

Dr. Andreas Knüpfer, Prof. Dr. Wolfgang E. Nagel
ZIH / TU Dresden

Forschungsgebiete des ZIH und der Professur für Rechnerarchitektur

Vorlesung Forschungslinie - Einführung in die Forschung
2019-07-08

Overview

- Overview about ZIH
- Research Fields
- Selected Research Topics, Projects, and Results

Overview about ZIH

ZIH Structure

- Director Prof. Dr. Wolfgang E. Nagel
- Deputy directors Dr. Björn Gehlsen and Dr. Andreas Knüpfer
- 7 departments

- Networking and Communication Services (NK)
- Operational Processes and Systems (OPS)
- Systems Design and Engineering (SDE)
- Service Desk

Service
oriented

- Interdisciplinary Application Development and Coordination (IAK)
- Distributed and Data Intensive Computing (VDR)
- Innovative Methods of Computing (IMC)

Research
oriented

- In total over 160 staff including apprentices at 6 locations on the campus

Topics overview

IT Services for TU Dresden

- Campus network, internet uplink, phone network
- E-Mail, groupware, data exchange, backup etc.
- Software procurement

Computational Science Services

- Virtual machines, hosting
- Supercomputing
- Big Data methods

Research and Development

- Parallel programming and algorithms
- Performance optimization and scaling, energy efficiency, ...
- Data Analytics applications, machine learning, ...

Topics overview

IT Services for TU Dresden

- Campus network, internet uplink, phone network
- E-Mail, groupware, data exchange, backup etc.
- Software procurement

Computational Science Services

- Virtual machines, hosting
- Supercomputing
- Big Data methods

Research and Development

- Parallel programming and algorithms
- Performance optimization and scaling, energy efficiency, ...
- Data Analytics applications, machine learning, ...

IT Services and the Service Catalog

- <https://tu-dresden.de/zih/dienste/service-katalog>
- Topic lists and service descriptions

ZENTRUM FÜR INFORMATIONSDIENSTE UND HOCHLEISTUNGSRECHNEN (ZIH)

Service Katalog

ANGEBOTENE DIENSTE FÜR EINZELNE NUTZER

Dienst	Studierende	Mitarbeiter	Gast
Zugangsvoraussetzung			
>ZIH Benutzer-Login	+	+	+



Forschungslinie – Einführung in die Forschung, ZIH, 2019-07-08

Service Katalog

ANGEBOTENE DIENSTE FÜR EINZELNE NUTZER

Dienst	Studierende	Mitarbeiter	Gast
Zugangsvoraussetzung			
>ZIH Benutzer-Login	+	+	+
Arbeitsumgebung			
>PC-Arbeitsplätze	IPC-Prüfung	IPC-Service	IPC-Service
>Software	MS-Betrieb	+	Spezial-Guest
>Laptop/Computer via WLAN & VPN	+	+	+
>Virtuelle	-	+	+
>Datenspeicher	-	+	+
>Backup & Archiv	Info über Datenspeicher oder VPC-Projekt	+	+
>E-Mail	-	+	+
>Druckerei	+	+	+
>Drucker / Scanner / Plotter	+	+	+
>Handbücher	-	+	+
Zusammenarbeiten & Forschen			
>Einarbeitung	Elaborat	+	+
>Vide / Videokonferenzen	Info als Subantrag	+	+
>Web-Lernen	-	+	+
>Software-Ermittlungsumgebung	+	+	+
>Eingeweinte (Skalierbare/Cloud)	-	+	+
Simulation / Hochleistungsrechnen			
>Zugriff auf die HPC-Ressourcen	Info über existierende Projekte	Info über existierende Projekte	Info über existierende Projekte
Service Desk			
>Dienstreue Ansprechpartner bei allen IT-Fragen			

ANGEBOTENE DIENSTE FÜR EINRICHTUNGEN DER TU DRESDEN

Dienst	Studierende	Mitarbeiter	Gast
Arbeitsumgebung			
>Zentrale Verwaltung von PC-Arbeitsplätzen / PC-Prüfung			
>Virtuelle Desktop (Active Directory) (VDI)			
>Kollaborationsumgebung (SharePoint)			
>Ermittlung der Nutzung von Software (License Management)			
>Anbindung von Einrichtungen aus externen Datenbanken			
>Zentrale IT-Abrechnung			
>MSI (Nameverwaltung)			
>Freizeit am Übergang zum Campusnetz			
>Remote Support			
Zusammenarbeiten & Forschen			
>Einarbeitung			
>Server-Hosting / Virtualisierung			
>Dienstleistung			

WEITERE IT-BEZUGENE BEREICHE DER TU DRESDEN

	Dienst	Links
Sachgebiet Informationsicherheit	Unterstützung bei allen Fragen rund um Datensicherheit und IT-Sicherheit	> Sachgebiet 3.5 Informationsicherheit
Medienzentrum	Angaben rund um audiovisuelle Medien, den Web-Auftritt der TU Dresden und die E-Learning-Plattform	> Medienzentrum > VHS > SFA
Zentrale Universitätsverwaltung	Lösungen zur zentralen Datenverwaltung für Angestellte und Studierende der TU Dresden	> ERP-System der TU Dresden > HIS2

IT Service Catalog

- <https://tu-dresden.de/zih/dienste/service-katalog>
- Topic lists and service descriptions
 - Target groups (students, staff, guests)
 - Specific service offerings
 - Details about
 - Scope of services and options
 - Application and permissions
 - Conditions and obligations
 - Accounting and costs
 - Service levels

Service Katalog

ANGEBOTENE DIENSTE FÜR EINZELNE NUTZER

Dienst	Studierende	Mitarbeiter	Guest
Zugangsverwaltung			
» ZIM Benutzer Login	+	+	+
Arbeitsumgebung			
» PC-Arbeitsplatz	IPC-Print	IPC-Service	IPC-Service
» Software	MS Office	+	Spezial Desktop
» Logging Component via LDAP & VPN	+	+	+
» Virenschutz	-	+	+
» Datenspeicher	-	+	+
» Backup & Archiv	Inter über Datenspeicher oder VPC-Projekt	+	+
» E-Mail	-	+	+
» Drucken	+	+	+
» Drucken / Scannen / Plotten	+	+	+
» Handbücher	-	+	+
Zusammenarbeiten & Forschen			
» Datenabruf	Elisabreit	+	+
» Virenschutz / Antivirenschutz	Inter als Subsystem	+	+
» E-Mail-Listen	+	+	+
» Software-Ermittlungsumgebung	+	+	+
» Einzelnutzer (Mitarbeiter/Guest)	-	+	+
Simulation / Hochleistungsrechnen			
» Zugriff auf die HPC-Ressourcen	Inter über existierende Projekte	Inter über existierende Projekte	Inter über existierende Projekte
Service Desk			
» Zentraler Ansprechpartner bei allen IT-Fragen			

ANGEBOTENE DIENSTE FÜR EINRICHTUNGEN DER TU DRESDEN

Dienst
Arbeitsumgebung
» Zentrale Verwaltung von PC-Arbeitsplätzen / PC-Print
» Microsoft Active Directory (AD)
» Authentifizierung (Identity Provider)
» Enterprise Mailung von Software (LicenseCentral)
» Verbindung von Einrichtungen aus zentrale Datenbank
» Zentrale IP-Adressverwaltung
» DNS (Namenauflösung)
» Firewall am Übergang zum Campusnetz
» Remote Support
Zusammenarbeiten & Forschen
» Datenabruf
» Server-Hosting / Virtualisierung
» Dienstleistung

WEITERE IT-BEZUGENE BEREICHE DER TU DRESDEN

	Dienst	Links
Servicecenter Informationsicherheit	Unterstützung bei allen Fragen rund um Datensicherheit und IT-Sicherheit	» Servicecenter 3.0 Informationsicherheit
Medienzentrum	Angaben rund um audiovisuelle Medien, den Web-Auftritt der TU Dresden und die E-Learning-Plattform	» Medienzentrum » VHS/ML » SFA
Zentrale Universitätsverwaltung	Lösungen zur zentralen Datenverwaltung für Angestellte und Studierende der TU Dresden	» ERP-System der TU Dresden » HIS/IS

IT Service Katalog

Services for individuals

- PC work environment
- Software
- Campus network, wireless, VPN
- Phones
- Data storage
- Backup & archiving
- E-Mail
- Time sync. services
- Printing, scanning, plotting
- Manuals
- Data exchange
- Video and phone conferences
- Mailing lists
- Groupware
- HPC access
- Service Desk

Services for departments

- Central Administration of PCs and PC pools
- Directory services (AD/LDAP)
- Authentication
- Software licenses
- Campus network connections
- IP address mgmt., DNS
- Firewalls
- Remote Support
- Data exchange
- Server hosting, VMs
- Service monitoring

Further service providers

- Information security unit
- Media center
- Central University mgmt.

Dienst	Studierende	Mitarbeiter	Guest
Zugangsverwaltung			
▪ DM Benutzer Login	+	+	+
Arbeitsumgebung			
▪ PC-Arbeitsplatz	IPC-Pool	IPC-Server	IPC-Server
▪ Software	MS-Betrieb	+	Spezial-Guest
▪ Logging (Computer via WLAN & VPN)	+	+	+
▪ Station	-	-	-
▪ Dienstleister	-	-	-
▪ Backup & Archiv	Man über Datenpartner oder VPC-Projekt	+	+
▪ E-Mail	-	+	+
▪ Adressen	-	+	+
▪ Drucker / Scanner / Plotter	+	+	+
▪ Handbücher	-	-	-
Zusammenarbeiten & Foren			
▪ Datentausch	(Cloud) oder	+	+
▪ Vire- / Scannerlösungen	Man als Subnetze	+	+
▪ E-Mail-Listen	-	+	+
▪ Software-Ermittlungsmöglichkeit	-	+	+
▪ Einzige (Skalierbar/Agil)	-	+	+
Simulation / Hochleistungsrechnen			
▪ Zugriff auf die HPC-Ressourcen	Man über existierende Projekte	Man über existierende Projekte	Man über existierende Projekte
Service Desk			
▪ Zentraler Ansprechpartner bei allen IT-Fragen			

Dienst
Arbeitsumgebung
▪ Zentrale Verwaltung von PC-Arbeitsplätzen / PC-Pools
▪ Webverzeichnis (Active Directory) (LDAP)
▪ Authentifizierung (Identity Provider)
▪ Enterprise Mailbox von Software & Contentmail
▪ Verbindung von Einrichtungen aus externen Datenbanken
▪ Zentrale IP-Adressverwaltung
▪ DNS (Namenauflösung)
▪ Firewall am Übergang zum Campusnetz
▪ Remote Support
Zusammenarbeiten & Foren
▪ Daten-Tausch
▪ Server-Hosting / Virtualisierung
▪ Dienstleistung

	Dienst	Links
Servicecenter Informationsicherheit	Unterstützung bei allen Fragen rund um Datensicherheit und IT-Sicherheit	▪ Servicecenter Informationsicherheit
Medienzentrum	Angaben rund um audiovisuelle Medien, den Web-Auftritt der TU Dresden und die E-Learning-Plattform	▪ Medienzentrum TU Dresden ▪ SFA
Zentrale Universitätsverwaltung	Lösungen zur zentralen Datenverwaltung für Angestellte und Studierende der TU Dresden	▪ ERP-System der TU Dresden ▪ HISQS

Self Service Portal

Self Service Portal

<https://selfservice.zih.tu-dresden.de/>

- Booking IT services
- Fully automated or with approval
- Requires ZIH login

	Dienst	Portal	Beschreibung
Zugangsvoraussetzung	ZIH-Benutzer-Login	Login-Antrag	Antrag für ein ZIH-Login generieren
		ZIH-Login über IDM verlängern	Im IDM können Sie ab sofort Gäste und Funktions-Logins verlängern, für die Sie als Kontaktperson zugeordnet sind.
		ZIH-Passwort ändern	ZIH-Passwort im IDM-Portal ändern
Arbeitsumgebung	Telefon	Sammlung Telekommunikations-Daten	Telefondaten ansehen und ggf. notwendige

Topics overview

IT Services for TU Dresden

- Campus network, internet uplink, phone network
- E-Mail, groupware, data exchange, backup etc.
- Software procurement

Computational Science Services

- Virtual machines, hosting
- Supercomputing
- Big Data methods

Research and Development

- Parallel programming and algorithms
- Performance optimization and scaling, energy efficiency, ...
- Data Analytics applications, machine learning, ...

HPC Resources at ZIH

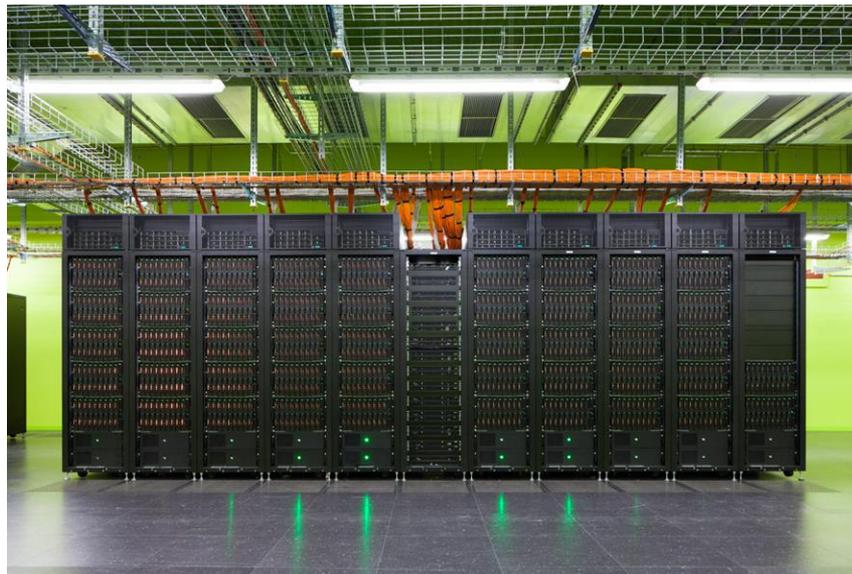
HRSK II 2015

- ~ 44,000 cores Intel (mostly Haswell)
- 256 GPUs Nvidia Tesla K80 +
- 44 GPUs: Nvidia Tesla K20
- 136 TB RAM, >5 PB scratch file system

HPC-DA extension 2018

- 22 Machine Learning nodes IBM AC922
- 2 PB NVME storage (90 nodes, NVMeoF) with 2 TB/s bandwidth in total
- 10 PB Object Storage

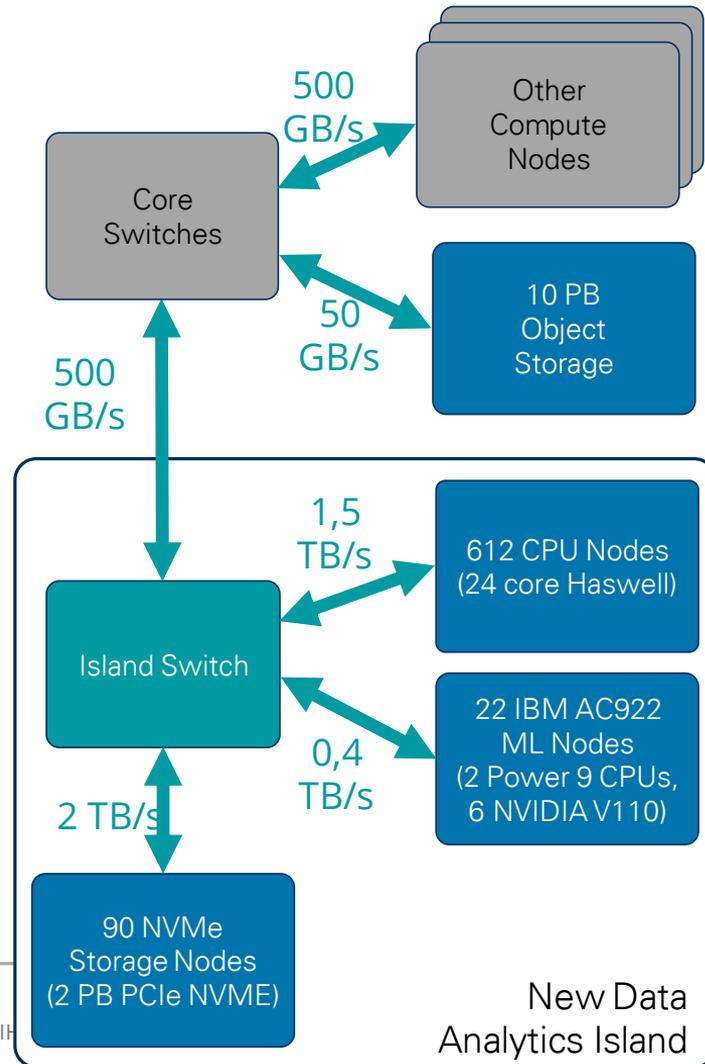
Follow-up extension in 2019 (approx. 4 M€)



HPC-DA Extension 2018

HPC-DA extension towards extremely fast I/O

- Redesigned one compute island of HRSK II
- Strong focus on highest bandwidth and low latency
- 612 existing CPU compute nodes
- 22 new Machine Learning Knoten IBM AC922
 - 2x Power-9 CPUs, 6x NVIDIA V100 GPUs, NVLink
- 90 NVME storage nodes
 - Each node with 8 3,2 TB PCIe x4 NVME cards
 - Dual-link EDR IB, NVME over fabric
- 10 TB Object Storage with 50 GB/s



HPC-DA ML Nodes

Hardware

- 22 IBM AC922 nodes
- 2x POWER9 CPU, 22 core, 4-way HT (176 threads per node in total)
- 2.80 GHz, 3.10 GHz boost
- 256 GB RAM DDR4 2666MHz
- 6x NVIDIA VOLTA V100 with 32GB HBM2
- NVLINK with 150 GB/s between GPUs and between host and GPUs
- CPUs and GPUs direct water cooled
- 0.4 TB/s aggregated bandwidth to NVME nodes

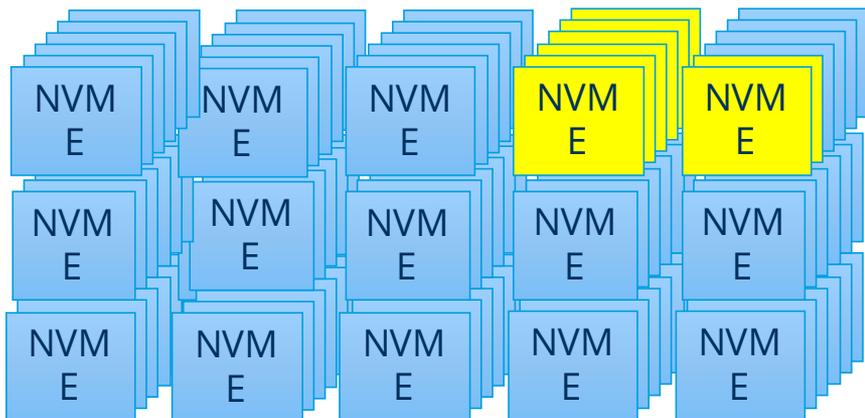


Image: <https://www.ibm.com/it-infrastructure/power/accelerated-computing>

HPC-DA NVME Usage Models

Allocation strategy

- NVME shares allocated as long-term “NVME leases” (weeks to months)
- Granularity of 1/8th node (1 NVME card) or full NVME nodes



Separate BeeGFSes in own NVME lease

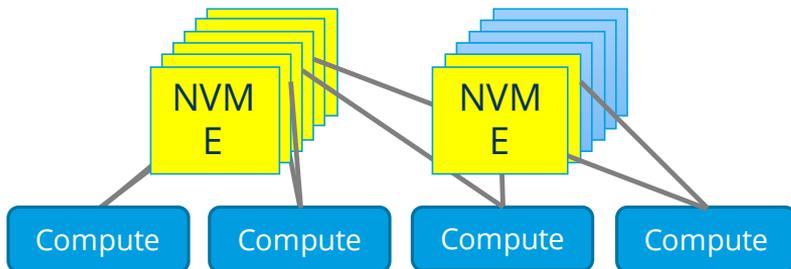
- Instantiate separate BeeGFS
- Granularity of ½ or 1 NVME nodes, including the EDR links and the CPU cores
- Separate MDSs and OSTs for this FS
- No meta-data interference with everyone else
- Full nominal bandwidth per NVME node*
- Have it mounted automatically to compute jobs of your HPC project

* Currently aggregated speed still limited due to IB routing issues.

HPC-DA NVME Storage Nodes

NVME over fabric to compute nodes

- Assign to compute nodes in n:m way
- NVMEof provides local block devices with $< 300 \mu\text{s}$ random access latency
- Need to manage (exclusive/parallel) accesses on your own! (through mmap or local FS)



Own services

- Run own services on NVME nodes, preferably data or storage services
- Reduced total IB bandwidth between the NVME nodes
- Consult HPC admin team in case this seems promising to you

Software for HPC-DA: Modules and Containers

Software modules

- Long list of software packages, multiple versions each, dependency management
- Open Source SW, scientific community software packages, commercial SW
- Application software, libraries, software tools

Singularity containers

- Tailored software environments that you can take with you or share with others
- Can be defined/built on top of each other
- Cannot combine two existing containers

- Challenge to build containers for Power9 because few have Power9 laptops ☹ yet.



Topics overview

IT Services for TU Dresden

- Campus network, internet uplink, phone network
- E-Mail, groupware, data exchange, backup etc.
- Software procurement

Computational Science Services

- Virtual machines, hosting
- Supercomputing
- Big Data methods

Research and Development

- Parallel programming and algorithms
- Performance optimization and scaling, energy efficiency, ...
- Data Analytics applications, machine learning, ...

Research Fields

ZIH Research Topics

- Scalable software tools to support the optimization of applications for HPC systems
- Data Intensive Computing and Data Life Cycle
- Performance and energy efficiency analysis for innovative computer architectures
- Distributed Computing and Cloud Computing
- Data analysis, methods and modeling in life sciences
- Parallel programming, algorithms and methods



Topic: Data Intensive Computing and Data Life Cycle

- ADA-FS - Advanced Data Placement via Ad-hoc File Systems at Extreme Scales
- EMuDIG 4.0 - Factory sensor monitoring
- EXPLOIDS - Monitoring tools for IT security
- GeRDi - Generic Research Data Infrastructure
- High Performance Deep Learning Framework
- MASi - Metadata Management for Applied Sciences
- IT support projects for SFB/TRR 205 and SFB 940
- ScaDS Dresden/Leipzig - National Big Data competence center
- VAVID - Comparative Analysis of Engineering Measurements and Simulation Data



Topic: Performance and Energy Efficiency Analysis for Innovative Computer Architectures

- READEX: Runtime Exploitation of Application Dynamism for Energy Efficient Exascale Computing
- HAEC: Highly Adaptive Energy Efficient Computing
- FIRESTARTER: A processor stress test utility
- HDEEM: High Definition Energy Efficiency Monitoring
- SPEC: Standard Performance Evaluation Cooperation



Topic: Distributed Computing and Cloud Computing

- UNICORE - Middleware for distributed computing and data
- VAVID - Comparative Analysis of Engineering Measurements and Simulation Data
- Chemomentum: Grid based software for complex chemistry workflows with a focus on data and knowledge management
- D-Grid Integration project (DGI 1 and 2), European Middleware
- FutureGrid: Experimental test and development environment
- GeneCloud: Secure semantic high performance cloud computing
- HEPCG: High Energy Physics Community Grid



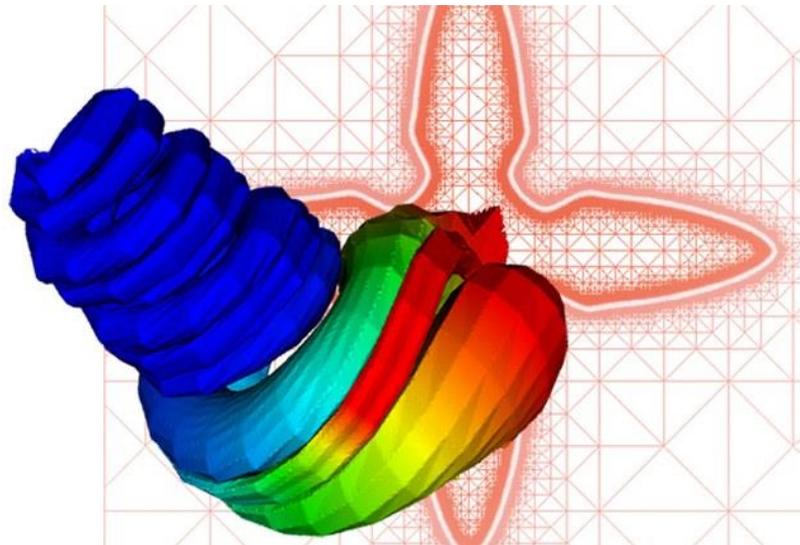
Topic: Data analysis, methods and modeling in life sciences

- Standardising the exchange of multicellular models in computational systems medicine
- Collective motion and swarming
- Stochastic processes, interacting cell systems and cellular automata
- Tumour development
- Endocytosis and systems biology
- Spatio-temporal pattern formation in cells and tissues
- Regeneration
- Bone remodelling



Parallel programming, algorithms and methods

- DASH: Hierarchical Arrays for Efficient and Productive Data-Intensive Exascale Computing
- High Performance Deep Learning Framework
- IPCC: Intel® Parallel Computing Center TU Dresden
- GCoE: NVIDIA GPU Center of Excellence
- MEPHISTO - Metaprogramming for Heterogeneous Distributed Systems
- PARADOM - Parallel Algorithmic Differentiation in OpenModelica
- ScaFES: Scalable Framework for Explicit Solvers
- OpenMP and OpenACC Standardization



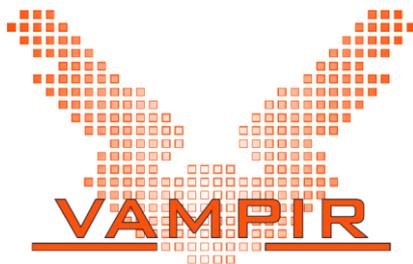
Selected Research Topics and Results

Parallel Performance Analysis Tools

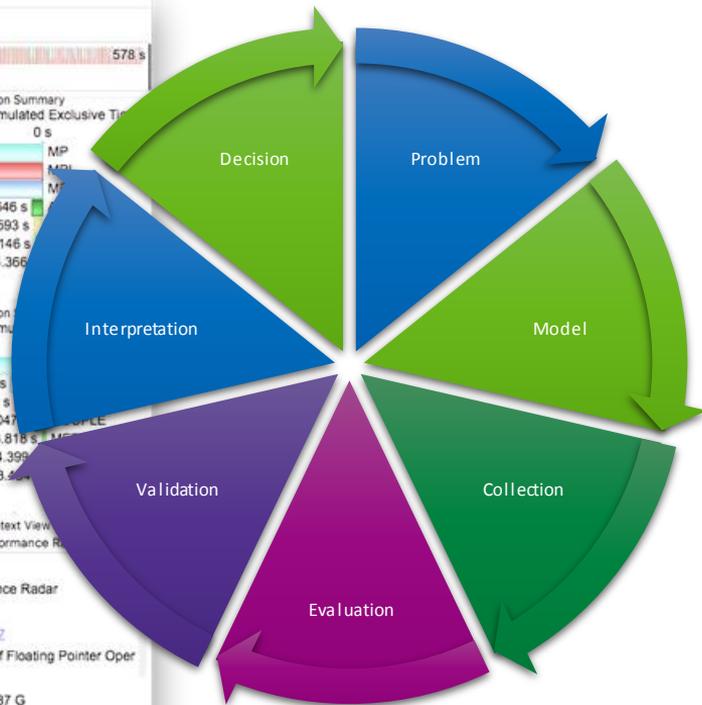
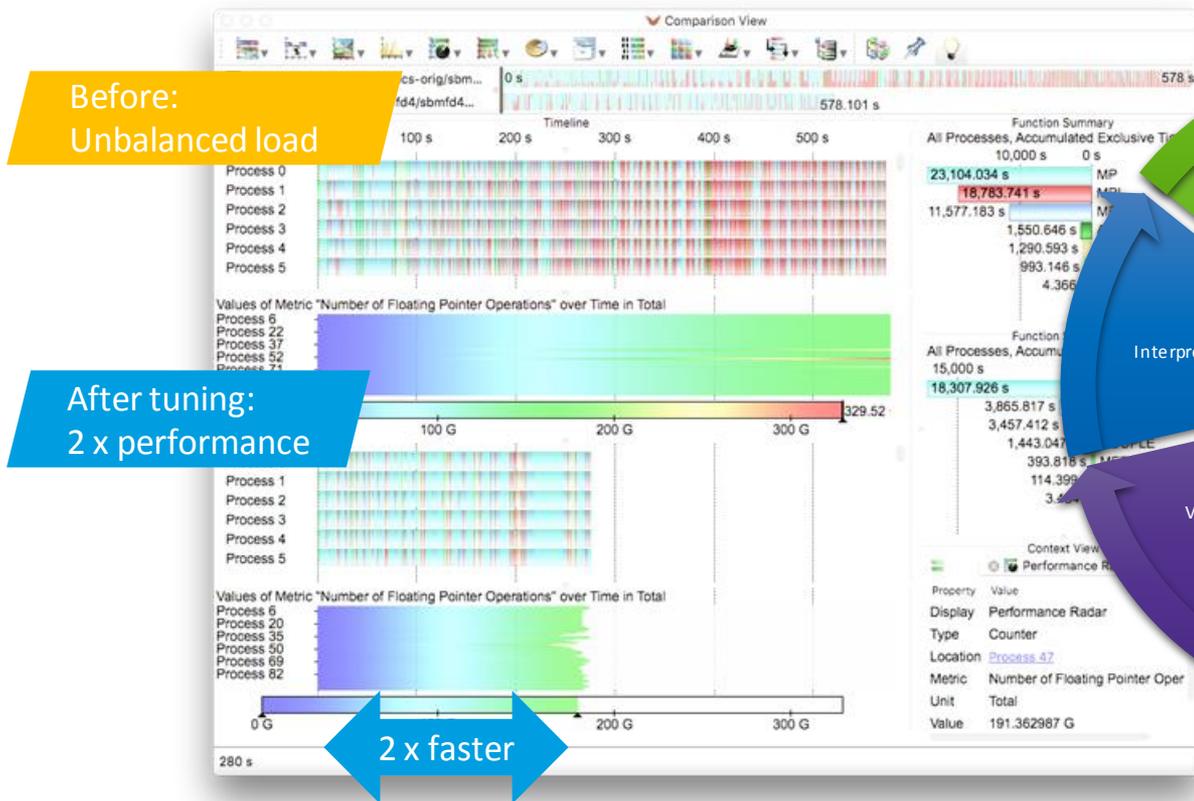
Parallel Performance Analysis

Enable or improve computational sciences

- Throughput
- Response time
- Scaling
- Quality
- Additional functionality
- Reliability
- Development Cycle

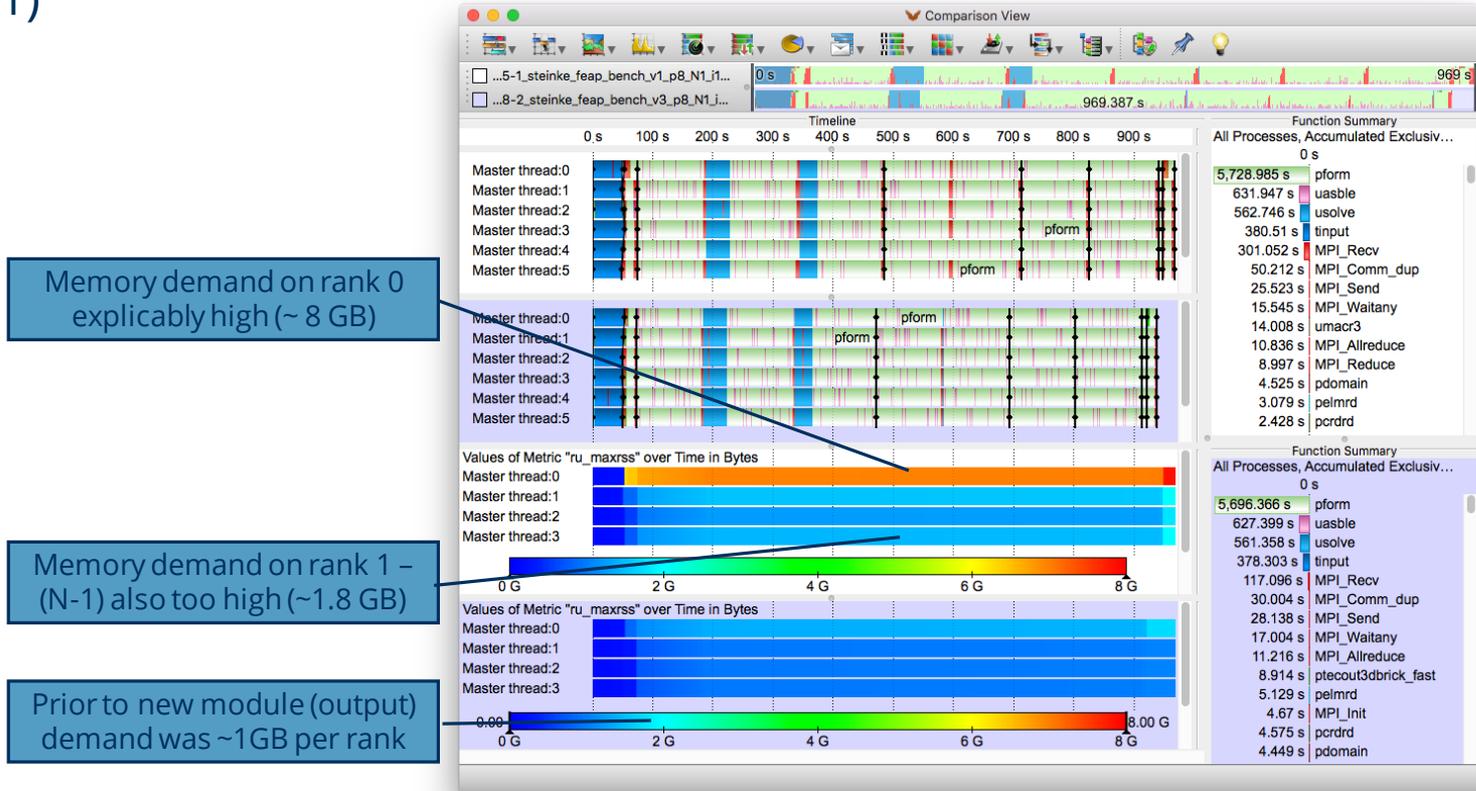


Performance Consulting is a Cyclic Process



Analyzing Unexpected Memory Demand

Analysis (1)



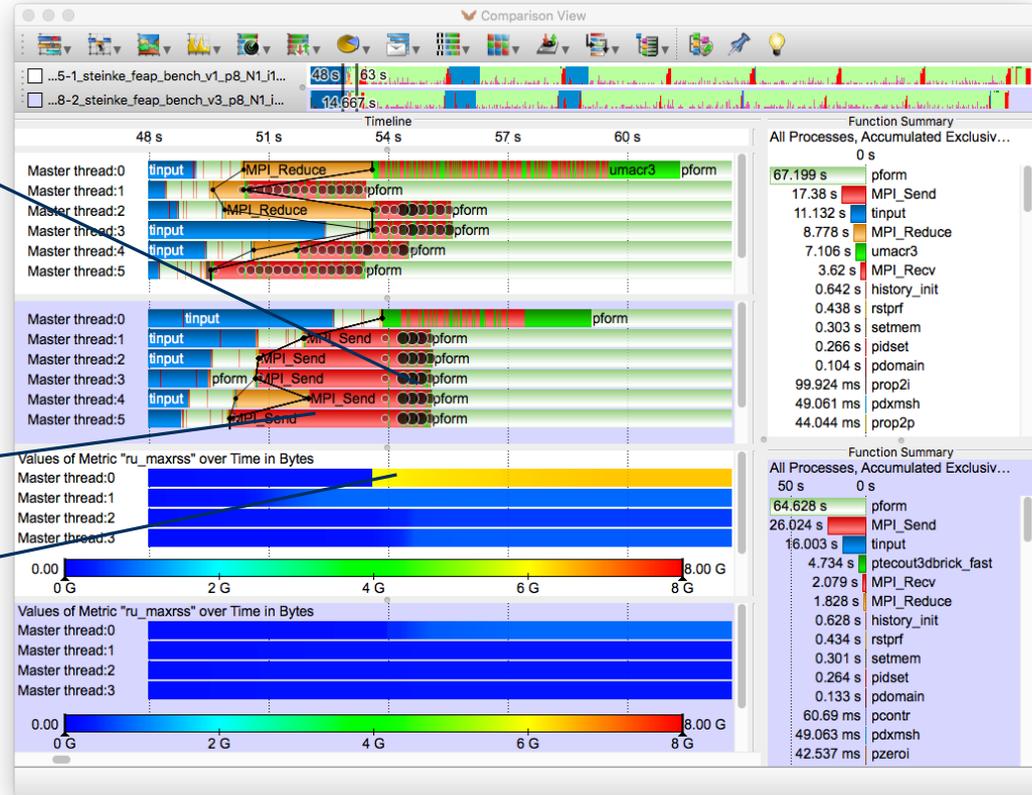
Analyzing Unexpected Memory Demand

Analysis (2)

Lots of messages are sent to rank 0

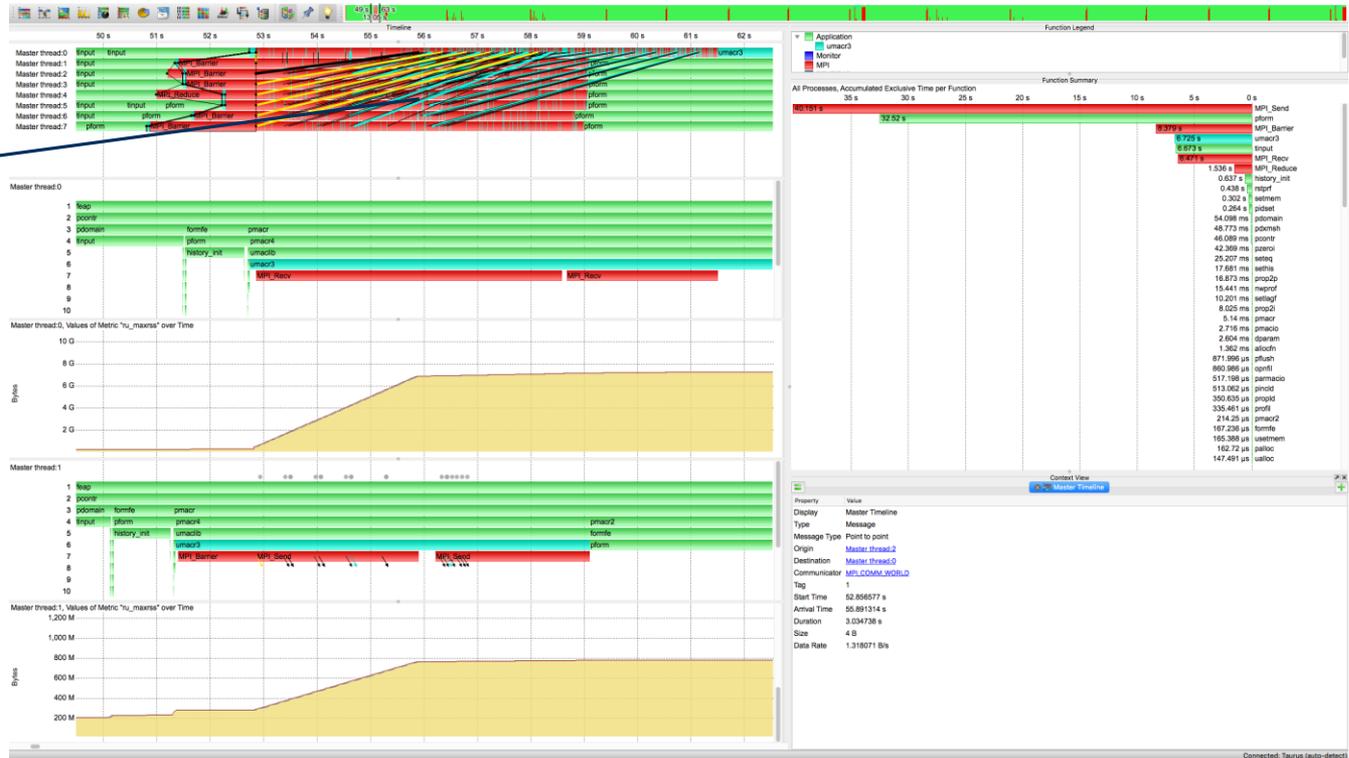
Educated guess: MPI needs to buffer data...

Memory demand explodes during MPI_Reduce on rank 0

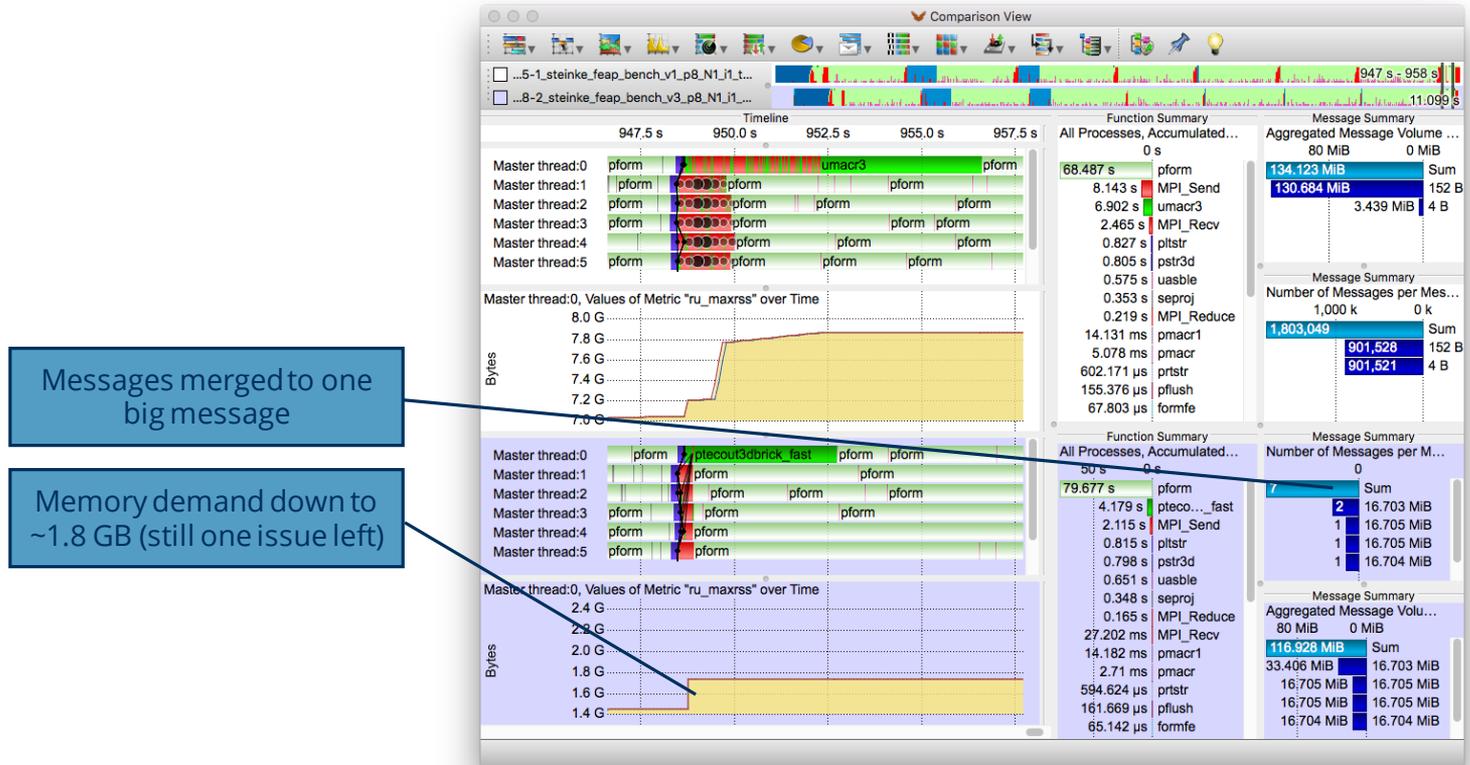


Analyzing Unexpected Memory Demand Analysis (3)

Cause: data is received rank by rank

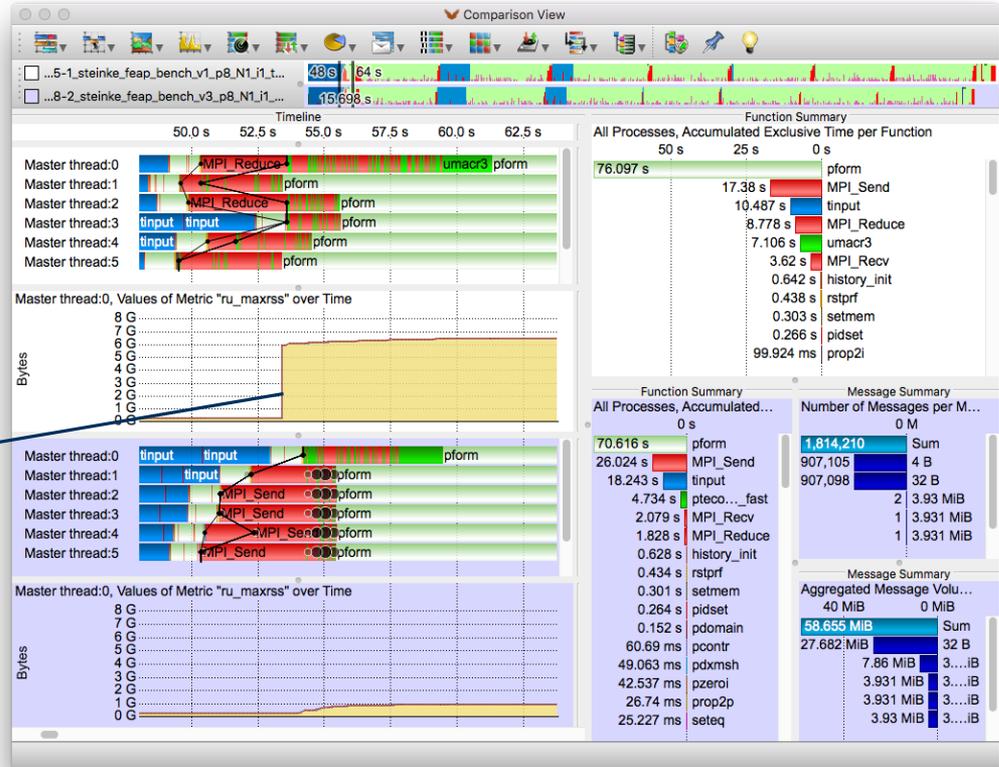


Analyzing Unexpected Memory Demand Analysis (4)



Analyzing Unexpected Memory Demand Analysis (5)

Many small messages received late, causing a buffer issue

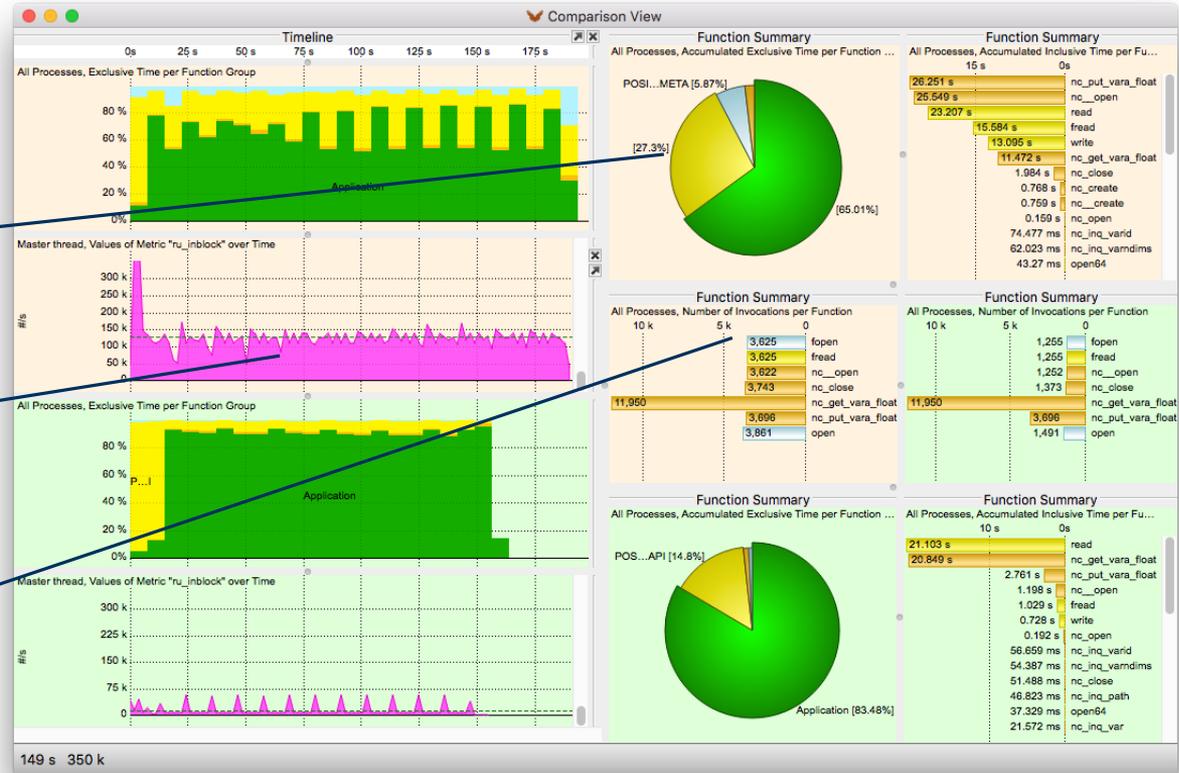


Analyzing Bottleneck in Multi-Stage-I/O Analysis (1)

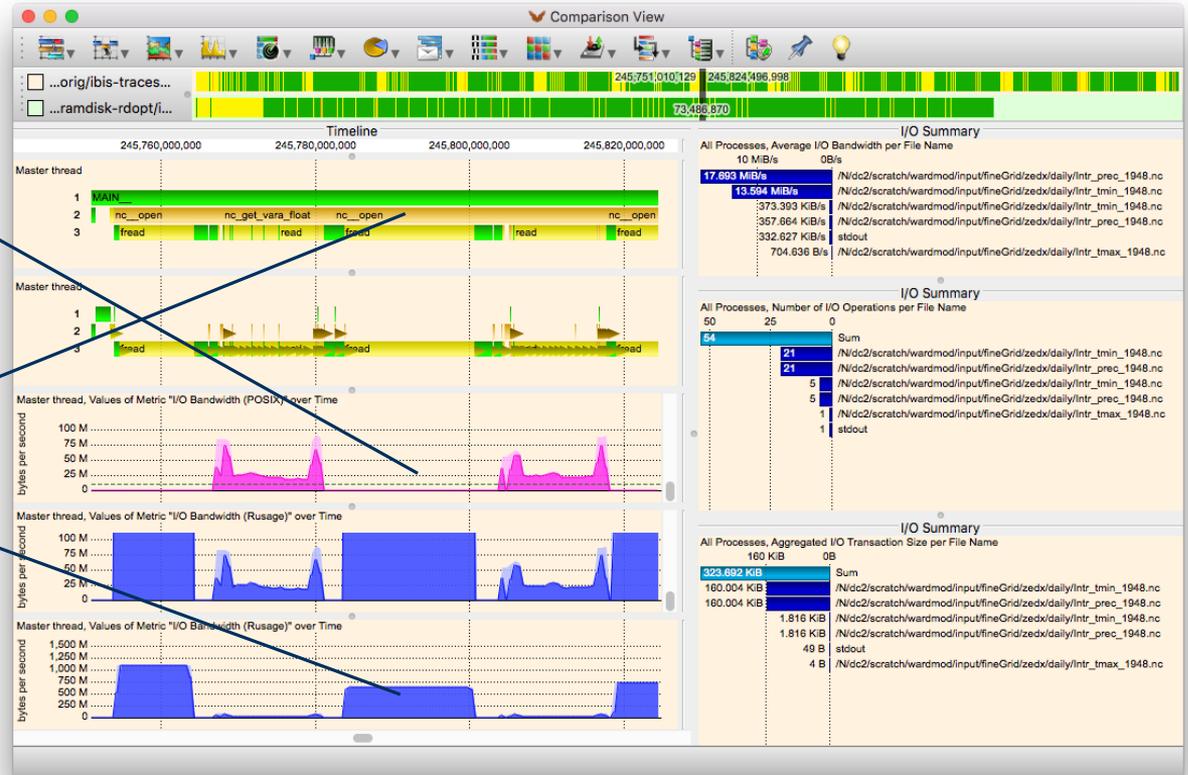
High I/O load on a single rank

~ 260 MB/s per rank !

High number of I/O open AND read operations



Analyzing Bottleneck in Multi-Stage-I/O Analysis (2)



I/O load measured from inside the application: 10 MB/s

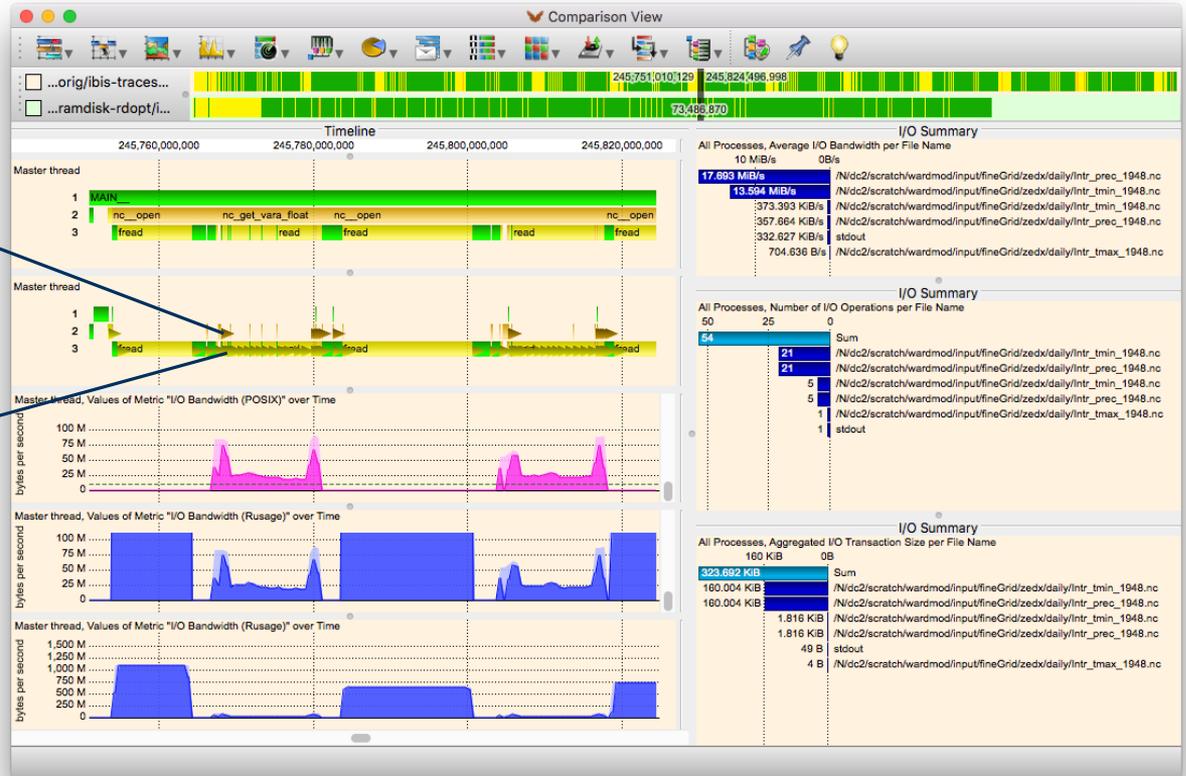
Data pre-fetching caused by open!?

I/O load measured on the kernel level: 260 MB/s !!!

Analyzing Bottleneck in Multi-Stage-I/O Analysis (3)

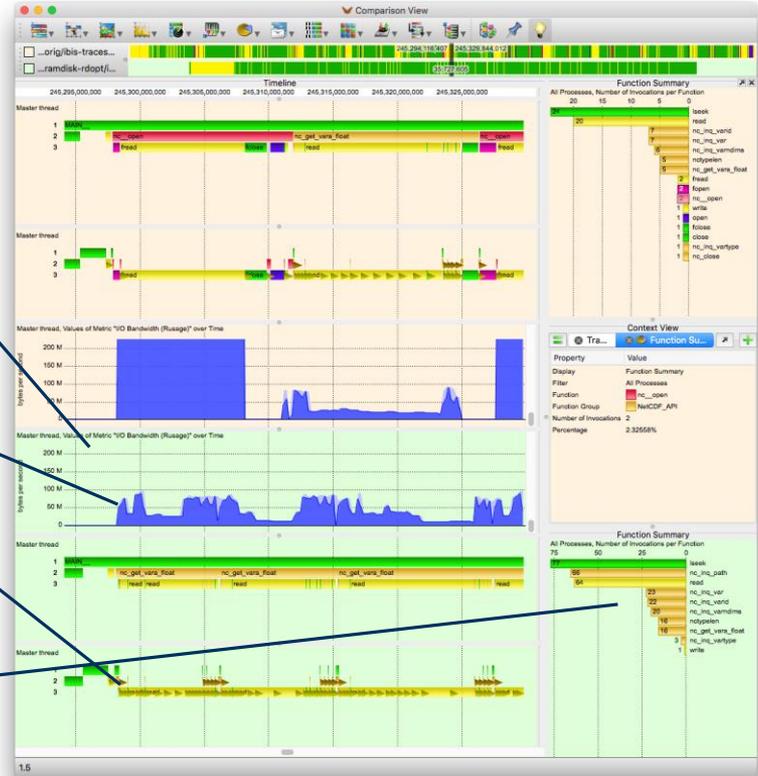
Single NetCDF call triggers a burst of POSIX calls

Strides in file access pattern !?



Analyzing Bottleneck in Multi-Stage-I/O Analysis (4)

- After tuning...
- I/O data rate at 10 MB/s per Rank
- I/O latency reduced by means of SSDs
- File open requests cached



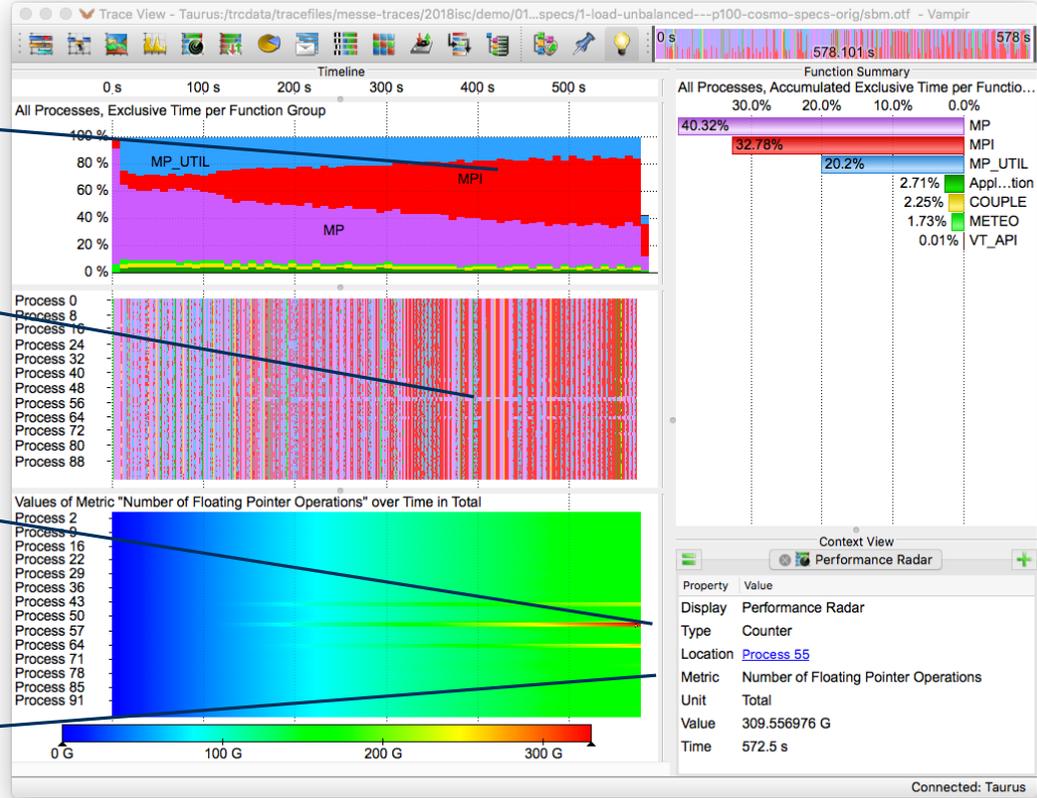
Load Balance Analysis of Weather Forecast Model Analysis

Message Passing share increases over time. Happens uniformly?

No, MPI Share does not increase on Process 56. Why not?

Process 56 is loaded heavily with FLOPs starting at t=130s

99 Processes wait for one process to finish

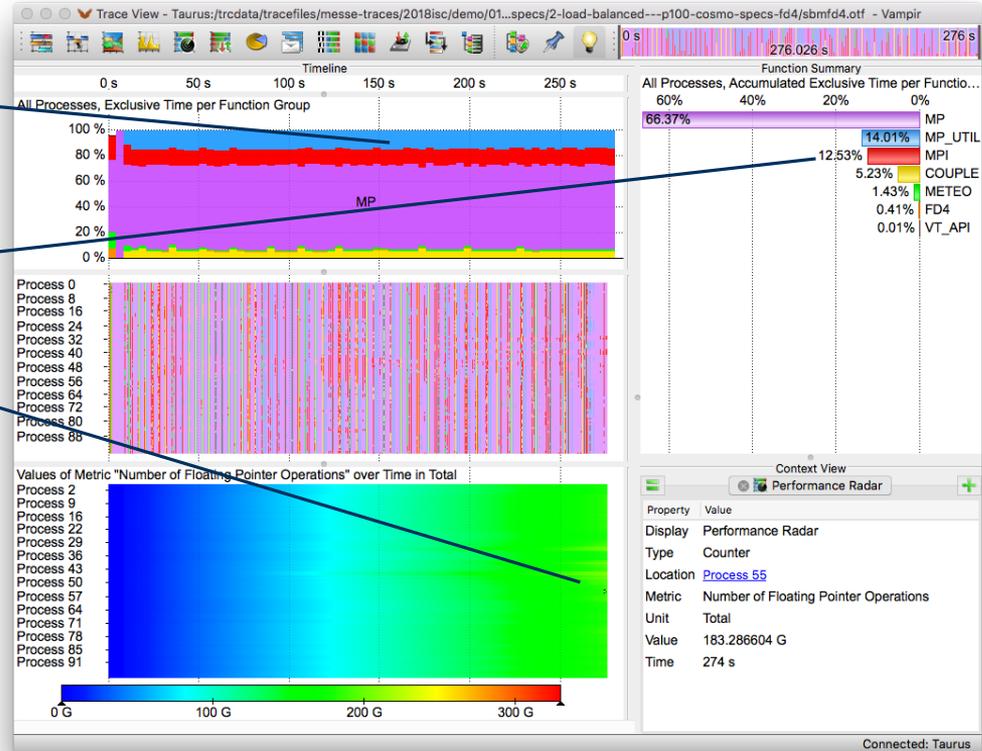


Load Balance Analysis of Weather Forecast Model After Tuning

No increase in MPI share anymore

MPI consumes 12.5 % of the total time

FLOP load is equally distributed now

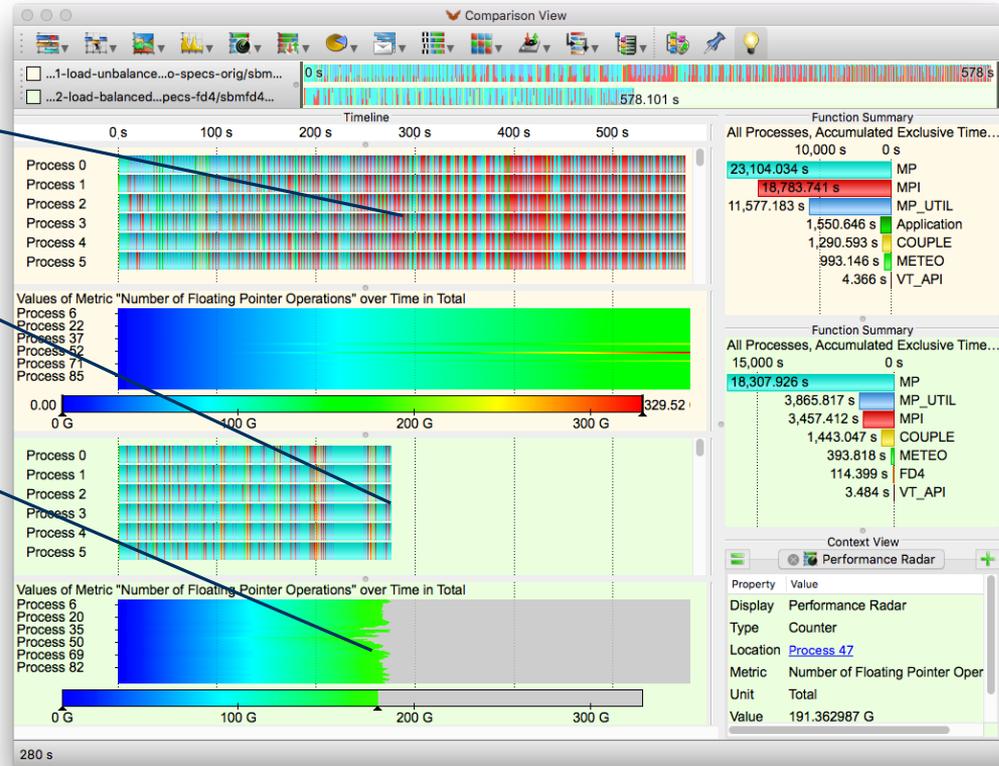


Load Balance Analysis of Weather Forecast Model Comparison

Before Load Optimization

After Load Optimization

Performance doubled!

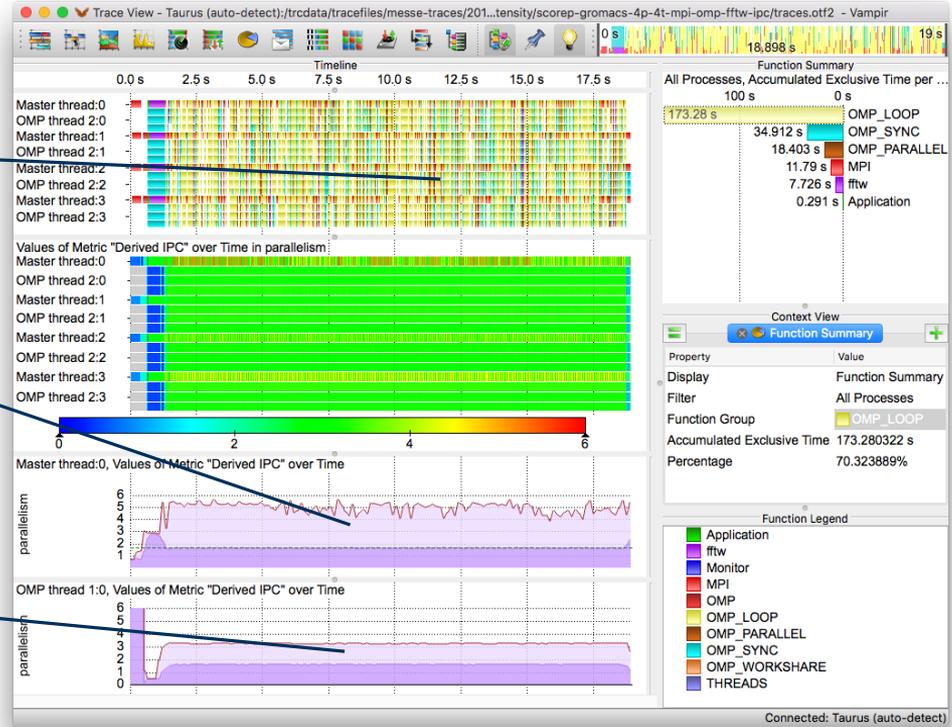


Analyzing Instructions per Cycle with Custom Metrics Analysis (1)

~12s time spent in OpenMP
(62% util.)

Master: Avg. IPC is 1.8,
Peak IPC is 5.5

OMP Threads: Avg.
IPC is 1.8, Peak IPC is 3.5

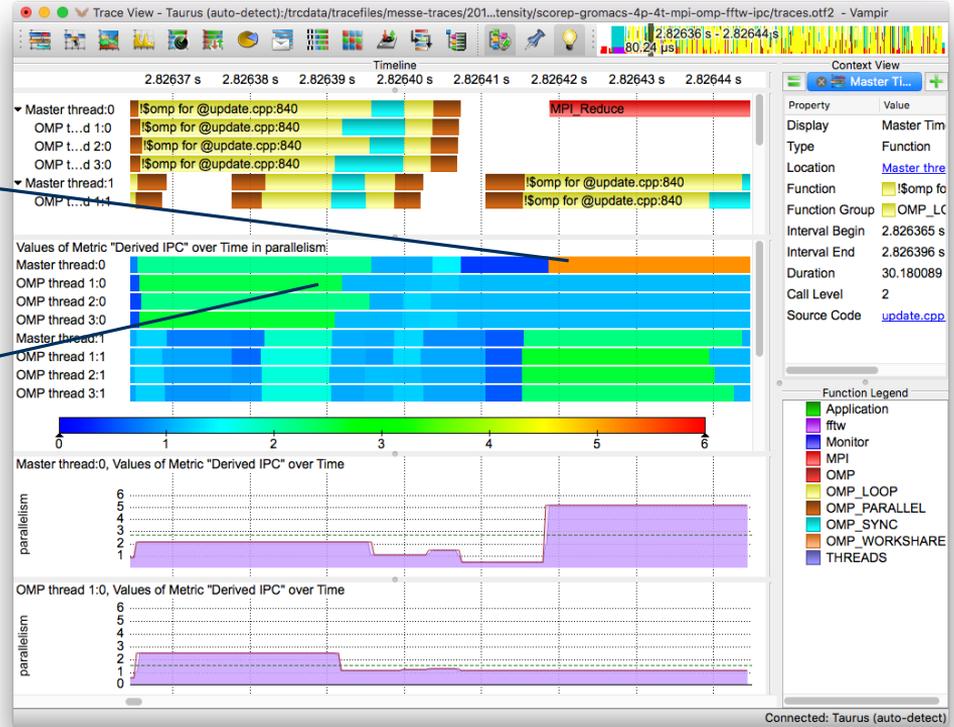


Analyzing Instructions per Cycle with Custom Metrics

Analysis (2)

MPI responsible for peak IPC. Irrelevant!

OpenMP loops have fair IPC of ~2.5



Analyzing Instructions per Cycle with Custom Metrics

Analysis (3)

Compute IPC for OpenMP regions only

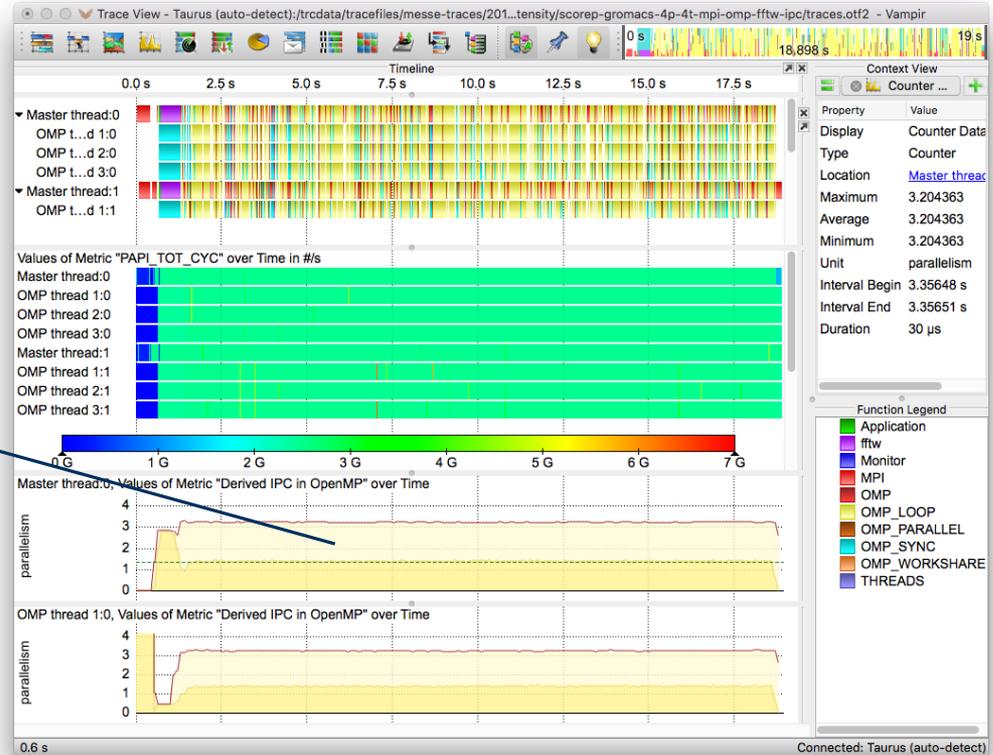
The screenshot shows the 'Custom Metrics' configuration window. The description is 'Derived IPC in OpenMP' and the unit is 'parallelism'. The workflow consists of the following steps:

- Metric 1:** Trace Counter, PAPI_TOT_INS, Increments per Second.
- Metric 2:** Trace Counter, PAPI_TOT_CYC, Increments per Second.
- Metric 3:** Function is Active, > parallel @vsite.cpp:644, Inclusive.
- Operation 1:** Divide (Metric 1 / Metric 2).
- Operation 2:** Multiply (Operation 1 * Metric 3).

An 'Include/Exclude All' dialog box is open, showing a list of OpenMP-related categories with checkboxes. The 'OMP_LOOP' and 'OMP_PARALLEL' categories are checked, indicating that the metric calculation is restricted to these regions.

Analyzing Instructions per Cycle with Custom Metrics Analysis (4)

In OpenMP: Avg. IPC is 1.4,
Peak IPC ~3.0

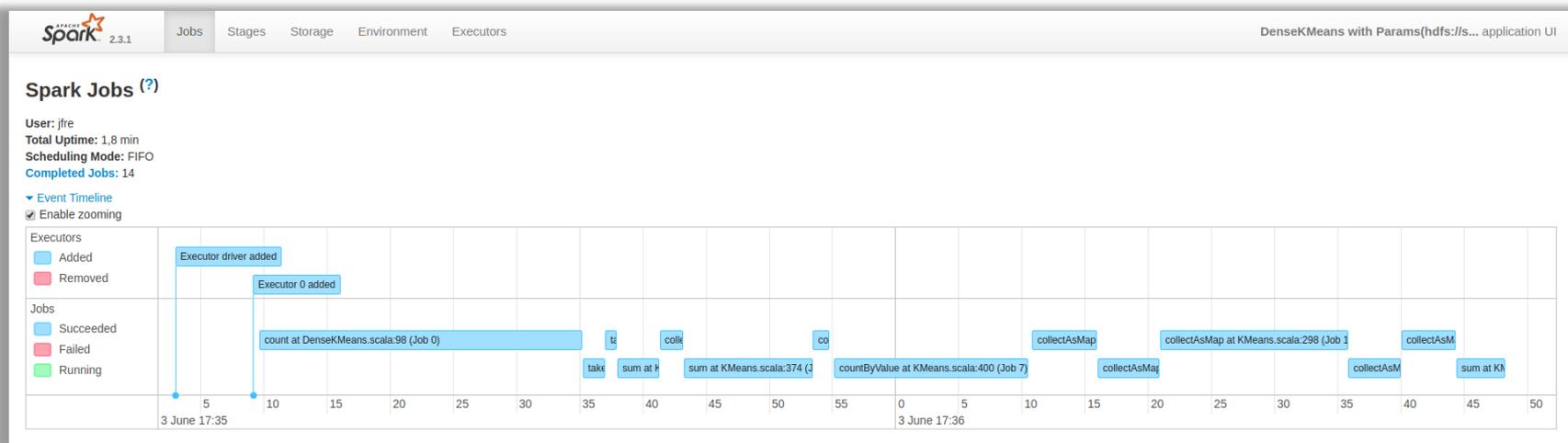


Jan Frenzel

Performance Analysis for Big-Data Frameworks

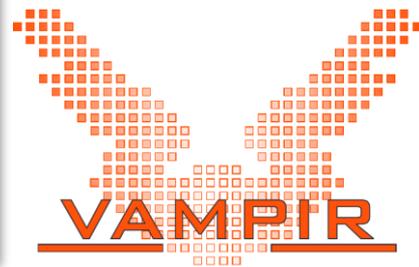
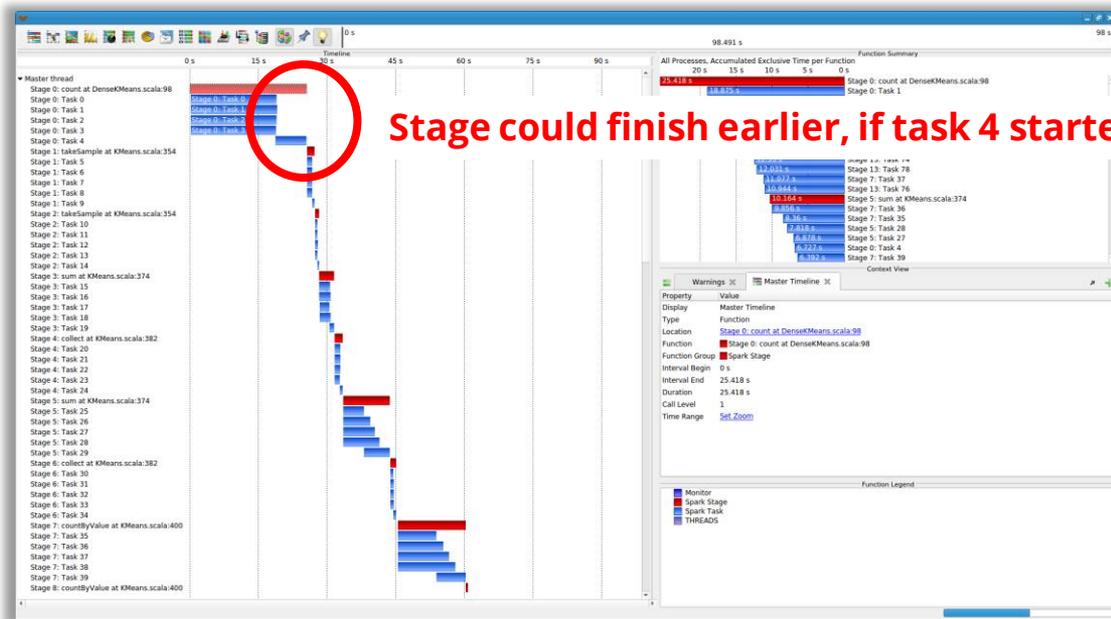
Information in the Spark Dashboard

- Easily usable (no extra tools needed)
- Overview over general information
- Limited usability for performance analysis (“Why and where is it slow?” remains unanswered)



Information of the Spark Dashboard in Vampir

- Overview over general information, easier to access information about tasks
- Limited usability for performance analysis (“Why and where is it slow?” only partially answered)



Further Topics

- Run-time monitoring with little overhead
- Support for upcoming hardware architectures (new CPU type --> small changes, introducing GPU computing --> major changes, close cooperation with vendors)
- Code instrumentation
- Sampling
- Profiling and Event Tracing data formats
- Automatic detection of typical performance issues

- ProPE: Continuous parallel performance monitoring, looking for abnormal performance behavior
- LO2S: Performance monitoring via sampling with Linux perf

Dr. René Jäkel

Supporting data intensive applications @HPC

Motivation

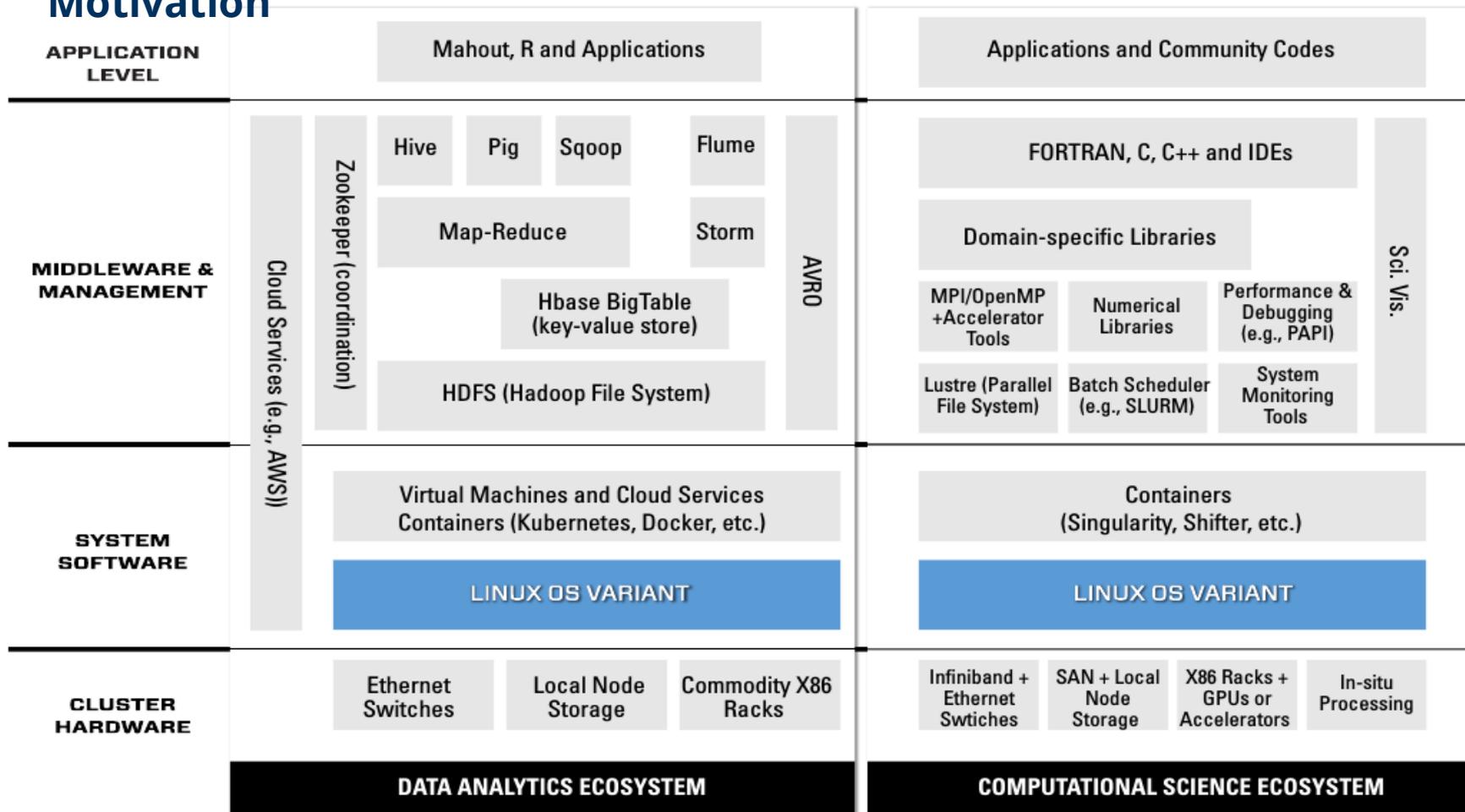
Domain Scientist

- Specifics of data/information: formats, content, error handling
- Combine theory-driven models with experimental data (e.g. simulation vs. exp.)
- Often knowledge not well formalized (“in the experts head”)
- Little or no HPC background

HPC Expertise

- Adoption of workloads to larger and more powerful infrastructures
- Optimization of workloads / parallel applications to infrastructure
- Support for use of hardware/software layers (parallel programming, filesystems, communication), but not on content
- Little or no domain knowledge

Motivation



Motivation



- Most important: bring experts from both sides together to investigate requirements of data-intensive applications and derive solutions
- Connect experts and application domain scientists

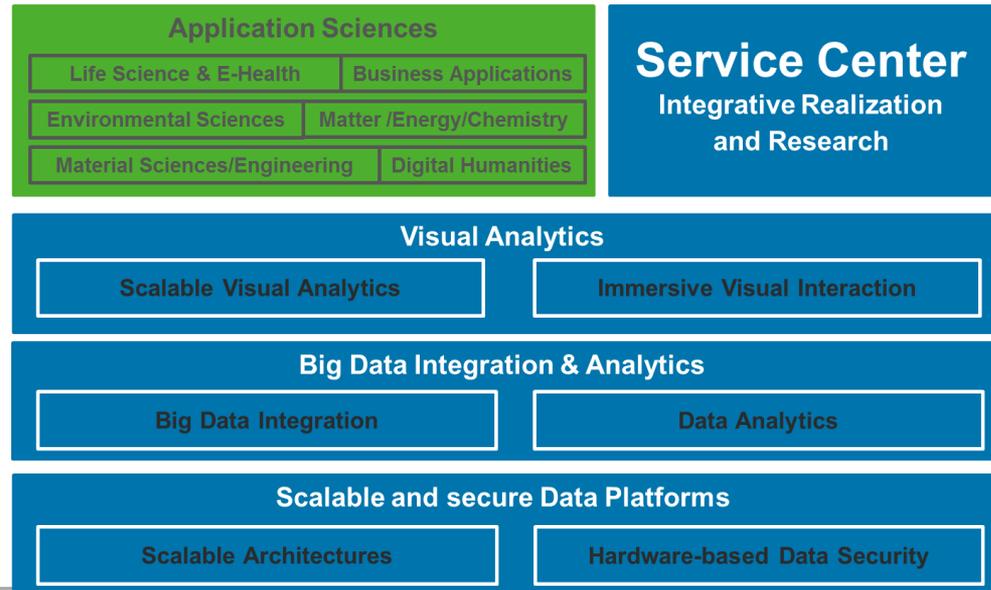


Motivation



ScaDS Dresden/Leipzig

- National Competence Center for collaborative Big Data driven research
- Established 2014 in Saxony: TU Dresden, Univ. Leipzig, MPI-CBG, IÖR, HZDR, UZF



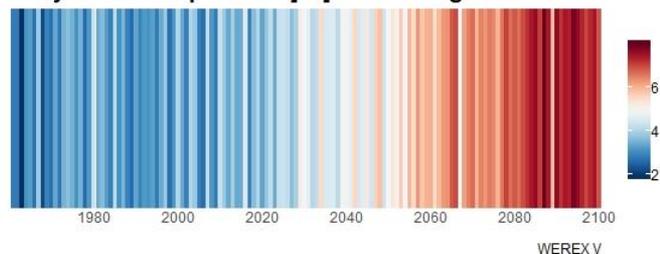
Flood risk analysis due to climate impact

Visualisation of Climate Data

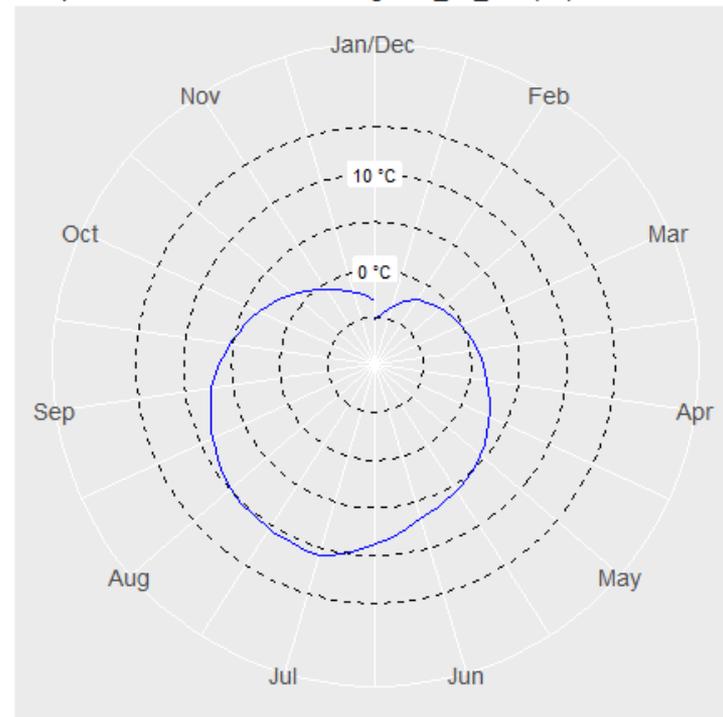
WEREX V Statistical Downscaling+ IPCC projection

- Rapid economic and population growth
- Quick spread of new & efficient technologies
- Convergent world - income and way of life converge between regions, extensive social and cultural interactions worldwide
- A balanced emphasis on all energy sources

Daily Mean Temperature [°C] Fichtelberg 1961-2100



Temperature at Station Fichtelberg EH5_L1_A1B(00) 1961-2100

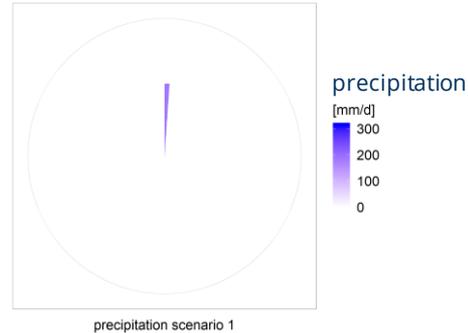


Spatial and temporal dynamics of flood risks

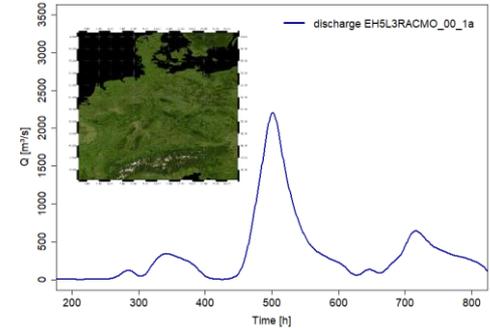
Estimation of flood risk due to climate change and the expected damages

- Climate scenarios
- Hydrologic modelling
- Hydrodynamic modelling
- Damage Modelling

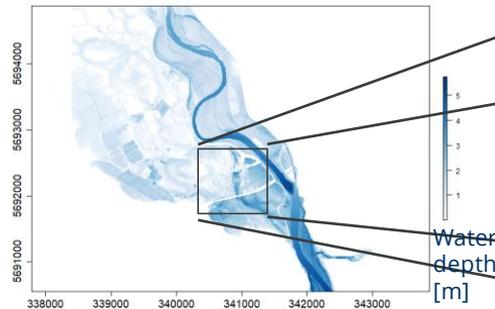
Climate Ensembles



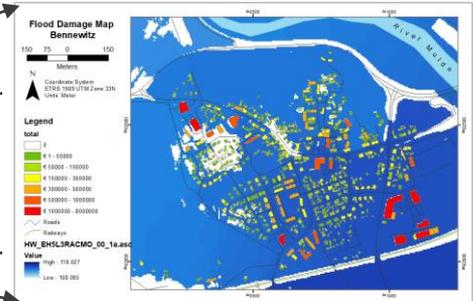
Hydrologic Modelling



Hydrodynamic Modelling



Damage Modelling



Genome assembly pipeline

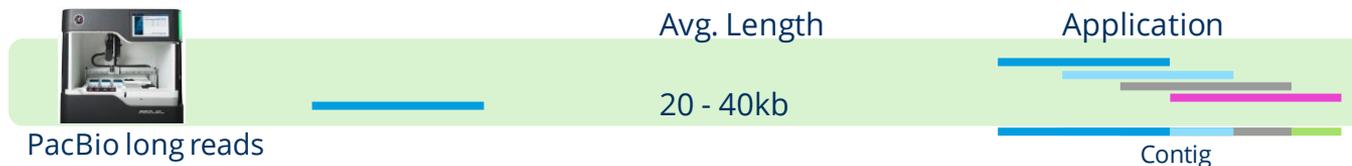
Platinum genome assemblies – The Bat1K project

- Goal: Catalog the unique genetic endowment and diversity present in all living bats
- understand the molecular basis of their unique adaptations
 - link genotype with phenotype
 - uncover their evolutionary history
 - better understand, promote, and conserve bats.

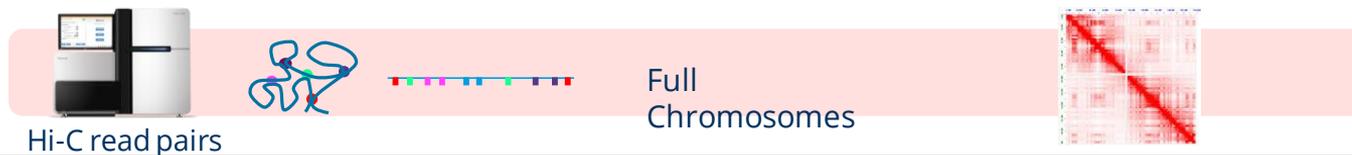
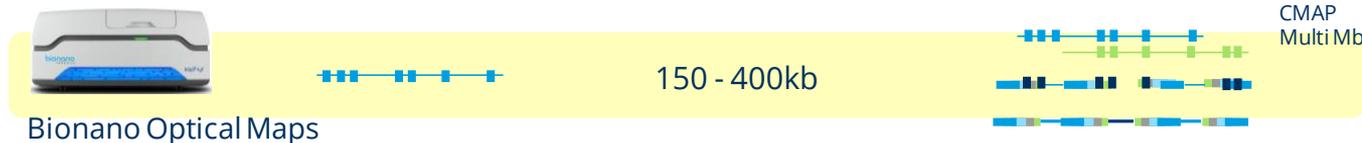


Multiple DNA sequencing technologies

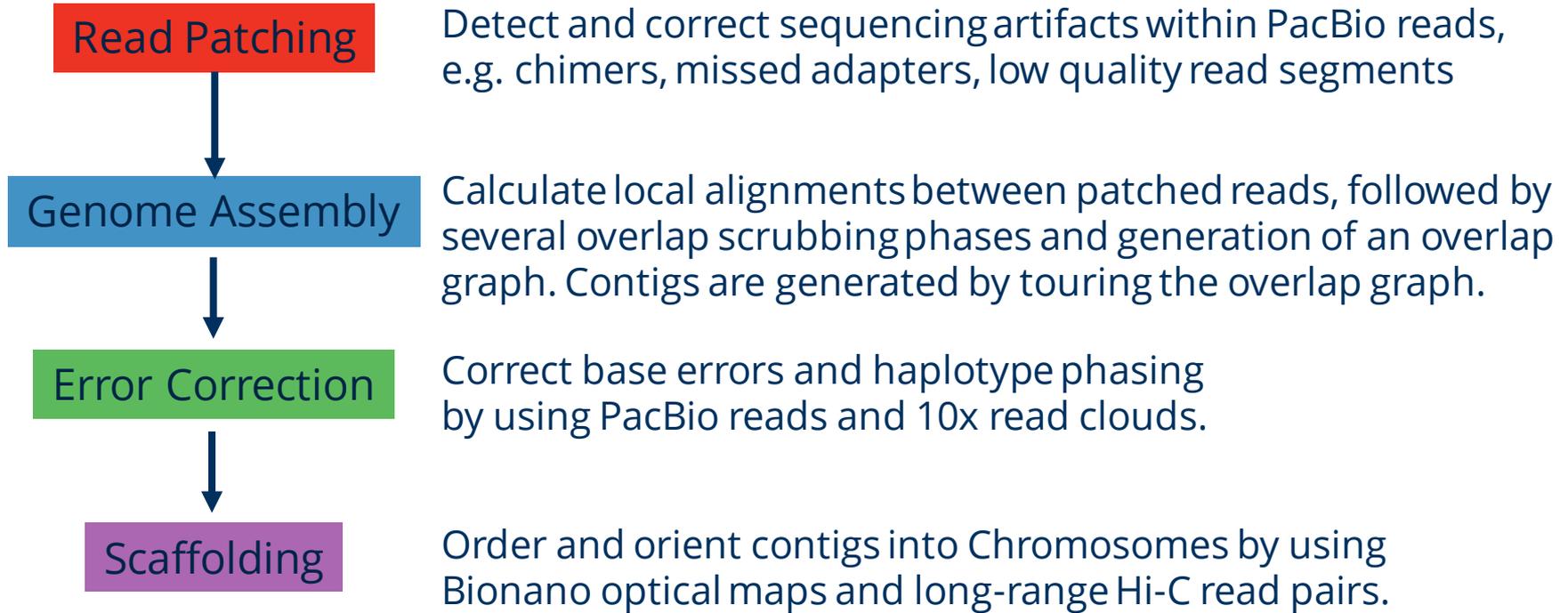
1. Genome Assembly based on noisy long reads



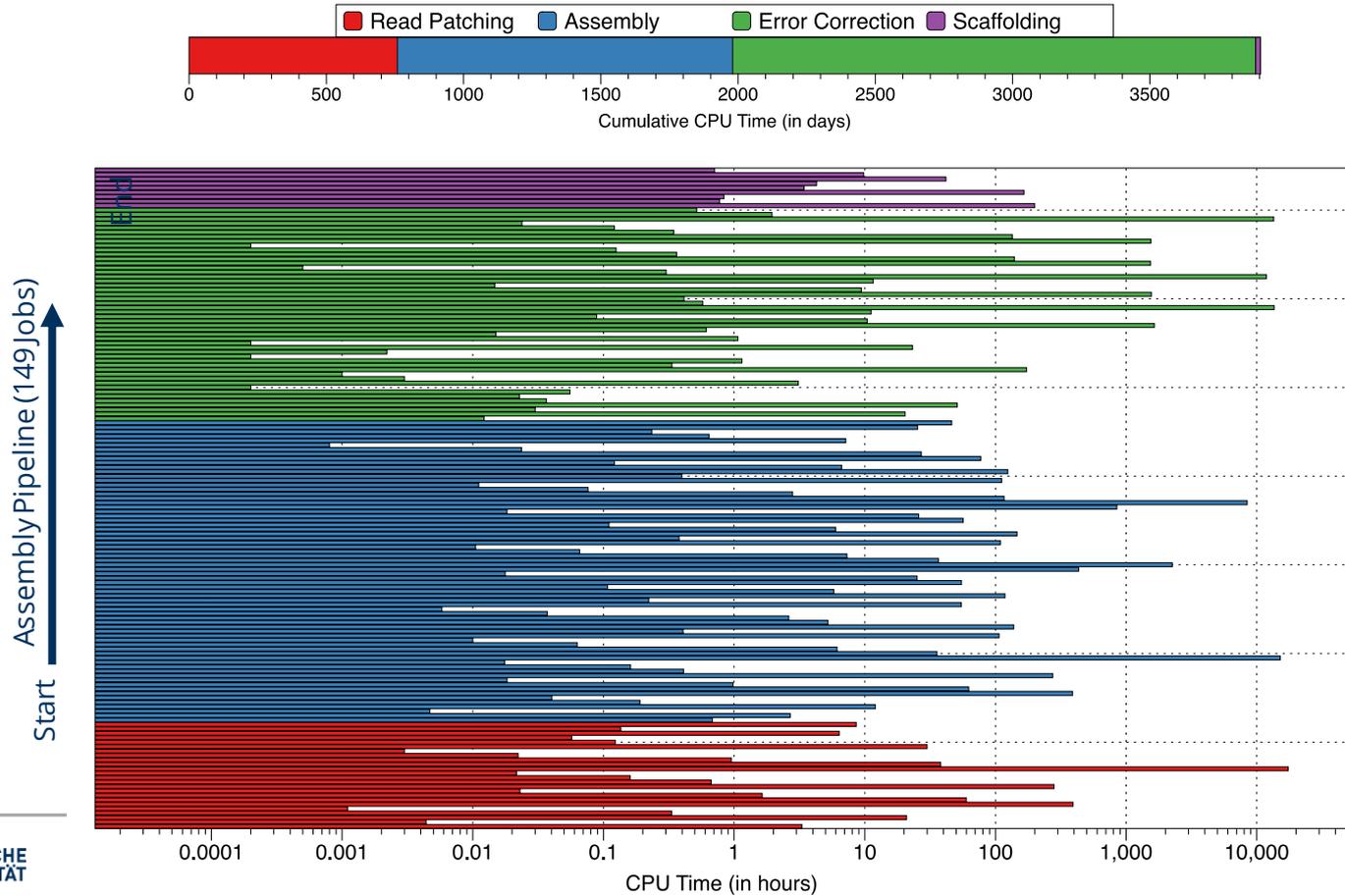
2. Scaffolding: order and orient contigs by using multiple sequencing technologies with increasing long-range information



Assembly pipeline



Assembly pipeline - runtime



Dr. Robert Schöne, Andreas Gocht

Energy Efficient HPC: READEX Project

Overview

- Tools aided methodology
- Automated energy efficiency tuning of parallel applications
- Dynamically adjust system parameters to actual resource requirements
- Co-design approach

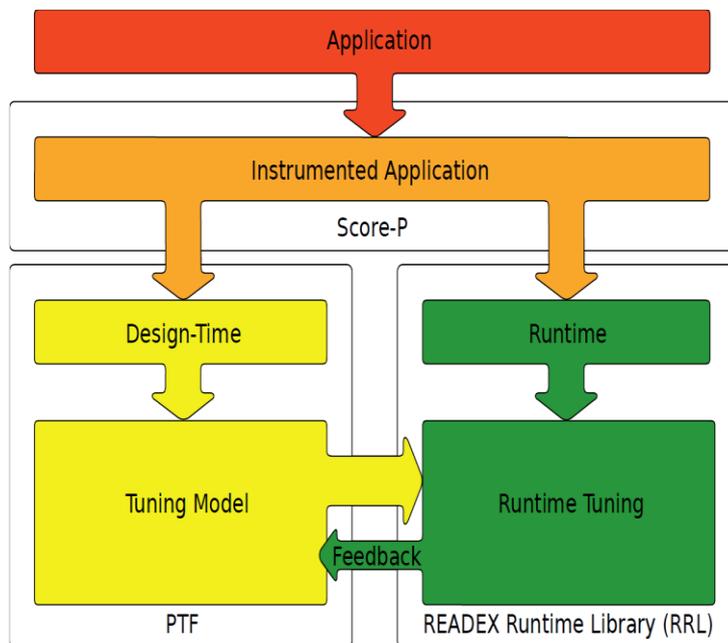


IT4Innovations
national01\$#&0
supercomputing
center@#01%101

Overview

Design-Time Analysis

- Periscope Tuning Framework
- Identify significant regions and runtime situations
- Test tuning parameters
- Detect optimal configuration
- Write tuning model



Runtime Tuning

- READEX Runtime Library
- Uses tuning model from DTA
- Detect runtime situations
- Adjust tuning parameters for energy efficiency
- Calibration of tuning model

Hardware Analysis and Control

- Background in computer and processor architecture
- Determine low-level side effects
- Determine influence of power saving mechanisms
- Provide access to low level hardware parameters

READEX

- Instrumentation
- Uses low-level power saving mechanisms

Linux Kernel Module

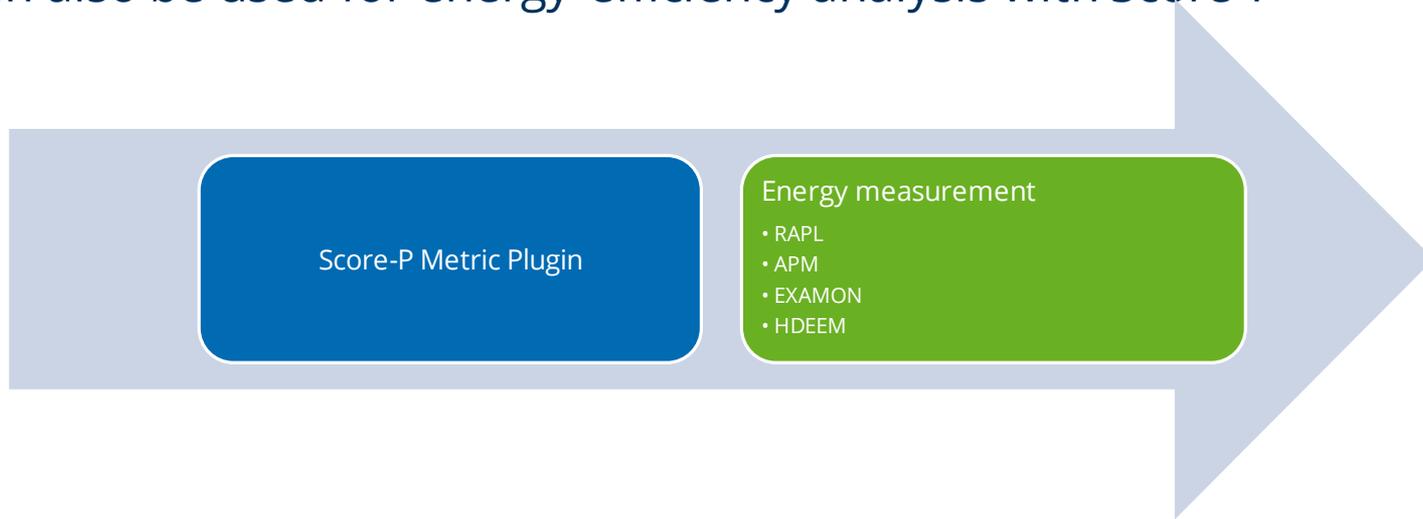
- Provide access to hardware
- Support for different processors

Processor

- Performance and Power knobs
- Energy counter

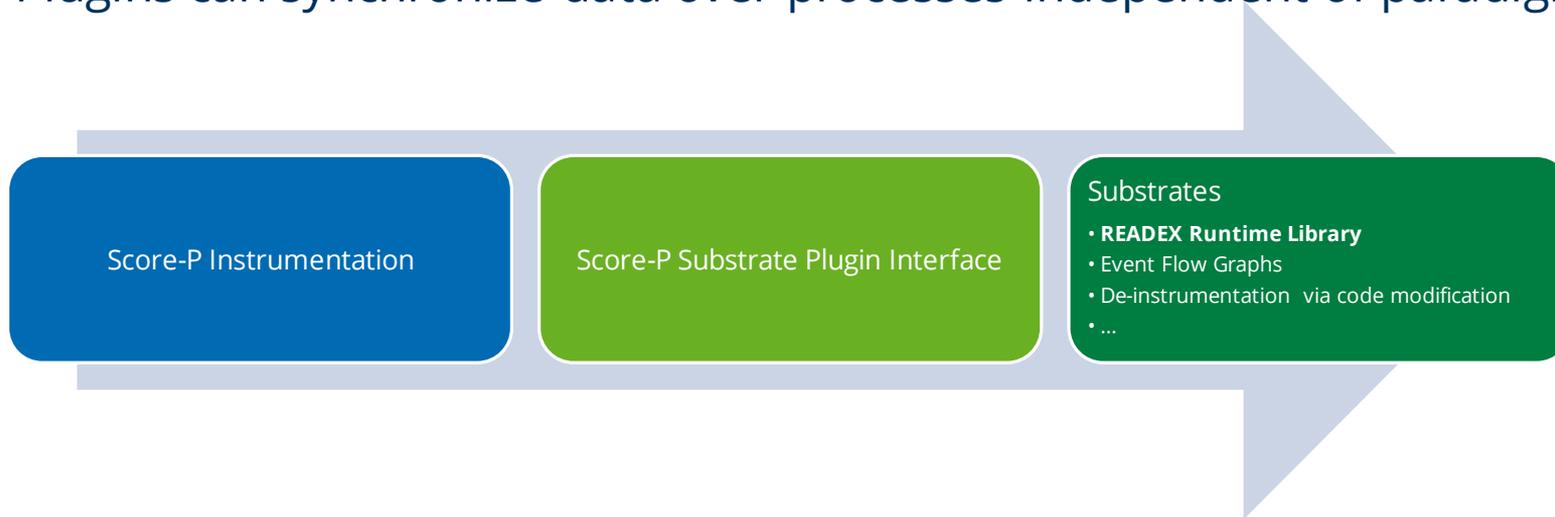
Energy Measurement

- Various back-ends for measuring energy
- Used by READEX during Design-Time or Runtime (online-tuning)
- Can also be used for energy-efficiency analysis with Score-P



Score-P Substrate Plugins

- Use Score-P instrumentation for other purposes (e.g., tuning)
- Consume instrumentation events and metrics
- Plugins can synchronize data over processes independent of paradigm

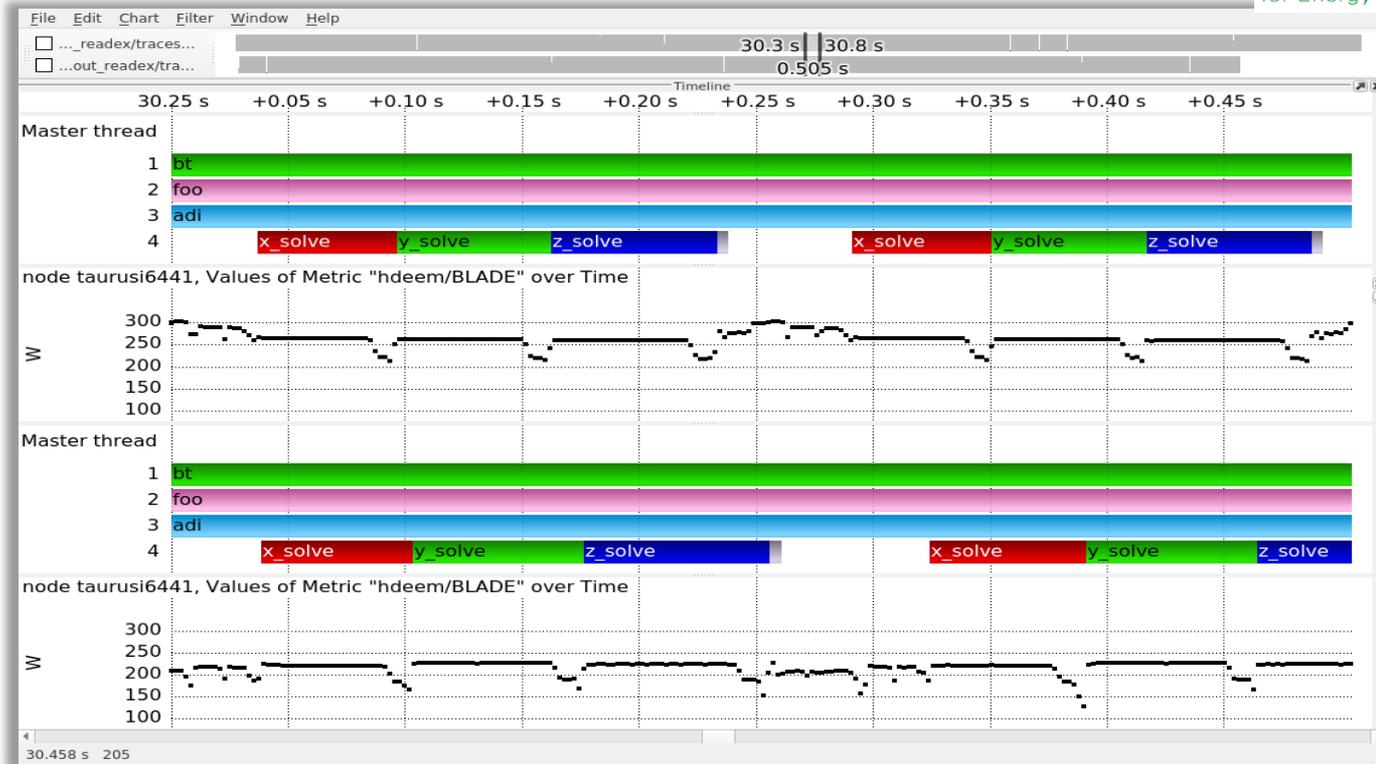


READEX Runtime Library

- Consumes Score-P Events
- Applies configuration changes during Design-Time and Runtime
- On-line tuning for unknown regions or standalone application without Design-Time

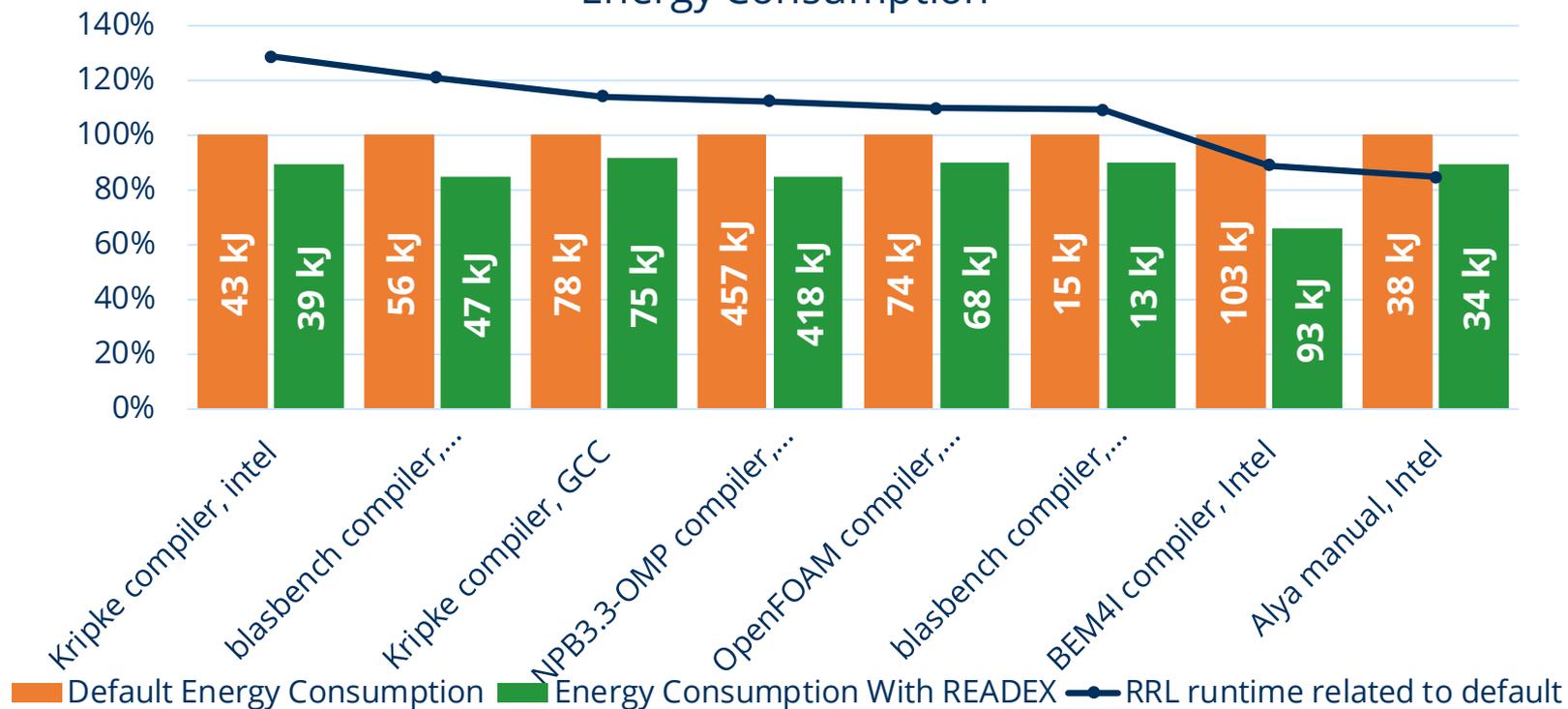


READEX Results



Power consumption of a untuned(top) and tuned(bottom) NAS BT.C benchmark run

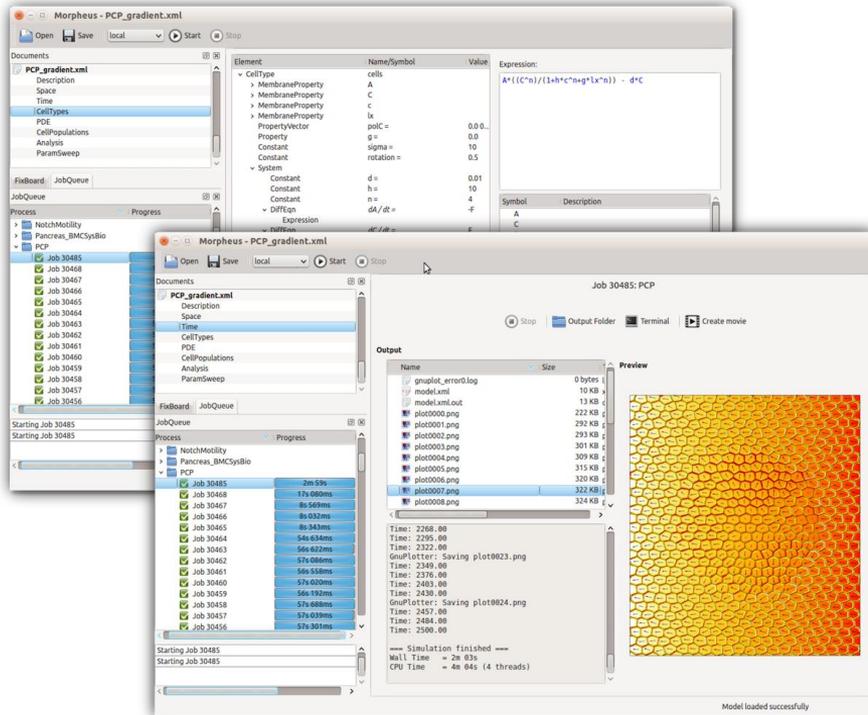
Energy Consumption



Dr. Lutz Brusch
Innovative Methods of Computing, ZIH

Development of Open Source Software to Enable Systems Biology and Systems Medicine

Open Source Software Morpheus - GUI-based Simulator for 3D Computational Biology



Homepage with installers for
Linux, Windows, Mac OSX:
<https://morpheus.gitlab.io>



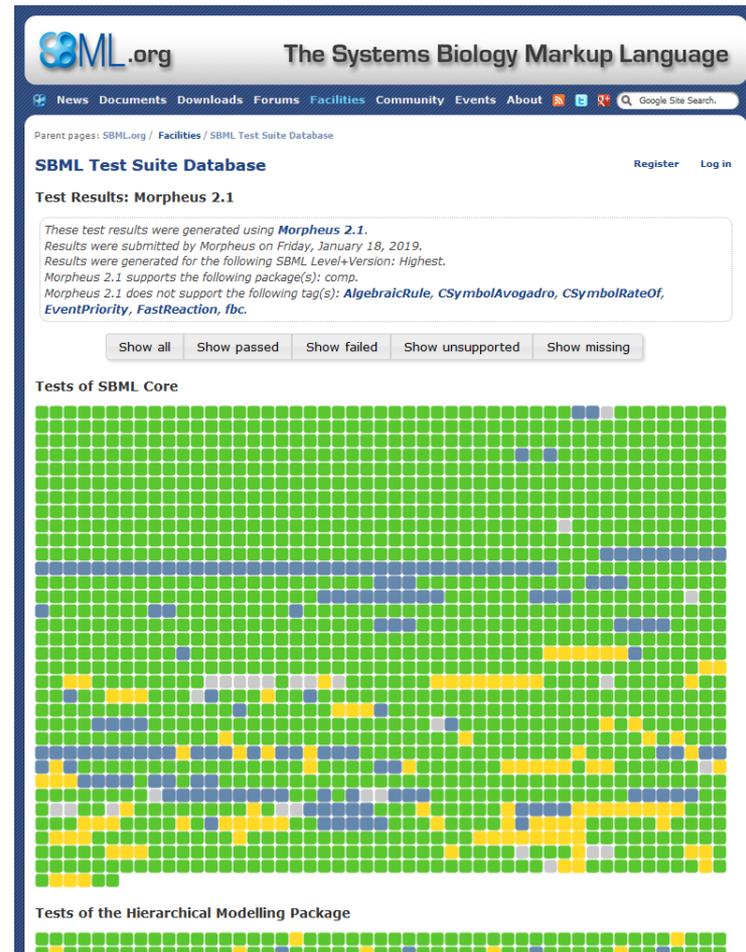
@morpheus_lab

Contact: Dr. L. Bruschi

Enabling Collaborative Research

MultiCellML – Standard Model Language

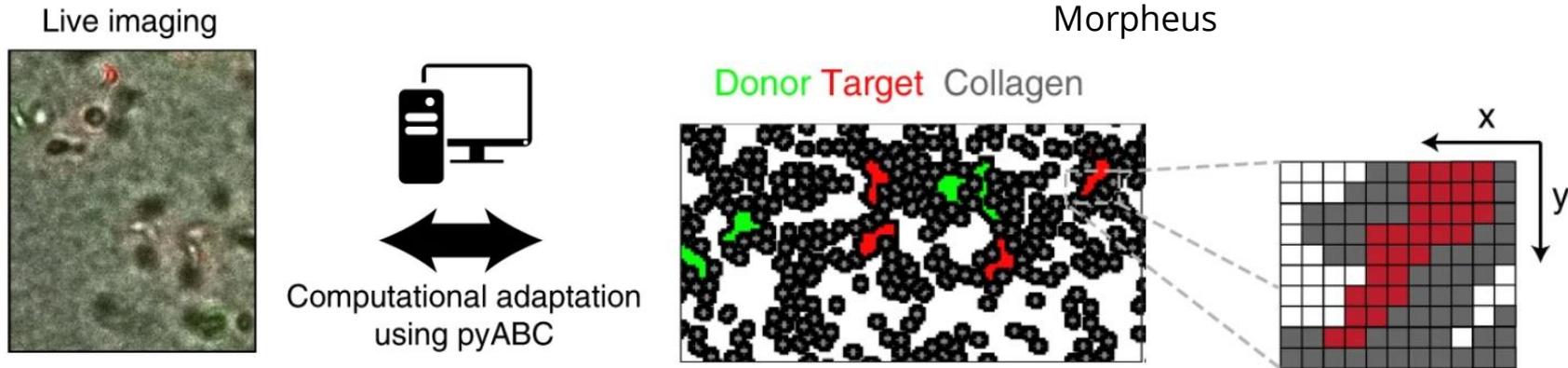
- Goal: Enable collaborative research and reproducibility of simulation studies through MultiCellML and model repository
- Build on success of SBML model language for biological networks
- Morpheus is 1 of world-wide 12 certified SBML simulators (Fig.: Morpheus' SBML-compliance certificate)
- Morpheus is world-wide the only of these simulators that also runs spatial tissue models
- Morpheus is world-wide the only simulator of spatial multicellular models that is able to completely define such models in a declarative language (XML) instead of execution code
- Project MultiCellML: Generalise Morpheus' solution such that different simulators can exchange models in the MultiCellML language
- Contact: Dr. L. Brusch



The screenshot shows the SBML Test Suite Database website. At the top, the SBML.org logo and the text "The Systems Biology Markup Language" are visible. Below the navigation bar, the page title is "SBML Test Suite Database". The main content area displays "Test Results: Morpheus 2.1". A text box contains the following information: "These test results were generated using **Morpheus 2.1**. Results were submitted by Morpheus on Friday, January 18, 2019. Results were generated for the following SBML Level+Version: Highest. Morpheus 2.1 supports the following package(s): comp. Morpheus 2.1 does not support the following tag(s): AlgebraicRule, CSymbolAvogadro, CSymbolRateOf, EventPriority, FastReaction, fbc." Below this text are five buttons: "Show all", "Show passed", "Show failed", "Show unsupported", and "Show missing". The "Show passed" button is highlighted. Below the buttons is a section titled "Tests of SBML Core" which contains a large grid of colored squares representing test results. The grid is mostly green, indicating passed tests, with some blue and yellow squares indicating failed or unsupported tests. Below the grid is a section titled "Tests of the Hierarchical Modelling Package" which also contains a grid of colored squares.

Data to Models

FitMultiCell – Parameter Estimation from Microscopy Data



- Goal: Enable parameter estimation for stochastic biological models
- data-intensive and compute-intensive
- Approximate Bayesian Computation in open source pyABC framework addresses data challenge,
- Parallel implementation of Morpheus addresses compute challenge
- Contact: Dr. L. Brusch

Gefördert vom



Bundesministerium
für Bildung
und Forschung

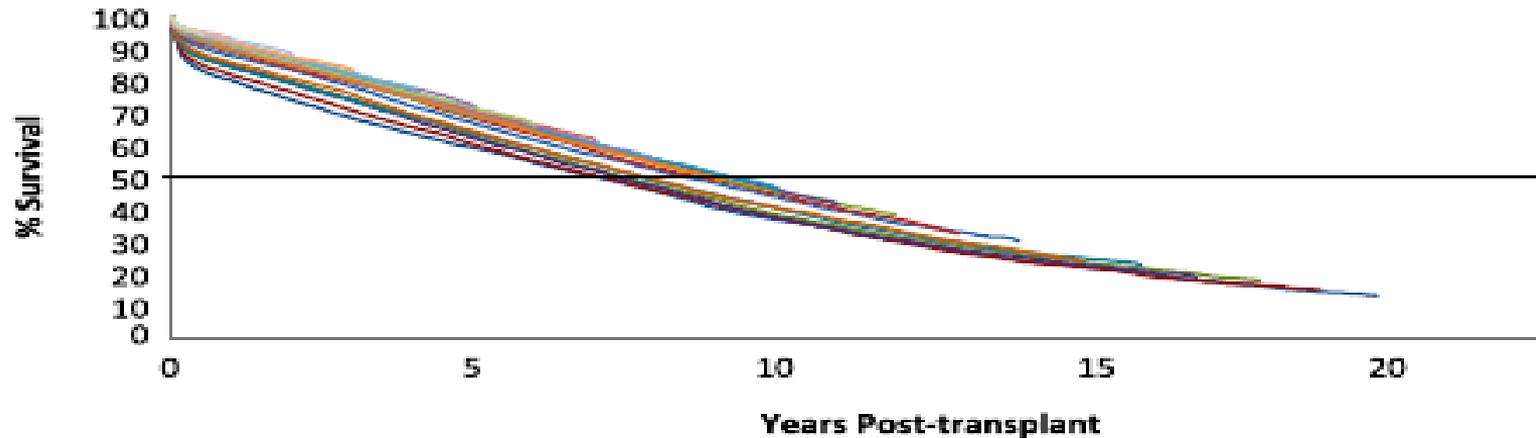
Prof. Dr. Andreas Deutsch

Reclassification using OmiCs integration in KidnEy Transplantation



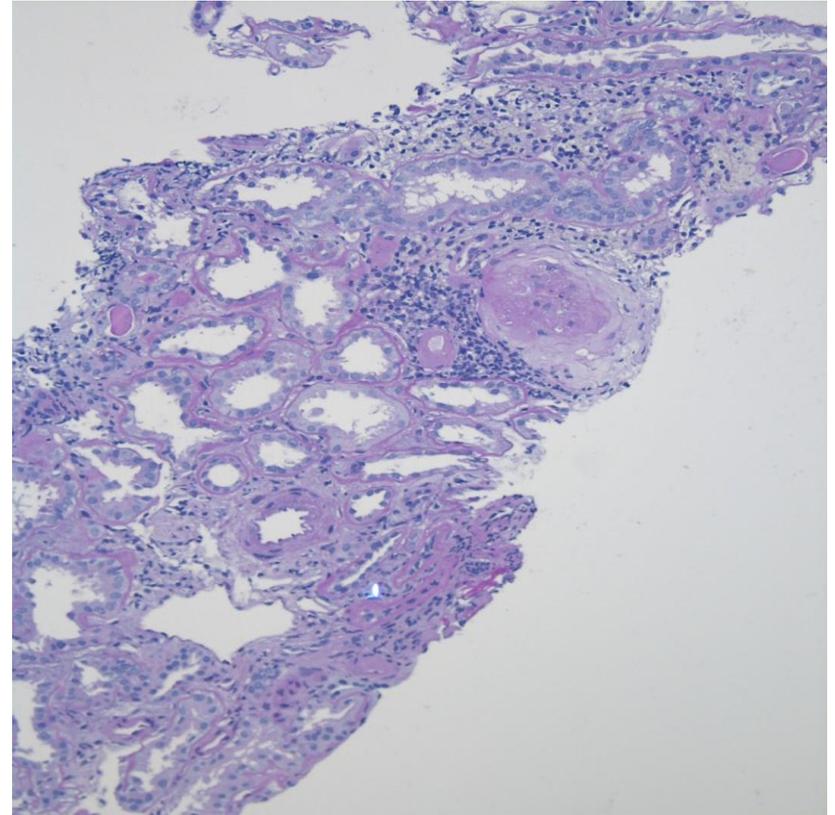
End stage renal failure

- 3.2 million patients with end stage renal failure world-wide
- Life-saving therapy: dialysis or kidney transplantation
- Dialysis: high morbidity and reduced life expectancy, low quality of live, high costs
- Transplantation: best therapy, but limited access and limited graft survival



Graft Biopsy as Gold Standard for Diagnosis

- “One-time stop” to diagnose all pathologies
- No continuous monitoring
- Invasive and costly
- Poor inter-observer concordance
- Despite an elaborate classification system: Very often diagnostic vagueness



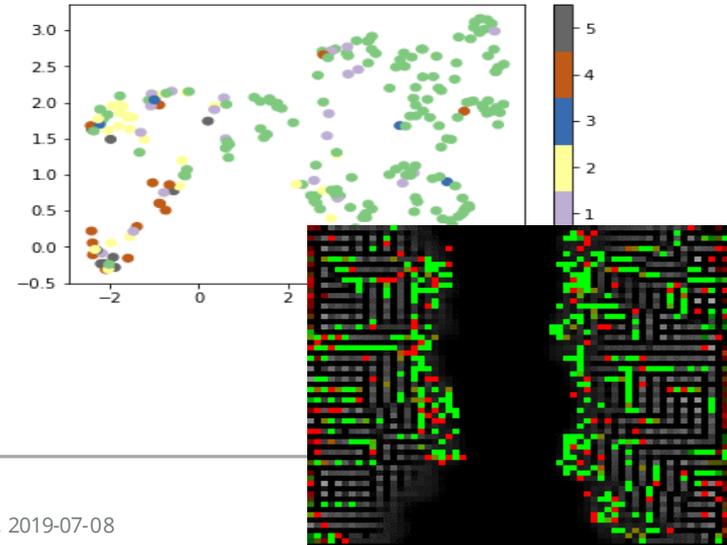
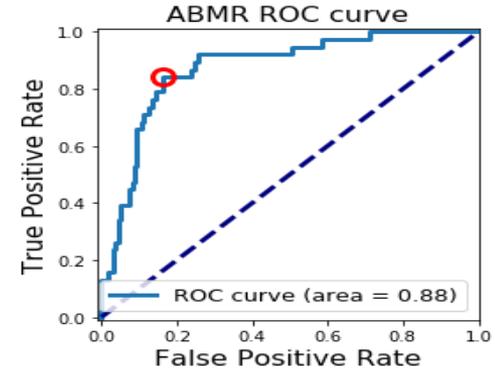
Machine learning and modeling approach for diagnosis and prognosis of kidney transplantation

Combined machine learning and modeling approach:

Gene expression data from rejection and no-rejection biopsies + machine learning = genetic rejection criteria

Reconstruction of hypothetical biological rejection mechanisms underlying two rejection phenotypes, based on machine learning-based dimensionality reduction

In silico testing of proposed mechanisms with non-spatial ODE modeling and spatial LGCA modeling and comparison with biopsy imaging data



Parallel Programming Abstractions with C++

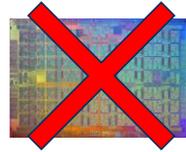
Spectrum of Computing Architectures



Supercomputer
100 000+ cores



Cluster
1000s of cores



Manycore



Server
10s of cores



Notebook
2-4 cores



Mobile
2-8 cores

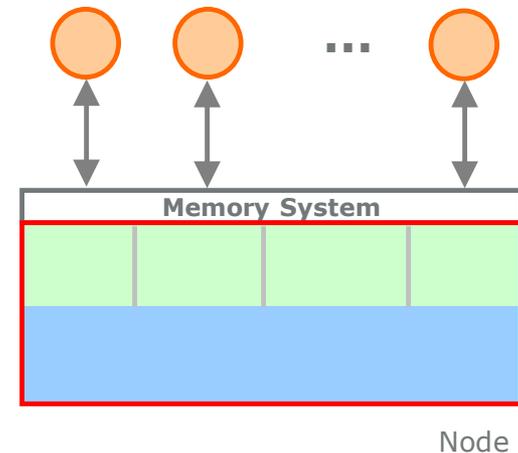
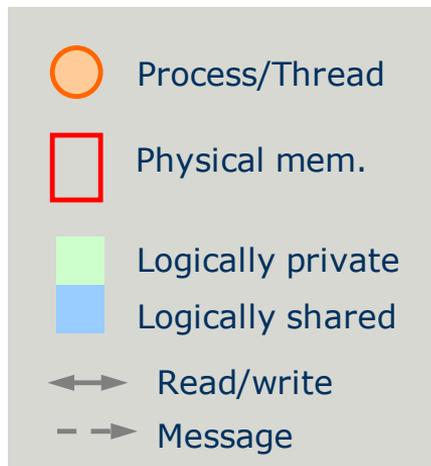
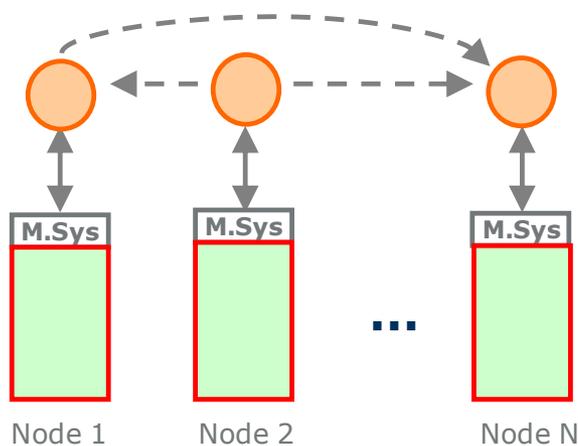
Distributed Memory (DM) 

DM programming: MPI,
Charm++, ...

 Shared Memory (SM)

SM programming: OpenMP, Pthreads,
Cilk, TBB, ...

Shared Memory vs. Distributed Memory Programming



Message Passing

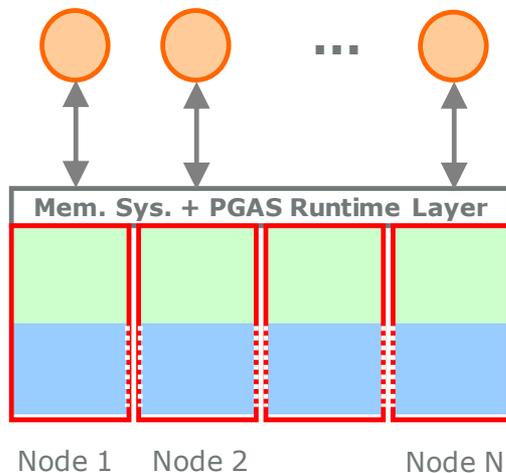
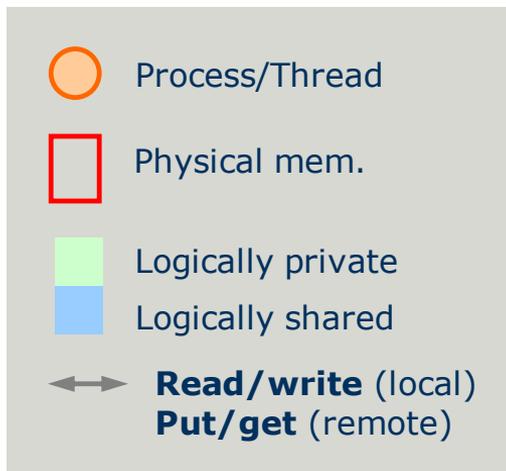
- + Performance, runs everywhere
- Productivity

Threading

- + Productivity
- Locality control, limited to SM hardware

PGAS – Combining the Advantages of both Approaches

PGAS: Partitioned Global Address Space



PGAS Languages

Chapel, CoArray
Fortran, UPC, ...

PGAS Libraries

Global Arrays (GA),
GASPI, OpenShmem,
MPI3.0 RMA



Locality control, runs everywhere,
performance and productivity

So you'd like to write parallel HPC codes in C++?

HPC programming today

- Large scale parallelism
- Heterogeneous architectures
- Hybrid parallelism --> multiple sources of complexity
- MPI+X as dominating parallel programming model
- Node-level model X strictly needed for portability and performance portability
- What if you bet on the wrong one?

MPI disregards C++

- Data distribution, data transfers, and synchronization deeply entangled
- The MPI C++ bindings deprecated in MPI 2.2 and removed in MPI 3.0*
- In C++ MPI codes you actually need to use MPI's C API
- C++ concepts like STL containers, iterators, and even basic data types are incompatible with MPI!

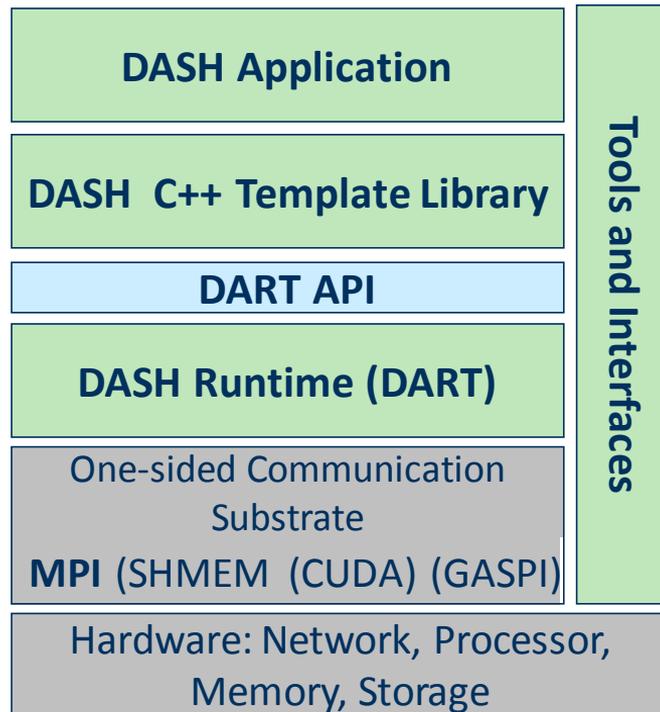
*<http://blogs.cisco.com/performance/the-mpi-c-bindings-what-happened-and-why>

DASH C++ Template Library for Parallel Programming

- C++ template library for application programmers
- Distributed data container classes
- Similar to the C++ STL container classes, compatible
- Built-in knowledge about distribution
- Algorithms similar to STL on distributed containers



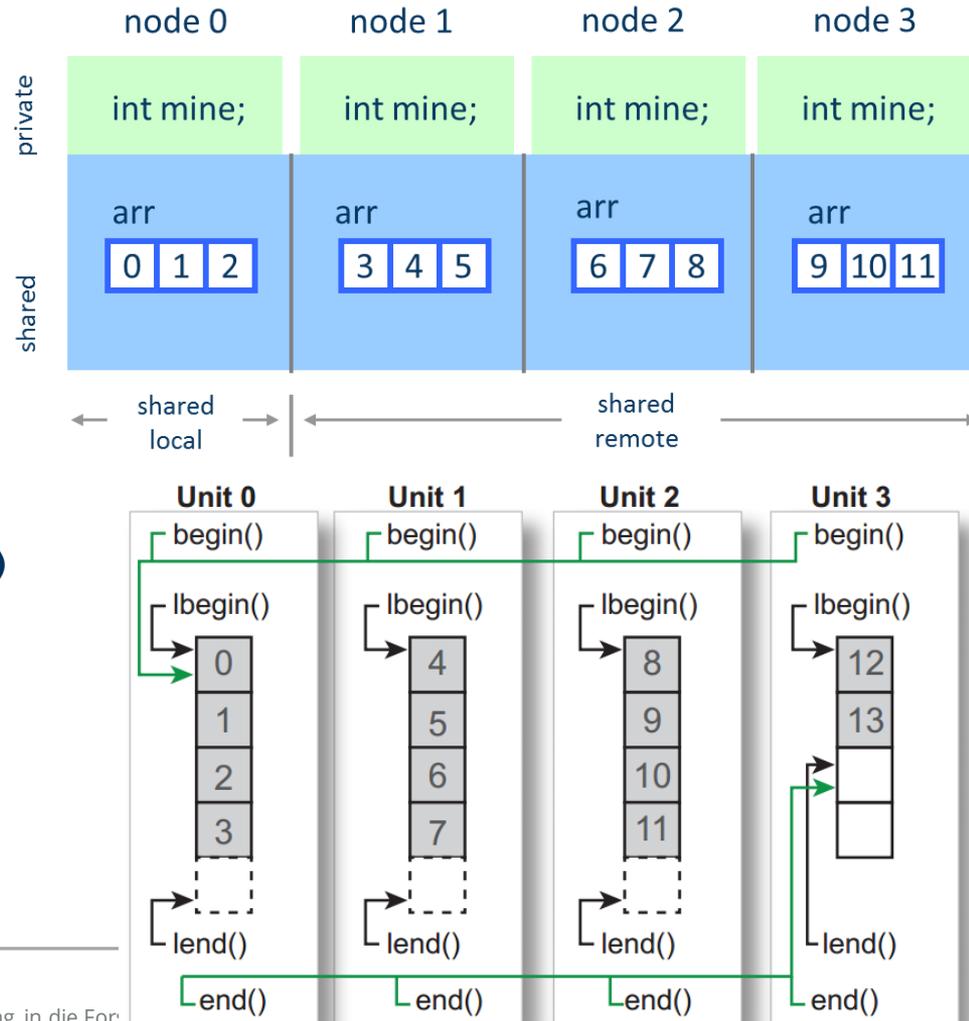
dash
www.dash-project.org
Distributed Data Structures
and Parallel Algorithms



DASH Array

DASH n-dimensional array

- Global random access with `begin()`, `end()` and `[]` ... via slow element-wise get
- Dedicated local access with `myarray.local.begin()` / `.end()` and `.local[]` ... direct and fast
- Configurable data distribution patterns in n dimensions
- STL-like algorithms considering actual data distribution patterns



Dr. Ralph Müller-Pfefferkorn

Data Intensive Computing and Research Data Management

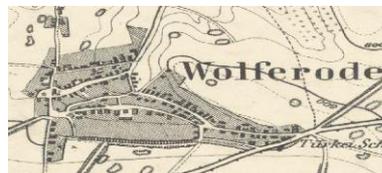
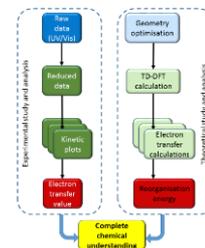
Research Data Management at ZIH: Generic Services

OpARA – Open Access Repository and Archive

- Research Data Repository for TUD and TUBA Freiberg, HTW Dresden to follow
- Open for all Researchers
- Funding by SMWK

MASI – Meta data management for Applied Sciences

- Scalable research data management
- Funded by DFG
- Use Cases from Chemistry, Environmental Sciences, and Humanities

A screenshot of a web-based data entry form titled 'Daten einreichen'. The form includes sections for 'Einarbeitung beschreiben', 'Einrichtung beschreiben', 'Titel des Datensatzes', 'Eigenschaften', 'Weitere Schlagwörter', 'Hintergrund', 'Informationsziele', 'Informationen zur Seite', and 'Angabe der Fachgebiete'. The right sidebar contains a search bar and a list of navigation options.

Research Data Management at ZIH: Together with Application Scientists

IT infrastructure project in SFB 940 (Neuroscience)

- Virtual research environments with focus on data management
- Image data (fMRT), EEG ...

IT infrastructure project in TRR 205 (Medicin)

- Virtual research environments with focus on data management
- Experiment data from tumor studies ...

EMuDIG 4.0 – IoT for heavy industry forging processes

- Sensor data in production environments, statistical analysis and machine learning
- Cloud environment for data analytics



EMuDig4.0

Kontaktstelle Forschungsdaten für die TU Dresden

Joint initiative by SLUB, ZIH, IGeweM (Institut für Geistiges Eigentum, Wettbewerbs- und Medienrecht), and TUD-CERT

Counseling and support for researchers:

- Organisation of RDM, data management concepts
- Metadata
- Tools
- Archiving
- Data publication
- Legal matters, ...



Contact

eMail: Kontaktstelle-Forschungsdaten@tu-dresden.de

<https://www.slub-dresden.de/en/service/knowledgebar/thema/gbList/34/>

<https://tu-dresden.de/forschung/services-fuer-forschende/kontaktstelle-forschungsdaten>

Further Topics

- NextGenIO: Programing for non-volatile memory (NV-DIMMs)
- HP-DLF: Highly Parallel Deep Learning Training
- HDEEM: High Definition Energy Efficiency Monitoring
- GCoE: Dresden GPU Center of Excellence
- IPCC: Intel Parallel Computing Center

Backup: Further Selected Results from Computer Science and Computational Science

Courtesy of Dr. Michael Bussmann et.al.

Novel Particle Accelerators and Highly Scalable GPU-based Simulation

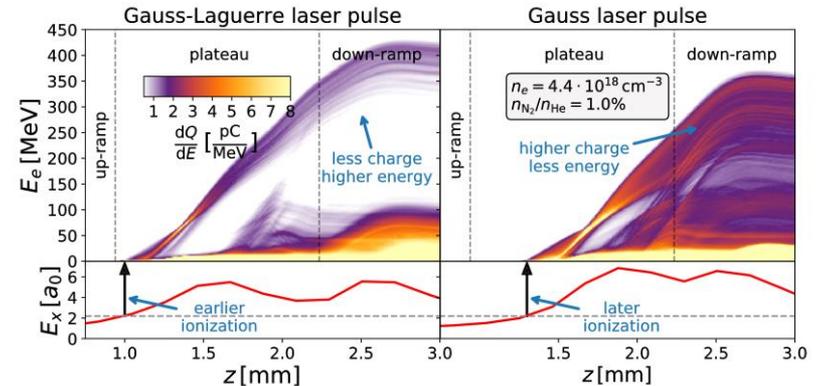
Laser Wakefield Acceleration Experiments at HZDR

J.P. Couperus et al.: Demonstration of a beam loaded nanocoulomb-class laser wakefield accelerator. **Nature Communications** 8.1 (2017)
A. Irman et al.: Improved performance of laser wakefield acceleration by tailored self-truncated ionization injection. **LPAW proceeding** pp.1-13 (2017)

- Advanced method of electron acceleration
- Based on highly non-linear laser plasma interaction
- Requires large scale particle-in-cell simulation for modeling
- Hundreds of simulations on up to 146 K80 GPUs performed with PIconGPU at Taurus/ZIH



Simulations accompanying experiment at HZDR



Studying the influence of higher order laser modes

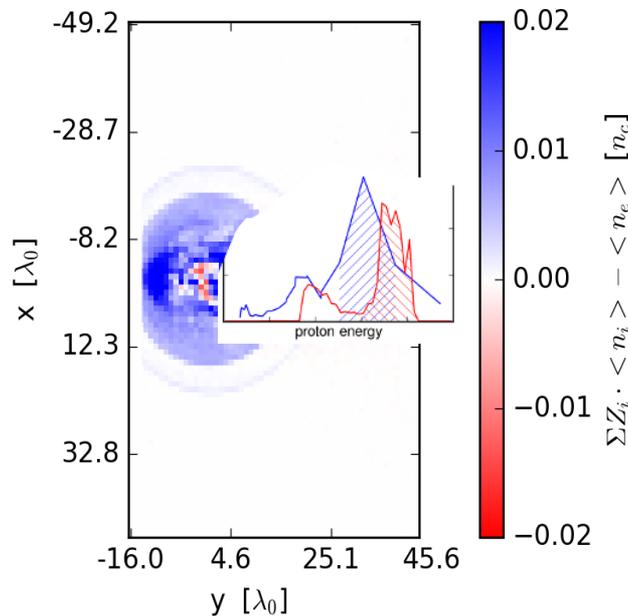
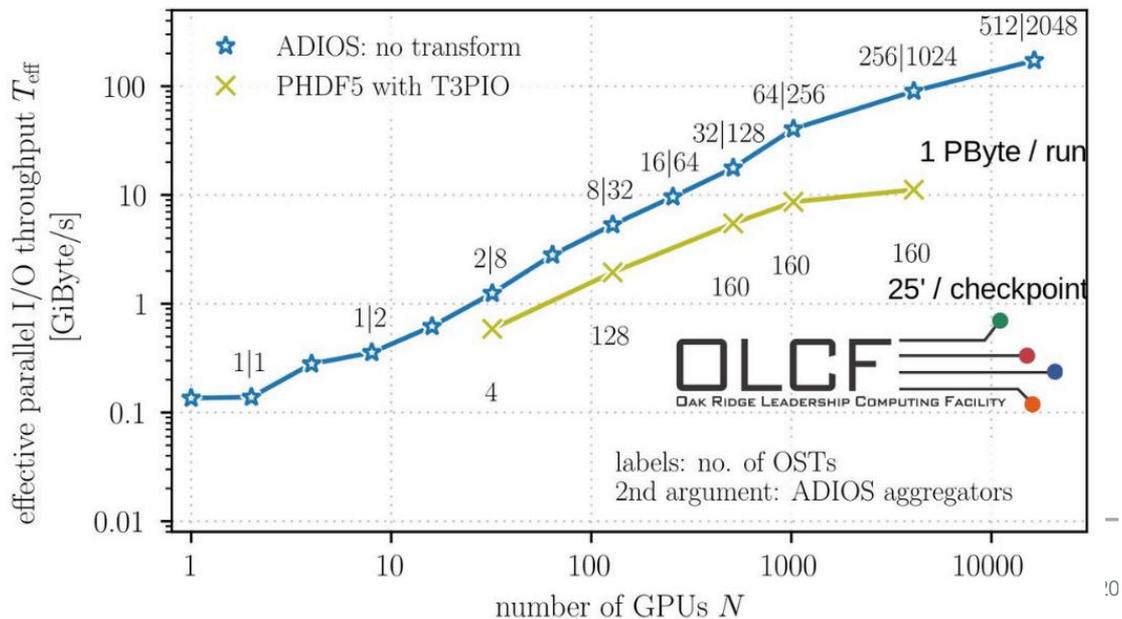


[Video](#)

Laser-Ion Acceleration with Mass-Limited Targets

P. Hilz, T.M. Ostermayr, A. Huebl et al.: Isolated proton bunch acceleration by a petawatt laser pulse. **Nature Communications** 9.423 (2018)
 A. Huebl et al.: On the Scalability of Data Reduction Techniques in Current and Upcoming HPC Systems from an Application Perspective, **ISC'17**, LNCS 10524

- 3D simulation of novel, fully isolated target for laser-ion acceleration
- 15 M CPUhrs (½ MGPUhrs), INCITE Award Highlight
- PByte-Scale I/O through **ADIOS** at Titan/OLCF



Performance Evaluation of GPU Applications

Kepler K40 | Pascal P100 | Volta V100



GPU

CENTER OF
EXCELLENCE

Bandwidth Improvements

Experiments done with gpumembench

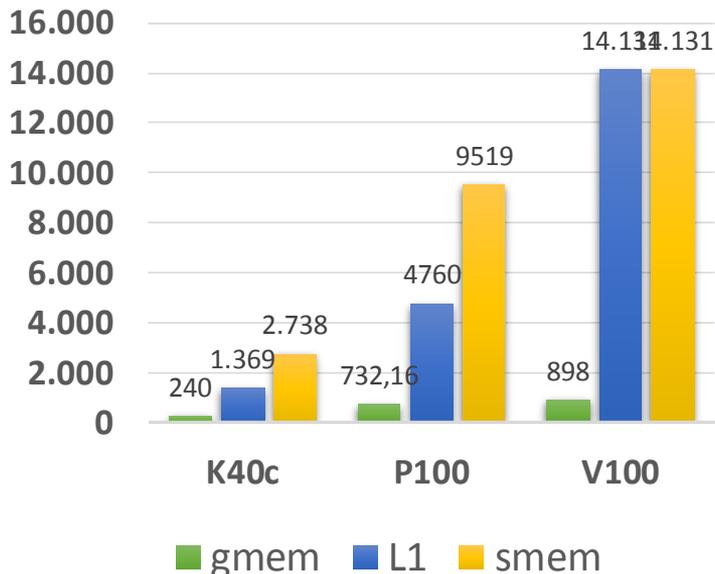


GPU

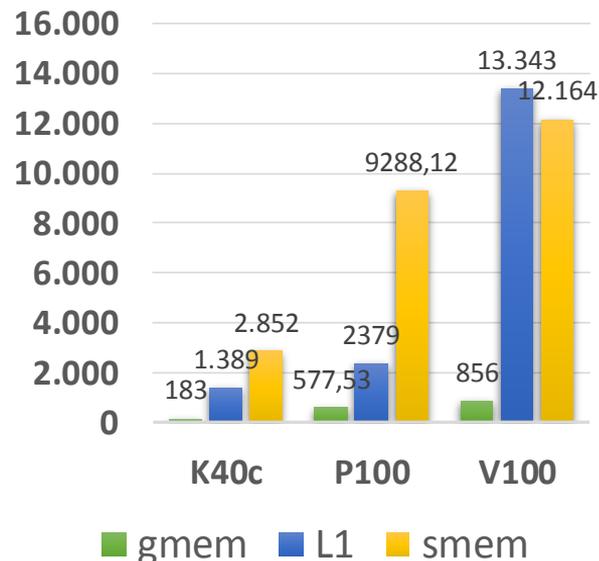
CENTER OF
EXCELLENCE

— throughput=measured, bandwidth=theoretical peak (load operations)

Bandwidth GB/s



Throughput GB/s



K40: L1+Smem unified
P100: L1+Tex Cache unified
V100: L1+Tex+Smem unified

Bandwidth Improvements

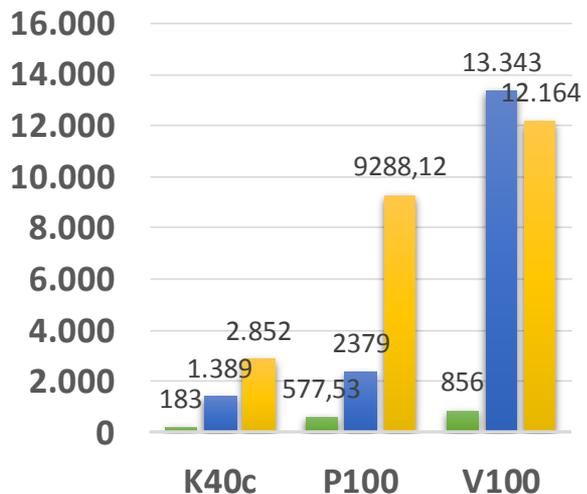
Experiments done with gpumembench

throughput=measured, bandwidth=theoretical peak (load operations)



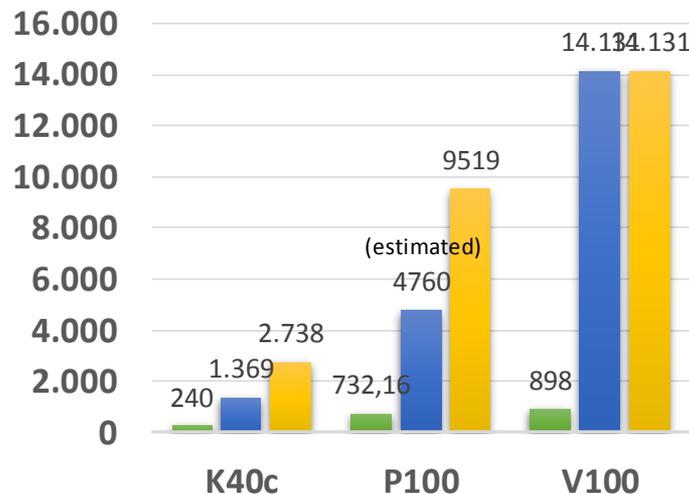
GPU
CENTER OF
EXCELLENCE

Throughput GB/s



■ gmem ■ L1 ■ smem

Bandwidth GB/s



■ gmem ■ L1 ■ smem

K40: L1+Smem unified
P100: L1+Tex Cache unified
V100: L1+Tex+Smem unified

AN-Coding Bruteforce Histogram Computations

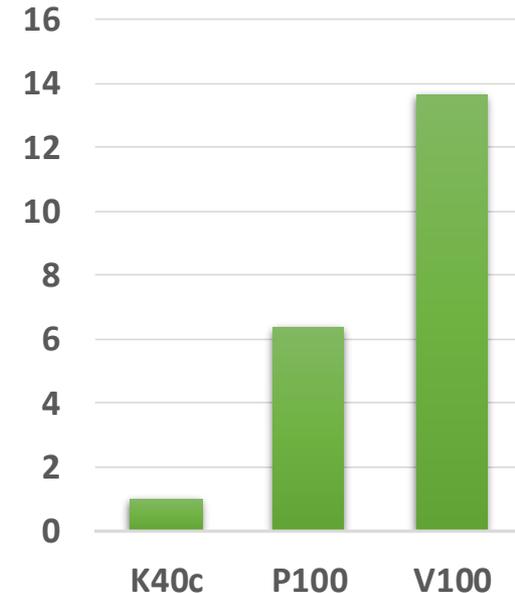


Bit flip Resilience for In-memory Column Stores

www.project-brics.de

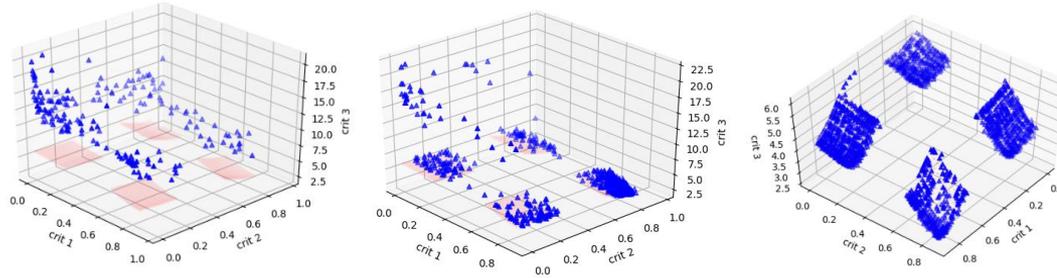
- AN-Coding is an arithmetic code for error detection, especially of multi-bit flips
- computes probability of silent data corruption by using the distance distribution of the code
- enumerates possible SDC bit patterns on GPU
- CUDA algorithm is mostly shared memory bound
- Almost 14x faster on V100 compared to K40

Speedup



Werner, M.; Kolditz, T.; Karnagel, T. et. al: Multi-GPU Approximation Methods for Silent Data Corruption of AN-Coding.

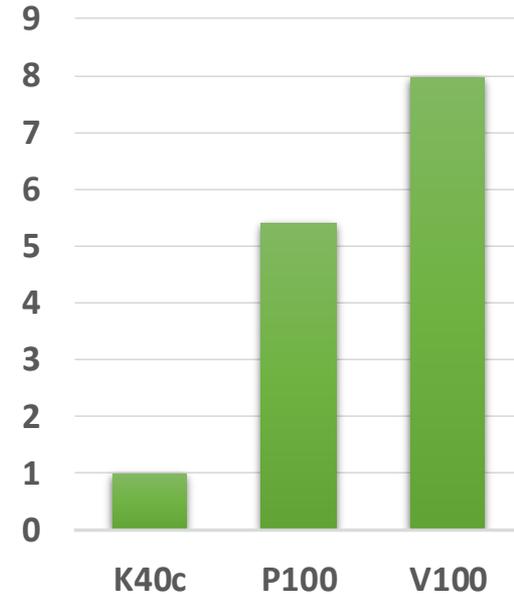
McEA Evolutionary Algorithm for Multi-Criteria Optimization



Pareto front (red=optimum): 10 generations 100 generations 1000 generations

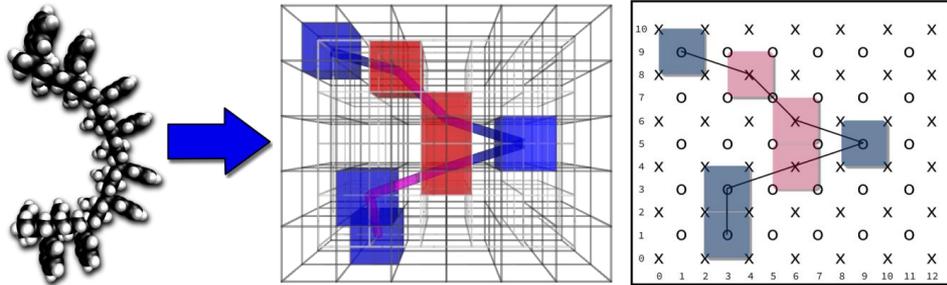
- multi-criteria optimization of production plans
- search heuristic: cellular evolutionary algorithm on GPUs (CUDA)
- population with over 1,000,000 individuals
- genetic evolution of 1000 generations
- for problem class see DTLZ-7*
- Almost 8x faster on V100 compared to K40c

Speedup

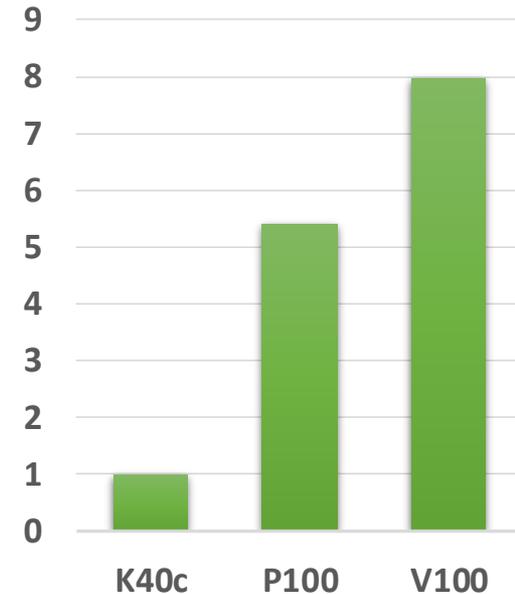


*Deb, Kalyanmoy, et al. "Scalable test problems for evolutionary multiobjective optimization." Evolutionary Multiobjective Optimization. Theoretical Advances and Applications (2005)

Polymer Science – Bond Fluctuation Model on GPUs



Speedup



- coarse-grained simulation model for polymers on GPU (CUDA)
- modified collision algorithm on a body-centered cubic grid
- exploration of new time and length scales using GPU
- almost 8x faster on V100 compared to K40

C. Jentzsch, R. Dockhorn, and J.-U. Sommer: A Highly Parallelizable Bond Fluctuation Model on the Body-Centered Cubic Lattice, in *Parallel Processing and Applied Mathematics*.

System Failure Analysis for HPC Clusters

Failures in HPC: Current Status

364 failures in year, 1990

2.33 failures per day, 2008

164,593 alerts per day, 2018

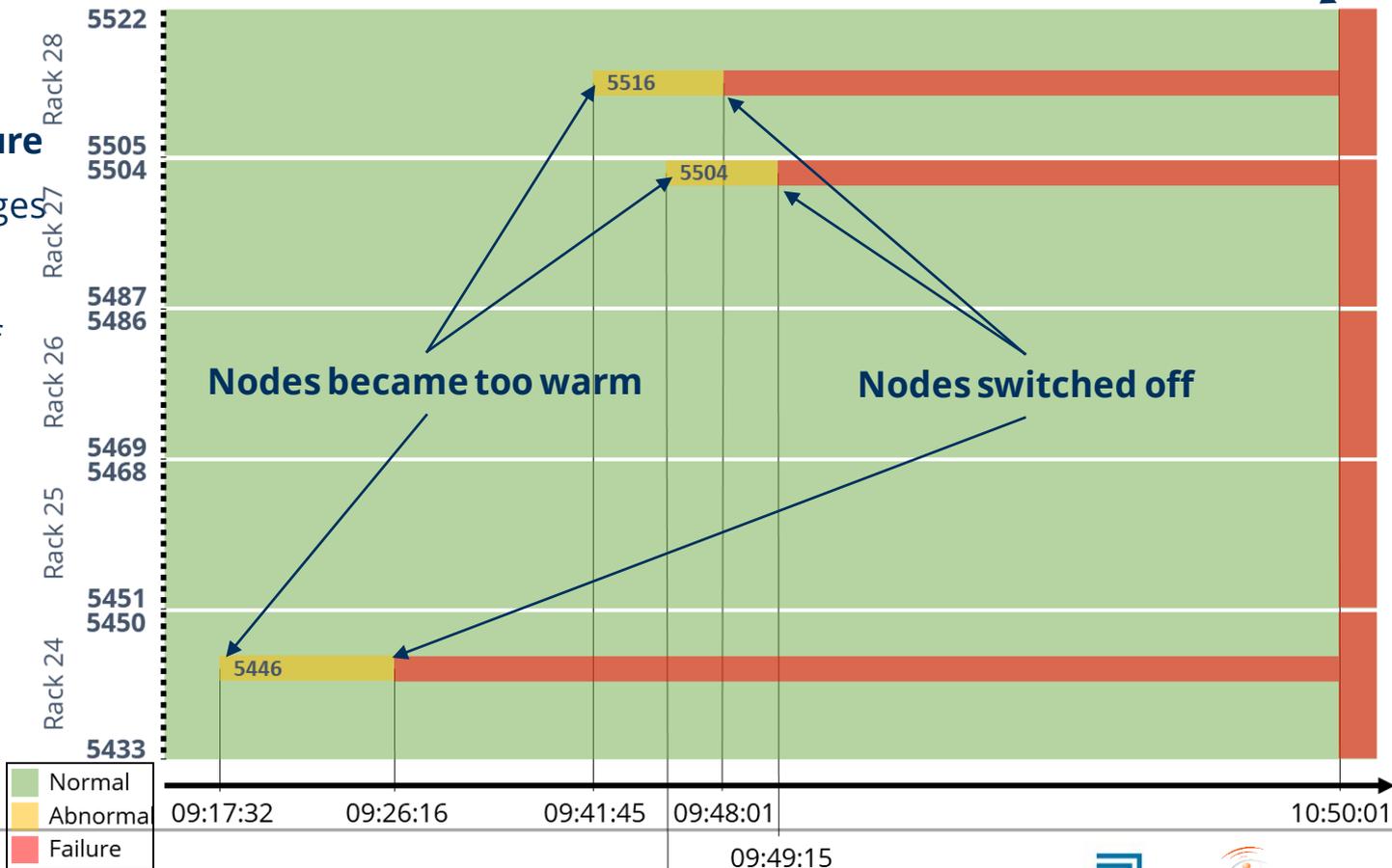
Date	Environment	Observation
1984-1986	IBM Mainframe	456 failures
1988-1990	Tandem	800 failures
1990	VAX	364 failures
1989-1990	VICE	300 failures
1999	70 Windows NT nodes	1100 failures
1999	503 nodes	2127 failures
2003	3000 nodes	501 failures
2003-2004	395 nodes	1285 failures
2004-2005	Liberty	7.8 alerts per day
2005-2006	Blue Gene/L	1,620 alerts per day
2005-2006	Thunderbird	13,312 alerts per day
2005-2007	Spirit (ICC2)	309,707 alerts per day
2005-2010	CENIC	16-302 failures per link
2006	Red Storm	16,016 alerts per day
2007	BlueGene/L Coastal	MTBF 7-10 days
2007-2009	Unknown	MTBF 3-37 minutes
2008-2010	Jaguar	2.33 failures per day
2008-2011	Jaguar XT4	MTBF 36.91 hours
2008-2011	Jaguar XT5	MTBF 22.67 hours
2008-2011	Jaguar XK6	MTBF 8.93 hours
2011-2017	Facebook network	MTBF 1.8 months
2012-2013	K Computer	Failure rate 1.6%
2013	Blue Waters	MTBF 4.2 hours
2014	Titan	317 HW and 270 SW failures
2014	Titan	9 failures per day
20013-2015	EOS XC 30	MTBF 189.04 hours
20013-2015	Titan XK7	MTBF 14.51 hours
2015	BlueGene/Q Mira	MTBF 5.5 hours
2015	Petascale systems	MTBF 7-10 hours
2015	Cielo	MTBF 24 hours
2015-2018	Titan	164,593 alerts per day
2017	Argonne FUSION	MTBF 3-52 minutes

Failures in HPC: Correlations

90 nodes switched off by overheating protection mechanism

Early detection of **failure chains** reduces damages

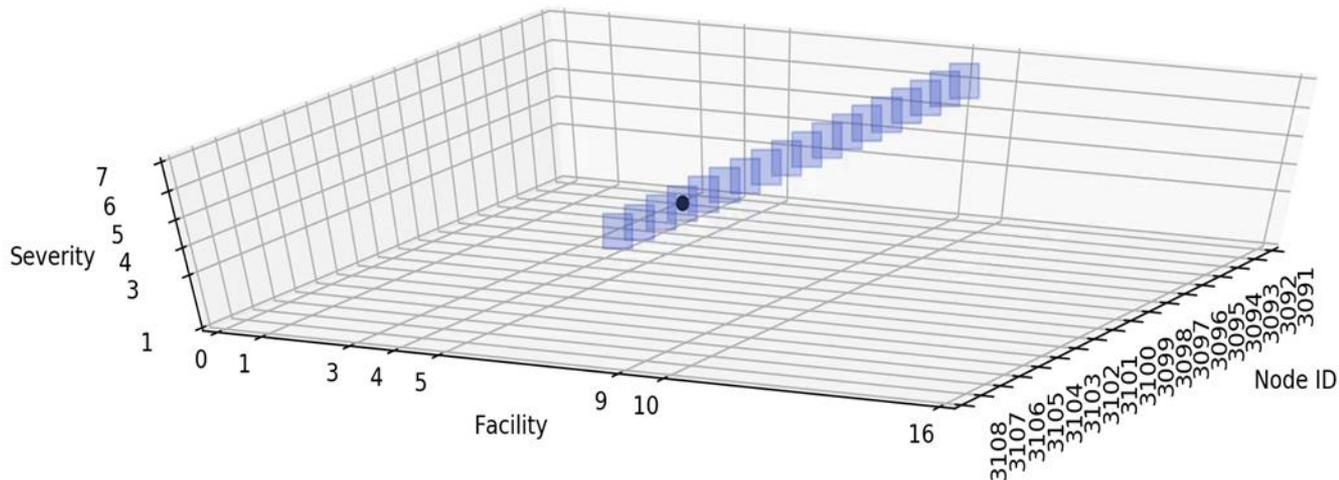
➤ **Proposed solution:** Statistical analysis of correlated failures in node vicinities



Failures in HPC: Timeline of a Failure in a Rack (18 nodes)

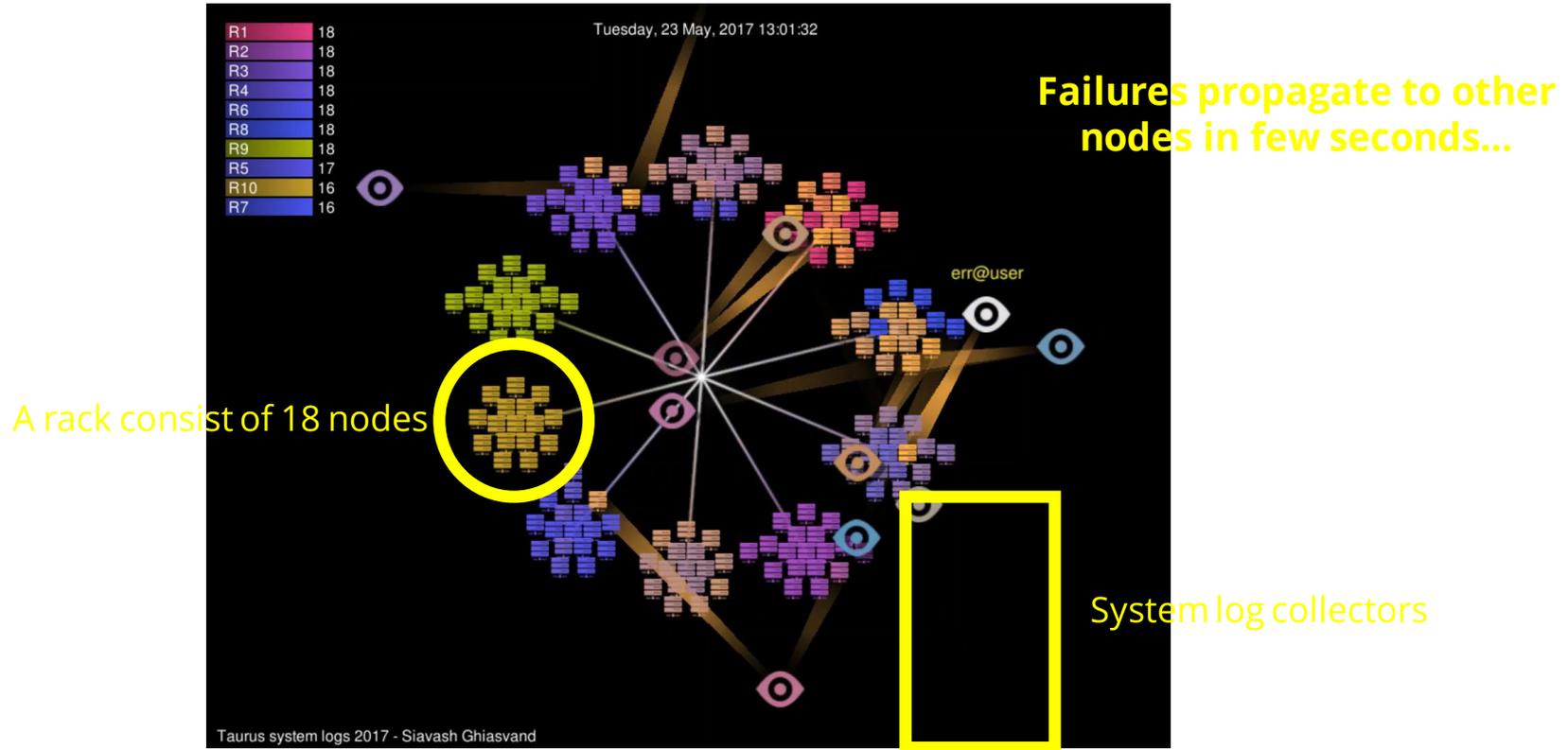


2017-05-21 23:50:00 to 2017-05-21 23:59:59
36000 seconds before 2017-05-22 09:49:09 node 3105



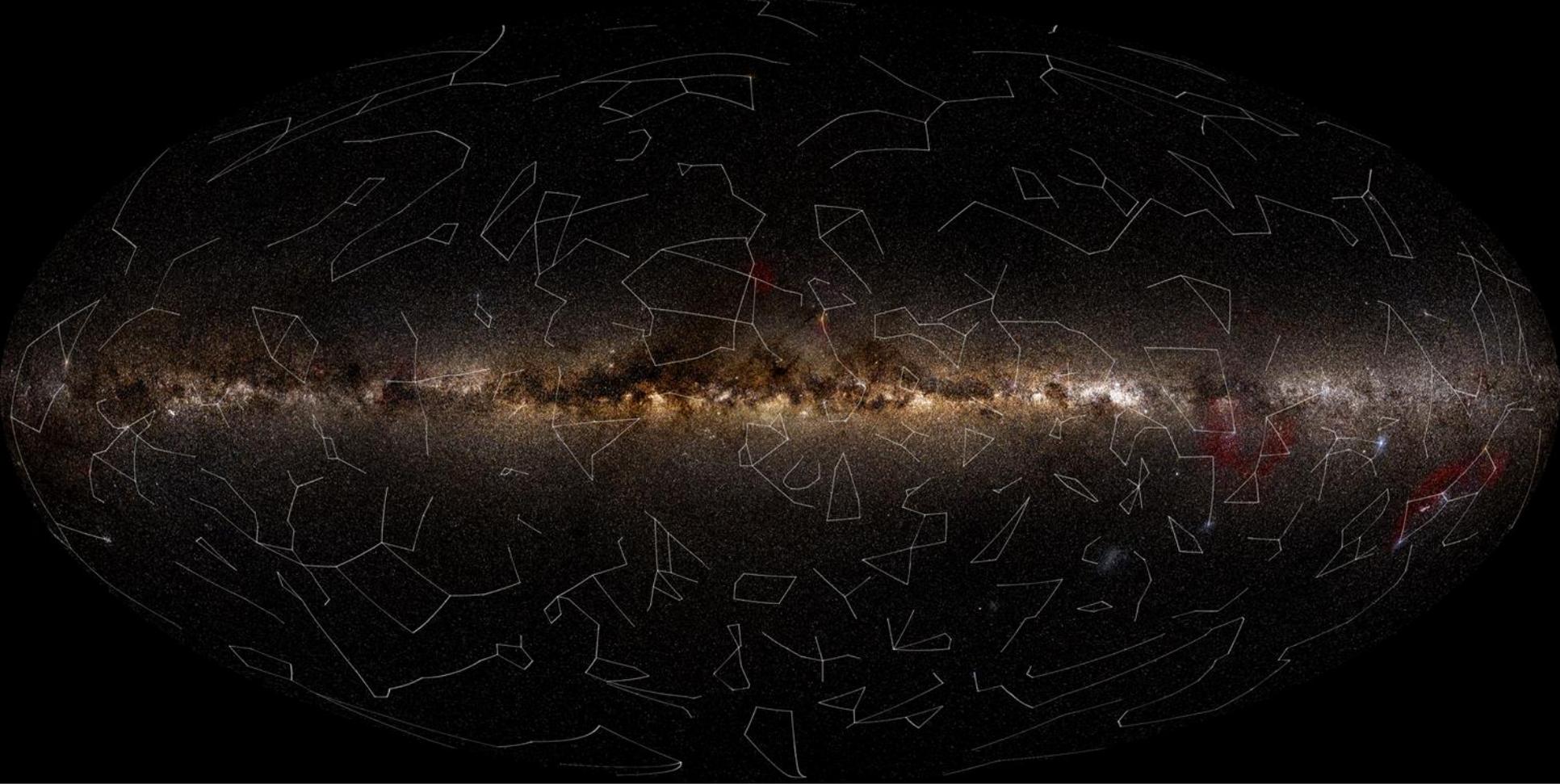
Stream of system logs, divided based on the facility and severity of entries

Failures in HPC: Timeline of a Major Failure in Island 3 (180 nodes)



Courtesy of Prof. Sergei A. Klioner et.al.

Astrophysics: The Gaia Project



One of the main problems of astronomy: distance

Without knowing how far the object is, physical understanding of that object is impossible...

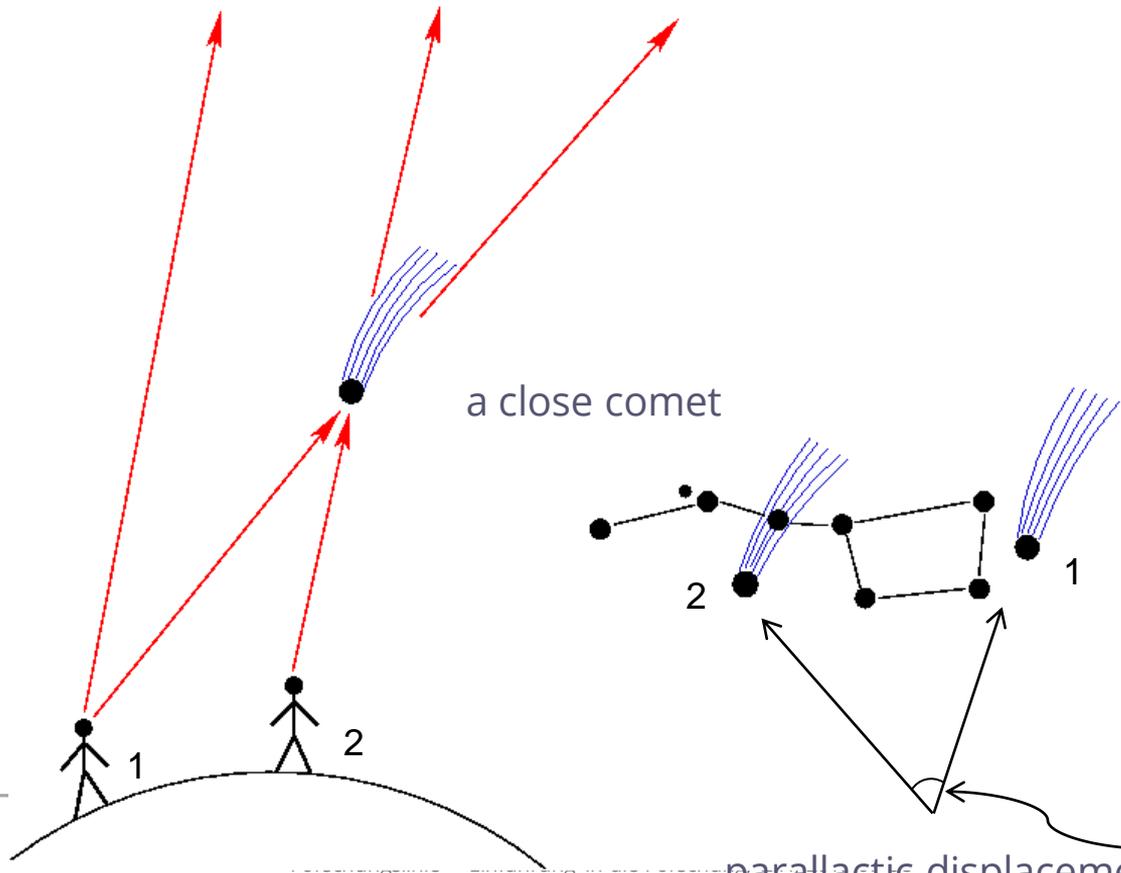


A comet: far away and very big or inside the Earth atmosphere and rather small?
Tycho Brahe, 1577

One of the main problems of astronomy: distance

a far away object

a close comet



Astrometry: the art of measuring stellar positions

Astronomy cannot touch
its objects!

Astronomy cannot make
experiments!

Astronomy analyses stellar light:

Astrometry	- direction
Photometry	- quantity
Spectroscopy	- colour and more
Polarimetry	- polarization

+ cosmic particles

++ gravitational waves



Why to bother?

- We need to understand stars.

(our Sun is a star!)

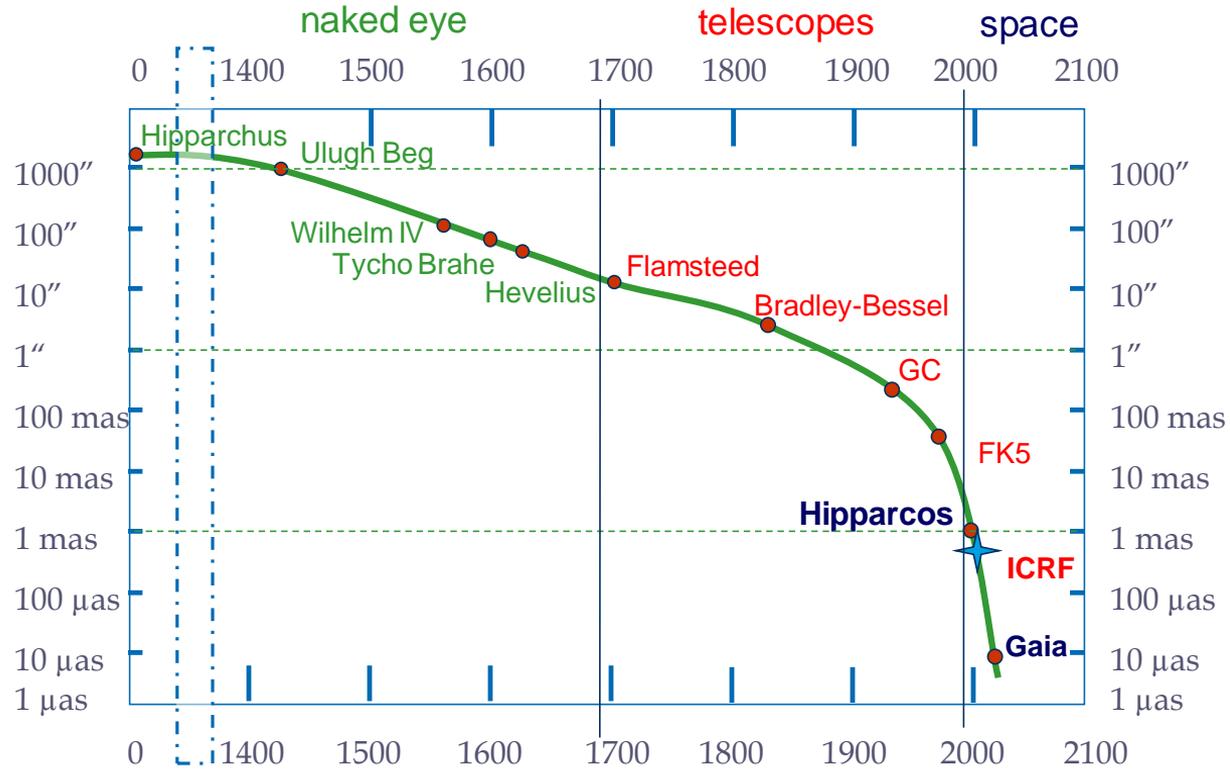
Without knowing the distance it is not possible to judge if a star is big or small, etc.

- We live in a galaxy.

We need to understand how our Galaxy was formed.



Accuracy of astrometric observations



What do we know about our Galaxy?

The Sun should be here:



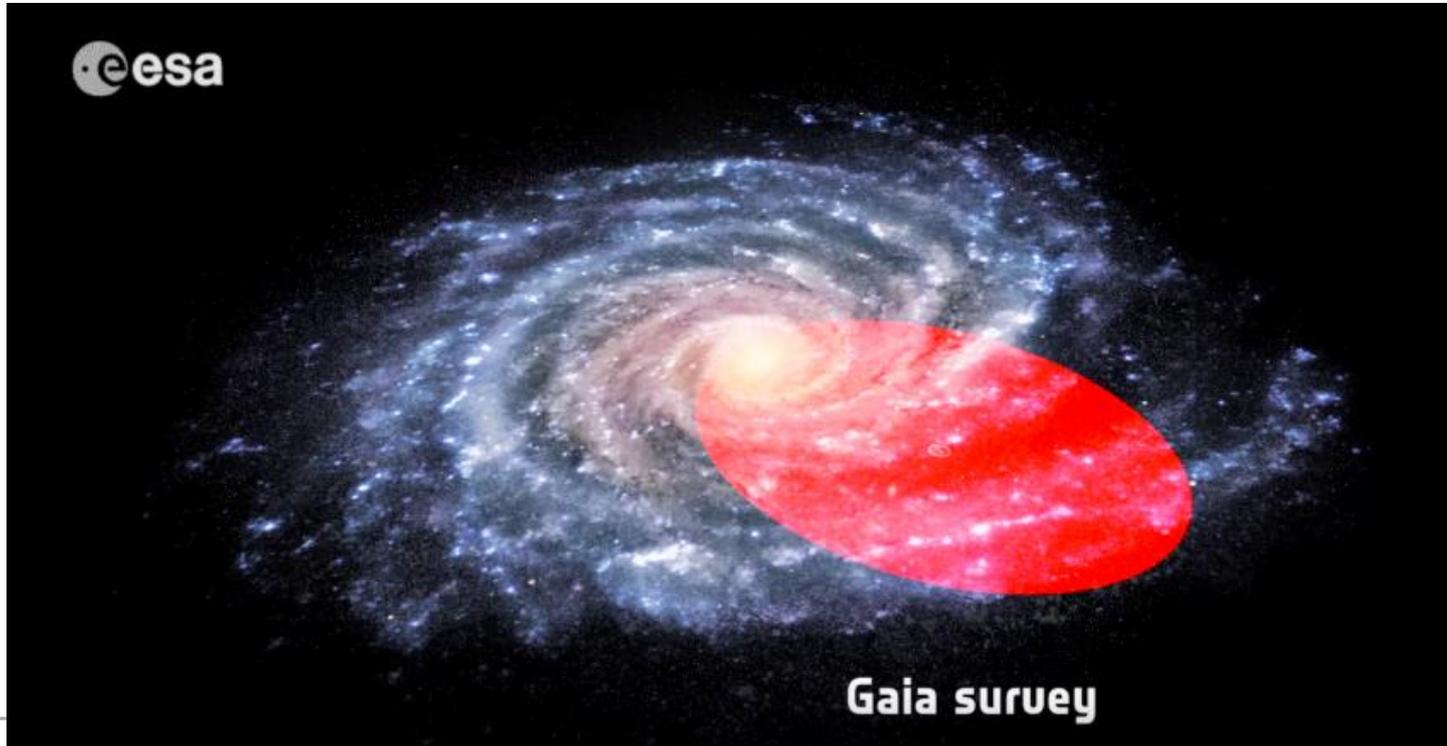
What do we know about our Galaxy?

The stars with distances known till 2016 are all in the small red spot:



What do we know about our Galaxy?

With Gaia we can explore a significant part of our Galaxy:

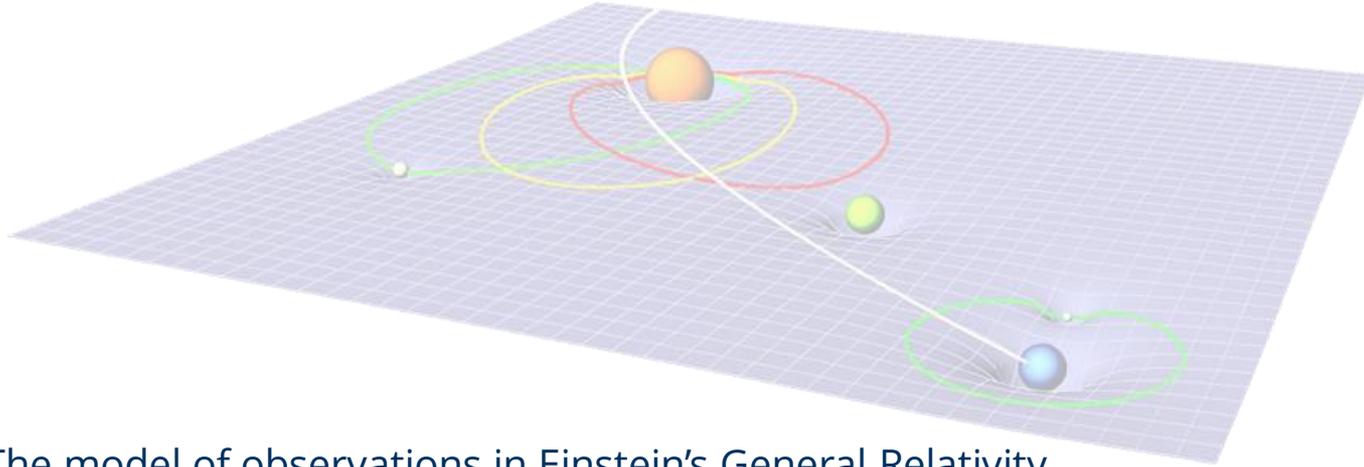


The challenge of data processing

- Parameters
 - At least 5 parameters for each star: $5 \times 1.7 \cdot 10^9$
 - 4 parameters of orientation each 15 seconds: 10^8
 - 2000 calibration parameters per day: $4 \cdot 10^6$
 - global parameters: $<10^4$
- Observations
 - about 1000 raw images for each star: 10^{12}
- Data volume: 1 PB (iteratively!)
- Computational efforts: $\sim 10^{22}$ flops
- Direct least squares solution is impossible



Gaia in Dresden



1. The model of observations in Einstein's General Relativity
2. Tests of fundamental physical laws with Gaia data
3. Analysis of the Gaia reference frame: quasars
4. Synchronization and monitoring of Gaia's atomic clock
5. Special astrometric solutions: stability and quality verification, special calibration of the instrument, relativistic tests

Special thanks to ZIH for about 3 Million CPU-hours by now!

[Video](#)

Summary

Computer Science and Computational Science at ZIH

- Exciting computer science research
- Broad spectrum of computational science topics together with application field scientists

Contact:

Dr. Andreas Knüpfer
Deputy Director / CTO of ZIH
andreas.knuepfer@tu-dresden.de

Tel. +49 351 463-38323

Willersbau A113

