# Data challenges in scientific computing

**Physics is overrated**

João Alvim, Anton Schwarz, Konrad Kockler
DESY Hamburg, 04.09.2024

HELMHOLTZ

DESY.

# Scientific Data Challenges
## Data processing

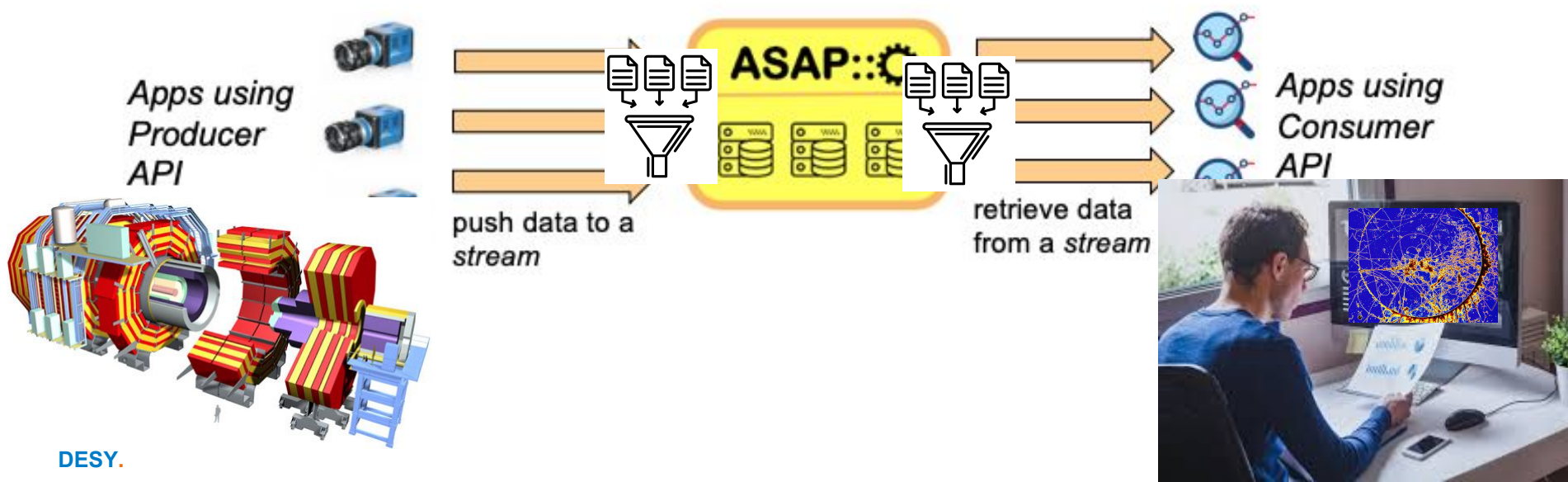| Ingest | Sharing & Exchange | Long Term Preservation | Analysis |
|---|---|---|---|
| - High data ingest rate<br>- Multiple parallel streams<br>- High durability<br>- Effective handling of large number of files | - 3$^{rd}$ party copy<br>- Effective WAN Access<br>- In-flight data protection<br>- Identity federation<br>- Access control | - High Reliability<br>- Self-healing<br>- Automatic technology migration<br>- Persistent identifier | - High CPU efficiency<br>- Unstructured access patterns<br>- Standard access protocols<br>- Access control<br>- Local user management |

# Improving ASAP::O Monitoring

João Alvim, Mikhail Karnevskiy

# ASAP::O

## What is it?

1.  **High performance distributed streaming platform**
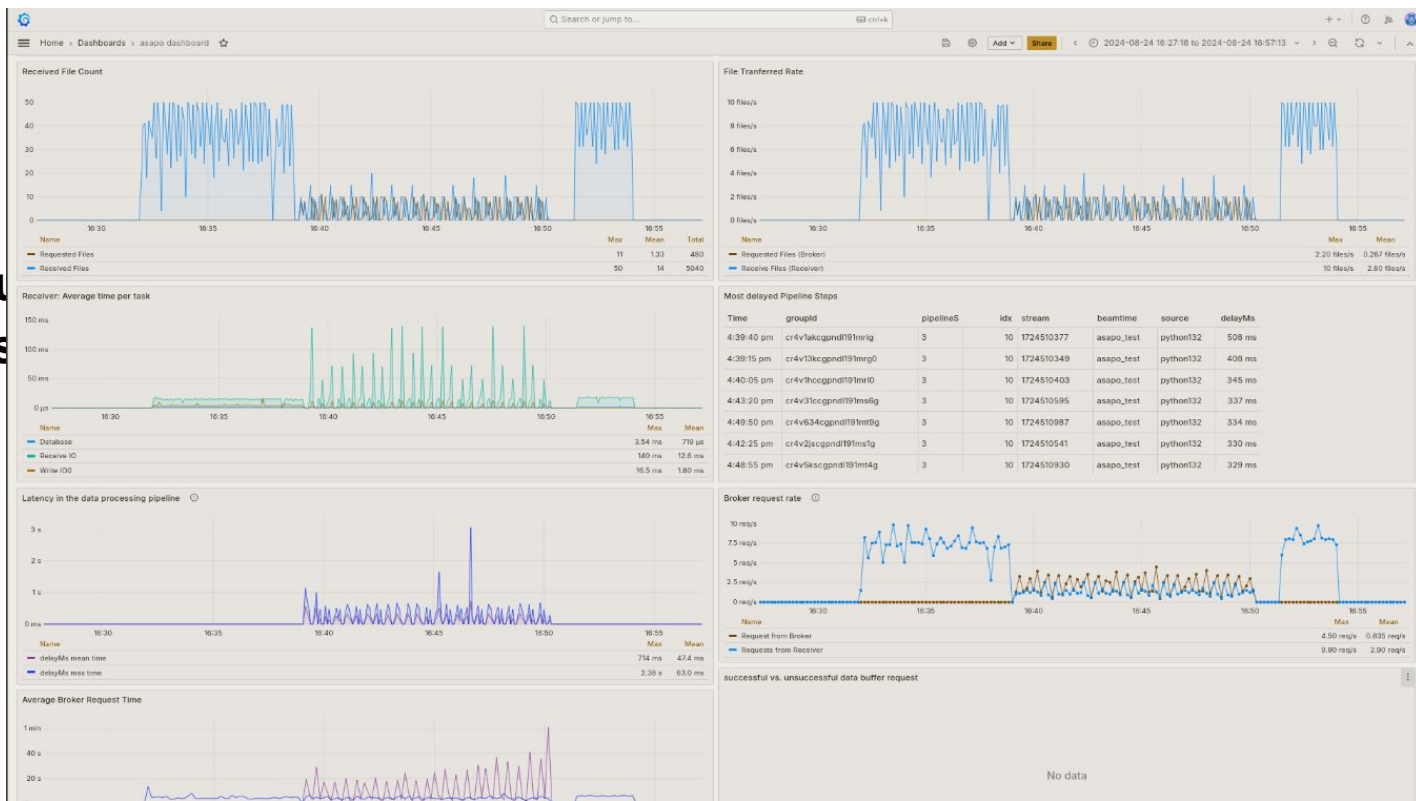2.  **Provides API to ingest/retrieve data to the system**

# Monitoring Improvements

**New Approach**



1. **Grafana**
   + **flexibility, bu**
   + **components**
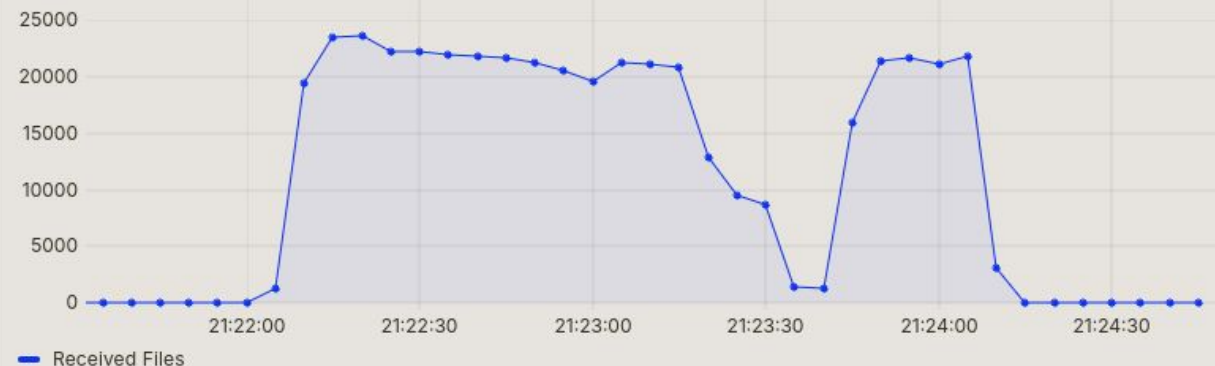   + **reliable**
   - **bugs**

2. **Collect new**

## Example case

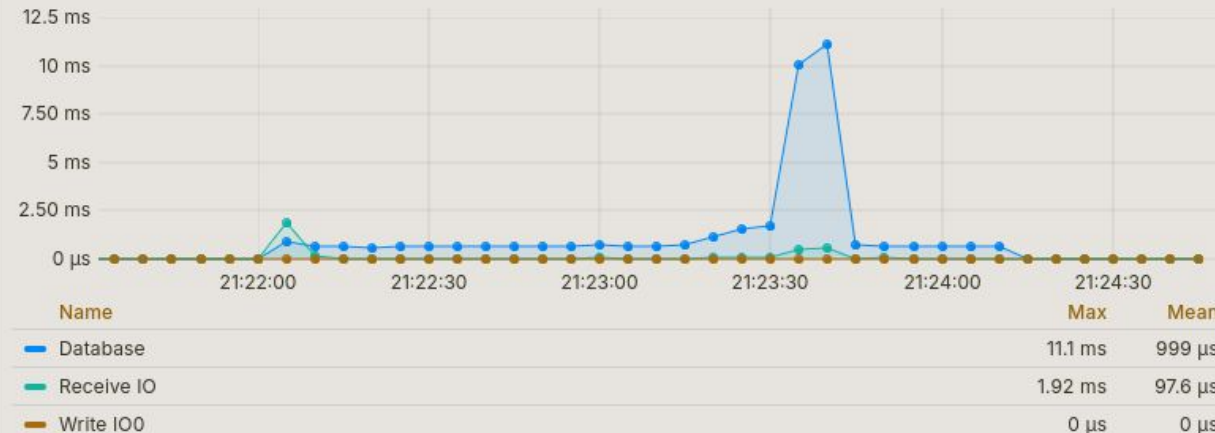**Amount of received files dropped significantly, because the time to access database increased**

**Raise Suspicion Potential Bottleneck**
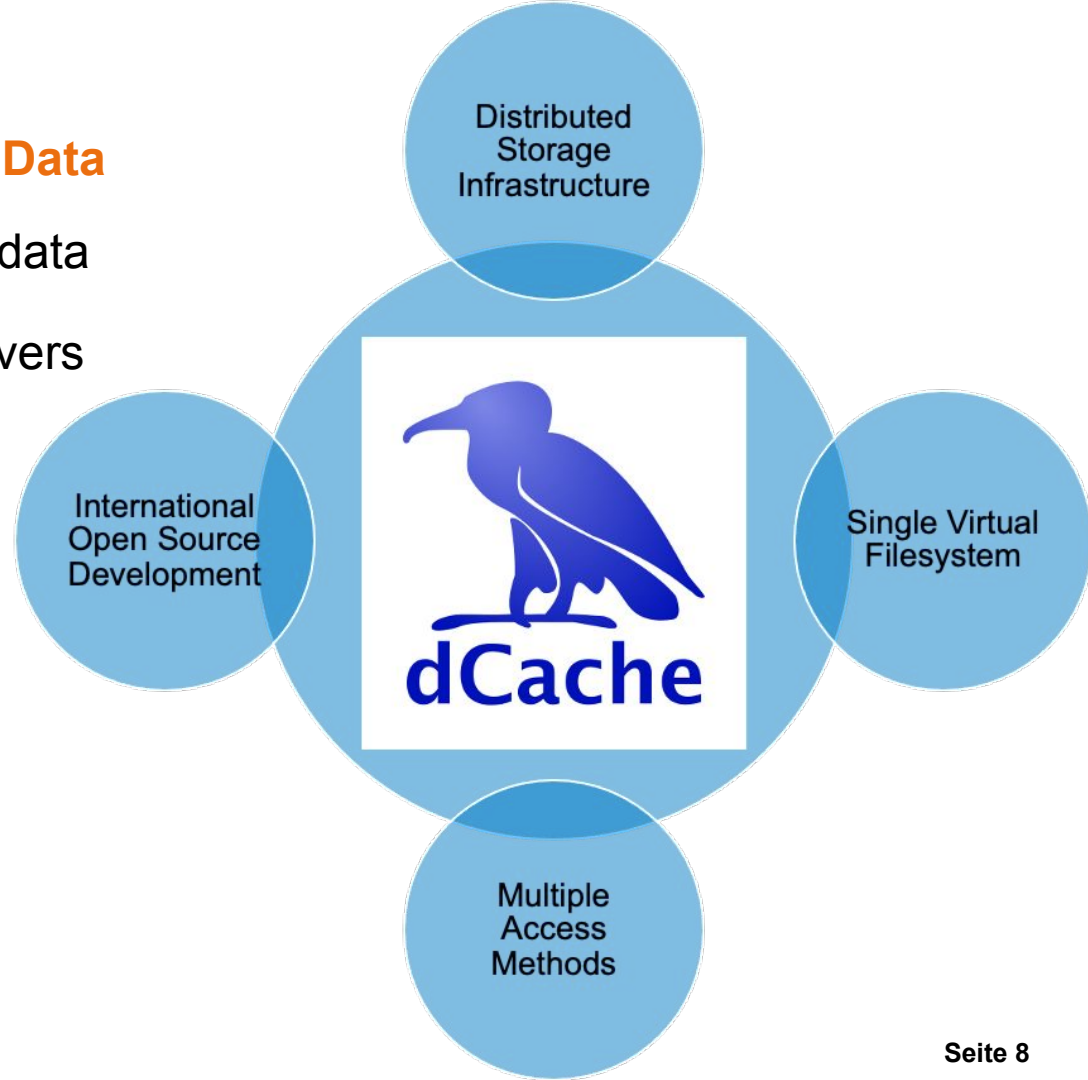
**=> More Reliable Software**

# Scientific Data Management with dCache
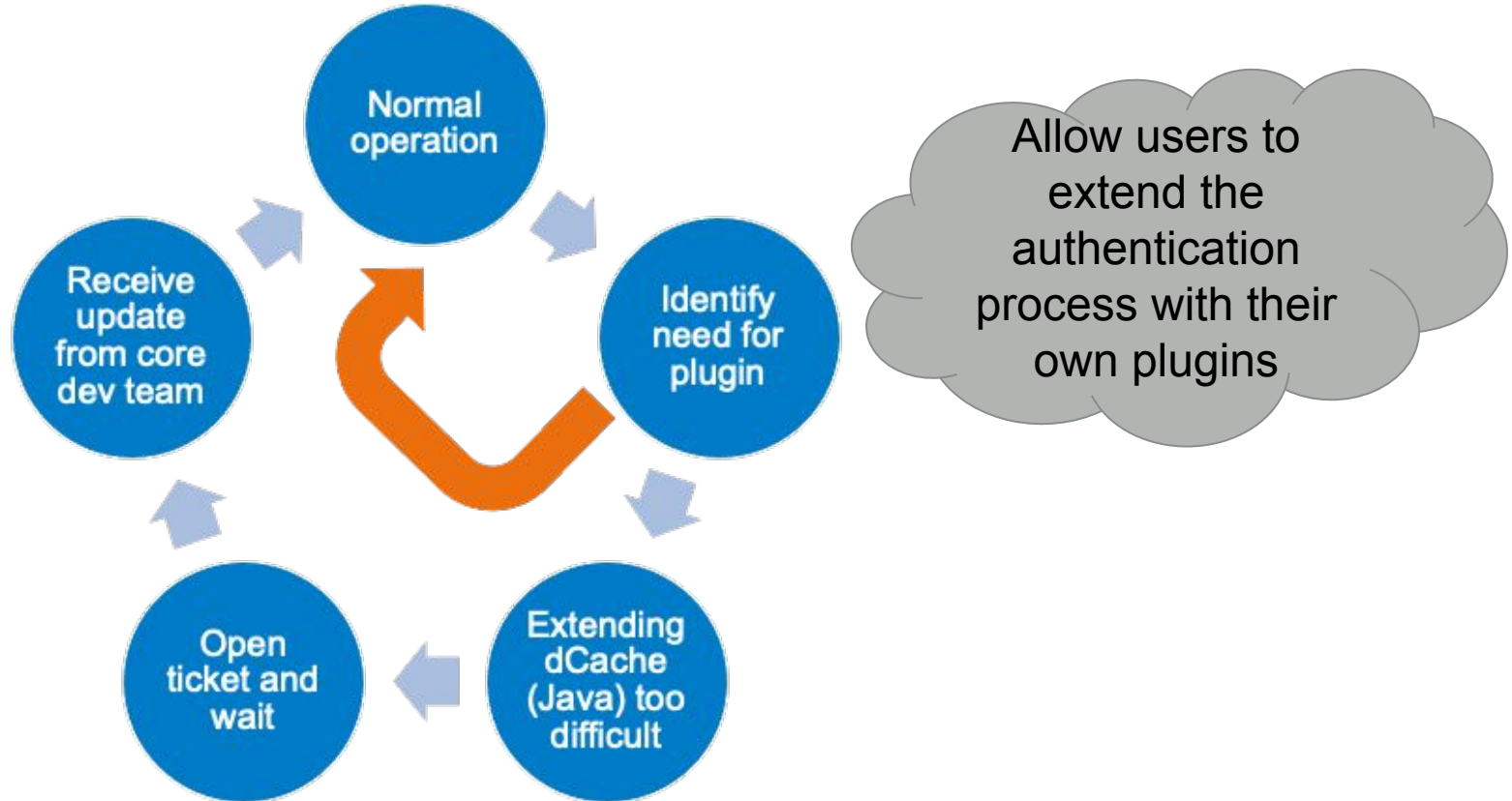
**Anton Schwarz, Tigran Mkrtchyan**

# dCache
## A System to Store and Retrieve Data

• System for storing and retrieving data

• Data is distributed over many servers

• In full production since 2001

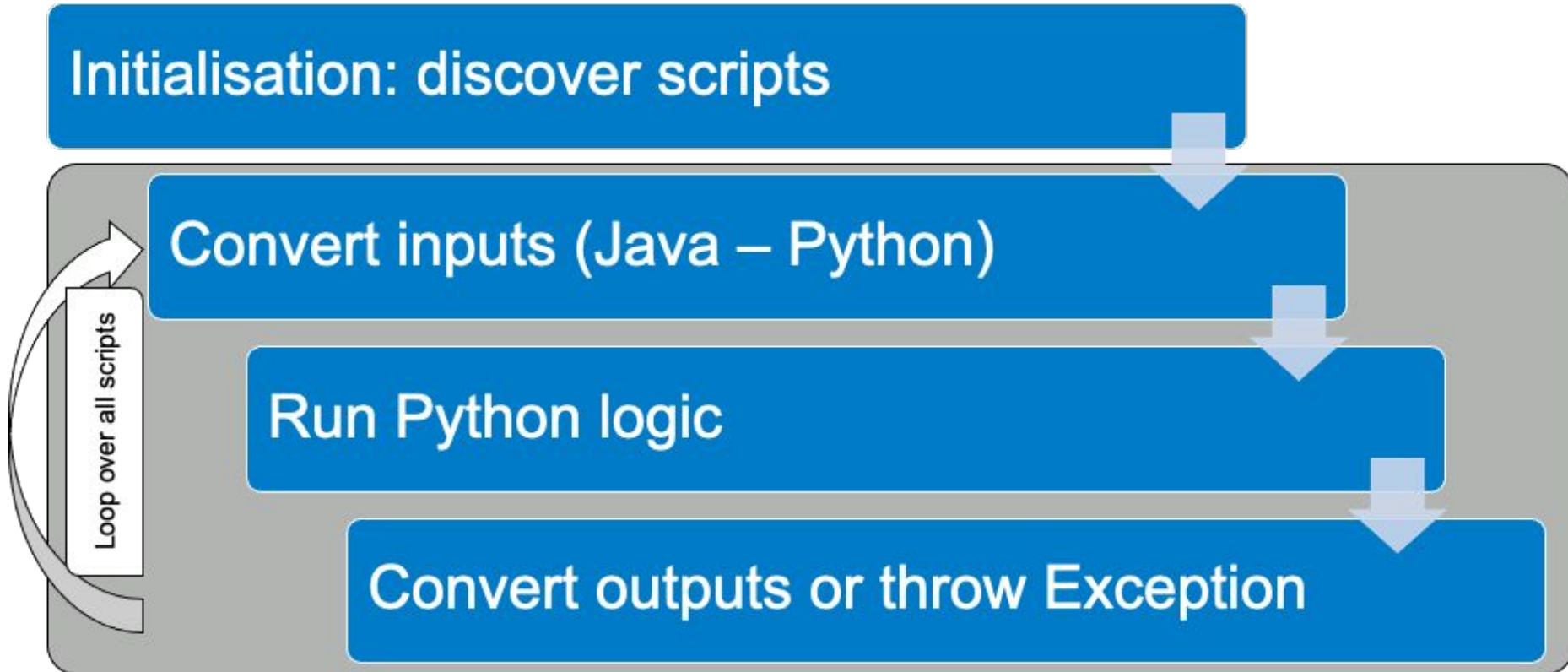• Written in Java

• Authentication and Authorisation

Distributed Storage Infrastructure

International Open Source Development

Single Virtual Filesystem

dCache

Multiple Access Methods

# Adapting the Authentication Process
## Motivation for my project



Normal operation

Identify need for plugin

Extending dCache (Java) too difficult

Open ticket and wait

Receive update from core dev team

Allow users to extend the authentication process with their own plugins

# Implementation

Initialisation: discover scripts

Loop over all scripts

Convert inputs (Java – Python)

Run Python logic

Convert outputs or throw Exception

# Timeline
## Where are we?



Make Jython Work → Write Plugin → Test, Document, Release → Find and satisfy users

# Research facility 2.0 Sustainable computing

**Konrad Kockler, Martin Gasthuber**

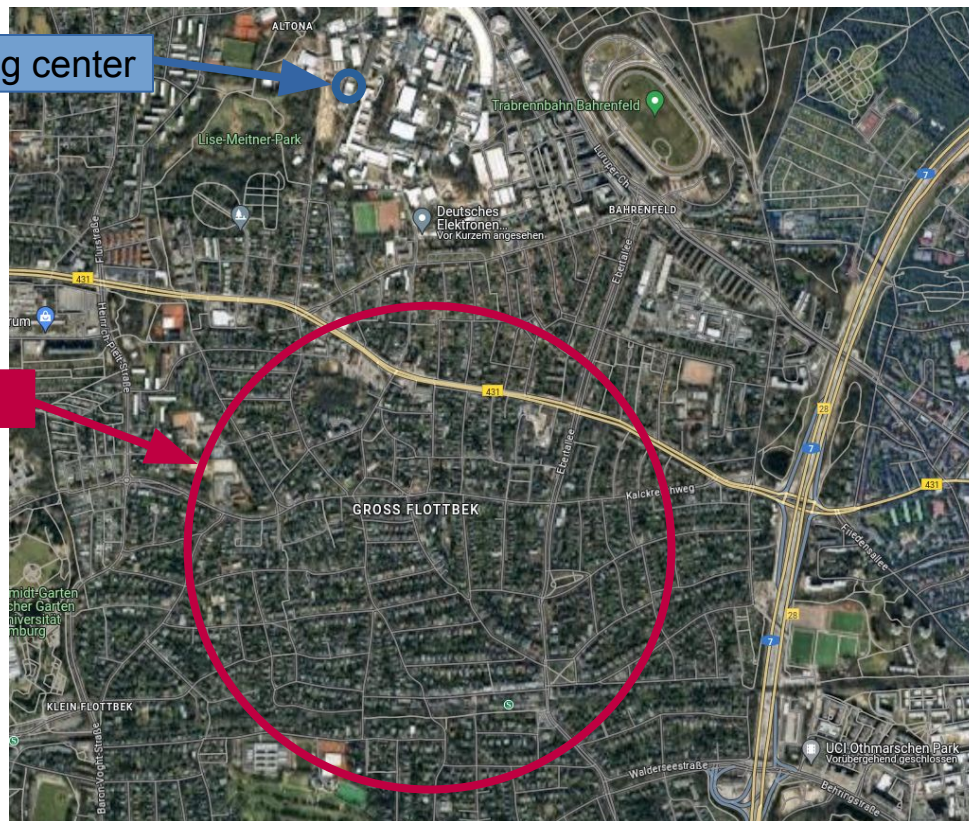# Current energy consumption

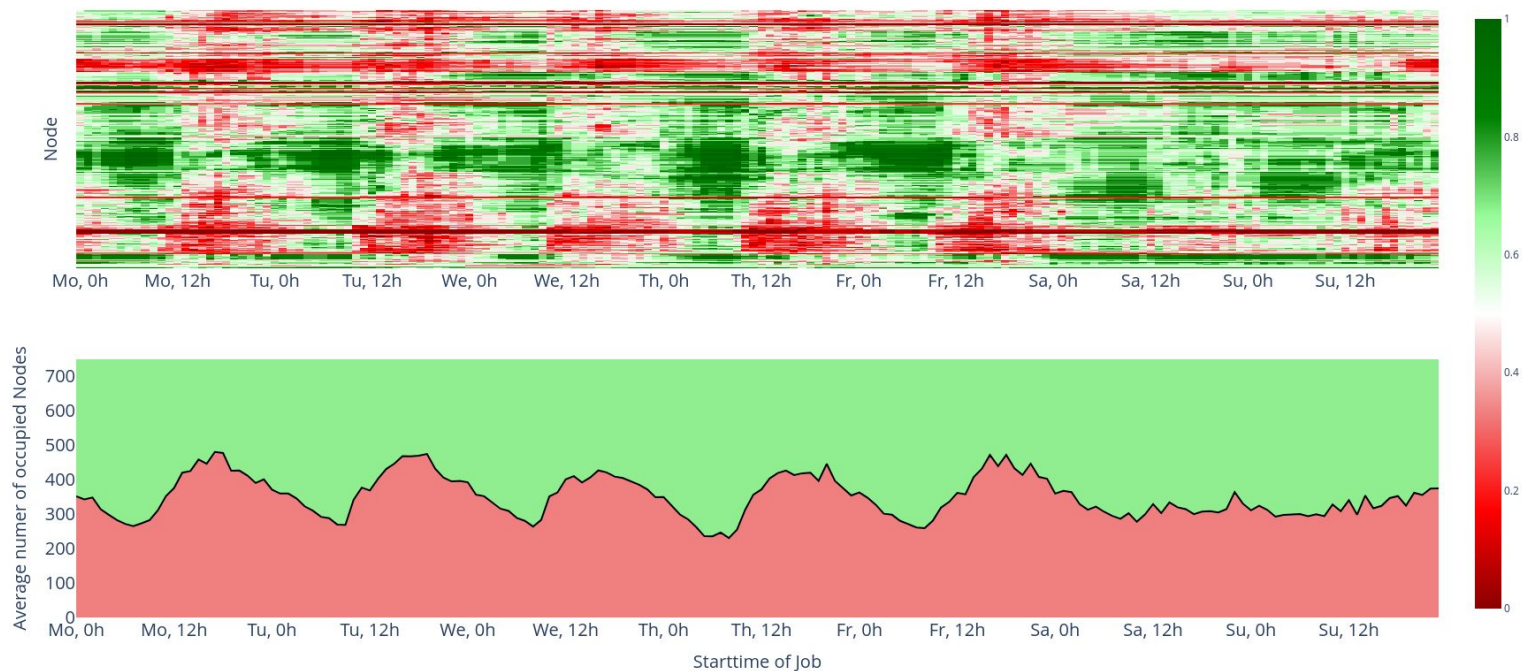**Per year (2023):**

## 12.9GWh

## 5160 tCO$_2$

**This is:**

- 250 000 trees

- 3 700 Single-family homes

- 1 Groß Flottbek



DESY computing center

Groß Flottbek

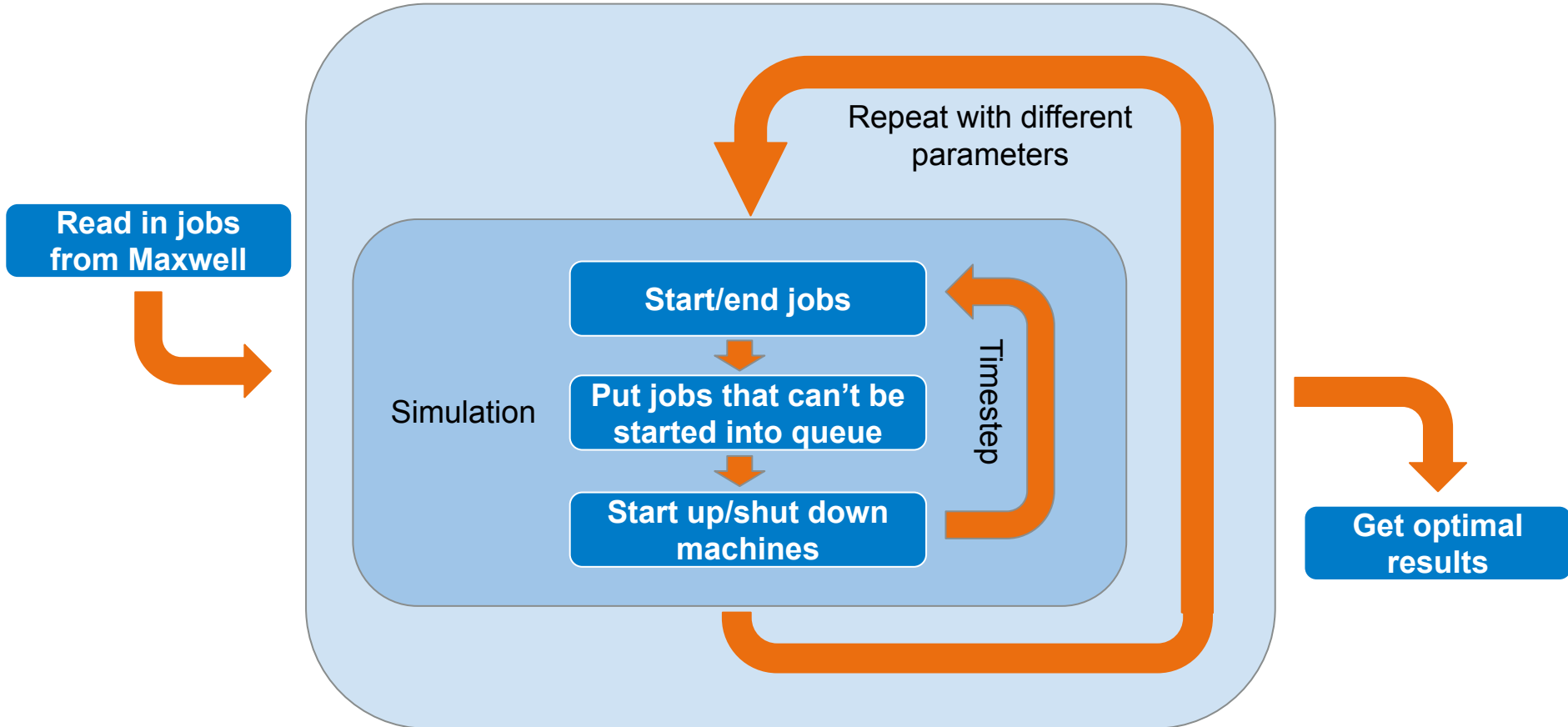Source: maps.google.com

# How much is actually used?

# Increasing the energy and CO$_2$ efficiency of the Maxwell cluster

1. **More efficient utilization**    ➤    **Put jobs from other grids on Maxwell**

2. **Switching off idle nodes**    ➤    **Up to 20% savings without the user noticing**

3. **Throttling CPUs**    ➤    **Save Electricity at night, when it has a high Carbon intensity**

**Huge savings potential!**
**Save ~30% in Carbon emissions**

# How my simulation worked



Read in jobs from Maxwell

Repeat with different parameters

Simulation

Start/end jobs

Put jobs that can't be started into queue

Start up/shut down machines

Timestep

Get optimal results

**Contact**

Deutsches Elektronen-
Synchrotron DESY

www.desy.de

João Alvim      - joao.alvim@desy.de
Anton Schwarz   - anton.schwarz@desy.de
Konrad Kockler   - konrad.kockler@desy.de

Scientific Computing