# Sustainability in the Digital Transformation of Basic Research on Universe & Matter

Ben Brueers · Marilyn Cruces · Guenter Duckeck · Michael Düren ·
Niclas Eich · Torsten Enßlin · Johannes Erdmann · Martin Erdmann ·
Peter Fackeldey · Christian Felder · Benjamin Fischer · Stefan
Fröse · Martin Gasthuber · Lukas Geiger · Andrew Grimshaw ·
Daniela Hadasch · Moritz Hannemann · Alexander Kappes · Raphael
Kleinemühl · Oleksiy Kozlov · Thomas Kuhr · Simon Neuhaus ·
Pardis Niknejadi · Özlem Özkan · Judith Reindl · Daniel Schindler ·
Astrid Schneidewind · Frank Schreiber · Markus Schumacher · Kilian
Schwarz · Jens Struckmeier · Florian von Cube · Rodney Walker ·
Cyrus Walther · Angela Warkentin · Sebastian Wozniewski · Kai Zhou

**Abstract** Sustainability

## 1 Introduction

The creative workshop develops strategic concepts for sustainability in the digital transformation of basic research on universe & matter - from the funding applications on. The workshop program focuses on working sessions with below-mentioned guiding questions. The 6 sub-groups of ca. 5 participants will gather relevant information for the final report which will be published timely after the workshop. Keynote presentations by high-profile experts will inform participants and stimulate discussions. Key measures will concern education, research and innovation, in line with the BMBFs sustainability goals. We encourage young scientists and experienced scientists to participate in order to enable broad discussions. Example of a citation [1].

## 2 Summary Andrew Grimshaw: Concrete Path to Sustainable Computing

Texas is essentially desert, very sunny and windy throughout the year. Variations of sunshine of course day/night. With wind, there are summer/winter

Physics Institute 3A, RWTH Aachen University, 52056 Aachen, Germany
Tel.: +49-241-80-27330
Fax: +49-241-80-22189
E-mail: erdmann@physik.rwth-aachen.de

differences as well as day/night differences. Northwest Texas leads the US in renewable energy. Capacity has expanded greatly over the past 20 years (most recently more than 8GW/year), reaching 70GW by the end of 2023. This is more than Germany needs. Often the wind turbines stand still because the electricity cannot be taken off at all. At the same time, electricity is priced negatively about 15% of the time. The prices are already low anyway. Cost parity with electricity from gas is not reached until 95%. Texas requires 40GW, which is split about half and half between renewables and gas. Germany has an electricity mix that still relies heavily on coal and requires additional purchases from abroad. The share of solar energy is significant only in summer, marginal in winter. Wind has a small share. He sees the possibilities for water-based power generation exhausted in the U.S. and Europe, so he considers solar and wind the technologies of the future. The cost of electricity for Germany is a factor of 10 greater (US$300/MWh compared to US$27/MWh for Lancium Company).

Transporting electricity is a major challenge. It is not financing that plays the major role, but licensing procedures. The cost of transmission lines is US$1-3 million per mile and per 400MW. There are only 16 transmission lines between Northwest Texas and Southeast Texas. As a result, electricity is significantly more expensive in the southeast than in the northeast. However, the permitting process for transmission lines is extremely costly and takes 2 decades, even in this desert country. In Germany, with its dense population, this is rather more difficult. In the medium term, he wants to

build a wind belt in the U.S. in a north-south direction that will supply large swaths of land. Germany has its wind potential on the coasts.

Two options in consequence are: Either wait until the lines to transport renewable energy arrive at the institution, or move the power consumers to the point of generation. This is difficult for most industries, but scientific computing is an ideal candidate. Computing centers require only power and fiber. Jobs are computed in batch mode. Scientists are patient.

The following criteria should be considered when setting up a computing center near the power generators. Since the price of electricity is extremely low, the performance of a processor per invested electricity (flops per watt) plays a subordinate role. Likewise, key figures such as PUE and density are not decisive. Therefore, hardware that does not belong to the latest generation can also be operated; it is also cheaper by a factor of 10. As a side effect, the footprint of newly produced hardware is saved. Communities should build their own data centers, the computing costs of the big companies (hyperscalers) are too expensive. He thinks a "stonesoup" approach makes sense (everyone contributes something to the soup): First you put together what you can find on a small scale, together it will become successful. Later you can scale up. A good network connection of the computing centers is important.

Stability of supply is of course a big issue to respond to. They work with predictions and shut down the CPUs within 15 seconds (GPUs are another issue). Restarting takes about 1 minute. They keep the servers and network switches running, their power consumption is minimal (a few watts), damage during shutdown and restart would be too great. The three Lancium Computing sites have 20 racks of 50kW per rack at a cost of US$ 1M. The network costs US$ 8000 per fiber per month. He sees a total cost of US$ 80/MWh. The footprint is zero CO2. The efficiency factor, what percentage of the time can be calculated at full power, is still being measured.

## 3 Summary Rod Walker: Measures for sustainable computing operation

In the field of data-intensive high-energy physics, several areas can be identified where improvements in sustainability are within reach.

In the area of software, there is a whole range of possibilities to reduce the use of resources. These include efficient algorithms, more parameterized or ML-based simulations, use of parallel structures (GPU, columnar analysis), shareable generator files for simulations, and

avoidance of duplications of all kinds. Comparing CPU processors, differences in power consumption per computational power measured in HS06 (HEPSpec2006) can be identified. Currently, ARM processors perform well here. Furthermore, the use of heat dissipation is a great opportunity for sustainability to improve the effectiveness of energy use measured in PUE (best value 1.0).

Grid and cloud computing have significantly reduced the need for institute-owned computing clusters. This greatly saves hardware and human resources, and localizations of computing centers can be chosen suitably.

One challenge is the responsiveness of computing centers to varying renewable energy power. The options range from

- Freezing computing jobs and CPUs in sleep mode, however still consuming power,
- Reducing CPU frequency to an optimum (e.g., from 3GHz to 1GHz), resulting in 50% power savings without significant additional loss of computational effectiveness (only the jobs run longer).
- Switch to battery power, which may become affordable in the long run.

Because of renewables, it makes sense to keep old hardware running longer (also because of the footprint of new hardware production). However, old hardware is substantially less efficient, so an optimal operating point must be sought with old hardware or its lifetime.

For data storage, a HEP GRID computing cluster at Tier2 level requires 40% of the total energy consumption. Lowering the rotation speed of the storage disks by 90% can save a factor of 10 energy and is thus at typical characteristics of tape units. This also allows one to react to variations in the energy supply of a data center.

## 4 Summary Salome Shokri-Kuehni: View of the BMBF

The BMBF framework program on sustainability is based on the UN's 17 goals on sustainability. The BMBF measures relate to the hierarchical levels of 1) German sustainability strategy, 2) action plan sustainability, 3) sustainability in science initiative and 4) training for sustainable developments. A clear distinction is noted between sustainable research and research for sustainability. There is already a statement from the Prisma Forum (high-level advisory board) that points to standards in evaluations, best practice, avoidance of climate killers, identification

## 4 Summary Özlem Özkan: Sustainability of Research Data Management

Estimates suggest that research data taken before 1990 are, with few exceptions, largely lost. The view of the value of data has changed considerably since then and is crystallized in national and international efforts to preserve research data and make it available as open data wherever possible. The term FAIR (findable, accessible, interoperable, reusable) has become the basic guideline. National efforts are condensed, for example, in the NFDI (National Research Data Infrastructure) and the HMC (Helmholtz Matter Collaboration).

Internationally, there are a wide variety of organizations (e.g., EOSC). It is noteworthy that the G7 countries are also clearly committed to open data. The realization that the financial benefit is quantifiable in billions of euros per year plays an important role here. Benefits of global collaboration based on open data became evident during the Covid pandemic in the decoding of genome sequences and the development of vaccines. Keeping data usable is possible today in terms of storage capacity. It requires energy consumption at the level of 1% of global energy production. In the UK, estimates for the cost per experiment project are 60h per hour. Further costs arise from the need for personnel to organize data preservation, training and education: Data Managers, Data Stewards, Librarians...

The concretization of FAIR is long-term availability, convenience of access, ease of use, and integrity of research data. This requires the commitment of the communities, which have the respective domain knowledge about the data and their value and can formulate policies, which data exactly and in which form they are preserved. One example is the NEXUS format, which is used in the Neutron, Muon and X-ray communities and has the goal of retaining research data for at least 10 years in accordance with FAIR. For this, there is an expert committee that brings together representatives from all research infrastructures and further stakeholders.

Overall, the preservation of research data requires on the community side well-defined data policies with foresight in terms of data value and on the institutional side budgets for staff and data facilities.

## 5 Summary Jens Struckmeier: Building sustainable systems

Computing centers already account for a large part of the worldwide energy consumption (20%). According to estimates, this proportion will rise to 50% by 2030. The thermal output of the devices per area is very high and exceeds that of a conventional induction stove plate. Utilization of the heat output is imperative in terms of sustainability and can be dissipated and made usable, for example, by transporting water. Signal runtimes play a role in the distribution of computing centers. Here, a response time of 1 ms requires a distance from the center of about 30km. In Germany, households consume about 30% of the total energy demand, with 75% going into indoor heating. Worldwide, the energy consumption for heating and cooling amounts to 48%. Utilizing the thermal output of computing centers has high sustainability potential. Exemplary examples of successful concrete implementations already exist, with efficiency measured in Power Usage Effectiveness (PUE). Such efforts exist not only at research infrastructures such as CERN in Geneva (Prevessin Site PUE=1.1) or at GSI in Darmstadt, but also for residential units. Best values are achieved for new residential buildings (PUE=1.024), but also for conversions of e.g. high-rise buildings far-reaching improvements from PUE 2 to PUE=1.27 could be achieved (Cloud&Heat). Various forms of such sustainable computing centers are being tested, including highly modularized units.

## 6 Summary Salome Shokri-Kuehni: View of the BMBF

The BMBF framework program on sustainability is based on the UN's 17 goals on sustainability. The BMBF measures relate to the hierarchical levels of 1) German sustainability strategy, 2) action plan sustainability, 3) sustainability in science initiative and 4) training for sustainable developments. A clear distinction is noted between sustainable research and research for sustainability. There is already a statement from the Prisma Forum (high-level advisory board) that points to standards in evaluations, best practice, avoidance of climate killers, identification and elimination of hindering bureaucratic regulations, strengthening of transfer, and more.

To develop guidelines and recommendations for funding in ErUM, 6 working groups have been formed in the so-called Prisma Trialogue to present initial feedback by September 23 and final results in the first quarter of 2024. All interested parties are invited to actively participate in the working groups:

- WG 1: Research planning and organization
- WG 2: Research funding in ErUM
- WG 3: Data and Computing
- WG 4: Technologies at research infrastructures
- WG 5: Data collection, monitoring and accounting
- WG 6: Research for Sustainability

and elimination of hindering bureaucratic regulations, strengthening of transfer, and more.

To develop guidelines and recommendations for funding in ErUM, 6 working groups have been formed in the so-called Prisma Trialogue to present initial feedback by September 23 and final results in the first quarter of 2024. All interested parties are invited to actively participate in the working groups:

– WG 1: Research planning and organization
– WG 2: Research funding in ErUM
– WG 3: Data and Computing
– WG 4: Technologies at research infrastructures
– WG 5: Data collection, monitoring and accounting
– WG 6: Research for Sustainability

## 5 Footprint by G. Duckeck, N. Eich, J. Erdmann, S. Neuhaus, M. Schumacher

A quantitative understanding of the CO2 footprint varies in maturity across different domains. The experiments of the large hadron collider LHC at CERN have joined together to form GRID Computing (WLCG) and record a very accurate monitoring of the more than 170 computing centers worldwide. This information could be used to estimate the carbon footprint. Similarly, theory groups are computing on application at HPCs. Individual centers of experiments (e.g., ATLAS) have already determined the CO2 footprint (see Rod Walker's talk). The CERN WLCG center also gives details of the split between CPU, disk, services and network. There are other examples mainly from experimental groups at research infrastructures (XFEL) with at least partial information. In contrast, no information at all is available for data transfers or clusters of institutes. Likewise, with a few exceptions, it is difficult to obtain information on the production and disposal process (an exception, e.g., DELL).

The opportunities for savings are numerous and reside in

– HW Technology progress
– SW/algorithm optimization
– Different architectures (GPU, ARM, . . . )
– ML optimization
– Dynamic power provision
– Lifetime extension
– Minimize losses

First of all, it is important to develop sensitivity to resource consumption at all different levels and to determine it using standard methods. Further training for this concerns users, developers and experts, managers

and convenors of working groups. It is also worth considering adding information about resource consumption to hackathons/challenges and theses/journal publications.

## 6 Dynamic Energy Supply by K. Schwarz, B. Brüers, O. Kozlov, D. Hadasch, A. Kappes, M. Cruces

The substantial fraction of scientific computations are performed in batch mode, so delays in computations due to varying energy supplies seem acceptable to scientists. Several themes emerge in the use of renewable energy for scientific computing (inspired by the Andrew Grimshaw talk).

Large computing centers with all services should be located and operated at the seaside, where renewable energy is available without distant transfers. The location of a German cloud for science should be decided only according to sustainability criteria, and as of today only northern German states would be eligible. All associated scientific facilities can be connected via the DFN scientific fiber optic network. It is still conceivable that the large centers will be supplemented by small computing clusters at universities with possibly dedicated tasks.

The energy supply of the large centers should happen by wind, sun, biomass and gas turbines as far as possible sustainably. Heat production should be largely put to beneficial utilization (see Jens Struckmeier's presentation). Possibly the centers will be supported by energy storage, which could be a to be explored mixture of batteries, compressed gas, flywheels. Energy supply is forecast driven with criteria such as weather, pricing, etc. Needed for this is the establishment of communication channels both with energy suppliers on energy demand and with users, who should be informed on the one hand about resource consumption and further on delays in computations.

The actions at the computing centers on the adjusted energy demand, such as energy reduction due to shortage or excessive consumption due to energy surplus, must be controlled with predictive algorithms. In case of shortage, it has to be decided whether jobs are frozen or compensation is needed and how this can be realized. This could rely on local energy storage and possibly non-renewable energy.

Challenges will also be found on the part of legislators and supplier standards to make the regulations of large computing centers and their dynamic supply price attractive and practical.

## 7 Hardware Lifetime by M. Gasthuber, R. Walker, C. Felder, B. Fischer, R. Kleinemühl

The life cycle of a computing center has three areas in terms of $CO_2$ footprint: 1) construction and maintenance of the building, its infrastructure and computing facilities (approx. 30%), 2) use for computations (approx. 70%), 3) dismantling and disposal.

The lifetimes of electronic components can vary widely and must be evaluated on an individual basis accordingly. Experience shows that servers with mainboards and CPU have typical lifetimes of 3-5 years, while the 1-10 Gigabit network ports have very long lifetimes (more than 10 years). Other components are even 30 years in service. Contrary to common assumptions, electronic failure rates or setup difficulties are mainly found in the first few months, while components run stably for years after that. Older components can be removed from the main center and deployed in computing clusters that run only at lower priority and, for example, exclusively with renewable energy on demand. Lifetimes of more than 7, sometimes even 12 years are achievable. Further developments of operating systems over the lifetime do not cause any problems even on old hardware. Firmware enhancements vary among manufacturers.

The successful and continuous operation of a computing center depends on the expertise and commitment of the local experts. This also applies to the compatibility of the various components, their energy consumption in combination with all other components and the provision of spare parts. Manufacturers are reluctant to provide on-site support. Therefore, the use of older hardware is less about financial savings because potential new acquisition costs are balanced by the human resources of the experts.

The obvious advantage of running clusters with older hardware lies in avoiding the carbon footprint of new hardware and the disposal of old hardware. Even if the computing clusters run slower with older hardware, by using renewable energy, simulations and data analysis can be performed here without $CO_2$ footprint. This could provide a special award as an example for final theses and journal publications.

## 8 Hardware & Algorithms by F. von Cube, P. Fackeldey, L. Geiger, D. Schindler, J. Struckmeier

The interaction of algorithms and hardware has important significance for the efficiency of the resources used. The assessment of efficient use is hardly manageable by
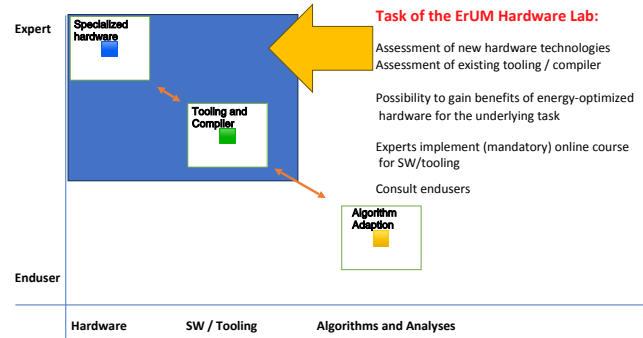


Fig. 1: ErUM-Data-Hardware-Lab

the normal user without the experts' evaluations and knowledge.

The experts know both the variants to specialized hardware, their functionalities and benchmarks as well as their resource efficiency. These include 1) GPU, TPU, FPGA, 2) CPU, SIMD ..., 3) ASICs, ARM64, 4) Neuromorphic and Quantum Computers. Furthermore, the experts know the tooling and compilers and thus the possibilities of efficient use of this hardware. Best practices for e.g. compiler flags, usage of current versions, JIT compilation, caching / checkpointing, alpaca, freezing.

In order for users to efficiently use the hardware for their needs and algorithms, they need knowledge about tools, vectorization, memory management, sustainable scheduling, which can be achieved through appropriate training.

To make this concert optimal, it makes sense to establish an ErUM Data Hardware Lab with expert staff and equipment for travel and hardware trials (Fig. 1. This could also be seen in conjunction with a second phase of the ErUM Data Hub.

The Lab's role would be to track and assess new hardware developments in relation to needs in ErUM communities. The Lab would be in close contact with the WLCG committees at CERN and would be specifically available as service to all ErUM communities.

The Lab educates users on the above-mentioned tools, promotes useful "middleware" to the communities, acts as contact for industry, and allows for prototyping with industry. Furthermore, a quick and unbureaucratic support for hardware procurement could be given. It benefits both research and industry, helping them to understand the needs. They could benchmark hardware, perform emission tracking of computing workflows, recommend energy-/CO2-optimized scheduling, and provide incentives for users to compute sustainably.

## 9 Smart Data by S. Wozniewski, M. Hannemann, C. Walther, K. Zhou

In order to extract information from large amounts of data, it is often necessary to reduce the amount of data to a level that is relatively user-friendly. On the one hand, the efficiency of the storage and, on the other hand, possible further uses of the data should be considered. This includes data and analysis preservation for the current scientific questions as well as possibly long-term storage of the data for later verification and further scientific exploitation in other contexts.

Policies are needed for these reductions, which can only be developed by experts with appropriate domain knowledge. The reduced data should be able to be produced in an automated way and should be self-explanatory, either through model variants such as pairs of declaration and value, or through metadata explaining the data structure, detector conditions, and their previous processings (ontology). Moreover, they should be easily machine-readable, especially if they are to be preserved for eternity.

## 10 Cultural Change by P. Niknejadi, M. Erdmann, S. Fröse, T. Kuhr, A. Schneidewind, Ö. Özkan, F. Schreiber

At the heart of a culture change is awareness of the sustainable use of resources (Fig. 2). The costs of data analysis, for example, consist not only in the immediate monetary energy costs (so far mostly invisible to the user), but also in the extensive footprints of hardware and human resources described above.

The focus of the individual scientists and working groups here is initially on concepts of sustainable working methods such as reliability, reproducibility (workflows), accessibility and portability of data analyses. In addition, there is also the conscious handling of the working time of directly or peripherally involved persons.

Many of these aspects are included in the variants of the term FAIR (finable, accessible, interoperable, reproducible) that apply to research data and now also to software and analyses. A change in awareness is to be supported on the one hand by means of education and furthermore in the form of prizes, certificates, citations. With all openness, recognition of intellectual properties is to be taken into account (see Embargo Periods for Open Data).

On the side of ministries and universities, structural support is needed in the form of funding for sustainability, which itself needs a sustainability seal. This should



Fig. 2: Aspects of cultural change

enable new career paths that have both attractiveness and sustainability as a focus of their activities.

Finally, the sustainable scientific exploration of data is about combining domain knowledge about the experiment and its questions with methodological skills from mathematics or computing science. Depending on the ErUM community, they may or may not come from a single hand. On the one side, a lot of communication is necessary to bring the domains together, on the other side, an improved and thus more sustainable use of the measurement data can be expected due to the combination of the different expertises.

## 11 Autonomization

## 12 Inquiries & Dynamics

## 13 Algorithmics & Software

## 14 Machine Models

## 15 Injected Intelligence

## 16 Workflow & Stakeholders

## Appendix

May be removed after writing: Twelve Guiding Questions: Sustainability in the Digital Transformation of Basic Research on Universe & Matter

Hardware & Research Data

1. Footprint: Constructing a comprehensive picture of the footprint of all ErUM-Data related activities. Where does quantitative knowledge exist, where is it lacking? What resource needs do you see, what opportunities for savings? What innovations are needed to keep sustainable use of resources in balance with demands? To what extent does continuing education play a role? How can feedback reduce a footprint through ML methods?
2. (Dynamic) Energy Supply: Where to locate & operate computing systems incl. storage? How could a dynamic energy supply look like, which largely covers the needs of ErUM-Data related activities with renewable energies? What information flows would be required for this? What mechanisms and what dynamics are required on a supra-regional basis to create compensation possibilities for windless/sunless periods?
3. Hardware Lifetime: How could prolonged / optimized usage of hardware resources in view of technology evolution be modeled beyond their usual lifetimes? What short- and medium-term monitoring would be required to signal indispensable replacements on the one hand, and to execute computing jobs matching their algorithmic requirements on prolonged or current hardware on the other?
4. Hardware & Algorithms: Which adaptive measures for hardware and algorithms could have a decisive impact on ErUM-Data? Which types of hardware (including e.g., GPU, TPU, FPGA, neuromorphic computing) could be considered and which automated mechanisms exist for adapting algorithms to non-specific or dedicated hardware?
5. Smart Data: Deciding when and how to discard information without losing scientific value, based on learning from nature and experiment. What mechanisms for transforming data to smart data can be envisioned, and how can evaluation and control of information gain or loss be accomplished? How can archiving and retrieving data be managed?
6. Cultural Change: What could a comprehensive educational area for rethinking, among other things, the use of computer hardware, actually required information (smart data), preparation of data packages (event loops versus event chunks), etc. look like?

How can we change to a culture of data reuse? Assessment of ethical implications and risk assessment.

Algorithms & Mindset

7. Autonomization: We witness the transfor-mation from the era of automation to an era of autonomization (e.g., unsupervised learning). Where will ErUM-Data benefit from autonomization, which innovations are necessary and how can the reliability of the autonomously obtained results be ensured?
8. Inquiries & Dynamics: How can input questions be posed to generate the best possible output from the machines? What relevance will dynamic learning algorithms and machines have for the field of ErUM-Data?
9. Algorithmics & Software: Our thinking in algorithms and software has a direct impact on resource requirements. What can sustainable algorithm & software engineering and an associated educational program in algorithm & software development look like to get ErUM-Data to the forefront of developers?
10. Machine Models: Pre-trained and gene-rative models have a high potential for energy savings in both their creation and usage of machine learning. What innovations are needed to achieve a reliable routine operation?
11. Injected Intelligence: How can reasoning by the physicist, mathematician, or any other kind of intelligence speed up the processes of learning or make them more energy efficient? What measures can we apply to avoid constantly reinventing the wheel? What can knowledge discovery of work already performed look like?
12. Workflow & Stakeholders: How can well-defined, reproducible workflows with high user dynamics (data analyses) be captured that remain functional in the long term? How can an overall picture be created with all stakeholders working together on a large-scale project for the benefit of sustainability across their departmental boundaries?

**Data Availability Statement**

No explicit data were used in this study.

**Literature**

[1] Shankha Banerjee et al. *Striving towards environmental sustainability in High Energy Physics, Cosmology and Astroparticle Physics (HECAP)*. 2022. URL: https://sustainable-hecap.github.io.