# DATA CURATION IN KCDC

V. Tokareva*, A. Haungs, D.Kang, F. Polgart, D. Wochele, J. Wochele

*victoria.tokareva@kit.edu

INSTITUT FÜR ASTROTEILCHENPHYSIK (IAP)
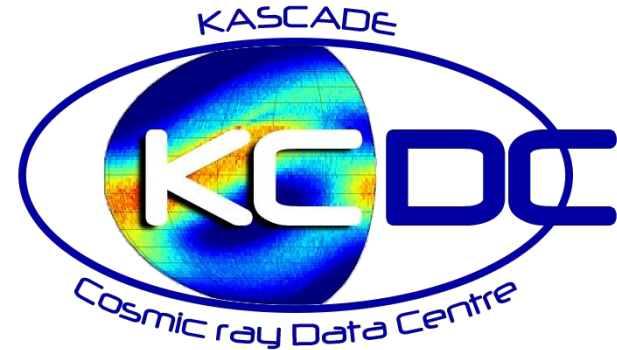
2 May 2022, TA4 WP2 meeting

# Content

- KCDC overview

- Data provided by KCDC

- Software architecture, data and metadata flows

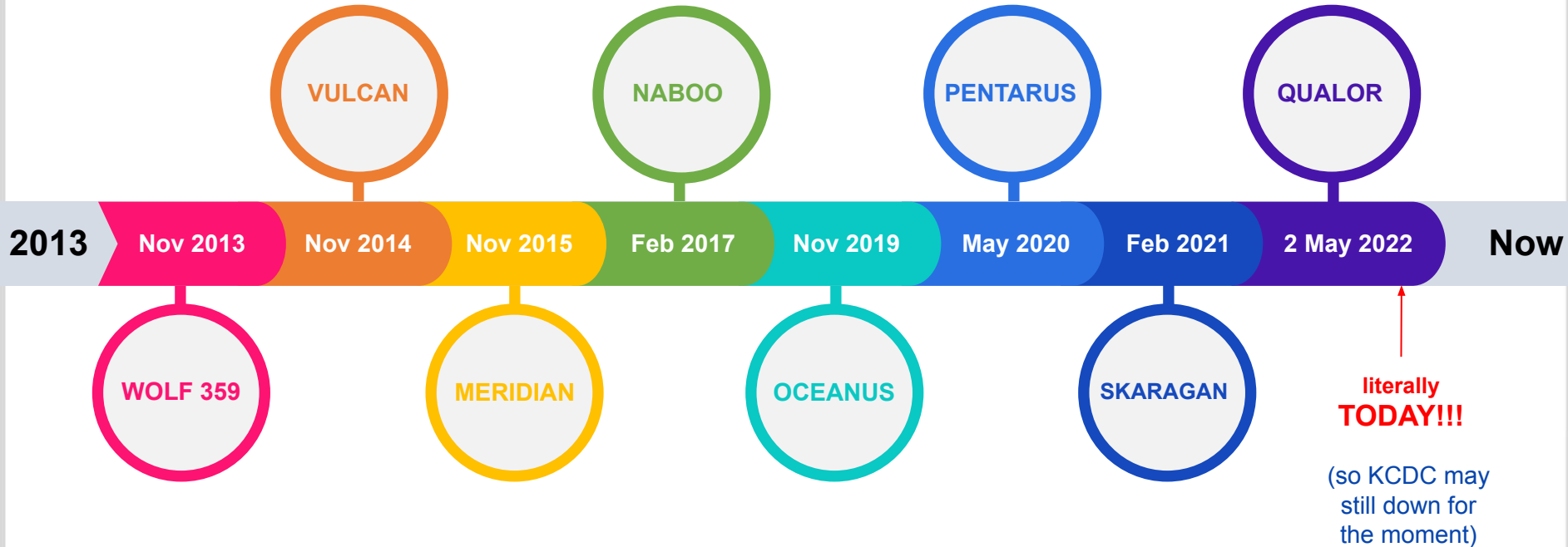- Jupyter Notebooks for data analysis

# KASCADE Cosmic-ray Data Center (KCDC)

- First released in 2013
- Aimed to provide free, unlimited, reliable open access to the data of various experiments measuring cosmic radiation by different methods and techniques both for scientists and the broad public

- Functions:
  - Data archive
  - Data analysis platform
  - Information center
  - Outreach platform

- Features:
  - Open data access
  - Allows custom data cuts
  - Ensures analysis reproducibility
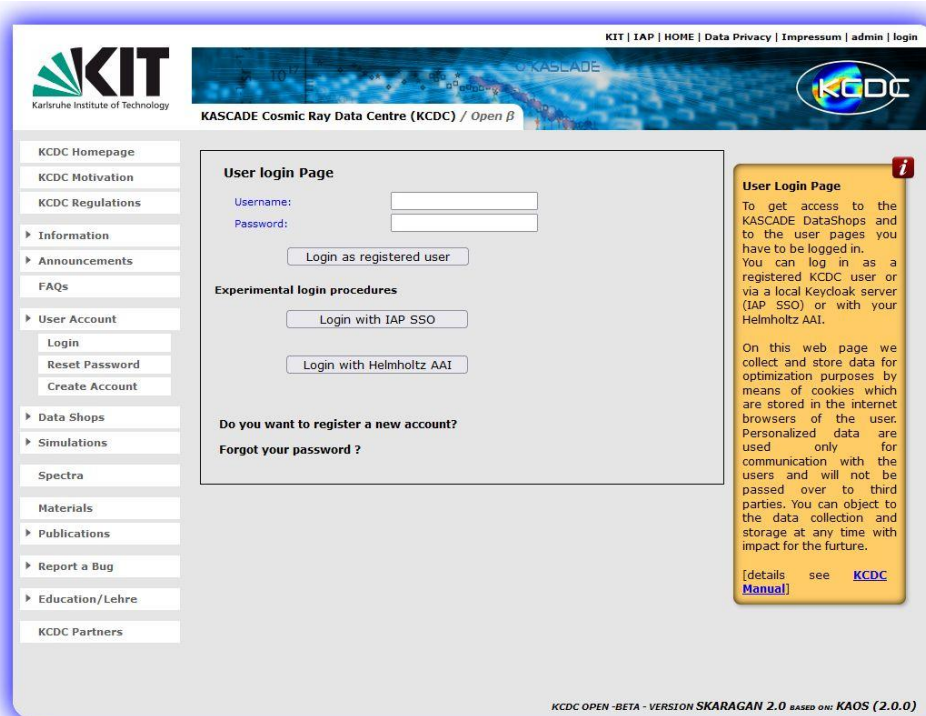  - Built on modern open source web technologies

https://kcdc.iap.kit.edu/

# KCDC timeline



2013

VULCAN

NABOO

PENTARUS

QUALOR

Nov 2013 | Nov 2014 | Nov 2015 | Feb 2017 | Nov 2019 | May 2020 | Feb 2021 | 2 May 2022 | Now

WOLF 359

MERIDIAN

OCEANUS

SKARAGAN

literally
TODAY!!!

(so KCDC may
still down for
the moment)

# What's new in QUALOR?



- Login via Helmholtz AAI and via a Keycloak server
- 4-shard database
- New QGSJet-II-04 Simulations

# Data overview

| Setup/ Detector component | Experimental data | | Simulations | |
|---|---|---|---|---|
| | Events | Size | Events | Size |
| KASCADE | 433 209 340 | 3 200 GB | 21 388 000 | 25 GB |
| GRANDE | 35 310 393 | 260 GB | 3 545 000 | 4 GB |
| COMBINED | 15 635 550 | 120 GB | 1 590 000 | 2 GB |
| LOPES | 3 058 | 25 MB | — | — |
| MAKET-ANI | 2 682 264 | 1 GB | — | — |

\* This table doesn't take into account the new QGSJet-II-04 simulations added in QUALOR

# How published: DOI



doi.org/10.17616/R3T
S4P

re3data.org

Repository details

## KASCADE Cosmic Ray Data Centre

General | Institutions | Terms | Standards

| | |
|---|---|
| Name of repository | **KASCADE Cosmic Ray Data Centre** |
| Additional name(s) | KCDC<br>Karlsruhe Shower Core and Array Detector |
| Repository URL | https://kcdc.ikp.kit.edu/ |
| Subject(s) | Particles, Nuclei and Fields  Astrophysics and Astronomy  Physics  Natural Sciences |
| Description | The aim of the project KCDC (KASCADE Cosmic Ray Data Centre) is the installation and establishment of a public data centre for high-energy astroparticle physics based on the data of the KASCADE experiment. KASCADE was a very successful large detector array which recorded data during more than 20 years on site of the KIT-Campus North, Karlsruhe, Germany (formerly Forschungszentrum, Karlsruhe) at 49,1°N, 8,4°O; 110m a.s.l. KASCADE collected within its lifetime more than 1.7 billion events of which some 425.000.000 survived all quality cuts. Initially about 160 million events are available here for public usage. |
| Contact | ikp-kcdc@lists.kit.edu |
| Content type(s) | Standard office documents  Plain text  Archived data  Scientific and statistical data formats |
| Keyword(s) | KASCADE GRANDE  air shower  astroparticle physics  cosmic rays  hadronic interactions  high-energy physics  large detector array  teaching materials |
| Repository type(s) | institutional |
| Mission statement for designated community | https://kcdc.ikp.kit.edu/static/pdf/kcdc_mainpage/kcdc-Manual.pdf |
| Research data repository language(s) | English German |
| Data and/or service provider | data provider |

# KCDC DataShops

The data sets are organised into so-called datashops:

- KASCADE - contains 'common data' and data from four detector components: KASCADE, GRANDE, CALORIMETER, LOPES
- COMBINED - includes 'common data', data from KASCADE and GRANDE detectors combined for joint analysis as well as data arrays from KASCADE and GRANDE and LOPES quntities
- Maket-Ani - provides quantities from the Maket-Ani setup

New data shops can be added.

# KCDC DataShops and data formats

They are supplied in the following file formats*:

- ASCII - plain text format
- ROOT - object oriented framework developed by CERN
- HDF5 - hierarchical data format

* Selectable by the user and depending on the quantities chosen

# Data quantities examples

**CALORIMETER Quantities**

| Var | Name | Available Data Range | Unit | Representation |
|---|---|---|---|---|
| Nhad | Nr of Hadrons | 0. - 511. | | |
| Ehad | Hadron Energy Sum | 0.; 1.e10 - 1.e16 | eV | log10 -> 10.0 - 16.0 |

**GRANDE Quantities**

| Var | Name | Available Data Range | Unit | Representation |
|---|---|---|---|---|
| Xc | X-Core Position | -500.0 - +100.0 | m | |
| Yc | Y-Core Position | -600.0 - +100.0 | m | |
| Ze | Zenith Angle | 0.0 - 40.0 | ° | |
| Az | Azimuth Angle G | 0.0 - 360.0 | ° | |
| Nch | Number of charged part | 11111. - 1,000,000,000. | | log10 -> 4.0 - 9.0 |
| Nmu | Number of Muons | 1500. - 100,000,000. | | log10 -> 3.2 - 8.0 |
| Age | Shower Age G | -0.385 - +1.485 | | |
| GDeposit | Energy Deposit charged | 0.0 - 100,000.0 | MeV | /station |
| GArrival | Arrival Time | 1000. - 10,000.0 | ns | /station |

# Usage of KCDC datashop GUI



The entry page of the KCDC DataShop pages

# Usage of KCDC datashop GUI



The confirmation page of the KCDC DataShop pages

# Usage of KCDC datashop GUI



The review page of the KCDC DataShop pages

# KCDC user's job data workflow

User GUI/API
- Data selection
- Meta information
- Tutorials
- Downloads

Job system
- Parallel processing
- Scalability

Server infrastructure
- CMS System
- User Management

Administrator interface
- Administration
- Monitoring

Databases
- Providing the data
- Providing selections

# Architecture and technology stack



Database - **MongoDB**; **REST**full API (starting from SKARAGAN release), **JupyterLab** for data analysis

# MongoDB data storage structure



*Wochele, D., Wochele, J., Polgart, F., Tokareva, V., Kang, D., & Haungs, A. Data Structure Adaption from Large-Scale Experiment for Public Re-Use. CEUR-WS (2019) 2406, 114*

# KAOS - KCDC's backend



- **K**arlsruhe **A**stroparticlephysics **O**pen data **S**oftware (KAOS)
- Implemented using a plugin based design with a focus on easy extensibility and modifiability
- Can work as well outside the context of KCDC

*Schoo, S. Energy Spectrum and Mass Composition of Cosmic Rays and How to Publish Air-Shower Data. PhD thesis, 2016, link: https://publikationen.bibliothek.kit.edu/1000055797*

# Django application structure



Django App

Browser

Web server

URLs

Views

Model

Templates

DB server

Static files

# Metadata on KCDC

JSON metadata schema, example of a record from KCDC

```json
392    {
393      "model": "kaos_datashop.quantity",
394      "fields": {
395        "quant_type": "num",
396        "allow_cuts": true,
397        "head_description": "<p class=dcInfoBoxHeaderDS>Zenith Angle Info</p>",
398        "descr_type": "HTML",
399        "composite_data_handler": "",
400        "unit": "\\u00B0",
401        "detector": [
402          "",
403          "grande"
404        ],
405        "quant_sub_type": "f64",
406        "display_format": "default",
407        "min_value": "0.0",
408        "display_name": "Zenith Angle",
409        "description": "<div>\r\n<span class=dcInfoBoxDetailsDS>\r\nThe reconstructed Zenith Angle of the
              KASCADE showers is derived from the arrival time distribution of the of the particles at the detector
              stations. The range is from <span class=dcMathFunct>0&deg;</span> to <span class=dcMathFunct>60&deg;</
              span> where <span class=math>0&deg;</span> corresponds to a vertical shower. The angular resolution is
              between <span class=dcMathFunct>0.4&deg;</span> and <span class=dcMathFunct>0.1&deg;</span> depending
              on the energy.\r\n<br>\r\n<b>We recommend to use data only up to 42&deg;</b><br><br></span>\r\n <span
              class=dcInfoBoxReference> [details see  <b>KCDC-Manual</b>]</span> \r\n</div>\r\n",
410        "name": "Ze",
411        "max_value": "60.0",
412        "selection_mode": "D",
413        "order": 2,
414        "descr_head_html": "<p class=dcInfoBoxHeaderDS>Zenith Angle Info"
415      }
416    },
```

# KCDC APPLICATION PROGRAMMING INTERFACE (API)

**Shell example:** Extraction of the all data with an energy range from 17-19eV[log10]

**Request:**

```
curl          --insecure          --request          POST          'https://kcdc-
dev.iap.kit.edu/datashop/api/submit' \
--header 'Authorization: Basic cG92dGVyOmhhcnJ5Kytxb3R0ZXI=' \
--header 'Content-Type: application/json' \
--data-raw '
{
    "reconstruction": "",
    "output_format": "ascii",
    "datasets": [
        {
            "name": "array",
            "quantities": [
                {
                    "name": "E",
                    "cuts": [[17, 19]]
                }
            ]
        }
    ]
}'
```
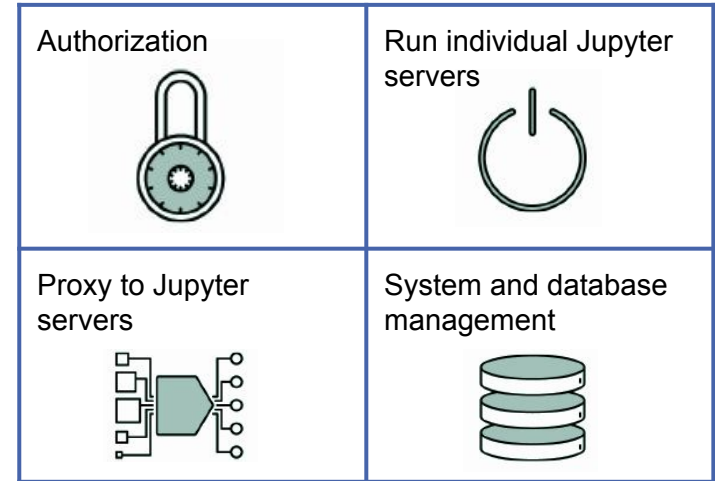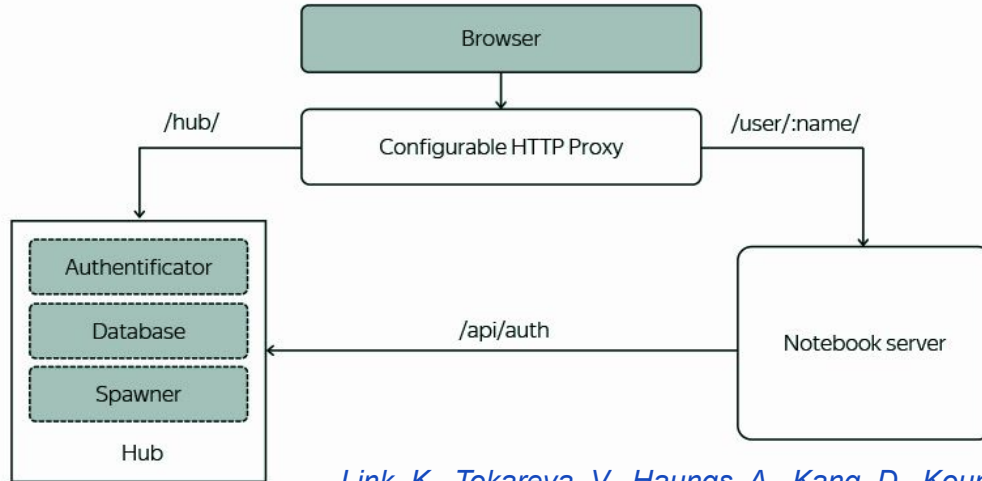
**Responce:**

job id:

{"id":"dbf1e608b6044223afe472125c020d88"}

or error message:

{"detail":"Invalid basic header. Credentials not correctly base64 encoded."}

- *Online API documentation: https://kcdc.iap.kit.edu/datashop/api/docs/index.html*
- *Wochele J. et al. KCDC User Manual: https://kcdc.iap.kit.edu/static/pdf/kcdc_mainpage/kcdc-Manual.pdf*

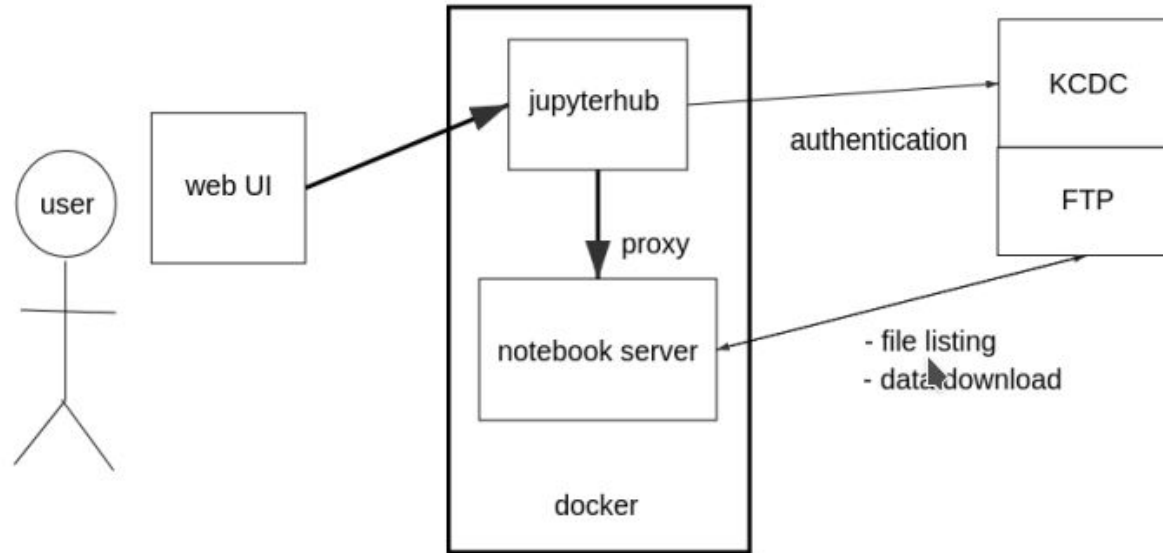# JupyterHub for data analysis

- Login via KCDC credentials
- Administration using Docker Swarm
- Tutorials by: KASCADE, IceCube, TRVO



| Authorization | Run individual Jupyter servers |
|---|---|
| Proxy to Jupyter servers | System and database management |

*Link, K., Tokareva, V., Haungs, A., Kang, D., Koundal, P., Polgart, F., Tkachenko, O.,Wochele, D., Wochele, J. Online masterclass built on the KASCADE cosmic ray data centre. In 37th International Cosmic Ray Conference (ICRC 2021), Online, 12.07. 2021–23.07.*

# JupyterHub integration in KCDC



*Polgart, F., Haungs, A., Kang, D., Wochele, D., Wochele, J., & Tokareva, V. (2020). An analysis framework for KCDC. In DLC 2020: Proceedings of the 4th International Workshop on Data Life Cycle in Physics. Ed.: A. Kryukov (p. 111).*

# Usage of KCDC's JupyterHub

https://jupyter.iap.kit.edu/

# Usage of KCDC's JupyterHub

# Thank you for your attention!

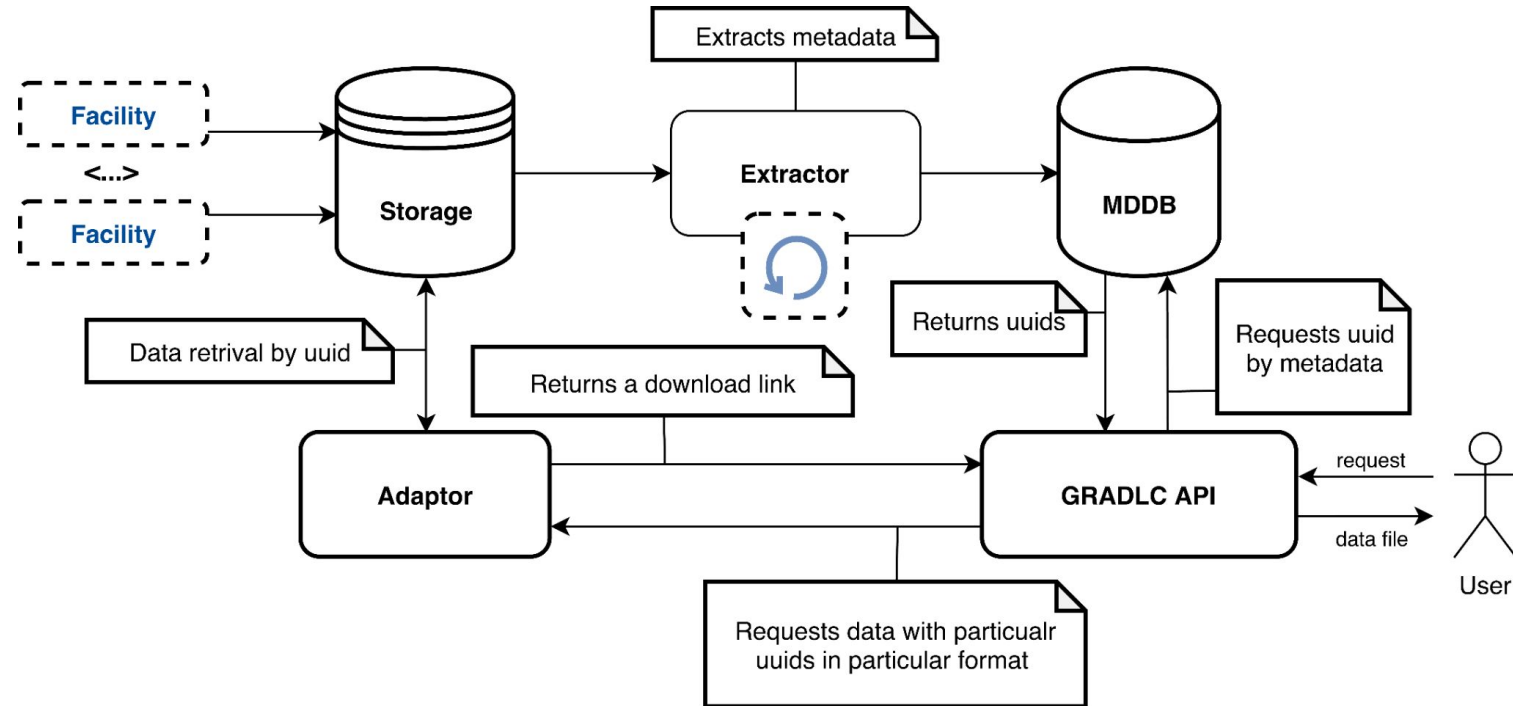**victoria.tokareva@kit.edu, iap-kcdc@lists.kit.edu**

Back up
# GRADLCI and KCDC

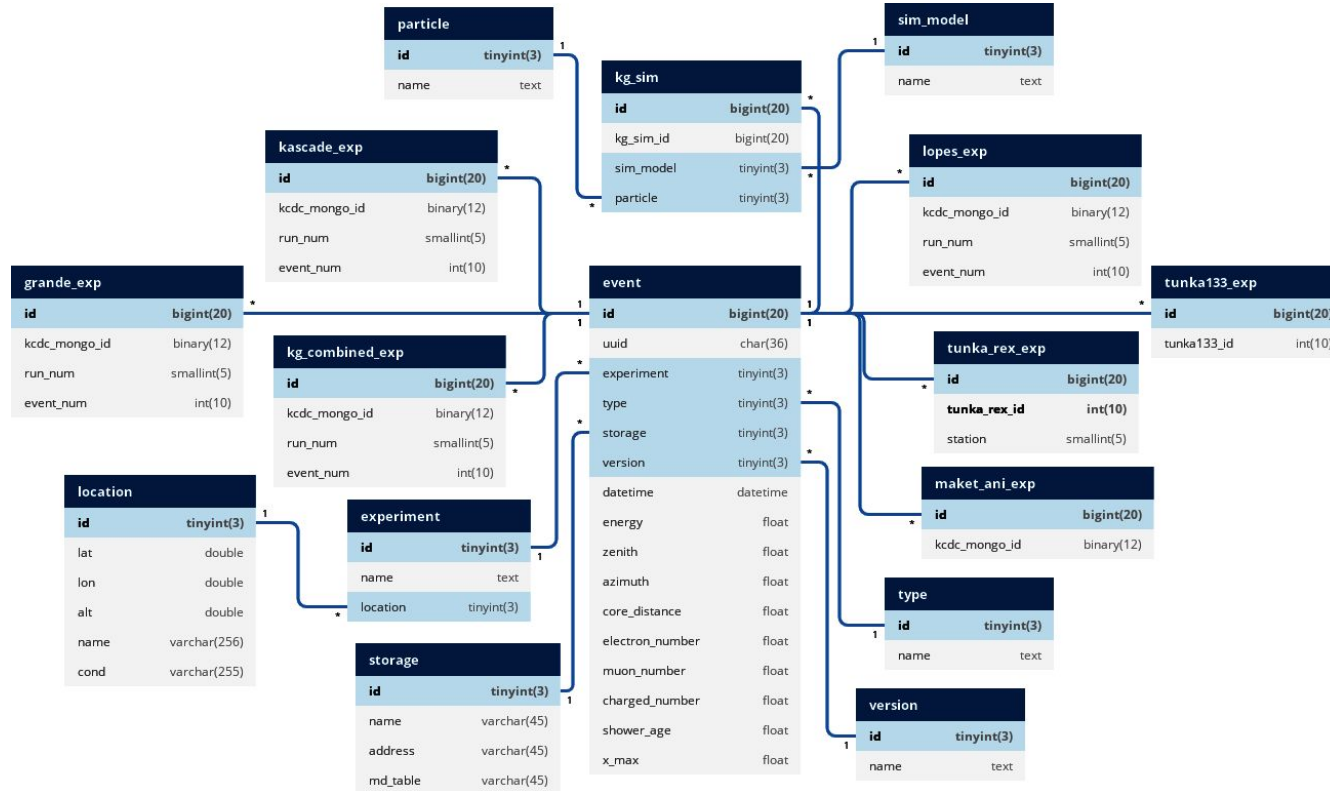# German-Russian Astroparticle Data Life Cycle Initiative (GRADLCI)

- The international initiative aiming at automatisation the maintenance of astroparticle-physics data throughout their entire life cycle
- 2018 - 2021
- Task areas:
  - KCDC extension
  - Prototype analysis and data center
  - Machine learning for astroparticle physics
  - Outreach and education
- Aggregated data by KCDC, Tunka-133, TAIGA and Tunka-Rex Virtual Observatory (TrVO)
- Data throughput:  4.5 TB
- Features:
  - Metadata database (MDDB) as SQL DB
  - 2 level metadata model:

    (1) file level metadata: file size, file type, last changed, etc.;

    (2) event parameter level: event id, datetime, setup, atmosphere, etc.

# GRADLCI (meta)-data flow
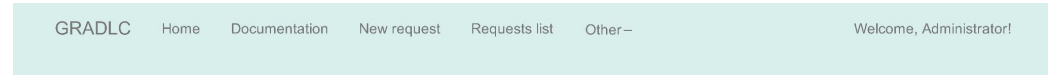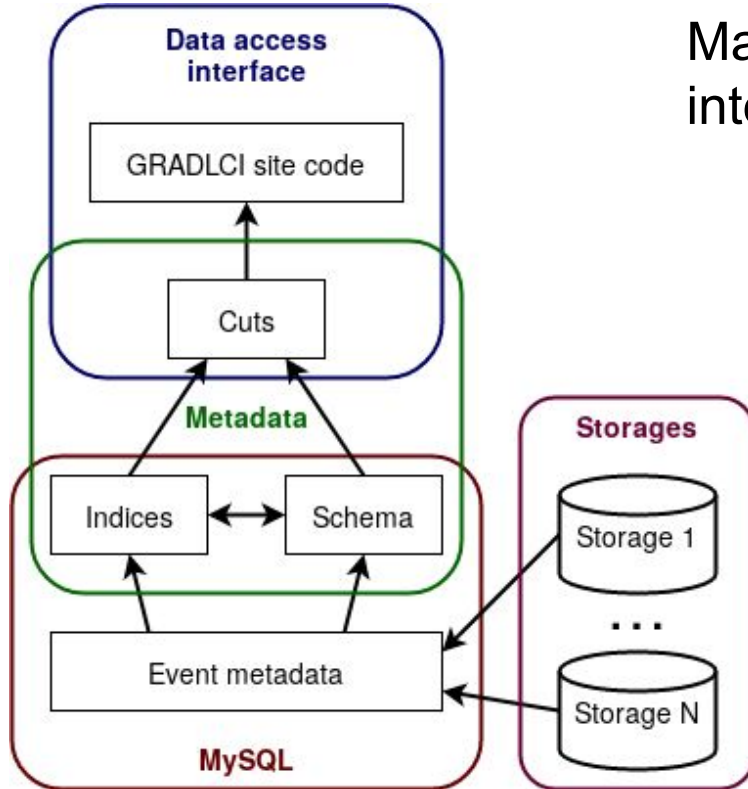


* Implemented on dedicated server at KIT. Considered to be integrated into PUNCH4NDFI data portal

# Metadata schemata GRADLCI

# GRADLCI metadata acquisition



Main idea: provide a common data selection interface for various data sources

# Aggregation server Web API[*]

**Domain name:** `gradlc-dc.iap.kit.edu`

**Request type:** JSON-RPC

**Protocol:** http

**Authentication:** HTTP Basic Auth

### Possible requests

- Data requests
- Request status
- List of requests
- Remove request from the list
- Download file

### Possible request status

- Running
- Scheduled
- Finished
- Failed
- Deleting
- Expired

## Example:

Request:

```
{"id": "4998715b-cd5d-4c17-80fb-8139a74d66ea", "jsonrpc": "2.0", "method": "new_task", "params":
{"kascade_exp": {"datetime_max": "2011-10-10 00:00:00", "datetime_min": "2010-10-10 00:00:00",
"zenith_max": 20.0, "zenith_min": 0.0}}}
```

Response:

```
{"id": "4998715b-cd5d-4c17-80fb-8139a74d66ea", "jsonrpc": "2.0", "result": {"url":
"http://gradlc-dc.ikp.kit.edu/download/c3feaa45-b654-44d3-83e7-671b1ac0499c.7z__", "uuid":
"c3feaa45-b654-44d3-83e7-671b1ac0499c"}}
```

*\* Application Programming Interface*