



# Erasmus Centre for Data Analytics

## Erasmus Data Collaboratory – House of AI Mixed Source Analytics

DWBISummit- March 27, 2024



**Erasmus Data  
Collaboratory**

AI, Data &  
Digitalisation

Erasmus  
University  
Rotterdam



# Erasmus Centre for Data Analytics

Flagship centre for cross-disciplinary insight on Artificial Intelligence, Data, Digitalization and Immersive Technologies at the Erasmus University Rotterdam.



# ECDA as supporting hub

Funding

AI, Data & Digitalisation

AI, Data & Digitalisation

Erasmus University Rotterdam



Erasmus University Rotterdam



Erasmus University Rotterdam



Supported Programs

**Erasmus Data Collaboratory House of AI**

*Strategic research programmes on Digitalization, AI & immersive tech*

**Convergence AIDD**

**Labs & Initiatives**

**AI-PACT**

**SSH**

influence of digitalization on work, prosperity and entrepreneurship

(under discussion)

**AI-MAPS**



**ALGOSOC**



**Smart Campus**

**AI @ EUR**

*Strategic EUR innovation programmes on Digitalization, AI & immersive tech*

**Energy Transition programme**

*Transdisciplinary interfaculty initiatives*

Offering

Facilitate Research & Innovation programs

Facilitate Hands on Education / (re)skilling

Share Thought Leadership XPs

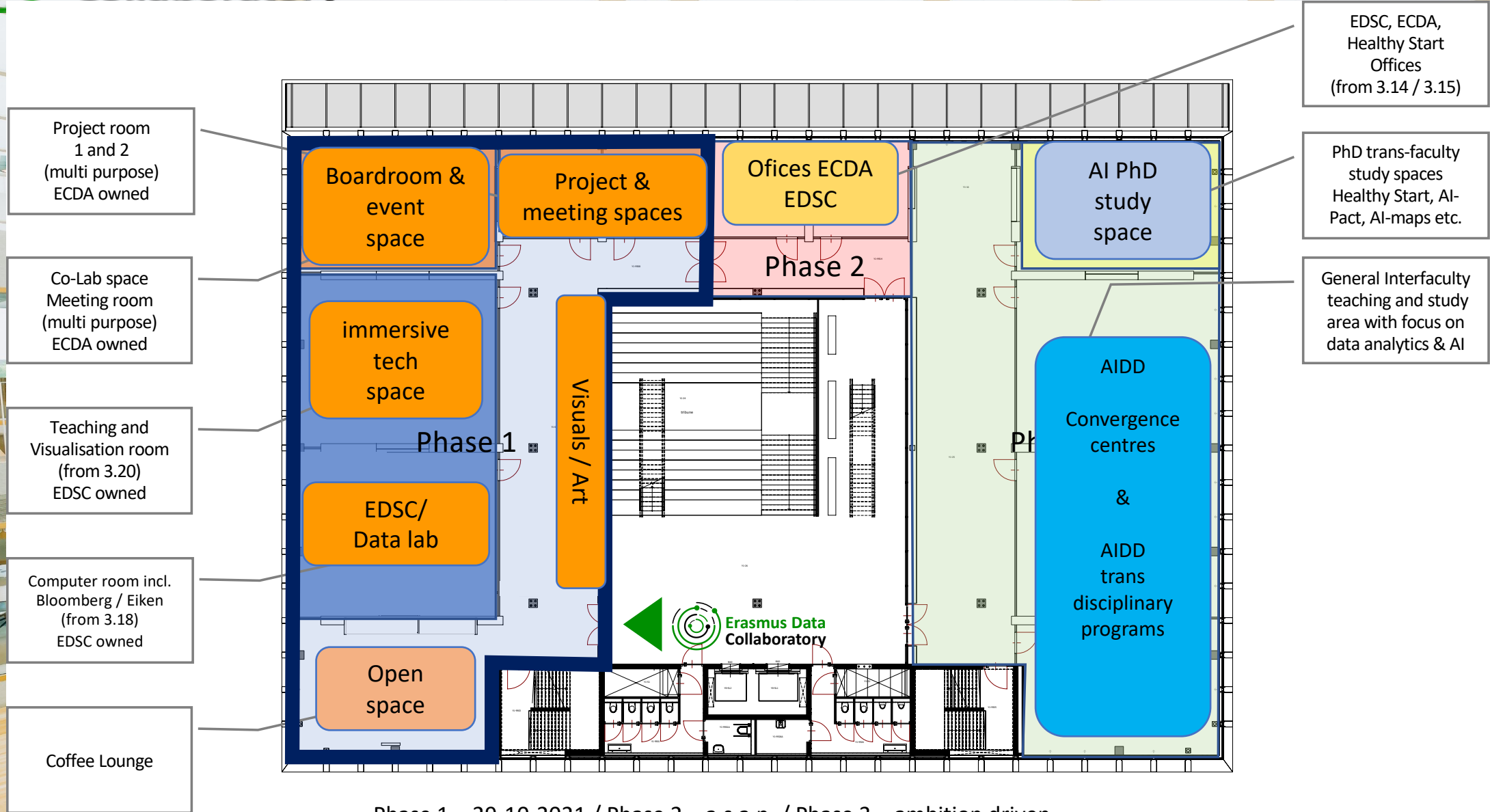
Community Building Connect & Inspire

Activate & Empower Partner organizations

Competencies

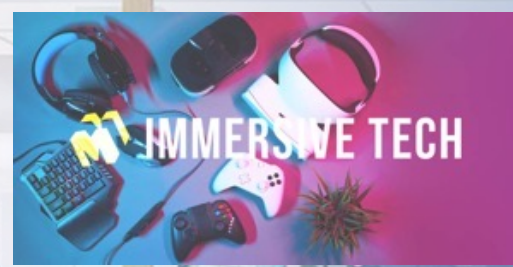
**ECDA**

AI & Data Competency / Marketing & Communications / Programs & Events / Business Development





# Visualization room Immersive Tech Space

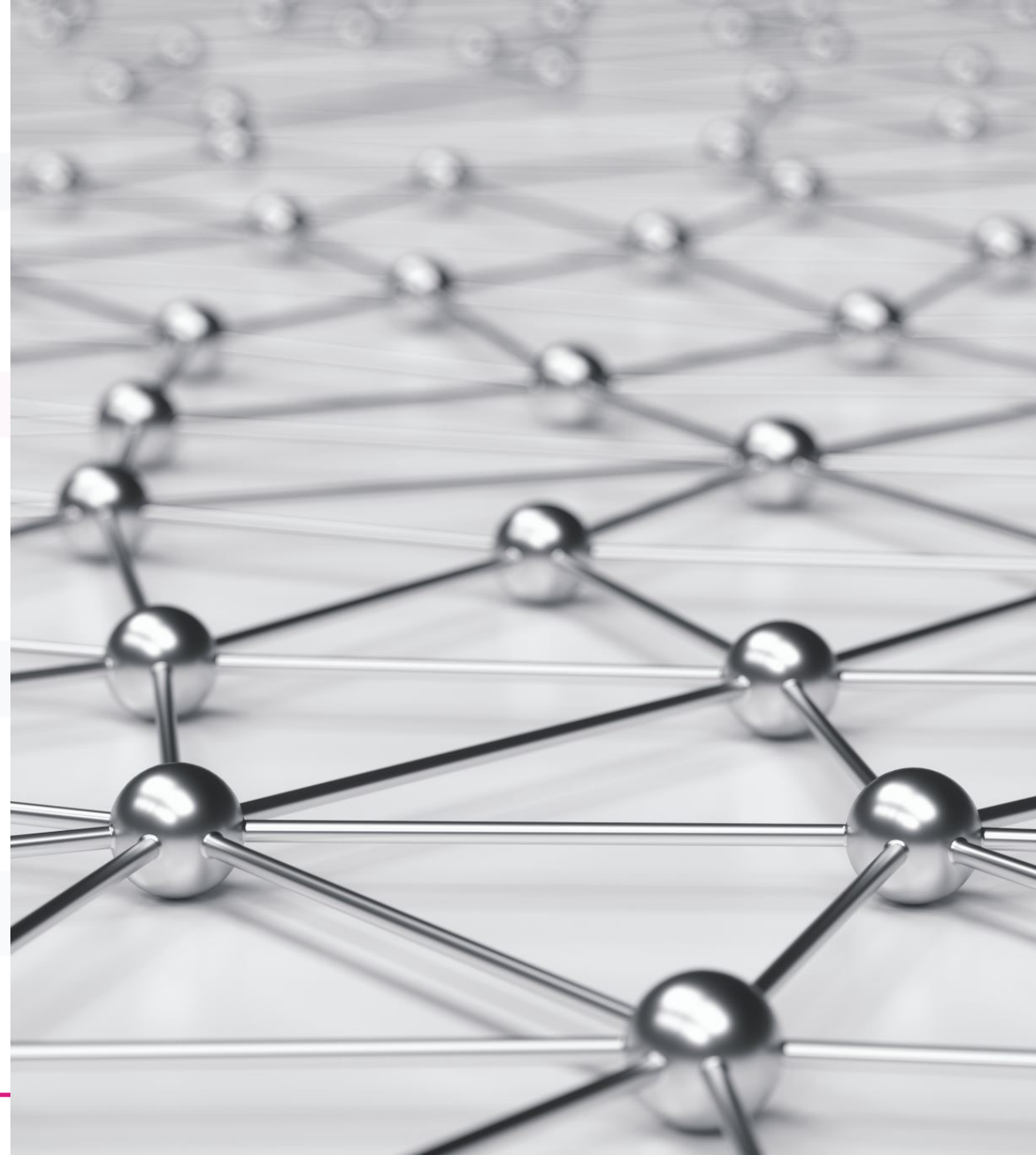




*“...support innovative data intensive projects with both internal and external parties at EUR ”*

If your project contains one or more of the following characteristics, please reach out to us:

- Innovative in nature
- Has a major data, AI or Immersive Technology component;
- Involves collaboration with External Parties with active exchange of data/expertise;
- Needs additional resources such as compute or engineering support.







Erasmus Data  
Collaboratory

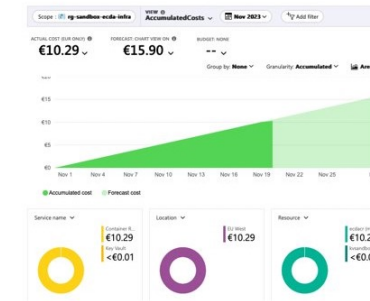
# Services for qualified research projects

- Access to external network (for data & funding)
- Access to internal resources (for skills & knowledge)
- Qualified data engineering support staff
- Safe and secure online collaboration environment
- Local compute & GPU (€0, limited)
- Access to cloud compute (€-€€€)



#### Local compute-storage-gpu:

- 104 vCpu
- 640 Gb Ram
- 3 \* Nvidia RTX 4090
- 12 TB Nvme SSD
- 25 TB SATA SSD
- 120 TB Harddisk



#### Microsoft Azure (EDIS)

- Virtually unlimited capacity
- we WILL charge you 😊



#### Lenovo Thinkstation P620:

- 1 \* Nvidia RTX 6000 Ada

#### HP Z8 Fury G5:

- 2 \* Nvidia RTX 6000 Ada

# Meet the Team!



**Data Team**



**Immersive Tech Team**

**Platform & Apps**



**Infrastructure**





# EDC Data Sandbox Core Principles



1. Never compromise on security
2. No to Personal Data (unless...)
3. Data is a strategic asset
4. Don't reinvent the wheel
5. Everything should be reusable
6. Embrace change
7. Stay independent
8. Don't (try to) boil the ocean
9. Maximize Utilization
10. Minimize Carbon Footprint

# EDC Sandbox: powered by Open Source



python™



PROXMOX



docker



kubernetes



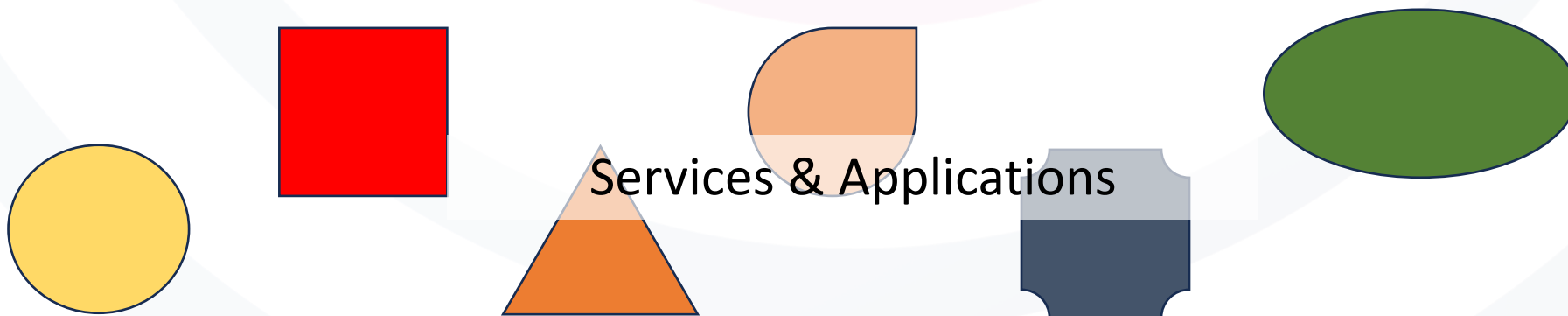
RANCHER®



# High Level Architecture

Identity & Access Management

SURF



kubernetes



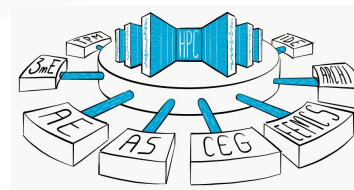
kubernetes



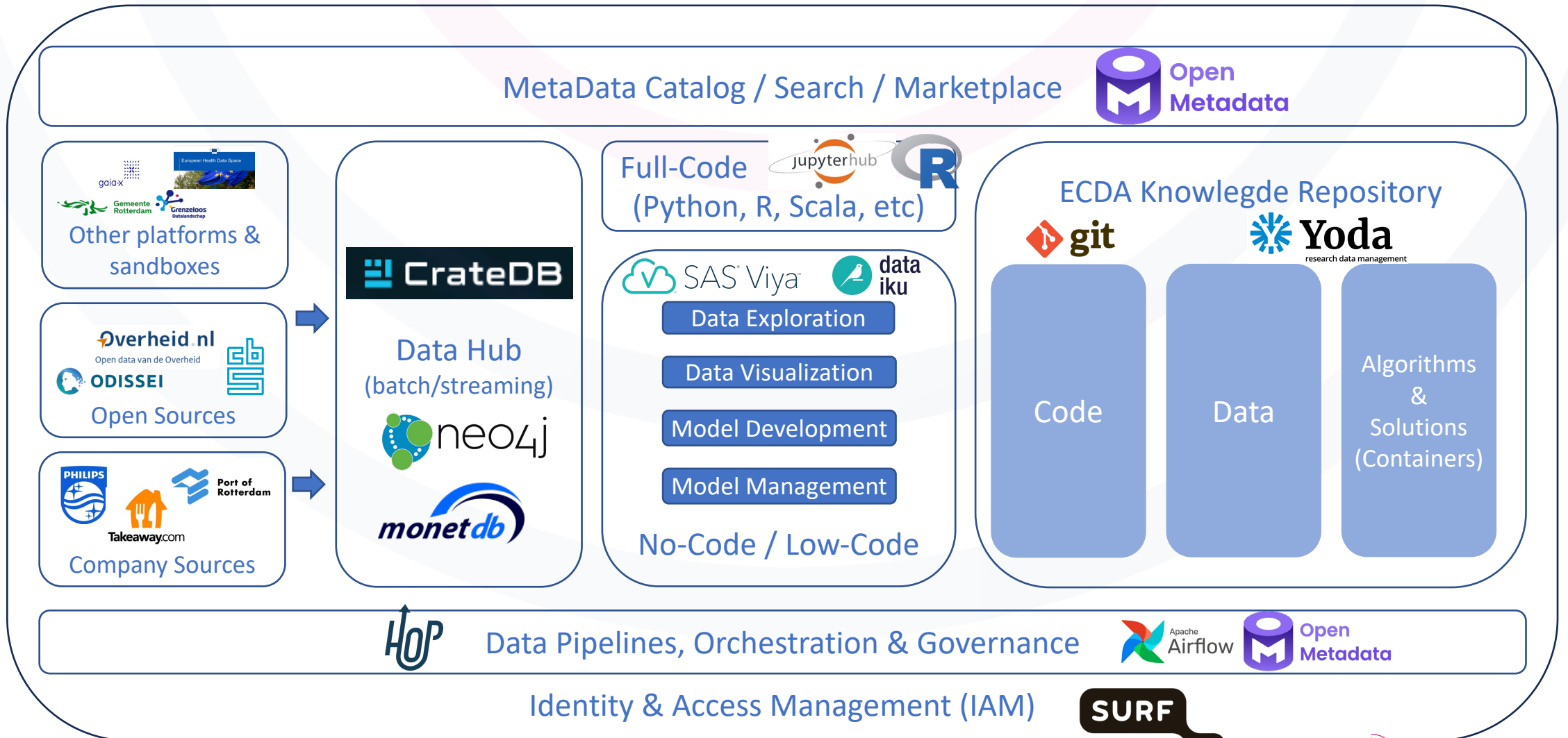
kubernetes



kubernetes



# Erasmus Data Collaboratory Sandbox Tech





# Sandbox = Hybrid 'Cloud'



**6 \* MinisForum UM 790 Pro**  
 16 vCpu, 64 GB Ram  
 internal GPU  
 2\*2TB Nvme SSD  
 10 Gbit interconnects



**4 \* GPU Workstation**  
 32 vCpu, 192 GB Ram  
 nVidia RTX 4090  
 2\*1TB + 2\*2TB Nvme SSD  
 25 Gbit interconnects



**Backup/ shared storage**  
 8 core cpu, 64 GB Ram  
 2\*2TB Nvme SSD, 2\*1.8TB SSD  
 6 \* 20 TB Sata  
 10 Gbit interconnect



*Monthly Cost on Azure would be > € 8.500, excl. GPU's!*

**EDC - House of AI Local Infrastructure**



**EDC – House of AI LeafCloud**

**EDC - House of AI EDIS-Azure**

**EDC - House of AI Surf Research Cloud**

**Total compute capacity**

- 112 cores/224 vCpu
- 1.15 TB Ram

**Total storage capacity (raw/net)**

- 50 / 34 TB Nvme SSD
- 120 / 60 TB Sata

**Total compute/storage capacity**

- Virtually unlimited

# The EDC Compute Cabinet 😊

✓ Fast

✓ Affordable

✓ Reliable

✓ Energy Efficient



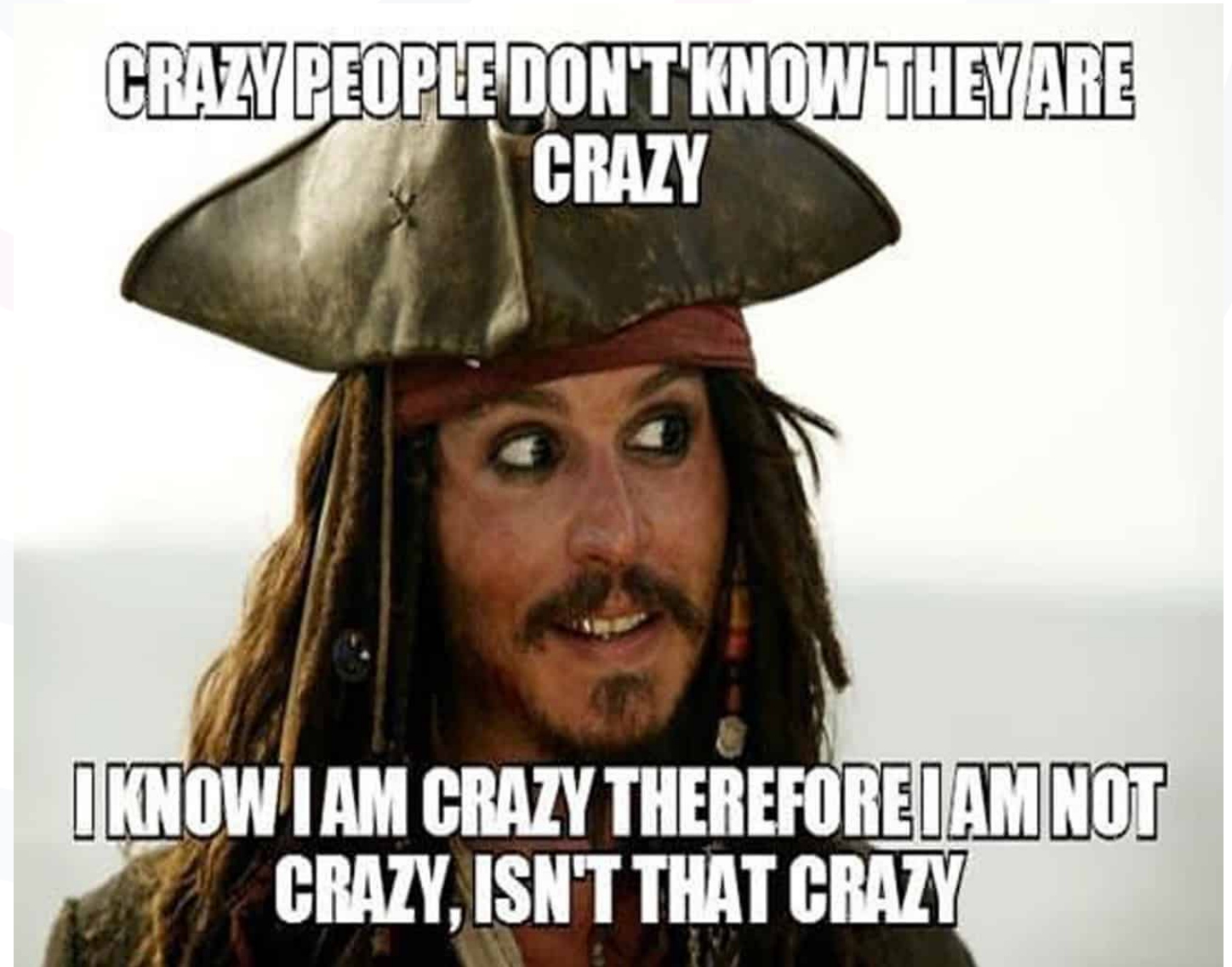


## Public Cloud vs On Prem

*"You're crazy if  
you don't start  
in the cloud;  
you're crazy if  
you stay on it"*

andreessen.  
horowitz

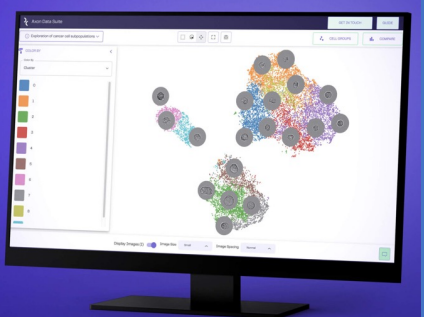
Source: [The Cost of Cloud, a Trillion Dollar Paradox](#)





# Sandbox Experiments / Projects

Deepcell 



Erasmus Centre for Data Analytics  
ECDA INITIATIVE  
**Smart Campus**  
Erasmus University Rotterdam  
Make your campus more sustainable with data

Erasmus University Rotterdam  
NEW INITIATIVE  
**Erasmian Language Model**

 **FeedbackFruits**



chat.ecda.ai

galax  
Gemeente Rotterdam  
Grenzeloos Ontwikkeling  
Other platform sandboxes

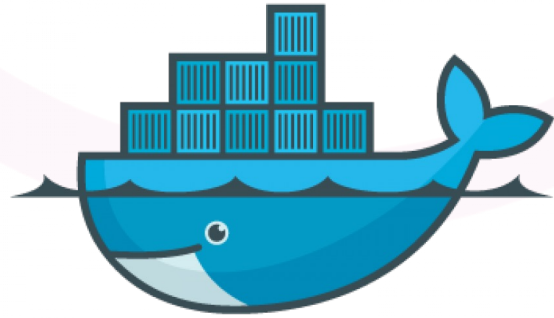
Overheid.nl  
Open data van de Overheid  
ODISSEI  
Open Sources

PHILIPS  
Takeaway.com  
Company Sources

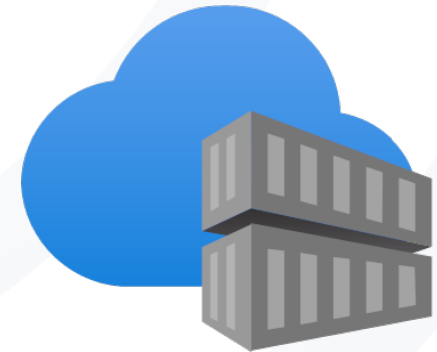
pository  
Algorithms & Solutions Containers)

# Data Engineering @ Erasmus Data Collaboratory

# Tools of the Trade



docker



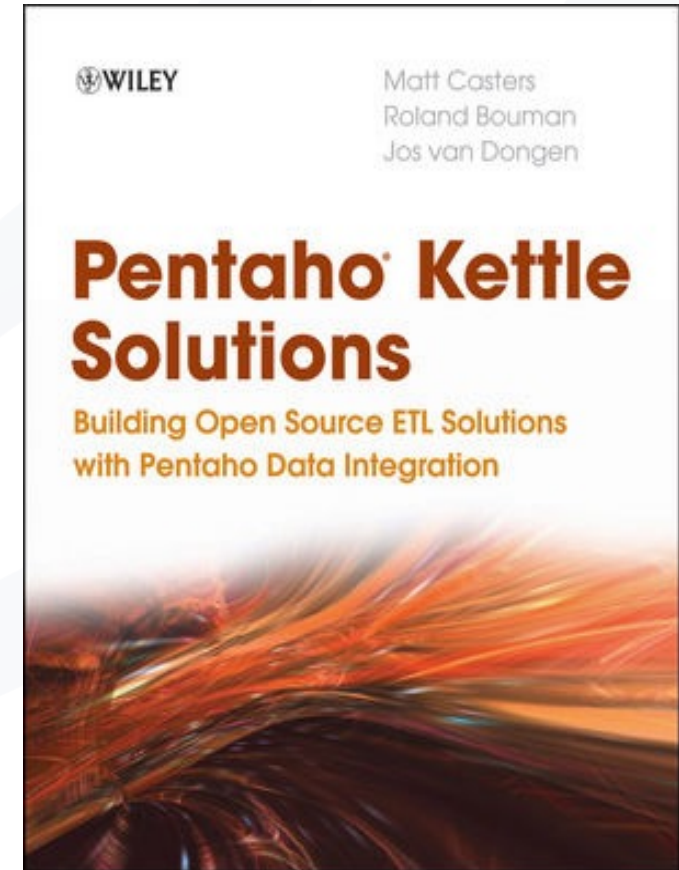
kubernetes



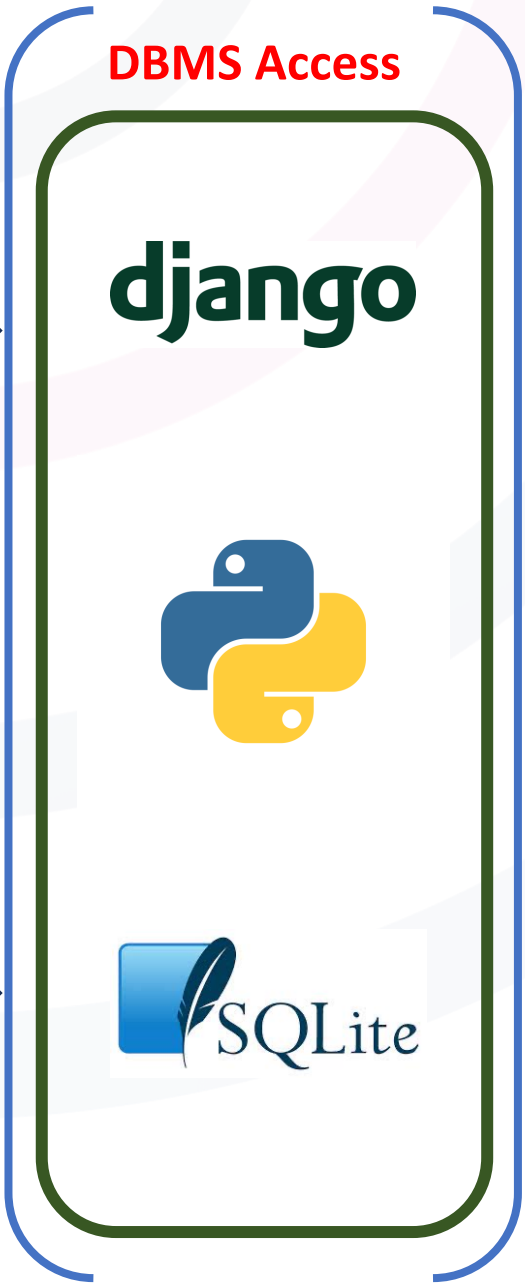
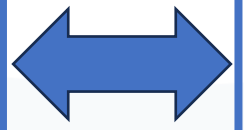
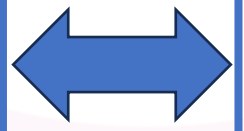
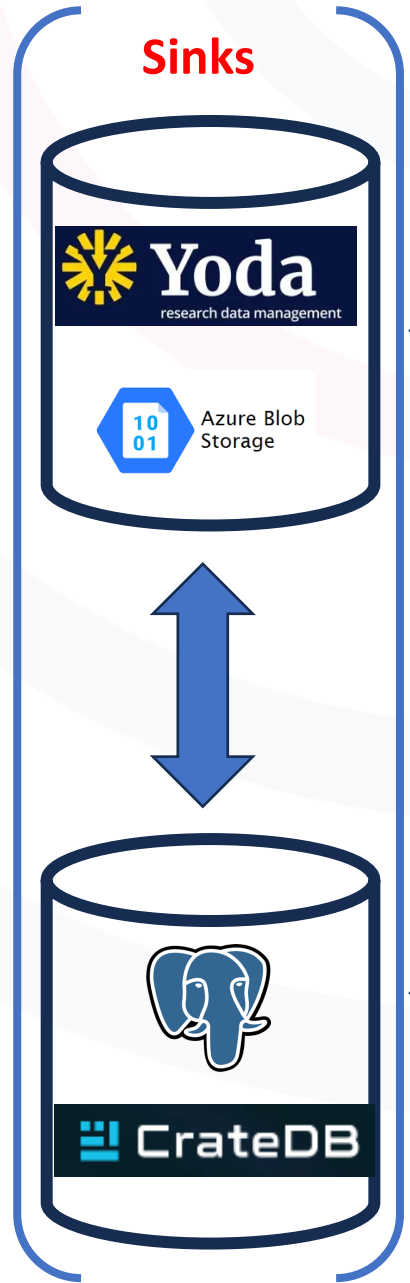
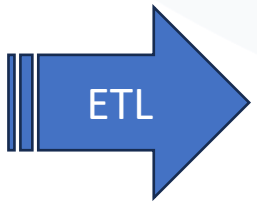
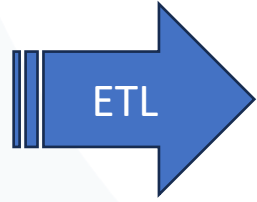
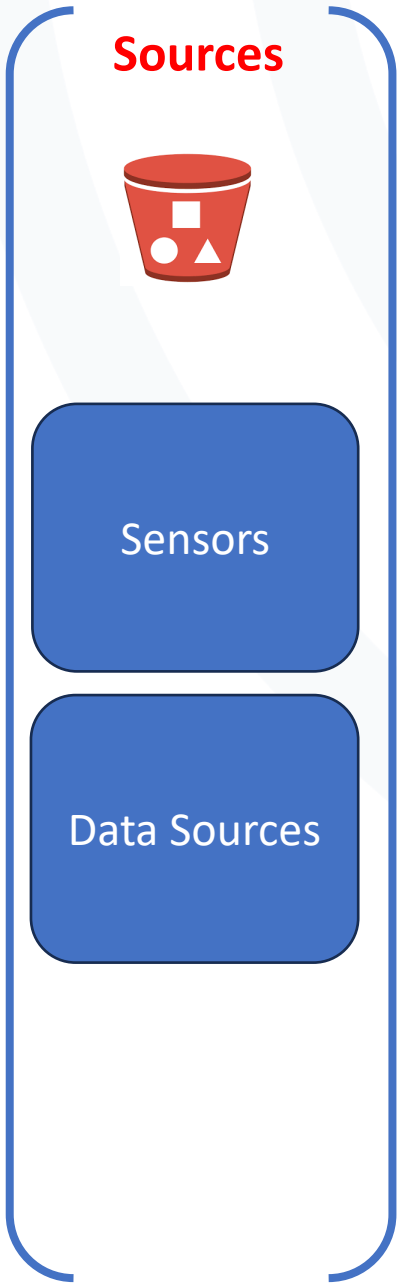
Apache  
Airflow

# Apache *HOP*

- A.k.a. Kettle a.k.a. Pentaho Data Integration
- Kettle founded by Matt Casters early 2000's
- 2005: Open Sourced & 'acquired' by Pentaho
- 2015: Pentaho acquired by Hitachi
- 2019: Forked by community
- 2020: Apache Incubator project
- 2022: Apache Top Level project (Jan 18)
- <https://hop.apache.org>





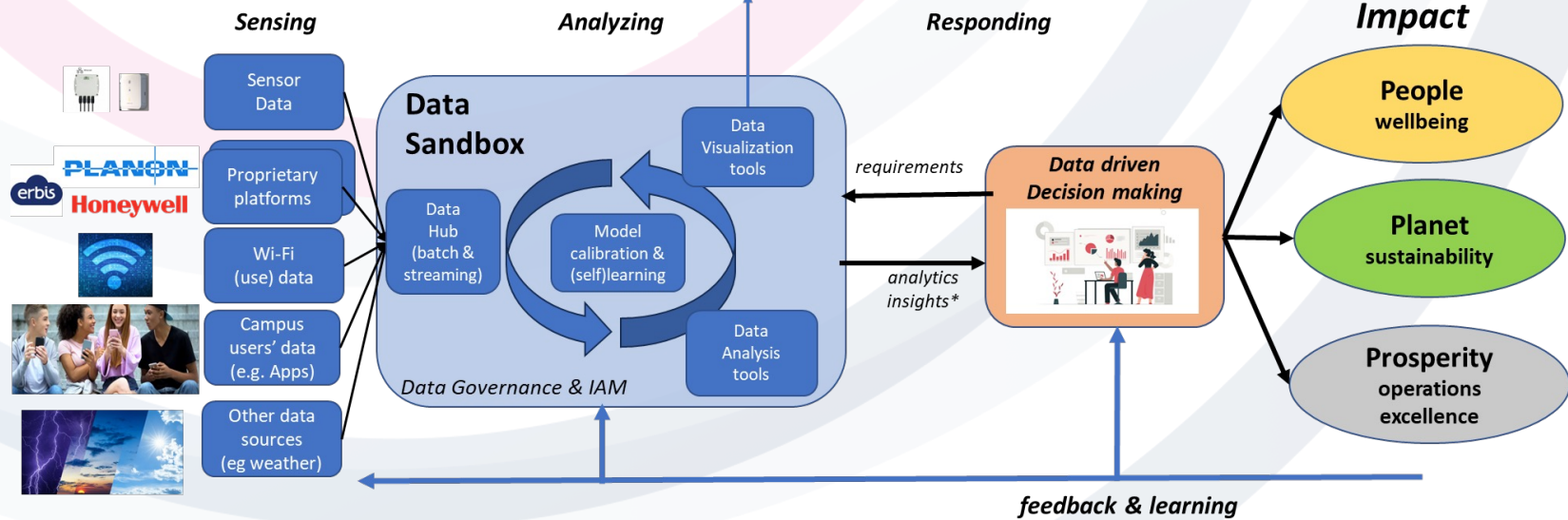


ECDA INITIATIVE

# Smart Campus

## Erasmus University Rotterdam

Make your campus more sustainable with data

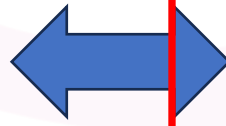
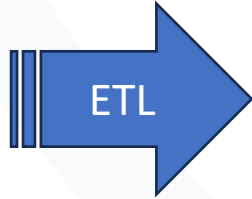


\* Descriptive, Diagnostic, Predictive, Prescriptive, Automated optimization

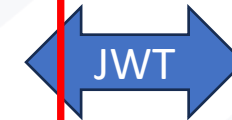
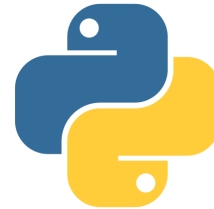


Sensor data  
From API  
System of the  
Measurement

HOP



django



## MEASUREMENTS

improving workplaces worldwide

/ Solutions

/ Our Expertise

/ Cases

[Get an estimate](#)

### Improving Office climate conditions Worldwide

By providing decision-making data, we contribute to better performing workplaces, positively impacting both the quality of work and life.

Over  
**+1.2 million**  
improved  
workspaces



We have  
**+36 million**  
datapoints

We work in  
**43 countries**





In Apache HOP, the data is extracted from API Swagger (Measurement Company) and transferred to CrateDB

The screenshot displays the Apache HOP interface with a data pipeline and a dialog box for the 'Get variables' transform.

**Data Pipeline:**

- Generate rows** (top left) feeds into **JSON input** (top middle).
- JSON input** feeds into **Filter rows** (top right).
- Filter rows** has two outgoing paths:
  - A green path to **Occupancy** (top right), which then feeds into **Occupancy sensor** (top right).
  - An orange path to **peopel counter** (bottom middle), which then feeds into **Peopel sensor** (bottom right).
- Get variables** (bottom left) feeds into **HTTP client** (bottom middle), which feeds into **JSON input**.

**Get variables Dialog Box:**

Transform name:

Fields:

	Name	Variable	Type	Format	Length	Precision	Currency	Decimal	Group	Trim type
1	x-token	\${x-token}	String							none
2	startTime	\${startTime}	String							none
3	endTime	\${endTime}	String							none

Buttons:



Project: sensors Environment: env\_sensors

bulk\_direct\_ecda\_sensors\_snow X

Zoom: 100% Unit test:

```
graph LR; GV[Get variables] --> HC[HTTP client]; HC --> JI[JSON input]; JI --> FR[Filter rows]; FR --> OCC[Occupancy]; OCC --> OS[Occupancy sensor]; GV --> JI; JI --> PC[peopel counter]; PC --> PS[Peopel sensor];
```

HTTP client

Transform name: HTTP client

General Fields

Settings

- URL: \${URL\_GET}
- Accept URL from field?
- Ignore SSL certificate check?
- URL field name: [dropdown]
- Encoding (empty means standard): UTF-8
- Connection timeout: 1000000
- Socket timeout: 1000000
- Connection close wait time: -1

Output fields

- Result field name: result
- HTTP status code field name: [dropdown]
- Response time (milliseconds) field name: [dropdown]
- Response header field name: [dropdown]

HTTP authentication

- Http Login: [input]
- HTTP Password: [input]

Proxy to use

- Proxy Host: [input]
- Proxy Port: [input]

Help OK Cancel



Project: sensors Environment: env\_sensors

bulk\_direct\_ecda\_sensors\_snow X

Zoom: 100% Unit test:

```
graph LR; GV[Get variables] --> JSON[JSON input]; HTTP[HTTP client] --> JSON; JSON --> FR[Filter rows]; FR --> OCC[Occupancy]; FR --> PC[peopel counter]; OCC --> OS[Occupancy sensor]; PC --> PS[Peopel sensor];
```

Snowflake Bulk Loader

Transform name: Occupancy sensor

Bulk Loader | Data type | Fields

Connection: ecda-snowflake

Schema: PUBLIC

Table name: TBL\_OCCUPANCY

Staging location type: User Location

Internal stage name:

Work Directory: \${java.io.tmpdir}

On error: Abort

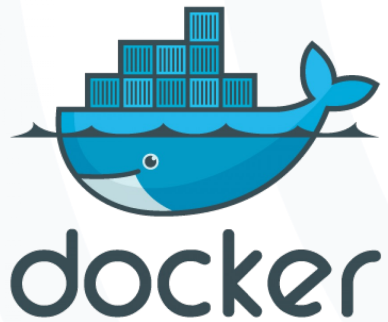
Error limit:

Split load files every ... rows: 20000

Remove files after load:

Help OK SQL Cancel





```
Dockerfile > ...
1  # Use the base Apache Hop image
2  FROM apache/hop:latest
3
4  # Set the working directory in the container
5  WORKDIR /project
6
7  # Copy your project files into the container.
8  # Adjust the paths according to where your Hop project and environment files are located on your build context.
9  COPY /sensors /project
10 COPY /sensors /project-config
11
12 # Set environment variables to match your original DockerOperator configuration
13 ENV HOP_LOG_LEVEL="Basic" \
14     HOP_FILE_PATH="/project/bulk_ecda_sensors_snow_hourly.hwf" \
15     HOP_PROJECT_DIRECTORY="/project" \
16     HOP_PROJECT_NAME="sensors" \
17     HOP_ENVIRONMENT_NAME="env_sensors" \
18     HOP_ENVIRONMENT_CONFIG_FILE_NAME_PATHS="/project-config/env_sensors.json" \
19     HOP_RUN_CONFIG="local"
```



Az login

```
Docker build -t ecdacr.azurecr.io/smartcampus-hop:2.1.0
```

```
Docker push ecdacr.azurecr.io/smartcampus-hop:2.1.0
```

Microsoft Azure Search resources, services, and docs (G+)

Home > ecdacr | Repositories > smartcampus-hop >

### smartcampus-hop

Repository

Refresh

**Essentials**

Repository: smartcampus-hop

Last updated date: 1/16/2024, 2:39 PM GMT+1

Tag count: 4

Manifest count: 4

Search to filter tags ...

Tags	Digest
552a198155650dd021...	sha256:7a8afc186c65
ea99b44673a1014daa...	sha256:18928e7d4d7
2.1.0	sha256:9617666a17b
1.1.1	sha256:93d5de06e10

### smartcampus-hop:2.1.0

sha256:9617666a17b2262b650545b444ecd70aeb8e84a0bd662d6438c11428511ccd88

Create streaming artifact Refresh

**Essentials**

Repository: smartcampus-hop Digest: sha256:9617666a17b2262b650545b444ecd70aeb8e84a0bd662d6438c11428511ccd88

Tag: 2.1.0 Manifest creation date: 12/18/2023, 1:14 PM GMT+1

Tag creation date: 12/18/2023, 1:14 PM GMT+1 Platform: linux / amd64

Tag last updated date: 12/18/2023, 1:14 PM GMT+1 Media type: application/vnd.docker.distribution

Manifest Referrers

Search to filter digests ...

Artifacts	Digest	Manifest type
No result		



# Apache Airflow

```
with DAG('smartcampus_k8s_hop_hourly',
         default_args=default_args,
         schedule_interval='0 * * * *',
         catchup=False,
         is_paused_upon_creation=False,
         max_active_runs=1) as dag:

    start_dag = EmptyOperator(
        task_id='start_dag'
    )

    end_dag = EmptyOperator(
        task_id='end_dag'
    )

    hop = KubernetesPodOperator(
        task_id='smartcampus_k8s_hop',
        image='ecdacr.azurecr.io/smartcampus-hop:2.1.0', #modify the name, repo, and version of the container accordingly
        namespace='airflow', # Change to your namespace if different
        name='smartcampus-hop-task',
        get_logs=True,
        image_pull_secrets='my-acr-secret',
    )

    start_dag >> hop >> end_dag
```





# kubernetes

```
! override-values.yaml
1   dags:
2     gitSync:
3       enabled: true
4       repo: git@github.com:Erasmus-Data-Collaboratory/eur-smart-campus.git
5       branch: main
6       subPath: "dags"
7       sshKeySecret: airflow-ssh-secret
8   extraSecrets:
9     airflow-ssh-secret:
10      data: |
11        | gitSshKey: ''
12
```



<> Code Issues Pull requests Actions Projects 1 Security Insights Settings

Files

main

Go to file

- .github
- airflow-helm
  - airflow-nginx
    - README.md
    - git-credentials.yaml
    - my-custom-values.yml
    - variable-configmap.yml
  - dags
    - .gitignore
    - README.md

Apache-hop / airflow-helm /

fhasanabadi add doc for nginx and cert-manager b6da1d4 · 4 months ago History

Name	Last commit message	Last commit date
..		
airflow-nginx	add doc for nginx and cert-manager	4 months ago
README.md	add doc for nginx and cert-manager	4 months ago
git-credentials.yaml	add helm airflow instructions	4 months ago
my-custom-values.yml	add helm airflow instructions	4 months ago
variable-configmap.yml	add helm airflow instructions	4 months ago

README.md

## Up and running Airflow on AKS cluster with gitSync

### Prepare and create the cluster

Apache-hop (Private)

main 1 Branch 0 Tags

Go to file Add file Code

actions-user commit from Erasmus-Data-Collaboratory/eur-smart-campus 698bca8 · 2 weeks ago 119 Commits

.github/workflows	Update image-build-push.yml	4 months ago
airflow-helm	add doc for nginx and cert-manager	4 months ago
dags	commit from Erasmus-Data-Collaboratory/eur-smart-campus	2 weeks ago
.gitignore	Added .DS_Store to .gitignore	5 months ago
README.md	edit readme	4 months ago

README

### Apache-hop

Repository for Hop data projects  
Including Apache HOP and Apache Airflow  
The repo contains the Airflow DAGs for syncing with a Airflow instance running on AKS. The DAGs are located in the dag folder according to the name of the project. For example, the DAGs belong to the project smartcampus are in /dags/smartcampus .

## ECDA Administration

## AUTHENTICATION AND AUTHORIZATION

[Groups](#) [+ Add](#) [Change](#)[Users](#) [+ Add](#) [Change](#)

## OTP\_TOTP

[TOTP devices](#) [+ Add](#) [Change](#)

## SMARTCAMPUS

[Databases](#) [+ Add](#) [Change](#)[Projects](#) [+ Add](#) [Change](#)[Rooms](#) [+ Add](#) [Change](#)[Sensor allocations](#) [+ Add](#) [Change](#)[Sensors](#) [+ Add](#) [Change](#)

## ADMIN

## Recent actions

## My actions

None available

## USERS

# API System for ECDA Projects

1.0.0 OAS 3.0

</authentication/api/schema/>

To get a Token API key, first request for an account. Your API key is a long alphanumeric string. Include this API key in the Authorization header with the Token key for all authentication endpoints.

[Authorize](#) 

## authentication

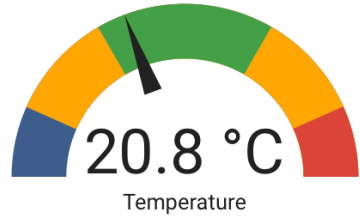
[GET](#) /authentication/api/schema/ [POST](#) /authentication/api/token/[POST](#) /authentication/api/token/refresh/

## smartcampus

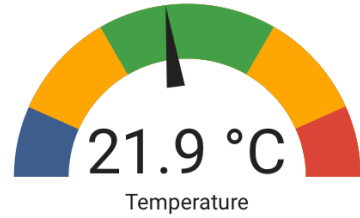
[GET](#) /smartcampus/roomsensors/ [GET](#) /smartcampus/sensors/ 



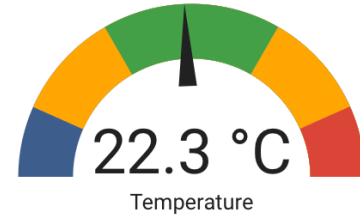
Y1-08



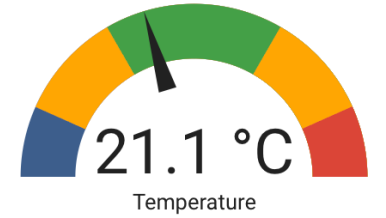
Y1-10



Y1-12



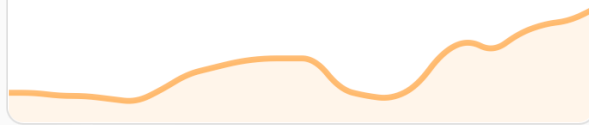
Y1-13



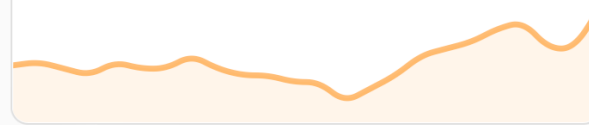
Temperature Trend  
20.8 °C



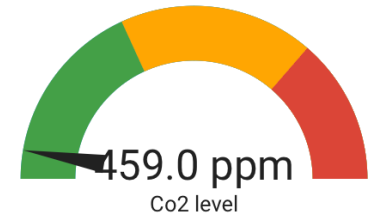
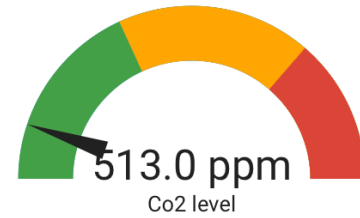
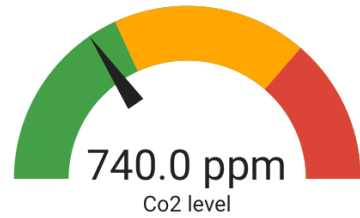
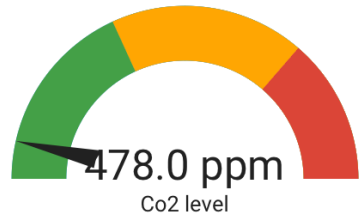
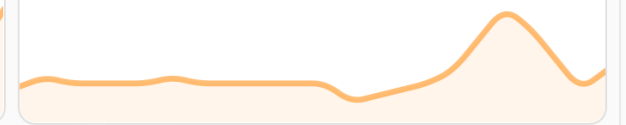
Temperature Trend  
21.9 °C



Temperature Trend  
22.3 °C



Temperature Trend  
21.1 °C



 Humidity 44.0 %

 Presence Detected


 Desk 1 38.8 W

 Desk 2 10.6 W

 Coffee Machine 0.0 W

 Humidity 47.0 %

 Presence Clear

 Humidity 42.0 %

 Presence Detected

 Desk 1 73.3 W

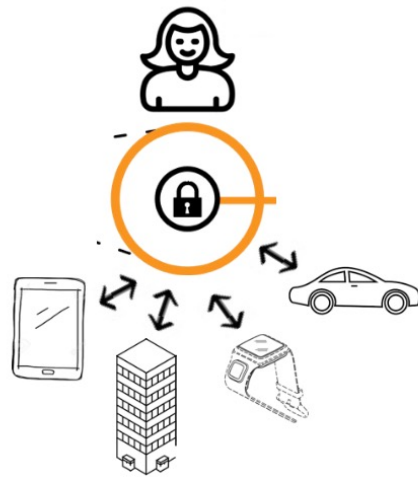
 Desk 2 18.1 W

 Humidity 42.0 %

 Presence Clear

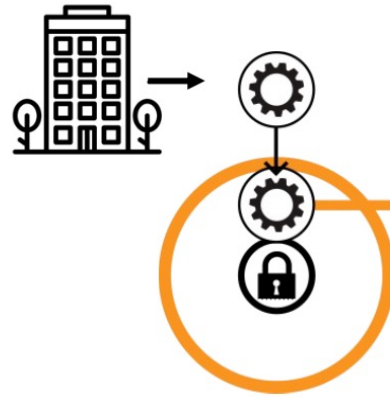
# New Pilot: **bubl cloud**

1 Demonstratable  
Security & Privacy



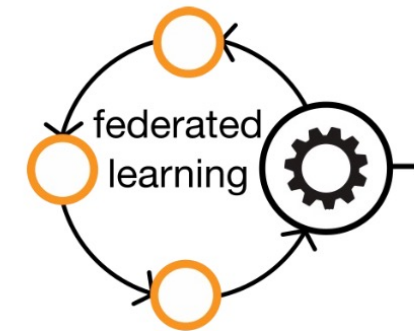
Personal encrypted  
data vault

2 Use AI on combined  
sensitive personal data



Sending the service  
to the data

3 Get detailed insights  
privacy safe



Privacy safe  
exchange



# Erasmus Centre for **Data Analytics**

## Questions?