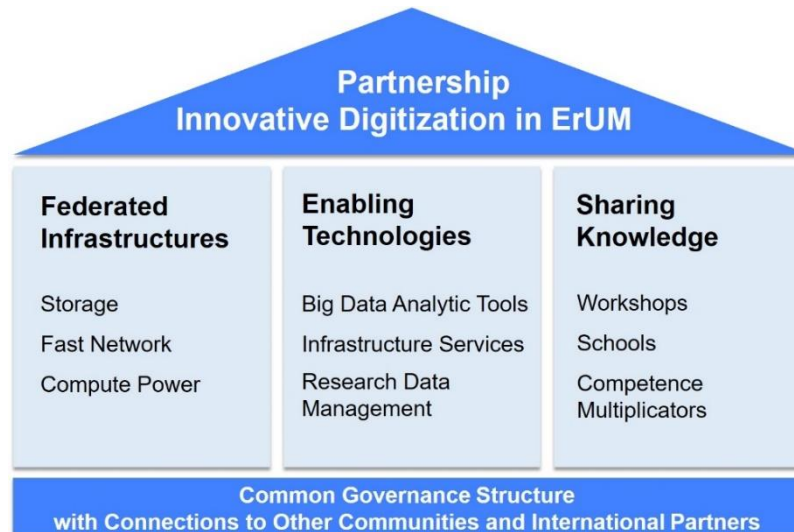


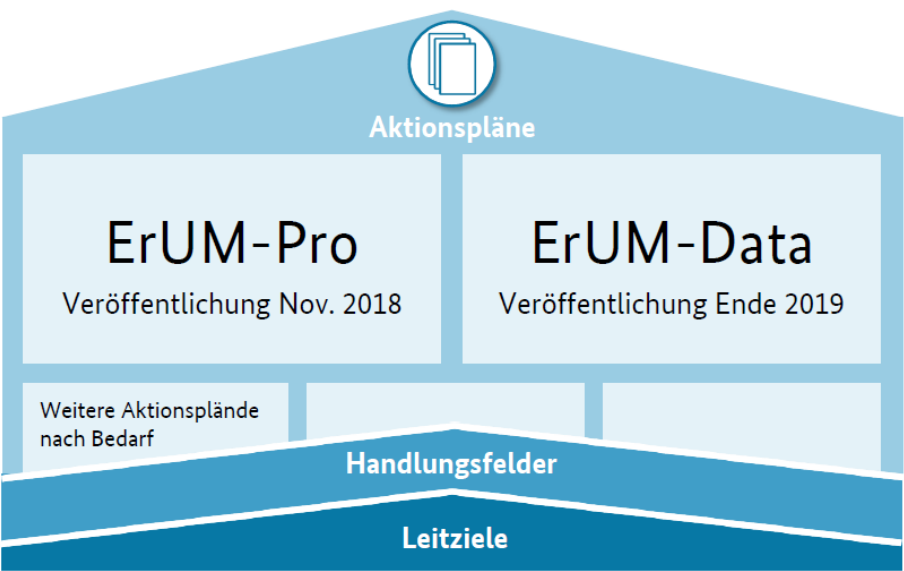
Challenges and Opportunities of Digital Transformation in Fundamental Research on Universe and Matter

HAP / AKPIK workshop | Big Data Science in Astroparticle Physics
Aachen, 18-20 February 2019

Andreas Haungs



action plan: 2017-2027



one plan of action:
ErUM-Data: Contributions
to the digital agenda

Committees related to
ErUM in Germany

Scientists
with doctoral
degree

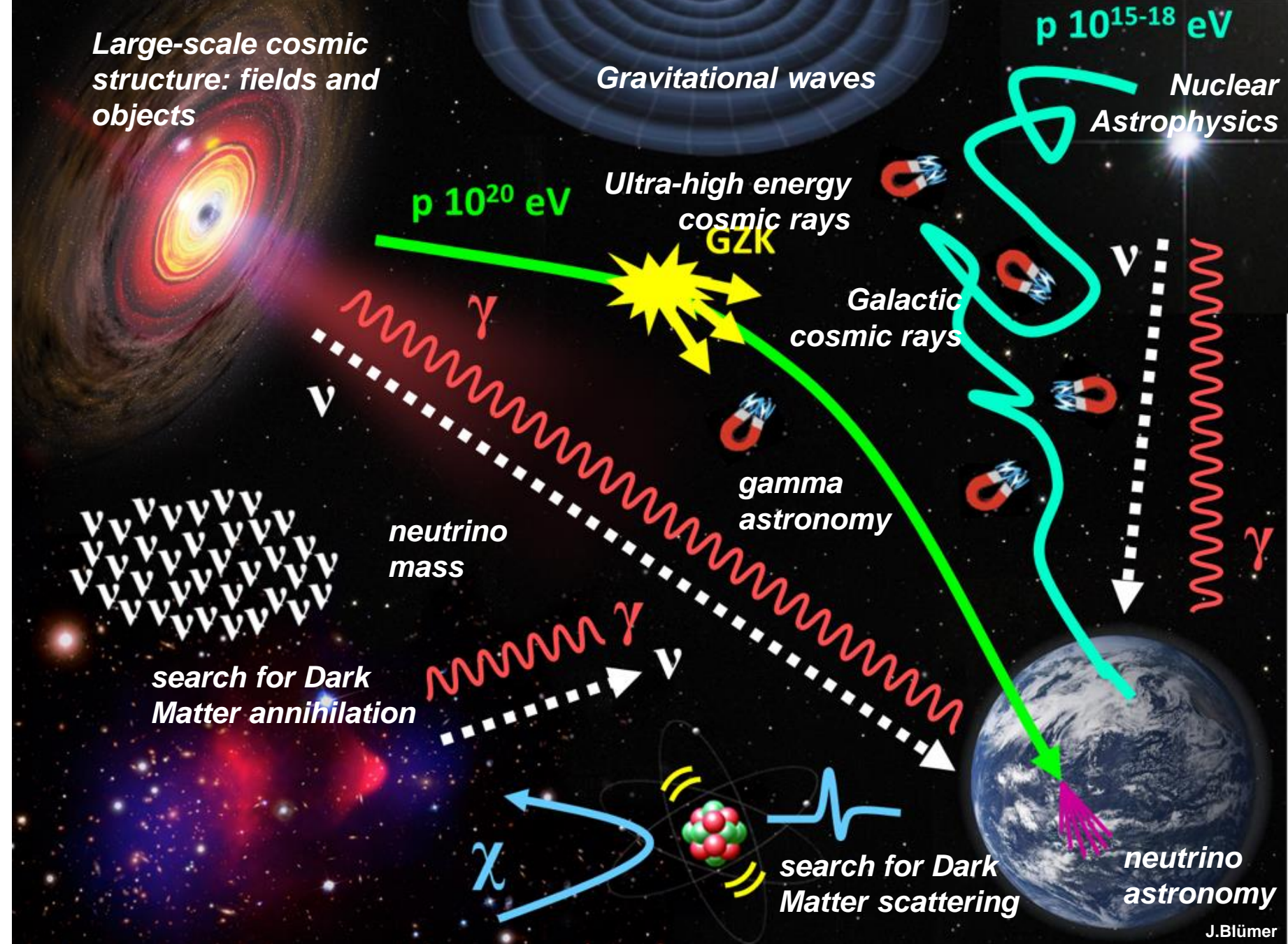
KFS	4.000
RDS	1.500
KHuK	1.500
KET	1.300
KFN	1.000
KAT	500
KfB	200
KFSI	100
	10.100

Initiative for a (global) Analysis & Data Center in Astroparticle Physics

Astroparticle Physics = Understanding the

- Multi-Messenger Universe
- Dark Universe

needs an **experiment-overarching** platform!



Initiative for a (global) Analysis & Data Center in Astroparticle Physics

- Astroparticle Physics requests for multi-messenger analyses - this needs an **experiment-overarching** platform!

■ Tasks

- Provide sustainable access to scientific data
- Archiving of Data and Meta-Data
- Providing analysis tools
- Education in Big Data Science
- Development area for multi-messenger analyses (e.g. Deep Learning)
- Platform for communication and exchange within Astroparticle Physics

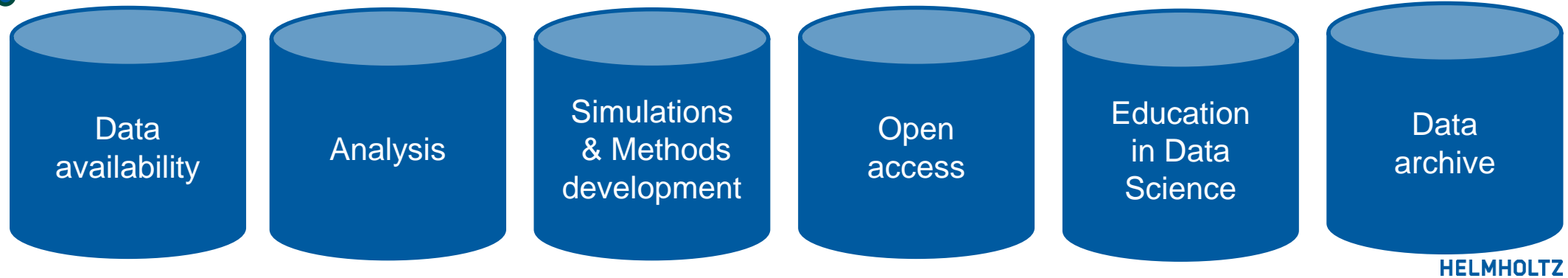
■ Elements

- Advancement, generalization of existing structures (like KCDC and others)
- In direction of a virtual Observatory (like in astronomy)
- In direction of Tier-systems and DPHEP (like in particle physics)
- „Digitale Agenda der Bundesregierung“
- OECD Principles and Guidelines for Access to Research Data from Public Funding
- Follow the FAIR principles of data handling

FINDABLE-ACCESSIBLE-INTEROPERABLE-REUSABLE

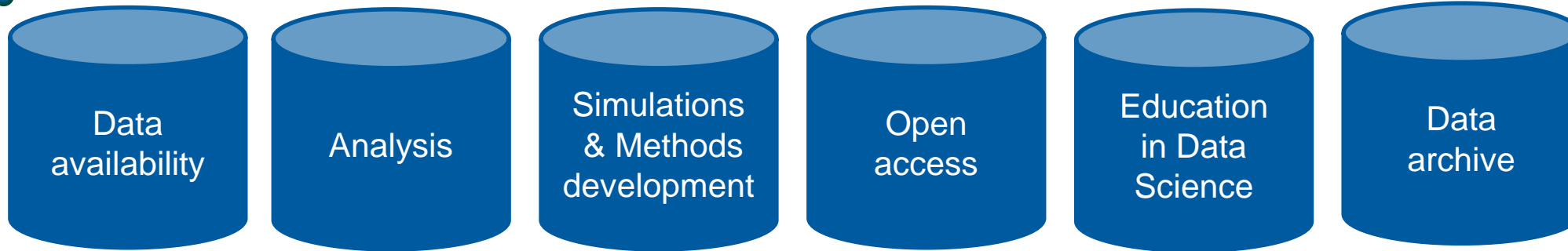


Analysis and Data Center in Astroparticle Physics



- **Develop a global analysis & data centre as user facility for multi-messenger studies in astroparticle physics**
- Motivation:
 - Needed, as experiments globally distributed and no worldwide centre like CERN exist
 - Implementation of 'Digital Agenda' and 'Big Data Science' in Astroparticle Physics
 - Apply 'FAIR' data handling in Astroparticle Physics
- Elements:
 - Data Preservation; virtual Observatory; distributed resources, data provider; outreach;
 - Based on experience of KCDC, GridKa, CTA data center, IceCube-Tier1, VISPA, NIFTY5
 - User-led facility (in Germany: 2 Helmholtz, 3 Max-Planck, 15 Universities)
- Realization as sustainable User Facility

Analysis and Data Center in Astroparticle Physics



Partly realized
in individual
experiments

➤ Data availability:

All researchers of the individual experiments or facilities require quick and easy access to the relevant data.

➤ Analysis:

Fast access to the generally distributed data from measurements and simulations is required. Corresponding computing capacities should also be available.

➤ Simulations and methods development:

The researchers need an environment for the production of relevant simulations and the development of new methods (machine learning).

➤ Open access:

More and more it is necessary to make the scientific data available not only to the internal research community, but also to the interested public: public data for public money!

➤ Education in data science:

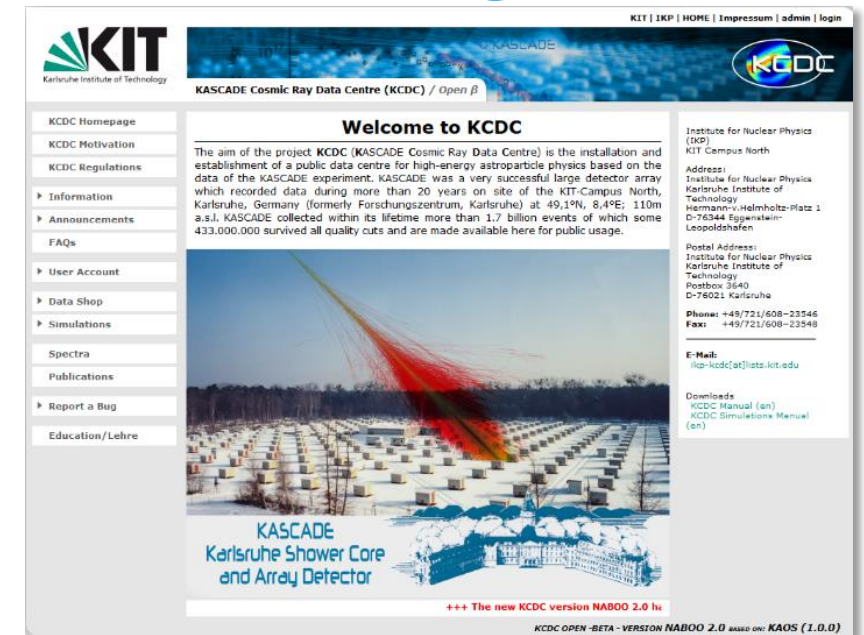
Not only data analysis itself, but also the efficient use of central data and computing infrastructures requires special training.

➤ Data archive:

The valuable scientific data and metadata must be preserved and remain interpretable for later use (data preservation).

KASCADE Cosmic ray Data Centre

- Motivation and Idea of KCDC:
 - public access to the data
 - data has to be preserved for future generations
- Web portal:
 - modern software solution
 - release the software as Open Source
 - educational courses
- Data access:
 - new release (Feb. 2017) with $4.3 \cdot 10^8$ EAS
 - simulation data
 - spectra
- Pioneering work in publishing research data in astroparticle physics



<https://kcdc.ikp.kit.edu/>

[J.Phys.Conf.Ser. 632 (2015) 012011]
[EPJ C78 (2018) no.9, 741]

Astroparticle Data Life Cycle Initiative

see talk Victoria Tokareva

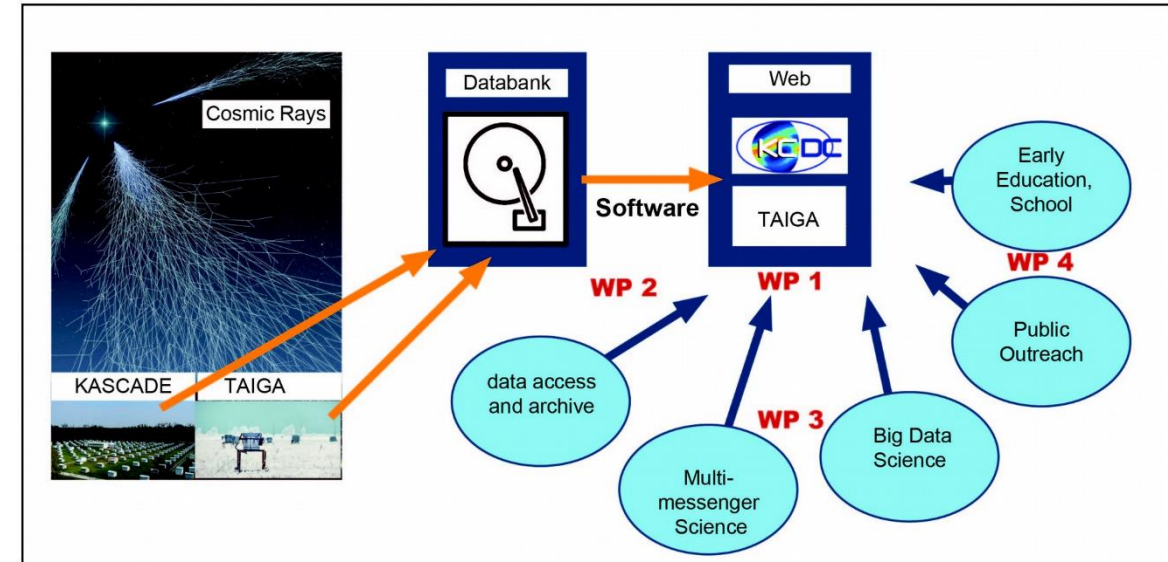
• Basics

- project period 2018-2020
- funded by Helmholtz and RSF
- Russia: SINP MSU, ISU, ISDCT SB RAS
Germany: KIT, DESY

• Main targets of the Project

- Extension: data from Tunka/TAIGA and KASCADE-Grande
- Developing solutions of distributed data storage techniques with a common metadata catalog
- Development of appropriate machine-learning workflows
- Perform experiment overarching multi-messenger astroparticle physics
- Learn to use GridKa environment
- Creation of an educational subsystem

<http://astroparticle.online>



Project is a first step in extension and generalization of KCDC

Analysis and Data Centre for Multi-messenger Astroparticle Physics

ADC-MAPP

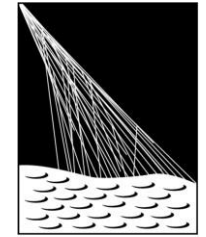
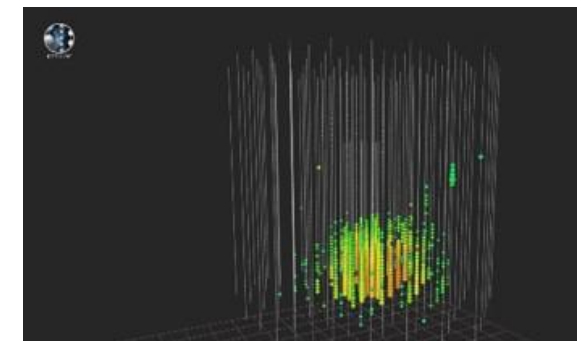
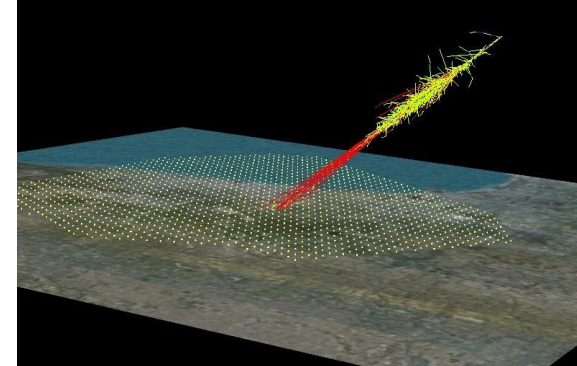
- **Basics**

- ADC-MAPP project period 2019-2020
- funded by Helmholtz

- **Main targets of the Project**

- Provide sustainable access to scientific data
- Archiving of Data and Meta-Data
- Providing analysis tools
- Education in Big Data Science
- Development area for multi-messenger analyses
(e.g. Deep Learning)
- Platform for communication and exchange within
Astroparticle Physics

Open positions for data
scientists - Contact me.... 😊



PIERRE
AUGER
OBSERVATORY



Support by the BMBF:

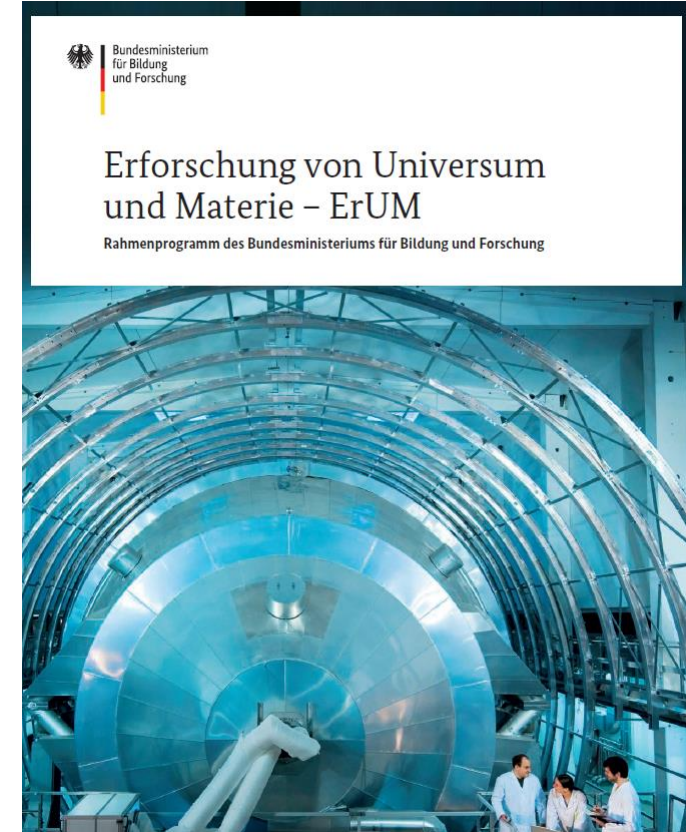
Innovative Digital Technologies for the Erforschung von Universum und Materie

Consortium KAT-KET-KHuK (pilot project ErUM-Data)

Proposal (11 Universities + 6 associated partners; coordinator T. Kuhr of LMU)

- Topic A: Development work for the provision of technologies to leverage heterogeneous computing resources
- Topic B: Application and testing of virtualized software components in the environment of heterogeneous computing resources
- Topic C: Deep learning, gaining knowledge through well-founded data-driven methods
- Topic D: Event reconstruction: cost and energy efficient use of computing resources

Approved for period 10/2018-9/2021



one plan of action:
**ErUM-Data: Contributions
to the digital agenda**

How do we organize our digitization era

Education/Schools



Deep Learning in Physics

- Master course, RWTH 17,18,19
- GRIDKA summer school, KIT 17,18
- Grad. Kolleg, Freiburg 18
- Belgium Dutch German summer school, Berlin 18
- ...

Workshops



Big Data Science

- made in Germany, Berlin 17
- Astroparticle Research, RWTH 17,18, [18.-20.2.19](#)
- Machine Learning, CERN 17,18,19
- ...

White papers

- KAT: Astroteilchenphysik im Licht der Digitalen Agenda
- KAT, KET, KfB, KFN, KFS, KHuK, RDS: Gemeinsames Strategiepapier
- A Roadmap for HEP Software and Computing R&D for the 2020s
- ...

Grants

Funded

- 12 Universities, 5 Research Centers: Innovative Digitale Technologien für die Erforschung von Universum & Materie
- ...

Whitepaper ErUM-Data

"Digitalisation in ErUM": BMBF-Workshop 4-5 October 2018

- Federated Infrastructures
 - Efficient usage
 - Services
- Research Data Management
 - Data life cycle;
 - Networking (NFDI, EOSC);
- Big Data Analytics
 - Deep Learning;
 - Provide sustainable algorithms and tools;

Challenges and Opportunities of Digital Transformation in Fundamental Research on *Universe and Matter*

Martin Erdmann¹, Christian Gutt², Andreas Haungs³,
Klaudia Hradil⁴, Thomas Kuhr⁵, Marcel Kunze⁶,
Anke-Susanne Müller⁷, Günter Quast⁸, and Matthias Steinmetz⁹

¹RWTH Aachen University, KAT

²University of Siegen, KFS

³Karlsruhe Institute of Technology, KAT

⁴Technische Universität Wien, KFN

⁵Ludwig Maximilians University Munich, KET

⁶Universität Heidelberg, KHuK

⁷Karlsruhe Institute of Technology, KfB

⁸Karlsruhe Institute of Technology, KET

⁹Leibniz-Institut für Astrophysik Potsdam, RDS

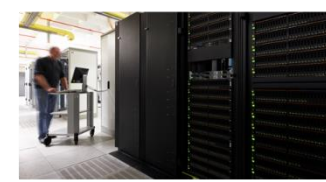
Our charge: write down concrete portfolio of measures → BMBF action plan → calls

Federated Infrastructures

- Increasingly heterogeneous computing infrastructures available and needed (HTC vs. HPC)
- Huge Storage: Multiple Exabytes
- Fast Networks: >100 Gb/s for entire ErUM
- Substantial large-scale experiences in all related aspects and connected to computer science, multiple domain specific aspects
- need large scale federated infrastructures from experienced providers (including commercial providers)
- Utilization needs sustained software development thus sustained positions
- Infrastructure in ErUM as building block of national (NFDI) and international (EOSC) initiatives



GridKa KIT



DESY Grid Centre & NAF



GreenCube GSI



ICT infrastructures at Univs

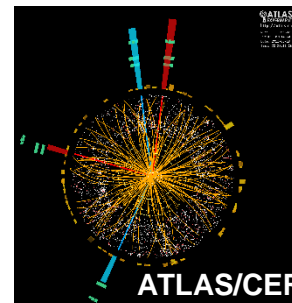


SuperMUC Garching

etc...

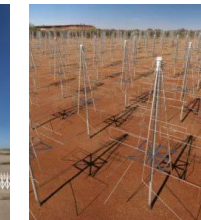


High luminosity LHC



ATLAS/CER

Square Kilometre Array



300 PetaBytes per year

Experiments at XFEL



Big Data Analytics

- **Utilization Big Data Analytics in national and international contexts:**
 - Development and implementation of tools for Big Data Analytics
 - Need for a collaborative effort in terms of Big Data Analytics including users, facilities, mathematics and computer science
 - A platform for sharing Big Data Analytics solutions (inside or even across communities).
 - Integration with data management (e.g. for efficient data access or mining archived data)
 - Integration with federated infrastructure (e.g. for utilizing resources optimized for Big Data Analytics tasks).
 - Training and education of the next generation of scientists in Big Data Analytics;
 - Ensure sustainable development and curation of algorithms and tools.



Scoogle: Scientist's data & algorithms

Scientists:
Questions

Web Interface

**Big Data
Analytic Tools**

Algorithms
Visualization
Machine Learning

Data

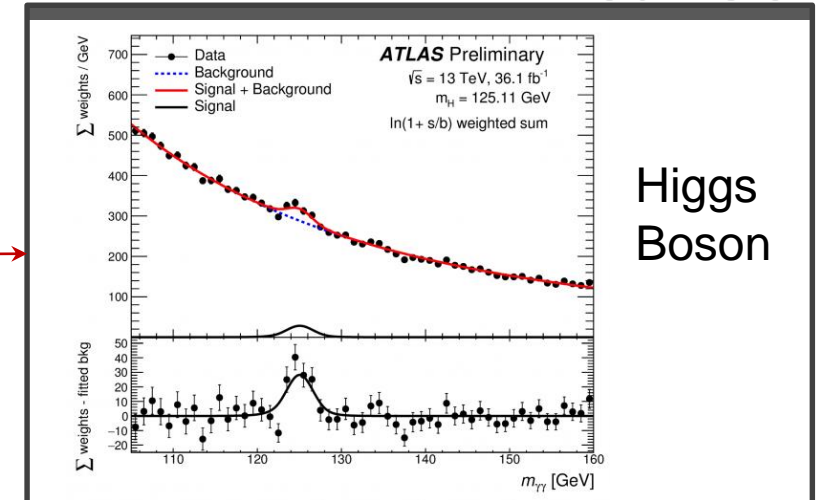
Experiment data
Metadata
Simulation

Web vision

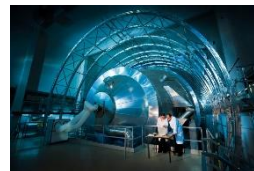
Access to

- Metadata & experimental data
- Computing resources
- Software libraries
- Algorithm development
- Workflow management system

Web vision



Higgs
Boson



Scoogle: Scientist's data & algorithms

Medium-scale prototypes exist, developed in our community (~5 years experience)

Scientists:
Questions

Web Interface

Big Data
Analytic Tools

Algorithms
Visualization
Machine Learning

Data

Experiment data
Metadata
Simulation

edit algorithms

execute program

inspect results

The screenshot shows a JupyterLab environment with a code editor on the left containing a Python script. The script imports libraries like argparse, logging, numpy, scipy, matplotlib, and stats, and defines a function to get data. On the right, there's a terminal window and a section for execution output. Below the code, three histograms are displayed, labeled _tbin_9.png, _tbin_8.png, and _tbin_7.png. Arrows point from the text labels to the corresponding parts of the interface.

The block features logos for VISPA (a blue sphere with a white 'L'), Jupyterhub (an orange circle with a white 'J'), Belle II (a blue swan), and SWAN (a blue cloud with a white bar chart). Surrounding these logos are images of various computing hardware: a laptop, a desktop monitor, a smartphone, server racks, and a Tesla accelerator card. Below the logos, three URLs are listed:

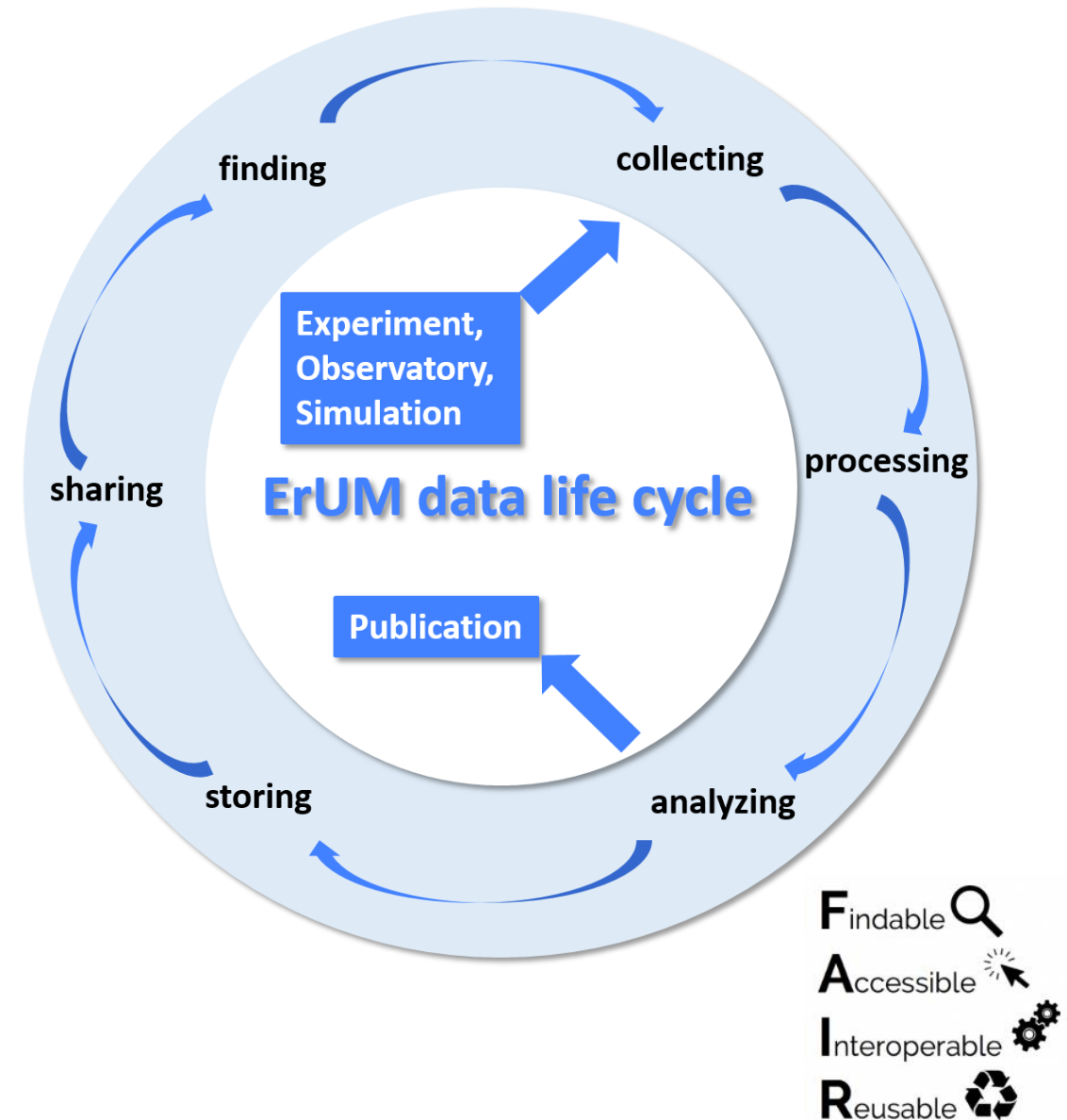
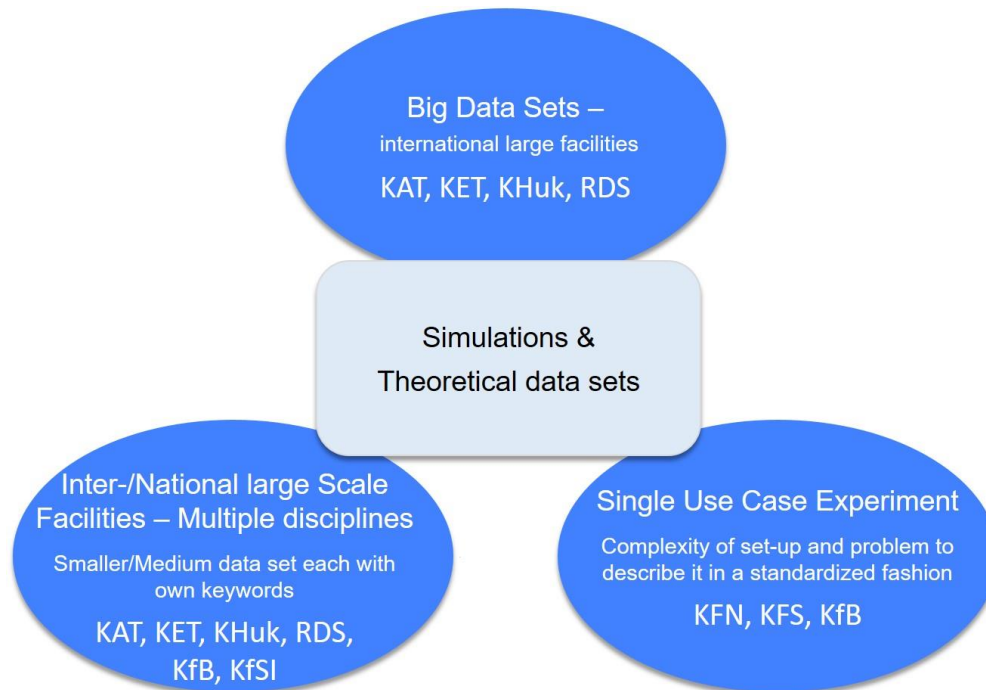
<https://swan.cern.ch>
<https://vispa.physik.rwth-aachen.de>
<https://stash.desy.de/projects/B2T/repos/b2-starterkit>

In ErUM: Substantial experiences in all related aspects, partly complementary, international context, also connected to mathematics, computer science, economy.

© Martin Erdmann, RWTH Aachen University

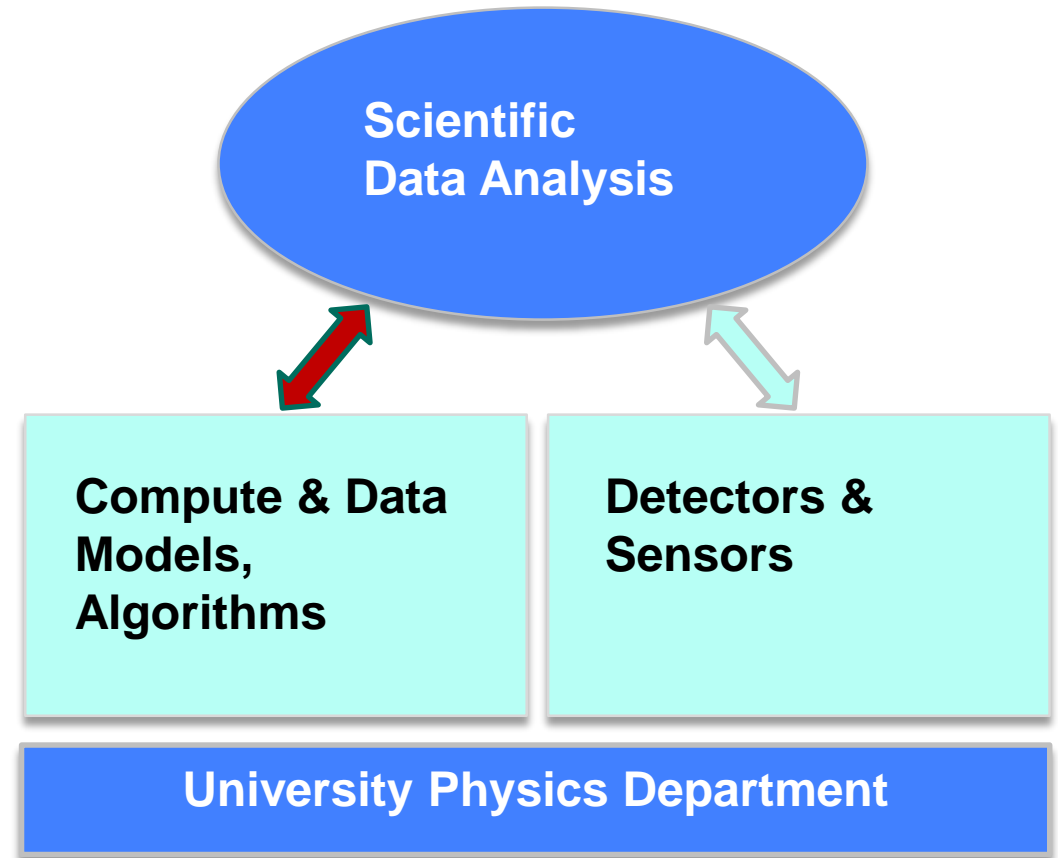
Research Data Management

- Where possible, common standards should be established to foster interoperability
- Importance of “data stewards” to manage the data life cycle and to act as a curator for metadata



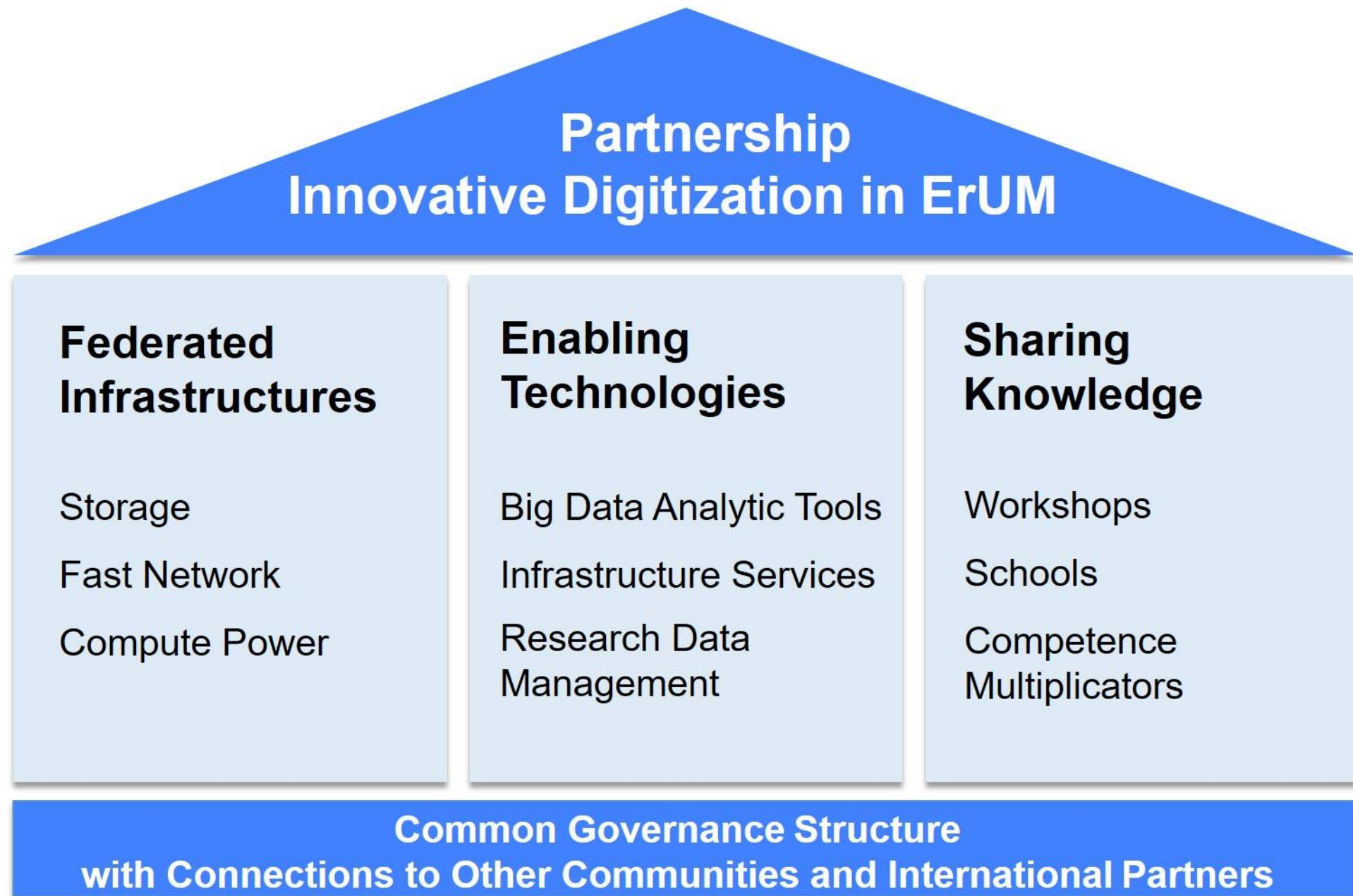
Tenure Track Programm

- **Education by scientific leaders: distribute and deepen knowledge in digitization**
- **Large tenure track programme for:**
 - Development of compute models for online & offline reconstruction, simulation, analytics
 - New algorithmic concepts, machine learning
 - Access to heterogeneous computing resources
 - New chairs will advance curricula



The ErUM House:

- user-led home to bundle and steer the activities



The needs and costs for digitization of the research area ErUM (estimate for a period of 10 years)

Manpower:

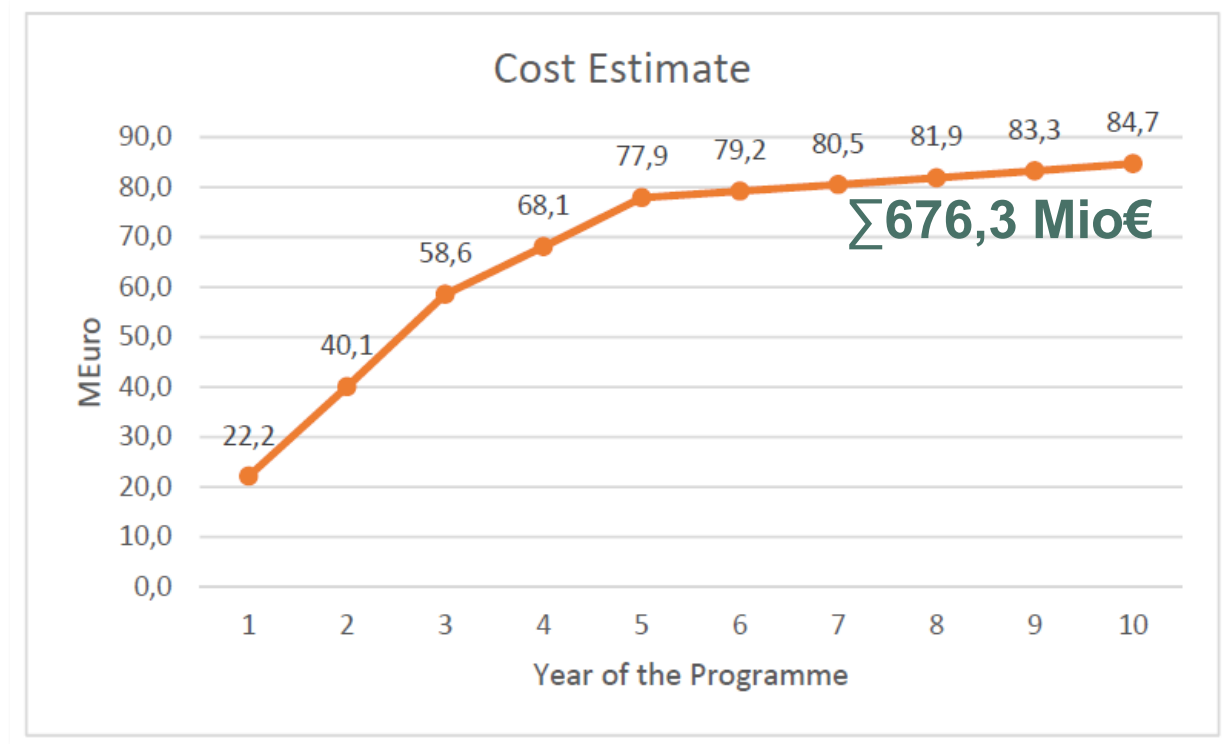
1. Workflows to exploit infrastructures: from 40/a to 100/a
2. Management of research data: from 40/a to 100/a
3. Big Data Analytics in physics research: from 40/a to 200/a
4. Scientist's web working environment: from 40/a to 100/a
5. Tenure track ErUM programme + 1 postdoc; from 40/a to 100/a

**No numbers →
No money**

Costs:

- **Full Time Equivalents:**
from 16M€/a to 59M€/a
- **Large-scale federated infrastructures:**
from 5M€/a to 25M€/a
- **Partnership for innovative digitization:**
for 1M€/a

➔ **10,100 ErUM scientists:**
0.6% FTE increase/year:
6.8k€/scientist/year!



DPG-AKPIK, Arbeitskreis Physik, moderne Informationstechnologie und Künstliche Intelligenz

Broad representation of interests for the relevant topics in physics



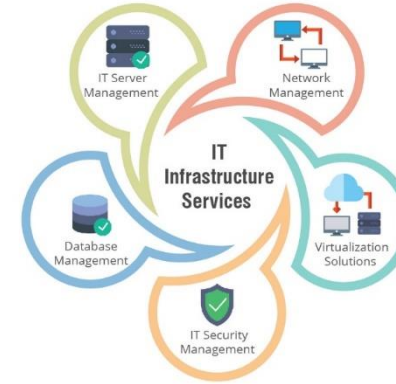
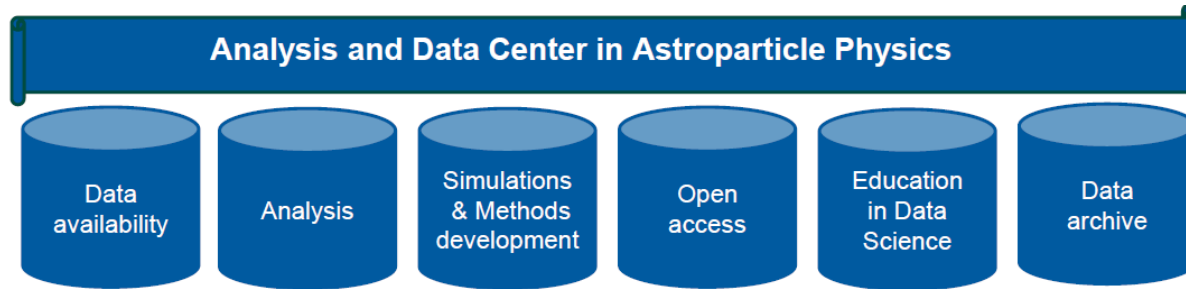
Topics:

1. **BIG DATA:** archiving, processing, management, analysis and simulation of complex data streams, HPC, information theory, statistical methods
2. **IT:** high-performance data readout systems and mass storage, visualization, smart sensors, bridge technologies for the next level of big data
3. **KI & ROBOTIK:** Data Driven Algorithms & Software, Autonomous Devices, Remote Control, Innovative Applications, Algorithms for Quantum Computers
4. **UNIVERSITY:** curricula and multi-disciplinary research centres, cooperation with the GI Task Force "Data Scientist", IT infrastructure
5. **INDUSTRY and SOCIETY:** Ethics, Technology Assessment, Sustainability, Business, Law, Start-Ups, Public

www.dpg-physik.de/dpg/gliederung/ak/akpik/index.html

Ask Karl Mannheim, Martin Erdmann, Wolfgang Rhode

Era of Digitization in Astroparticle Physics



Lots of activities in

- **Infrastructures**
- **Big Data Analytics**
- **Research Data Management**

Data Life Cycle

Next steps:

- **ErUM-Data is a big chance!!**
- **NFDI = Nationale ForschungsDatenInfrastruktur**
- **EOSC = European Open Science Cloud**

