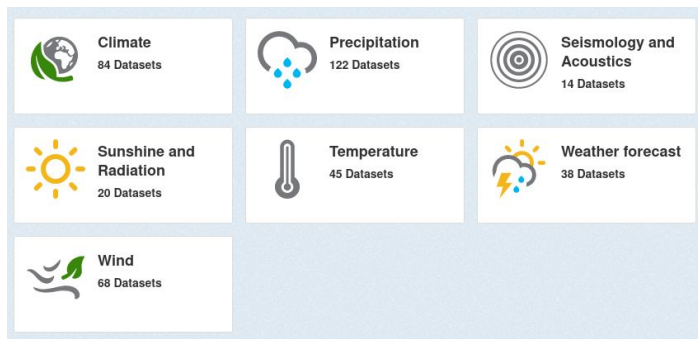


GOVERNMENT

Finding a dataset  
**KNMI.nl**

<https://www.knmi.nl/home>

Koninklijk Nederlands Meteorologisch Instituut.

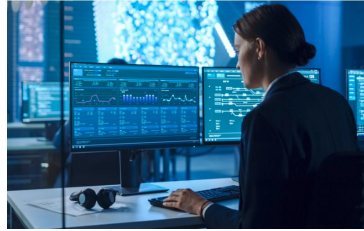


Koninklijk Nederlands  
Meteorologisch Instituut  
*Ministerie van Infrastructuur en Waterstaat*



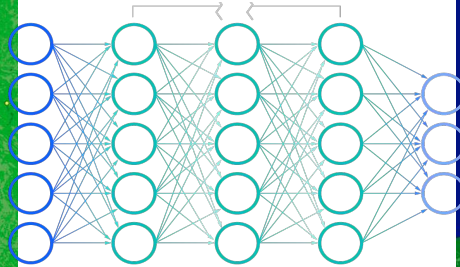
Use case

## Weather prediction models



$T = 1$

LSTM or a Transformer



$T = 2$

## Finding a dataset

# KNMI

<https://www.knmi.nl/home>

Koninklijk Nederlands Meteorologisch Instituut.



Koninklijk Nederlands

## API Catalog

### Notification Service

The Notification Service allows you to receive timely updates about events of the KNMI Data Platform.

[View Documentation](#)

### EDR API (ALPHA)

The *Environmental Data Retrieval* (EDR) REST API. This API allows users to query datasets on specific parameters in a spatio-temporal manner, reducing the need to download entire datasets. The EDR API adheres to [Open Geospatial Consortium EDR standard](#), facilitating integration with other geospatial tools and systems.

DISCLAIMER: This pre-release version is still in development. The API may contain bugs, and data integrity is not yet verified.

[View Documentation](#)

[API Reference \(Swagger\)](#)

### Web Map Service (WMS)

Web service for geospatial data integration in online mapping applications. The web service is compatible with OGC WMS specification versions 1.0.0, 1.1.1, and 1.3.0 (default, unless otherwise requested).

[View Documentation](#)

✓ Request an API key

✓ Request an API key

✓ Request an API key

## Finding a dataset KNMI

<https://www.knmi.nl/home>

Koninklijk Nederlands Meteorologisch Instituut.



Koninklijk Nederlands

## API Catalog

### Notification Service

The Notification Service allows you to receive timely updates about events of the KNMI Data Platform.

[View Documentation](#)

### EDR API (ALPHA)

The *Environmental Data Retrieval* (EDR) REST API. This API allows users to query datasets on specific parameters in a spatio-temporal manner, reducing the need to download entire datasets. The EDR API adheres to [Open Geospatial Consortium EDR standard](#), facilitating integration with other geospatial tools and systems.

DISCLAIMER: This pre-release version is still in development. The API may contain bugs, and data integrity is not yet verified.

[View Documentation](#)

[API Reference \(Swagger\)](#)

### Web Map Service (WMS)

Web service for geospatial data integration in online mapping applications. The web service is compatible with OGC WMS specification versions 1.0.0, 1.1.1, and 1.3.0 (default, unless otherwise requested).

[View Documentation](#)

✓ Request an API key

✓ Request an API key

✓ Request an API key



Accessing the data

## Requesting access to the mqtt stream



Royal Netherlands  
Meteorological Institute  
*Ministry of Infrastructure  
and Water Management*

KNMI Developer Portal

News

FAQ

Documentation ▾

API Catalog

Register

Login

Login

Email address

Password

Login

Don't have an account ? [Register](#) here

Forgot password? [Request password reset](#)



Accessing the data

## Requesting access to the mqtt stream



Royal Netherlands  
Meteorological Institute  
*Ministry of Infrastructure  
and Water Management*

KNMI Developer Portal

[News](#)

[FAQ](#)

[Documentation](#) ▼

[API Catalog](#)

[Register](#)

[Login](#)

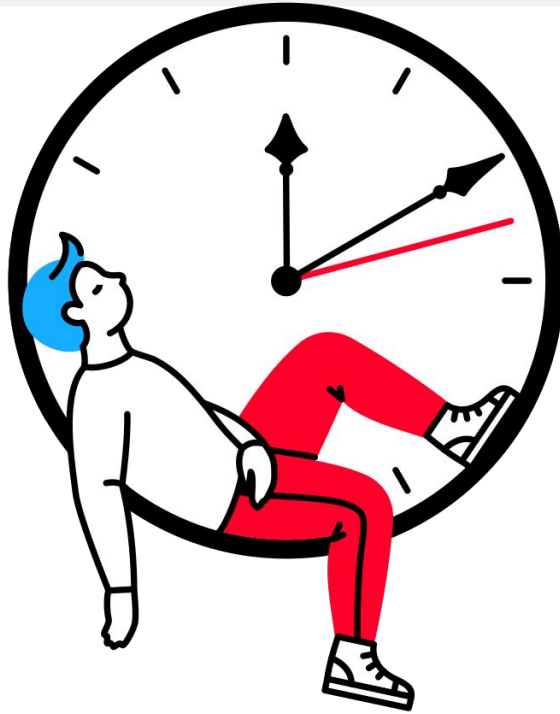
## Request a Key

Use case - To help KNMI understand how the API is used.

Request Key

Requesting access

**Waiting for a response**



Accessing the data

## Getting the notification token

### Your API Subscriptions

API Key Hash: 8f901bf12c7c456d856dc28f94c1c730

Accessing the data

## Subscribing to the mqtt stream

```
create database if not exists;
# Version 3.1.1 also supported
PROTOCOL = mqtt client.MQTTv5
logging.basicConfig(
logger = logging.getLogger( name )
logger.setLevel("INFO")
def connect mqtt() -> mqtt client:
    def on connect(c: mqtt client, userdata, flags, rc, reason code, props=None):
        logger.info(f"Connected using client ID:{str(c._client_id)}")
        logger.info(f"Session present: {str(flags['session present'])}")
        logger.info(f"Connection result: {str(rc)}")
        # Subscribe here so it is automatically done after disconnect
        subscribe(c, TOPIC)
    client = mqtt client.Client(paho.mqtt.enums.CallbackAPIVersion.VERSION1 client_id=CLIENT ID, protocol=PROTOCOL, transport="websockets")
    client.tls set(tls version=ssl.PROTOCOL_TLS
    connect properties= properties.Properties(properties.PacketTypes.CONNECT
    # Maximum is 3600
    connect properties.SessionExpiryInterval =3600
    # The MQTT username is not used for authentication, only the token
    username = "token"
    client.username pw set(username, TOKEN)
    client.on connect = on connect
    client.connect(host=BROKER DOMAIN, port=443, keepalive=60, clean start=False, properties=connect properties)
    return client

def run():
    client = connect mqtt()
    client.enable logger(logger=logger)
    client.loop forever()
```

Accessing the data

## Mqtt stream items

```
{
  "specversion": "1.0",
  "type": "nl.knmi.dataplatform.file.updated.v1",
  "source": "https://dataplatform.knmi.nl", "id": "458081e2-3842-8c0a-b912-85fb93e8c0ed",
  "time": "2024-06-05T14:19:36Z", "datacontenttype": "application/json",
  "data": {
    "datasetName": "Actuele10mindataKNMIstations",
    "datasetVersion": "2",
    "filename": "KMDS__OPER_P___10M_OBS_L2_202406051410.nc",
    "url": "https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS__OPER_P___10M_OBS_L2_202406051410.nc/url"
  }
}
```

Accessing the data

**Lets download the file**

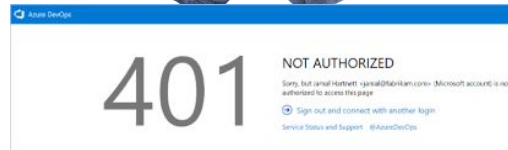
GET:

[https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS\\_OPER\\_P\\_10M\\_OBS\\_L2\\_202406051410.nc/url](https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS_OPER_P_10M_OBS_L2_202406051410.nc/url)

Accessing the data  
**Lets download the file**

GET:

[https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS\\_OPER\\_P\\_10M\\_OBS\\_L2\\_202406051410.nc/url](https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS_OPER_P_10M_OBS_L2_202406051410.nc/url)



Accessing the data

## Second KNMI API key

<https://www.knmi.nl/home>

Koninklijk Nederlands Meteorologisch Instituut.

### Open Data API

The Open Data REST API allows you to download files from datasets published on KNMI Data Platform - [Data Catalog](#).

[View Documentation](#)

[API Reference \(Swagger\)](#)

 [Request an API key](#)



Koninklijk Nederlands  
Meteorologisch Instituut  
*Ministerie van Infrastructuur en Waterstaat*



Accessing the data

## Requesting access - logging back in



Royal Netherlands  
Meteorological Institute  
*Ministry of Infrastructure  
and Water Management*

KNMI Developer Portal

News

FAQ

Documentation ▾

API Catalog

Register

Login

Login

Email address

Password

Login

Don't have an account ? [Register](#) here

Forgot password? [Request password reset](#)

Accessing the data

## Requesting access again



Royal Netherlands  
Meteorological Institute  
*Ministry of Infrastructure  
and Water Management*

KNMI Developer Portal

[News](#)

[FAQ](#)

[Documentation](#) ▼

[API Catalog](#)

[Register](#)

[Login](#)

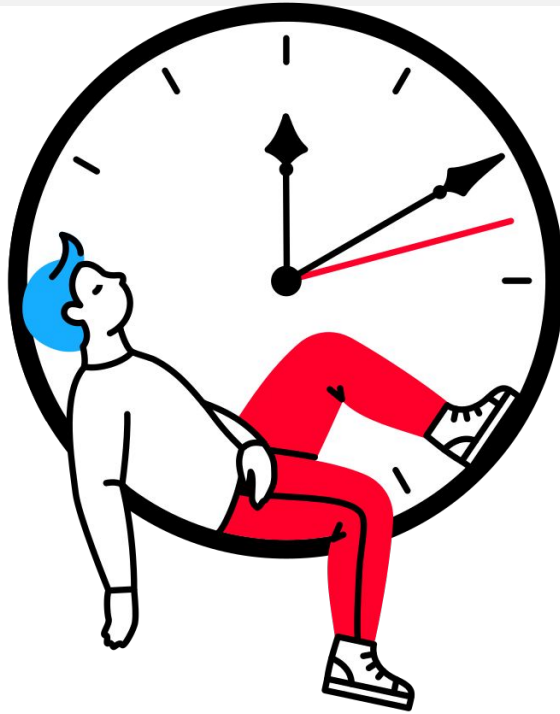
## Request a Key

Use case - To help KNMI understand how the API is used.

Request Key

Accessing the data

**Waiting for a response again**



Accessing the data

## Getting the second token

### Your API Subscriptions

API Key Hash: 8f901bf12c7c456d856dc28f94c1c730

Accessing the data

## Downloading the data second attempt

GET:

[https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS\\_OPER\\_P\\_10M\\_OBS\\_L2\\_202406051410.nc/url](https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS_OPER_P_10M_OBS_L2_202406051410.nc/url)




```
headers={"Authorization": f"Bearer {FILE_DOWNLOAD_TOKEN}"}
```

Accessing the data

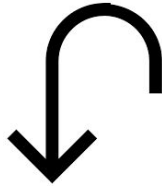
## Downloading the data second attempt

GET:

[https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS\\_OPER\\_P\\_10M\\_OBS\\_L2\\_202406051410.nc/url](https://api.dataplatform.knmi.nl/open-data/v1/datasets/Actuele10mindataKNMIstations/versions/2/files/KMDS_OPER_P_10M_OBS_L2_202406051410.nc/url)



```
headers={"Authorization": f"Bearer {FILE_DOWNLOAD_TOKEN}"}
```



```
{"temporaryDownloadUrl": "https://amazon.eu-west.supersuperlongurl.temp.kmni.nl"}
```

Accessing the data

## **Downloading the data third attempt**

GET: <https://amazon.eu-west.supersuperlongurl.temp.kmni.nl>

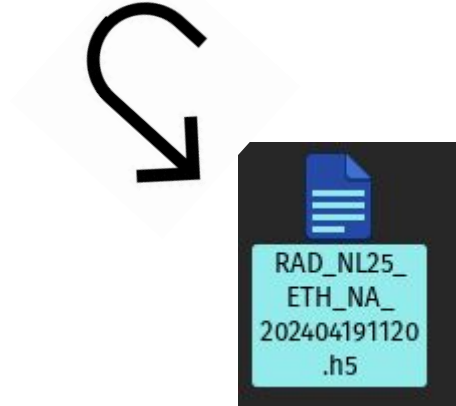
(without token)

Accessing the data

## Downloading the data third attempt

GET: <https://amazon.eu-west.supersuperlongurl.temp.kmni.nl>

(without token)





## Accessing the data

### Recap

Do get data from kmni you only need to

1. Request a token
2. Wait
3. Request another token
4. Wait
5. Listen to the stream with token 1
6. Get link from stream item
7. Request the download link from the stream item link with token 2
8. Download file from this link WITHOUT a token

Then you get an **h5** file...



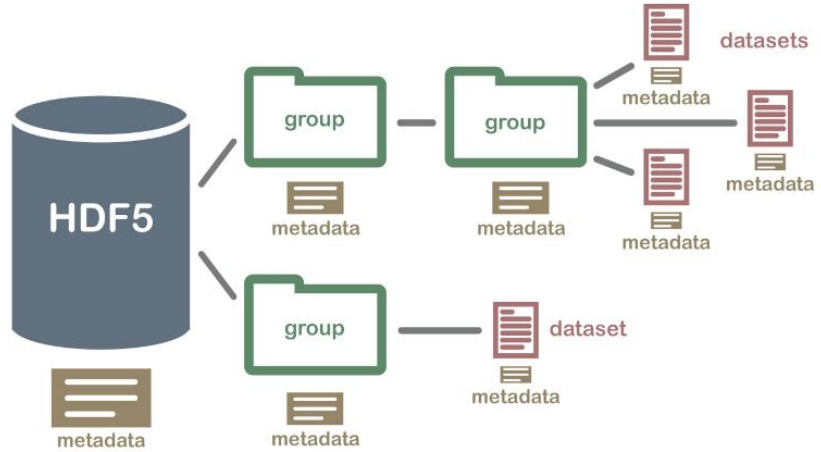
The data itself

**Wtf is an h5 file?**



## Wtf is an h5 file?

- Binary zip file with **groups** and **datasets**

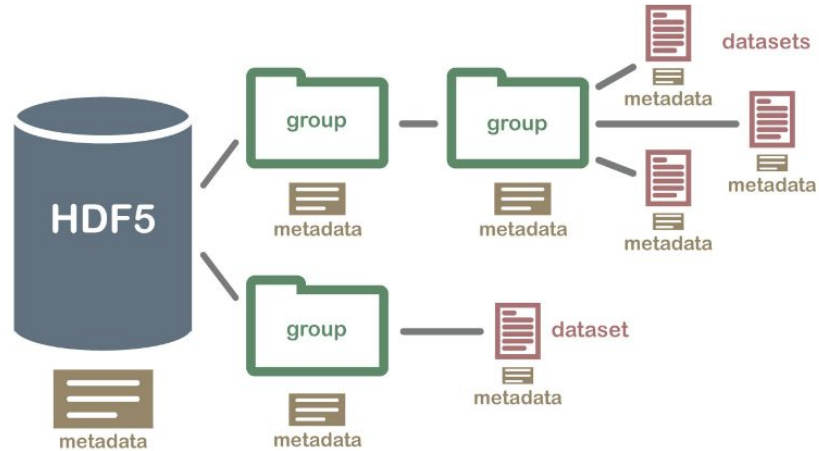


<https://www.neonscience.org/resources/learning-hub/tutorials/about-hdf5>

The data itself

## Wtf is an h5 file?

- Binary zip file with **groups** and **datasets**
- Datasets are homogeneous arrays
- Groups contain groups and data

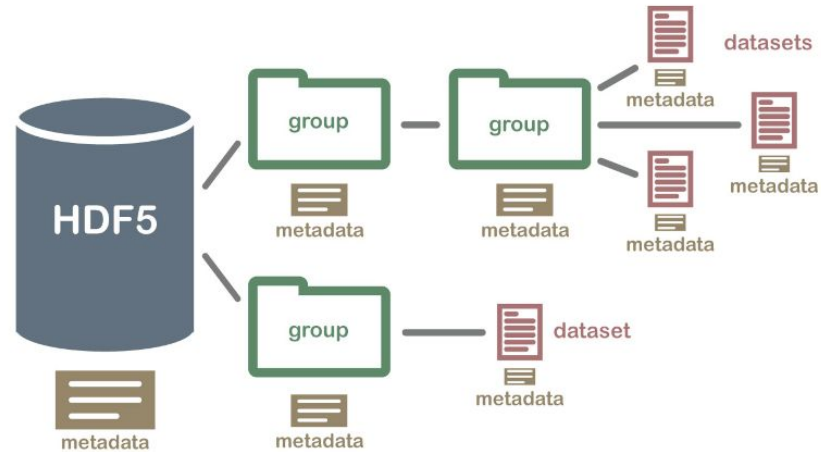


<https://www.neonscience.org/resources/learning-hub/tutorials/about-hdf5>

The data

## Wtf is an h5 file?

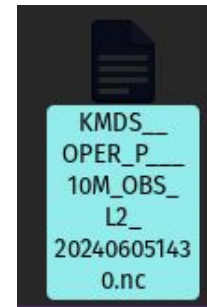
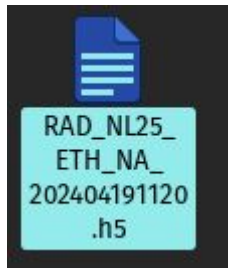
- Binary zip file with **groups** and **datasets**
- Datasets are homogeneous arrays
- Groups contain groups and dataset
- Groups and datasets have Metadata



<https://www.neonscience.org/resources/learning-hub/tutorials/about-hdf5>

The data

## KMNI API update

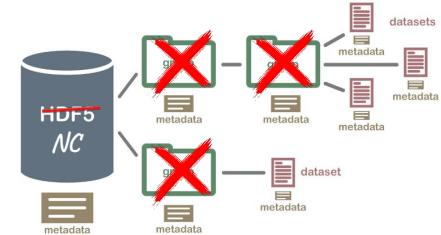


## What is an nc file



**nc** file instead gives you something like this:

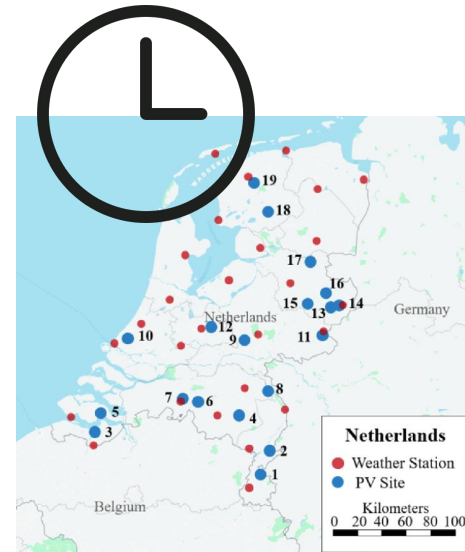
(No groups)



## What is IN each nc file?

Each nc file has the measurements of 94 variables for ~69 weather station at a single point in time.

Name	Long name	Shape	Unit	Type
dsd	Wind Direction 10 Min Std Dev with MD	(69, 1)	degree	float64
dx	Wind Direction Sensor 10 Min Maximum with MD	(69, 1)	degree	float64
ff	Wind Speed at 10m 10 Min Average with MD	(69, 1)	m s-1	float64
ffs	Wind Speed Sensor 10 Min Average with MD	(69, 1)	m s-1	float64
fed	Wind Speed 10 Min Std Dev with MD	(69, 1)	m s-1	float64





Quality

## The data quality is actually very good

- Detailed Comments,
- Iso standard ISO 19115
- Clear Units (c, %, ...)
- Missing values (but masked arrays)

Name	Long name	Shape	Unit	Type
dsd	Wind Direction 10 Min Std Dev with MD	(69, 1)	degree	float64
dx	Wind Direction Sensor 10 Min Maximum with MD	(69, 1)	degree	float64
ff	Wind Speed at 10m 10 Min Average with MD	(69, 1)	m s-1	float64
ffs	Wind Speed Sensor 10 Min Average with MD	(69, 1)	m s-1	float64
fed	Wind Speed 10 Min Std Dev with MD	(69, 1)	m s-1	float64

Quality

## The data quality is actually very good

- Detailed Comments,
- Iso standard ISO 19115
- Clear Units (c, %, ...)
- Missing values (but masked arrays)

But

Do get data from kmni you only need to

1. Request a token
2. Wait
3. Request another token
4. Wait
5. Listen to the stream with token 1
6. Get link from stream item
7. Request the download link from the stream item link with token 2
8. Download file from this link WITHOUT a token

Then you get an **h5** file...

h5

Name	Long name	Shape	Unit	Type
dsd	Wind Direction 10 Min Std Dev with MD	(69, 1)	degree	float64
dx	Wind Direction Sensor 10 Min Maximum with MD	(69, 1)	degree	float64
ff	Wind Speed at 10m 10 Min Average with MD	(69, 1)	m s-1	float64
ffs	Wind Speed Sensor 10 Min Average with MD	(69, 1)	m s-1	float64
ffd	Wind Speed 10 Min Std Dev with MD	(69, 1)	m s-1	float64

Easiest path to train test split?



Easiest path to train test split?  
Combining multiple nc files?



## Serving **Dataset interface**

```
•[1]:
```

```
KMNIdataset(  
    root="data",  
    train=True,  
    download=True,  
    after="2024-01-05T01:00:00.000Z",  
    until="2024-06-05T16:25:14.948Z",  
    stations=("06201", "06203", "06204", "06205", "06207", "06208", "06209", "06211")  
    variables=("fxs", "h", "n2", "p0", "Q1H", "rg")  
)
```



## Serving Dataset interface defaults

Just load everything defaults

```
•[1]:
```

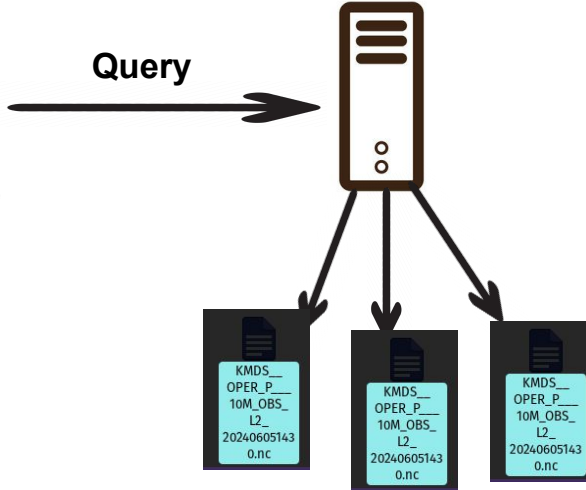
```
KMNIdataset(  
    root="data",  
    train=True,  
    download=True,  
)
```



## Serving Querying the nc files directly?

•[1]:

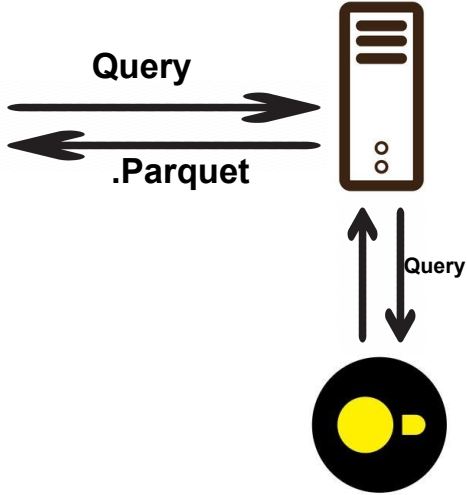
```
KMNIdataset(  
  root="data",  
  train=True,  
  download=True,  
  after="2024-01-05T01:00:00.000Z",  
  until="2024-06-05T16:25:14.948Z",  
  stations=("06201", "06203", "06204", "06205", "06207", "06208", "06209", "06211")  
  variables=("fxs", "h", "n2", "p0", "Q1H", "rg")  
)
```



## Querying the data from a database

•[1]:

```
KMNIdataset(  
    root="data",  
    train=True,  
    download=True,  
    after="2024-01-05T01:00:00.000Z",  
    until="2024-06-05T16:25:14.948Z",  
    stations=("06201", "06203", "06204", "06205", "06207", "06208", "06209", "06211")  
    variables=("fxs", "h", "n2", "p0", "Q1H", "rg")  
)
```





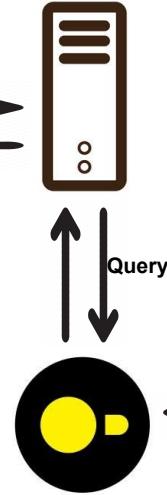
# Wrangling

## Pipeline to convert nc files into rows

•[1]:

```
KMNIdataset(  
  root="data",  
  train=True,  
  download=True,  
  after="2024-01-05T01:00:00.000Z",  
  until="2024-06-05T16:25:14.948Z",  
  stations=("06201", "06203", "06204", "06205", "06207", "06208", "06209", "06211")  
  variables=("fxs", "h", "n2", "p0", "Q1H", "rg")  
)
```

Query  
↓  
↑  
.Parquet

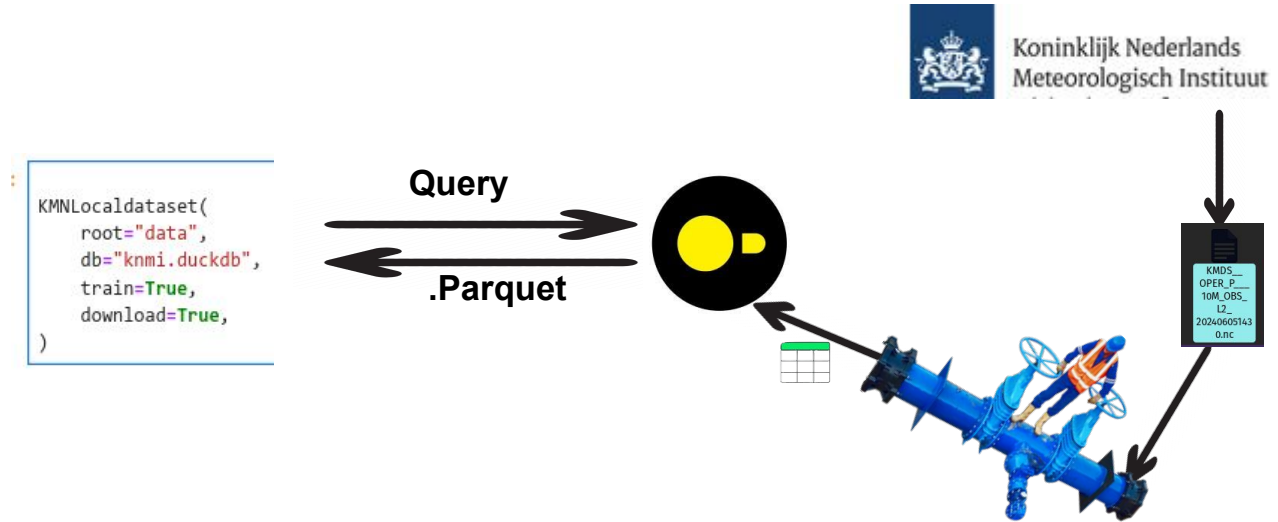


Koninklijk Nederlands  
Meteorologisch Instituut



Serving

## Pipeline to convert nc files into rows LOCALLY



# Airflow



## Cool but too complex

Airflow DAGs Cluster Activity Datasets Security Browse Admin Docs

DAGs

Active 2 Paused 27

Running 0 Failed 0

Filter DAGs by tag

DAG	Owner	Runs	Schedule	Last Run	Next Run
<a href="#">conditional_dataset_and_time_based_timetable</a> <small>dataset-time-based-timetable</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset or 0 1 * * *		2024-06-05
<a href="#">consume_1_and_2_with_dataset_expressions</a>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		0 of 2 datasets
<a href="#">consume_1_or_2_with_dataset_expressions</a>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		0 of 2 datasets
<a href="#">consume_1_or_both_2_and_3_with_dataset_expressions</a>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		0 of 3 datasets
<a href="#">dataset_consumes_1</a> <small>consumes dataset-scheduled</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		On a3/dag1
<a href="#">dataset_consumes_1_and_2</a> <small>consumes dataset-scheduled</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		0 of 2 datasets
<a href="#">dataset_consumes_1_never_scheduled</a> <small>consumes dataset-scheduled</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		0 of 2 datasets updated
<a href="#">dataset_consumes_unknown_never_scheduled</a> <small>dataset-scheduled</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		0 of 2 datasets updated
<a href="#">dataset_produces_1</a> <small>dataset-scheduled produces</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	@daily		2024-06-04, 00:00:00
<a href="#">dataset_produces_2</a> <small>dataset-scheduled produces</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	None		
<a href="#">example_branch_decorator</a>	airflow	<div><div></div><div></div><div></div><div></div></div>	None		
<a href="#">example_branch_operator</a> <small>example example2</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	0 0 * * *		2024-06-04, 00:00:00
<a href="#">example_branch_datetime_operator</a> <small>example</small>	airflow	<div><div></div><div></div><div></div><div></div></div>	@daily		2024-06-04, 00:00:00

Airflow DAGs Cluster Activity Datasets Security Browse Admin Docs

16:04 UTC

DAG: tutorial\_dag DAG tutorial

Schedule: None Next Run ID: None

DAG Docs

05-06-20 16:04:09 All Run Types All Run States Clear Filters

Auto-refresh 25

Press Shift + F for Shortcuts

delete failed queued removed resending retrying scheduled finished skipped success up\_for\_reschedule up\_for\_retry upstream\_failed no\_status

tutorial\_dag

Details Graph Gantt Code Audit Log Run Duration Calendar

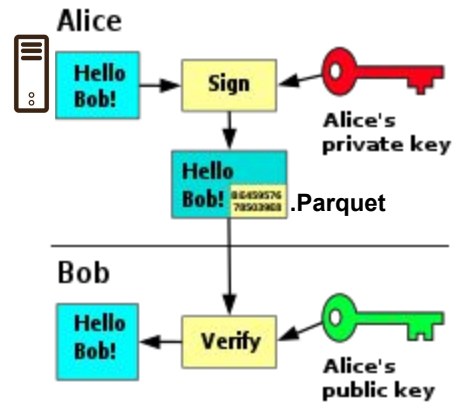
Layout: Left → Right

extract transform load

extract PythonOperator transform PythonOperator load PythonOperator



One more thing  
**Integrity**

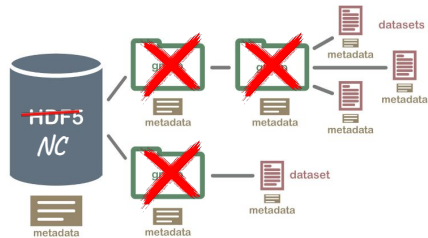


## Questions

Do get data from kmni you only need to

1. Request a token
2. Wait
3. Request another token
4. Wait
5. Listen to the stream with token 1
6. Get link from stream item
7. Request the download link from the stream item link with token 2
8. Download file from this link WITHOUT a token

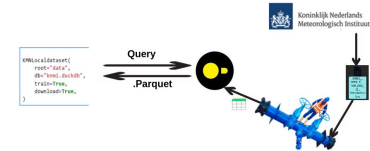
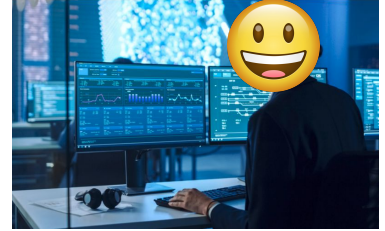
Then you get an nc file...



Koninklijk Nederlands  
Meteorologisch Instituut  
*Ministerie van Infrastructuur en Waterstaat*



data loader



## Wrangling Pipeline to convert nc files into rows

