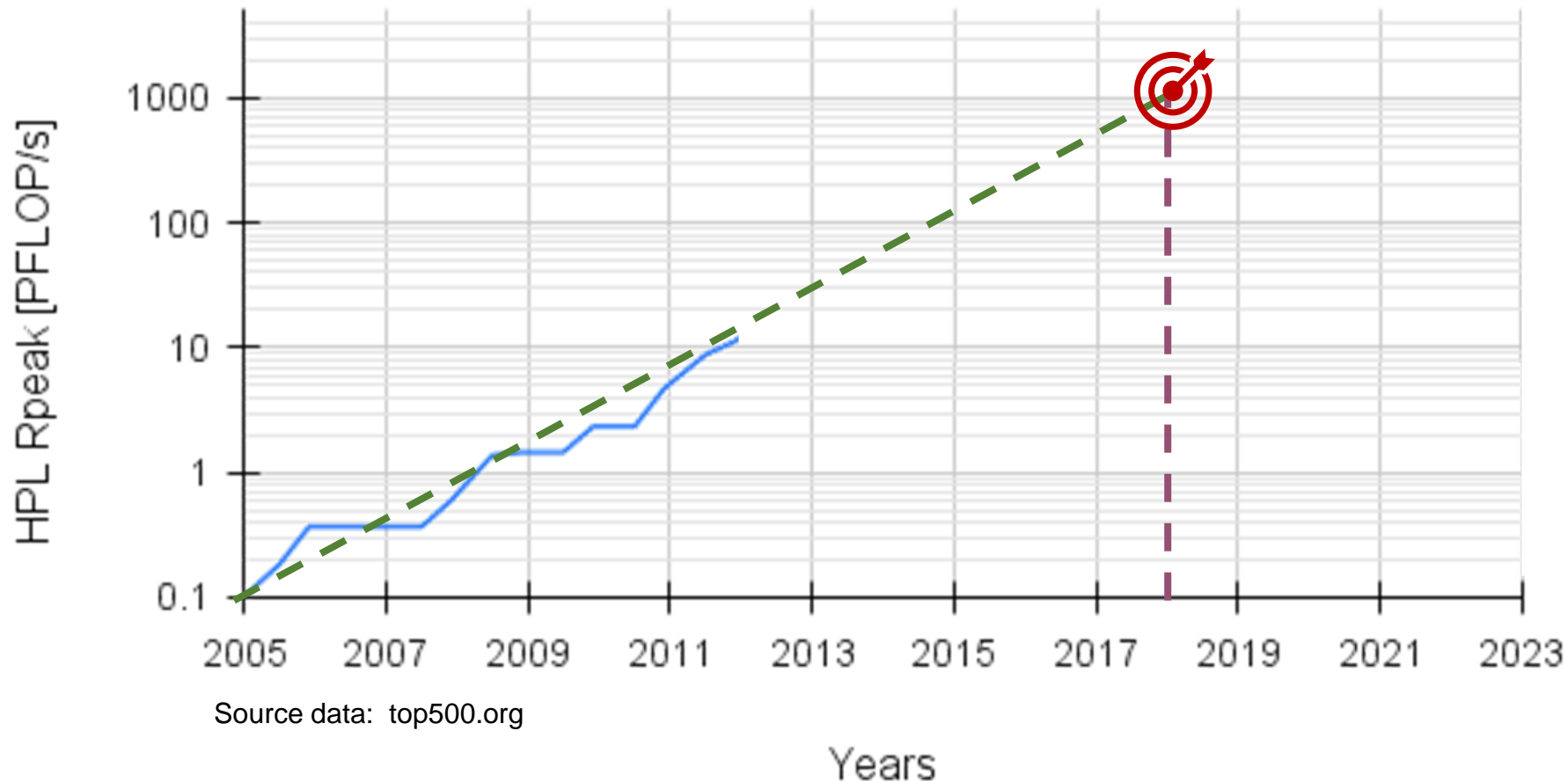


# The DEEP Series of Projects – Application-Driven Co-Design Towards Exascale

Hans-Christian Hoppe, Jülich Supercomputing Centre  
CECAM Flagship Workshop 2024,

# Destination Exascale – 2012 PoV

Top #1: HPL Rpeak [PFLOP/s]



1997: First **1TFlop/s** computer:  
(ASCI Red/9152)

2008: First **1 PFlop/s** computer: (Roadrunner)

There is a **solid trendline** – what could possibly go wrong?

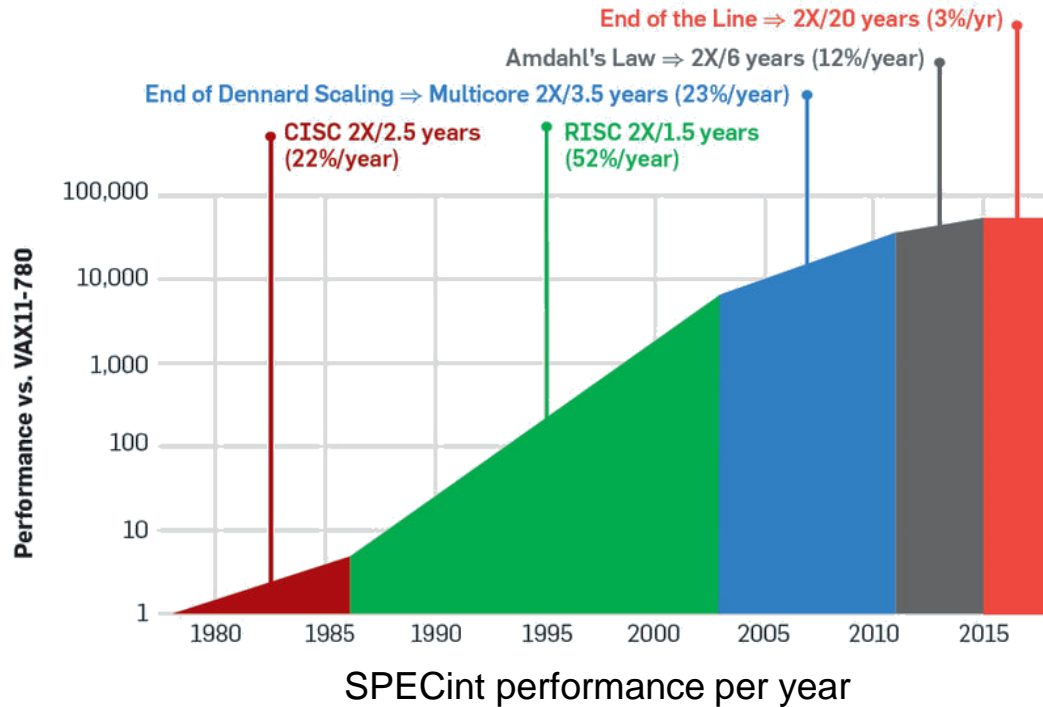
Of course, Top500 data reflects deployed, operational systems ...

- There were clear challenges ahead

# Moore's Law Slowing Down is not Your Friend

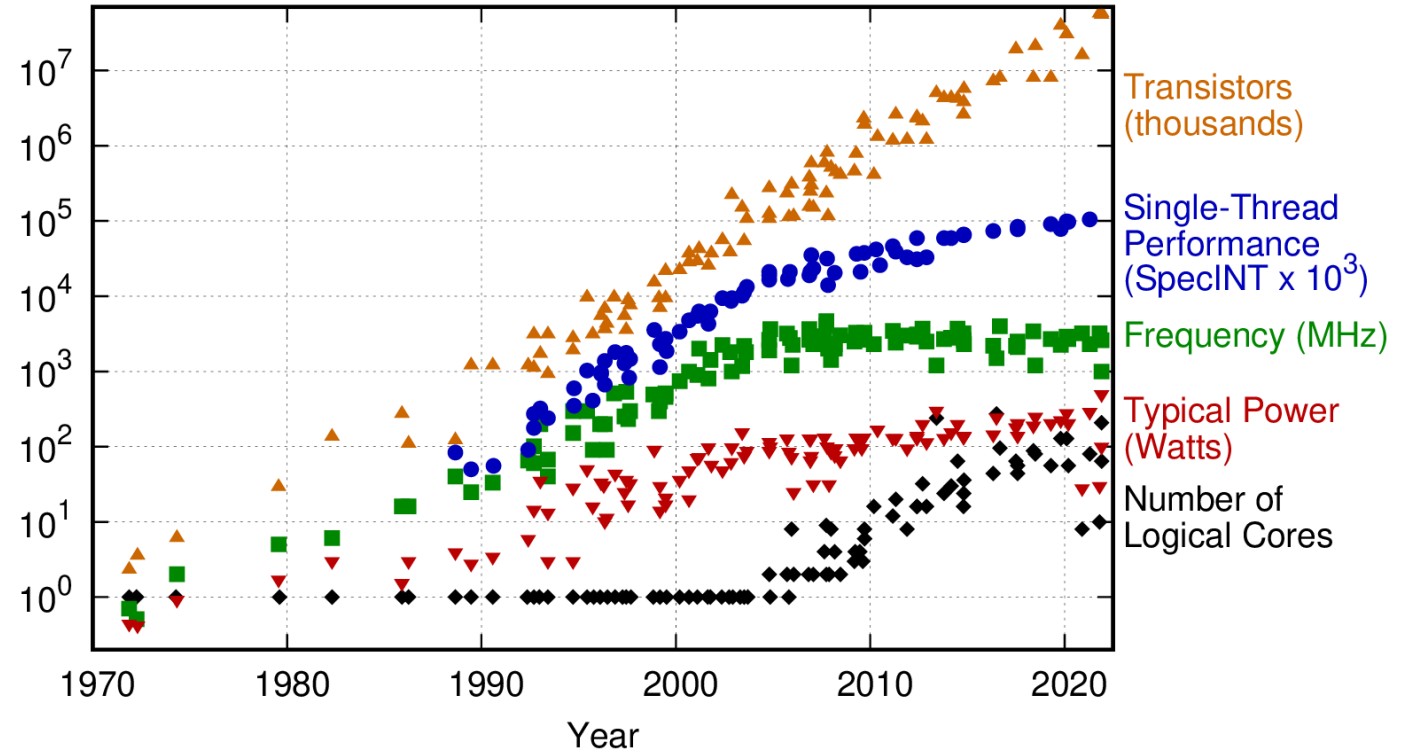
2012

2012



[https://www.researchgate.net/publication/343096513\\_Energy\\_Efficient\\_Computing\\_Systems\\_Architectures\\_Abstractions\\_and\\_Modeling\\_to\\_Techniques\\_and\\_Standards](https://www.researchgate.net/publication/343096513_Energy_Efficient_Computing_Systems_Architectures_Abstractions_and_Modeling_to_Techniques_and_Standards)

50 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten  
New plot and data collected for 2010-2021 by K. Rupp

<https://github.com/karlrupp/microprocessor-trend-data>

# And, There are More Exascale Challenges

## Application parallelism

- Applications must support billions of individual threads
- Lower-scaling applications / parts of applications should not run on a full Exascale system

DEEP Projects

## Truly scalable systems

- Huge numbers of devices need to exchange data with each other
- Collective communication operations are “slowing down” due to larger system sizes
- Network contention and reliability become worries

## Energy efficiency

- Accelerators clearly beat CPUs for many (most?) codes
- System heterogeneity is a must
- Yet – portable accelerator programming is hard

DEEP Projects

## Memory and storage

- Ever growing gap between compute throughput and memory bandwidth
- New technologies like HBM suffer from capacity limitations & high energy consumption

DEEP Projects

## Workload diversity

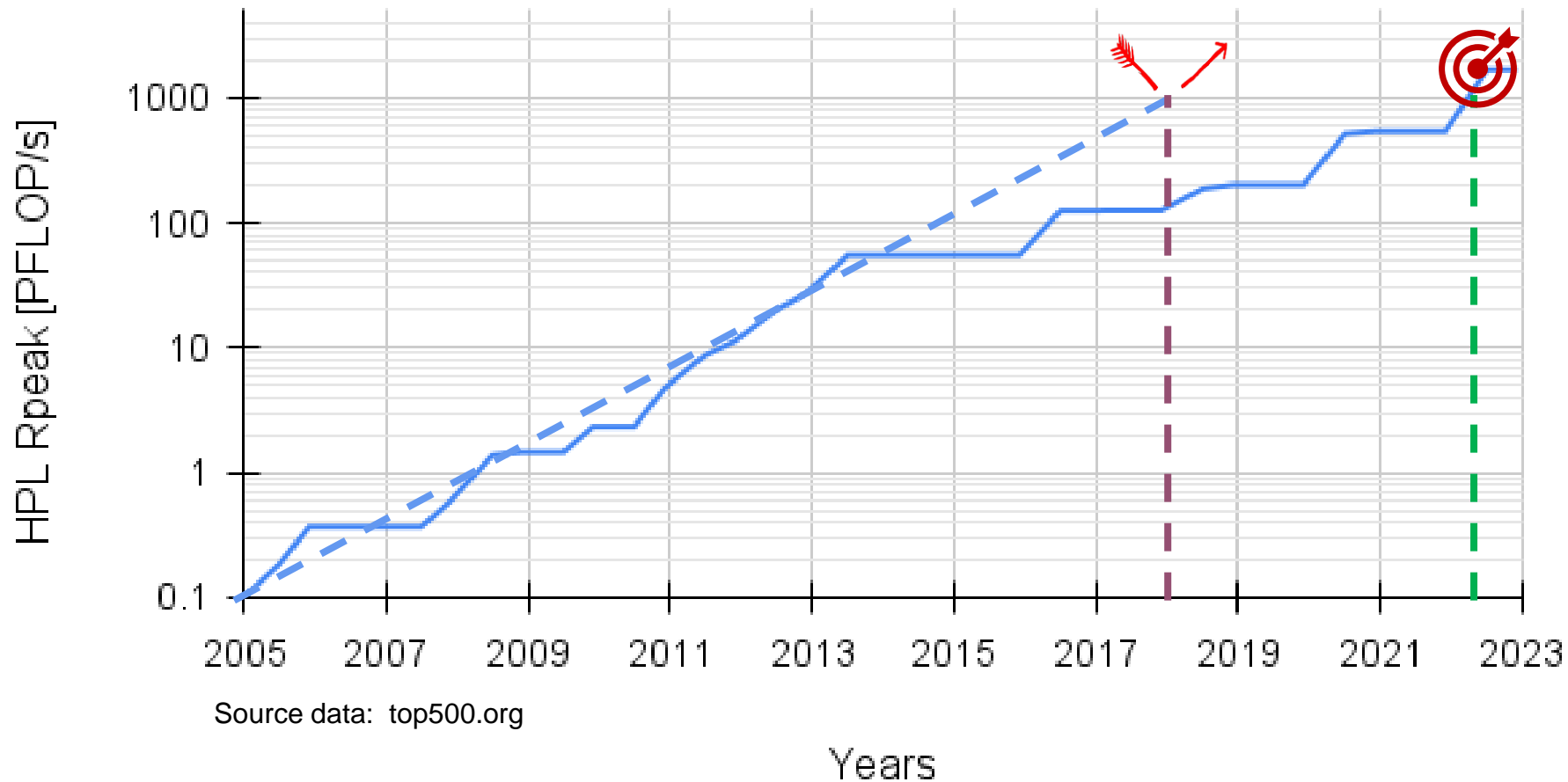
- Exascale centers must run a wide variety of HPC, AI and data analytics workloads with highest energy efficiency
- One size does not fit all

DEEP Projects



# Destination Exascale – Hindsight is 20/20

## Top #1: HPL Rpeak [PFLOP/s]



**1997:** First **1TFlop/s** computer:  
(ASCI Red/9152)

**2008:** First **1 PFlop/s** computer: (*Roadrunner*)

So.... First **1 EFlop/s** computer: **2018 !!**

– Well... not really

It took 4 years longer....

**2022**  
for *Frontier* to appear

# Co-Design to the Rescue?

## Wider interpretation

“Electronic and computer scientists speak of co-design when they want to describe the way hardware anticipates software and software adapts to hardware, both evolving towards a better integration”  
<https://codesignlab.wp.imt.fr>

## Sequential flow

- Pipelining could increase „throughput“
- Careful definition of APIs to overlap steps 2 & 3

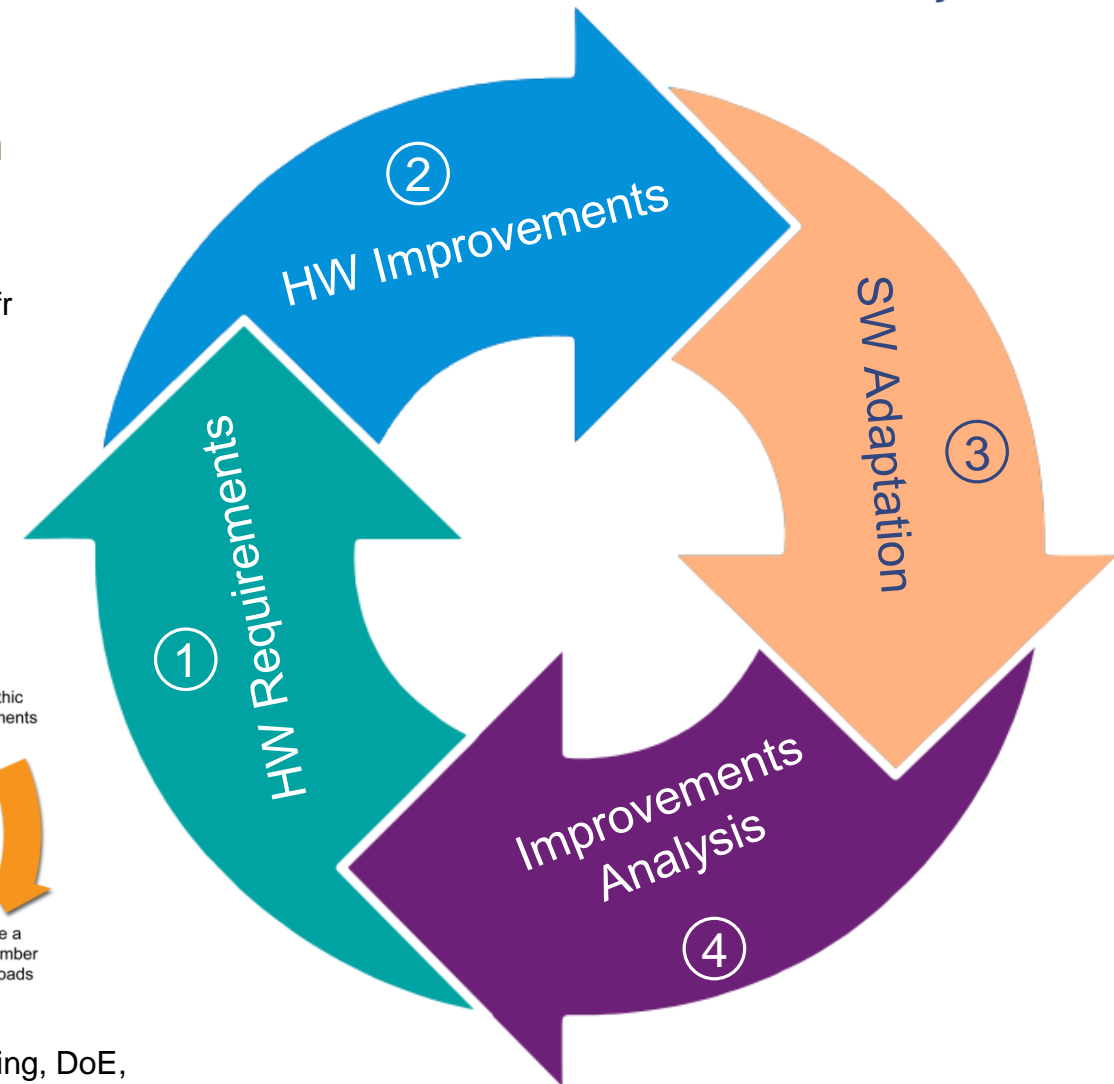
## Problematic steps

- Step 2 – HW takes (lots of) time & money
- Simulators/emulators can accelerate steps 2 & 3
- Step 4 – Adapting representative workloads takes time

## How many iterations

- Do we need
- Can we afford?

Reimagining Codesign for Advanced Scientific Computing, DoE,  
<https://doi.org/10.2172/1822199>



# Modified Co-Design Approach for DEEP

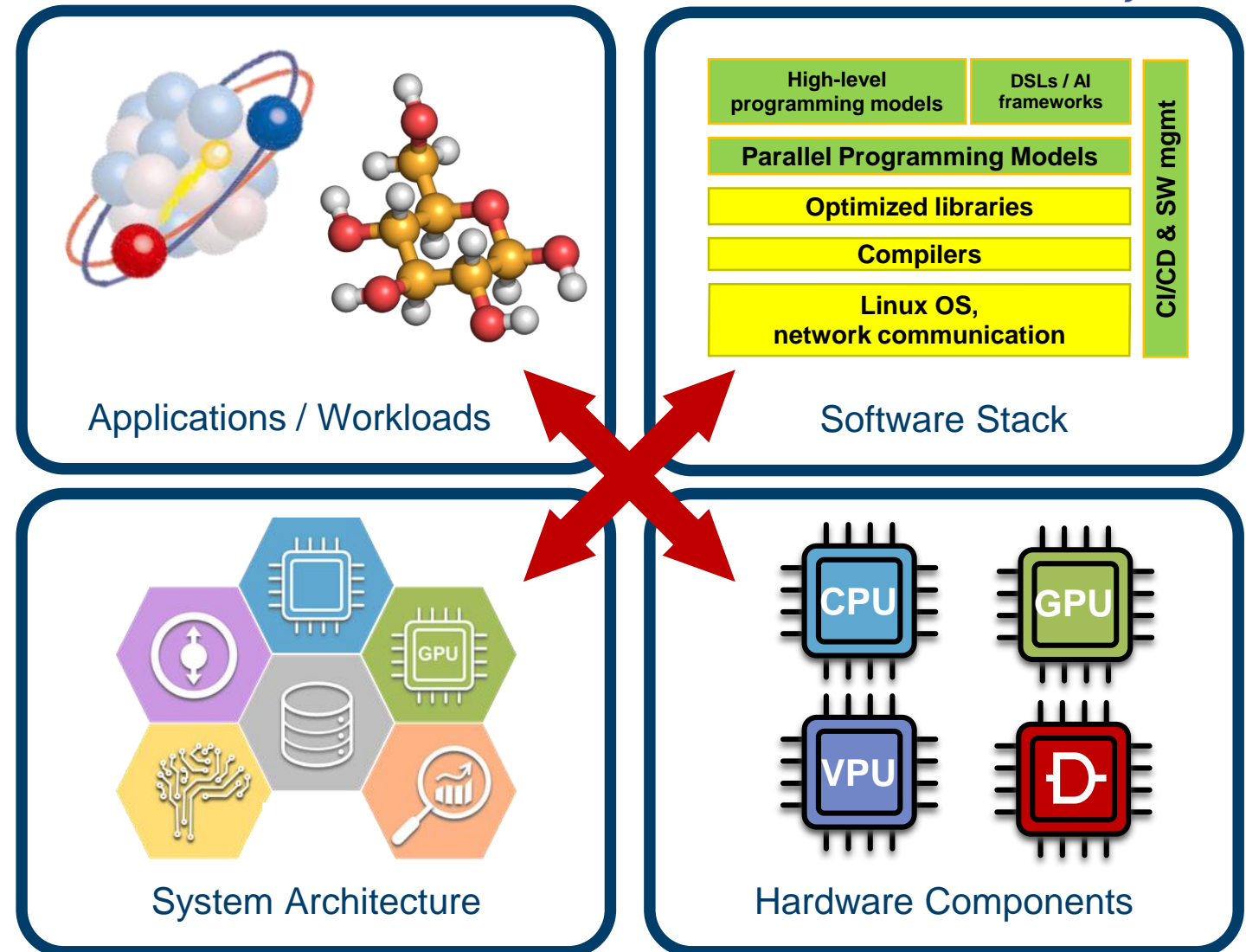
Optimise **all four** components to achieve **combined**

- **Performance** and
- **Energy efficiency**

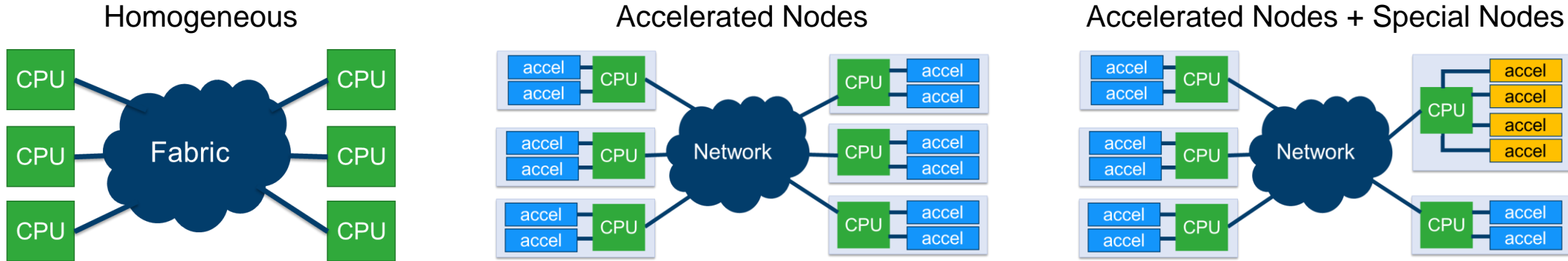
Identify **interactions** and **simulate/predict effects** of changes to other components

Ideally, **all components** evolve

- In reality, some are **less “malleable”** than others



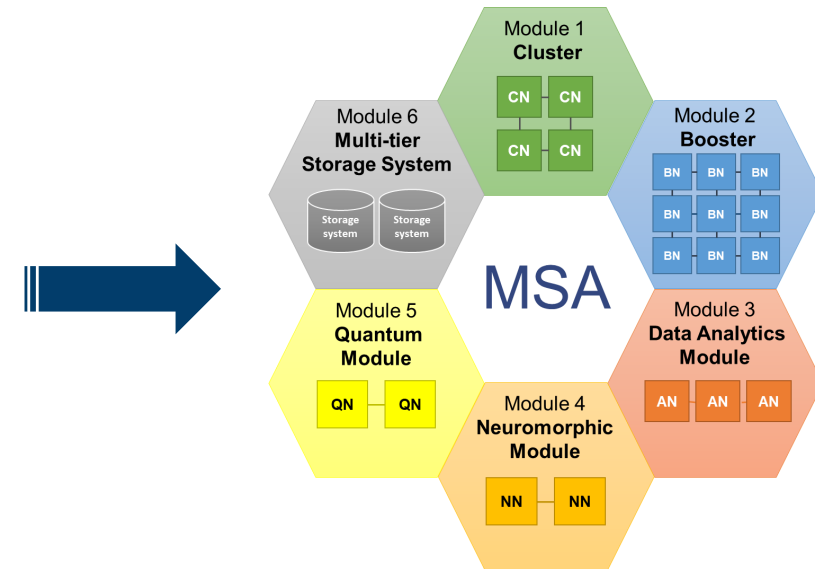
# Approaches to System Heterogeneity



Homogeneous systems lack efficiency\*

Accelerated nodes fix the ratio of CPUs vs. accelerators, complicate sharing resources across nodes

Adding „special nodes“ for certain tasks



\*: certainly for AI and dense linear algebra applications



# DEEP Prototype Systems

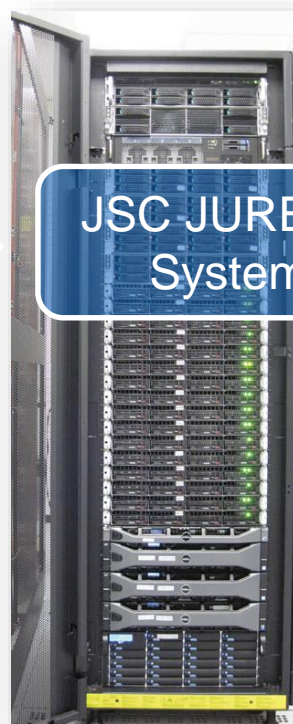


## DEEP Prototype

128 Xeon, 284 KNC nodes  
InfiniBand + Extoll  
550 TFlop/s

Directly connected accelerators

Hot water cooling



## JSC JURECA System

## DEEP-ER Prototype

16 Xeon, 8 KNL nodes  
100Gbit/s Extoll  
40 TFlop/s

European low-latency NW

Aggressive dense  
packaging



## JSC JUWELS System

## JSC JUPITER system

## DEEP-EST Prototype

55 Cluster (Xeon), 75 Booster (V100), 16 Data Analytics (FPGA) nodes  
100 Gbit Extoll, InfiniBand, Ethernet  
800 TFlop/s

Intel Optane persistent memory

First true multi-module system

# Modular Supercomputing Architecture

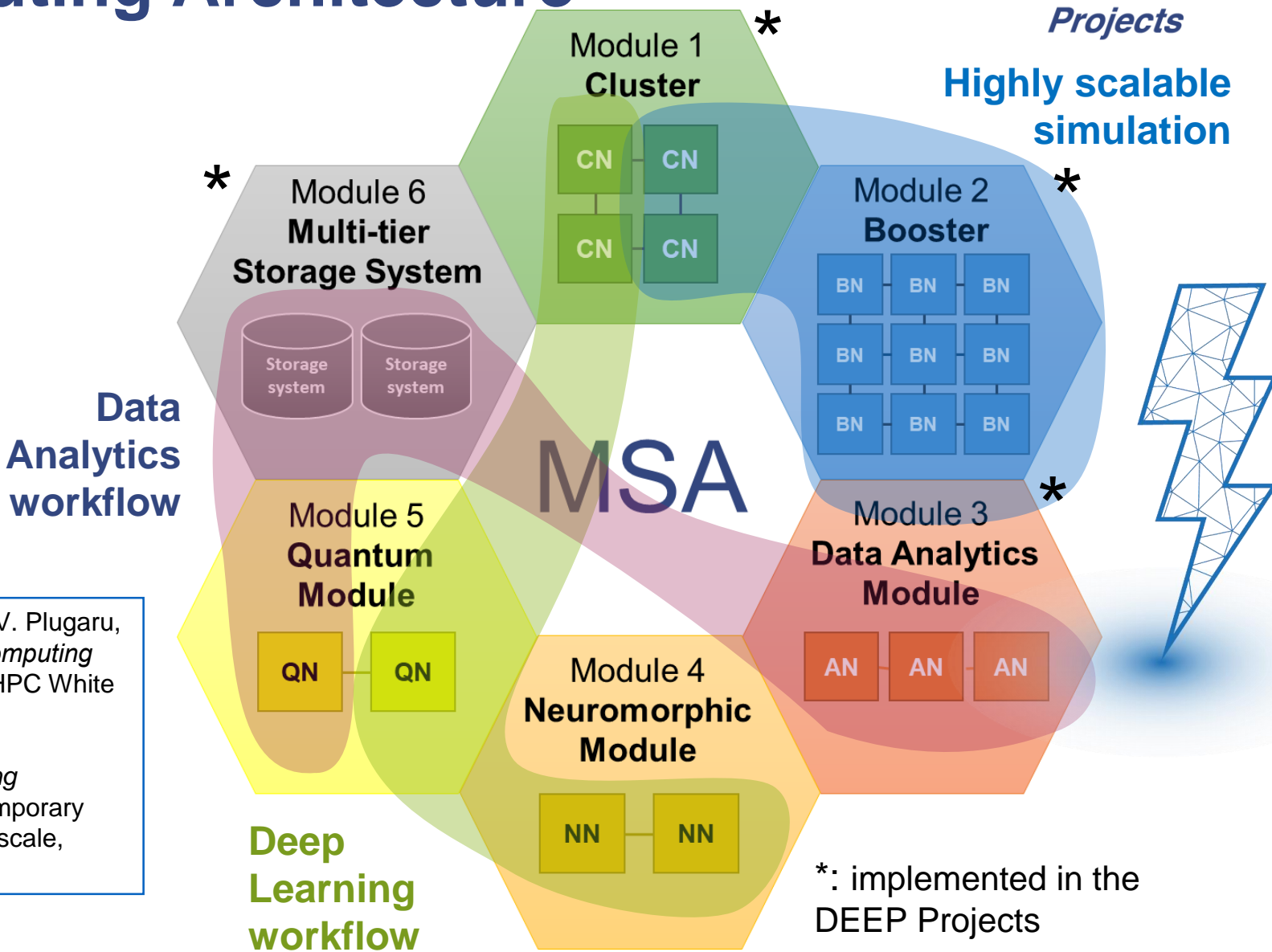
## Composability of heterogeneous resources

Cost-effective scaling

Effective resource-sharing

Match workload diversity

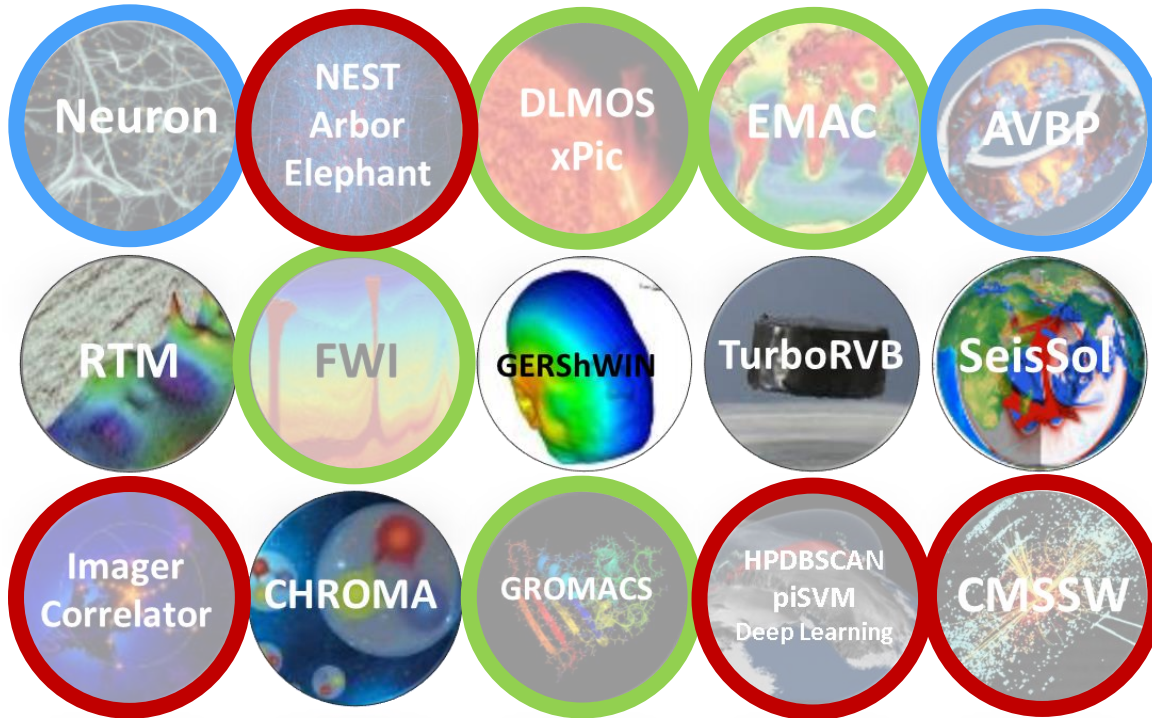
- Data analytics
- Machine- and Deep Learning
- Artificial Intelligence



- E. Suarez, N. Eicker, T. Moschny, S. Pickartz, C. Clauss, V. Plugaru, A. Herten, Kristel Michielsen, T. Lippert, "Modular Supercomputing Architecture – A Success Story of European R&D", ETP4HPC White Paper. (2022) Available at <https://www.etp4hpc.eu/white-papers.html#msa>.
- E. Suarez, N. Eicker, Th. Lippert, "Modular Supercomputing Architecture: from idea to production", Chapter 9 in Contemporary High Performance Computing: from Petascale toward Exascale, Volume 3, p 223-251, CRC Press. (2019)



# DEEP Projects Co-Design Applications

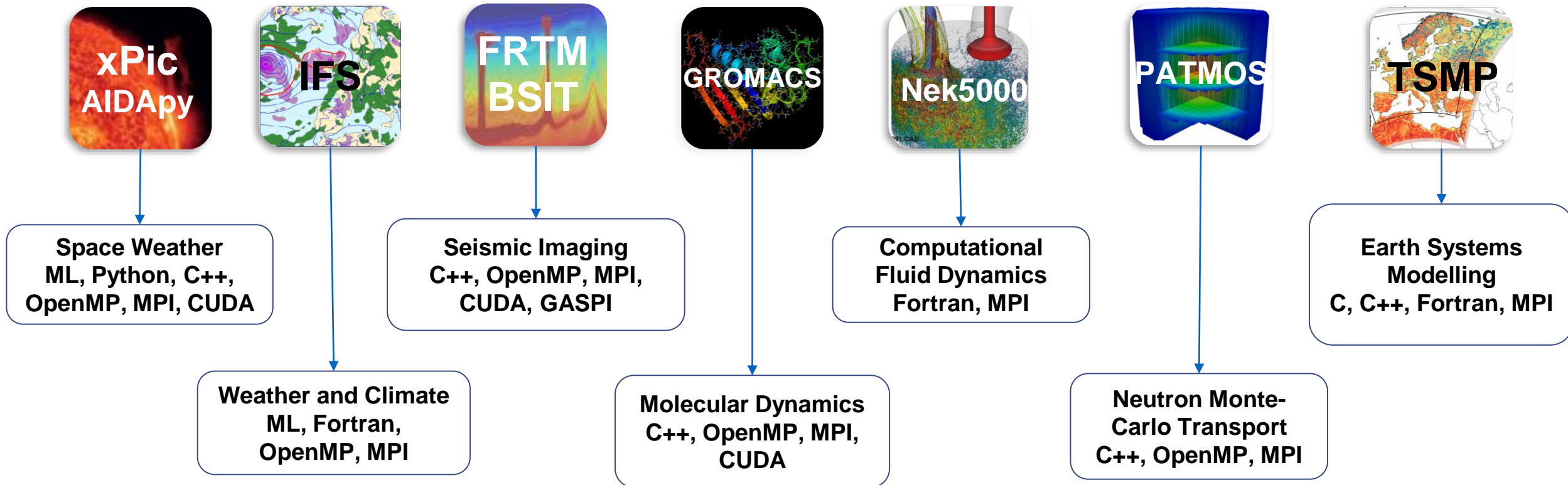


MSA usage modes

Code partition  
Workflow  
I/O forward

- **Kreuzer, et al.**, *Application Performance on a Cluster-Booster System*. IPDPSW – HCW (2018) [10.1109/IPDPSW.2018.00019]
- **Kreuzer et al.**, *The DEEP-ER project: I/O and resiliency extensions for the Cluster-Booster architecture*. HPCC'18 proceedings (2018) [10.1109/HPCC/SmartCity/DSS.2018.00046]
- Wolf et al., *PIC algorithms on DEEP: The iPic3D case study*. PARS-Mitteilungen 32, 38-48 (2015)
- Christou et al., *EMAC on DEEP*, Geoscientific model devel.(2016) [10.5194/gmd-9-3483-2016]
- Kumbhar et al., *Leveraging a Cluster-Booster Architecture for Brain-Scale Simulations*, Lecture Notes in Computer Science 9697 (2016) [10.1007/978-3-319-41321-1\_19]
- Leger et al., *Adapting a Finite-Element Type Solver for Bioelectromagnetics to the DEEP-ER Platform*. ParCo 2015, Advances in Parallel Computing, 27 (2016) [10.3233/978-1-61499-621-7-349]

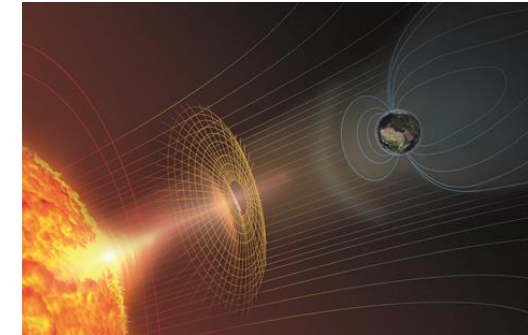
# DEEP-SEA Co-Design Applications



# Heterogenous Systems – Application View

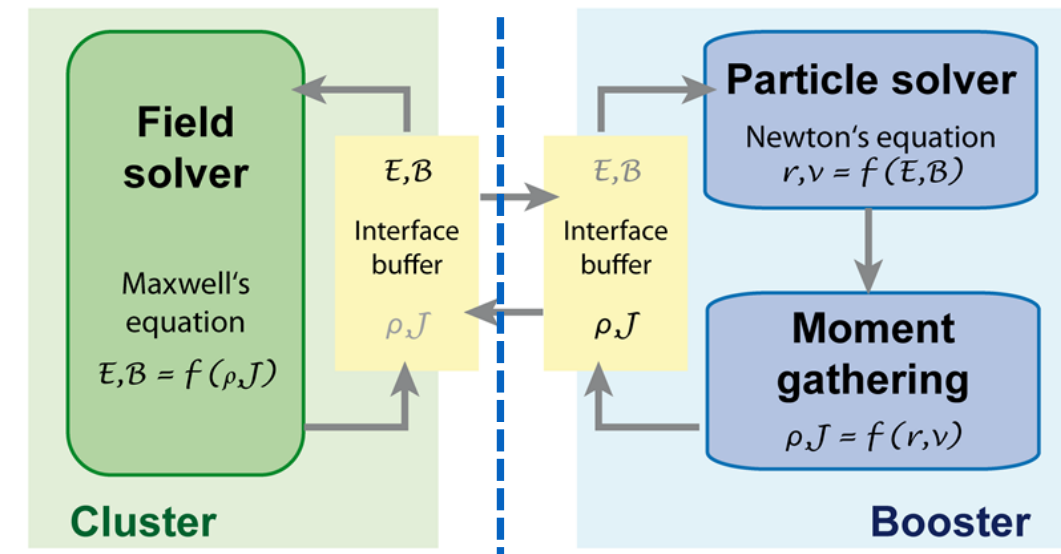
- **Space Weather simulation**

- Simulates plasma produced in solar eruptions and its interaction with the Earth magnetosphere
- Particle-in-Cell (PIC) code
- Authors: KU Leuven



- **Two solvers:**

- **Field solver:** Computes electromagnetic (EM) field evolution
  - *Limited code scalability*
  - *Frequent, global communication*
- **Particle solver:** Calculates motion of charged particles in EM-fields
  - *Highly parallel*
  - *Billions of particles*
  - *Long-range communication*



**A. Kreuzer**, J. Amaya, N. Eicker, E. Suarez, "Application performance on a Cluster-Booster system", 2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), HCW (20th International Heterogeneity in Computing Workshop), Vancouver (2018), p: 69 - 78. [doi: [10.1109/IPDPSW.2018.00019](https://doi.org/10.1109/IPDPSW.2018.00019)]



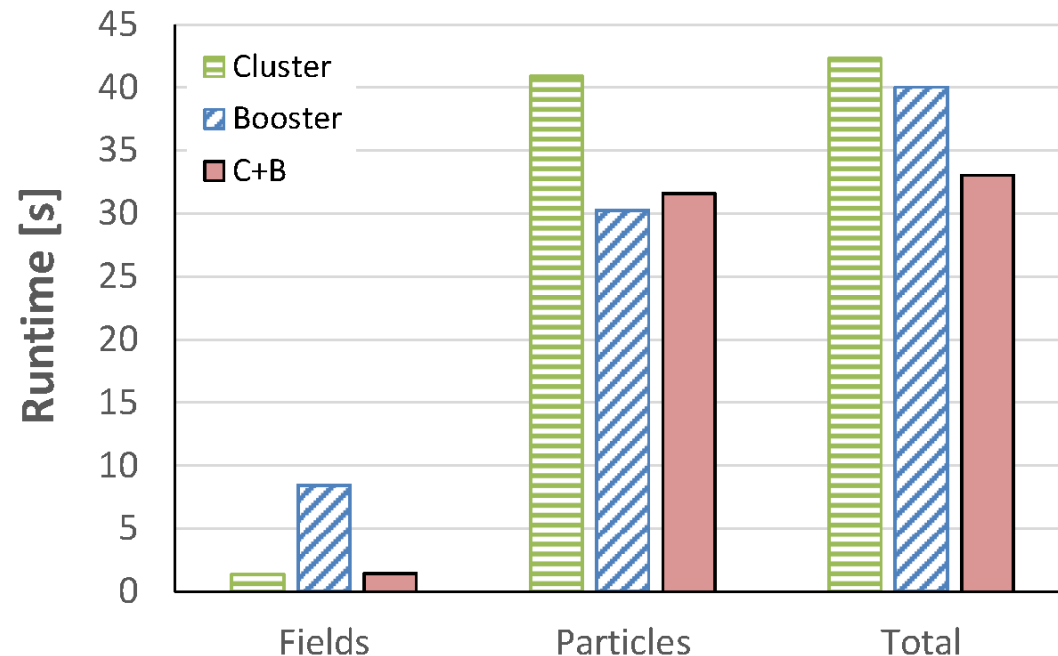
# xPic – Small Scale Performance Results

- **Field solver:** 6× faster on **Cluster**
- **Particle solver:** 1.35 × faster on **Booster**
- **Overall performance gain:**

**1× node**    **28% × gain** compared to Cluster alone  
**21% × gain** compared to Booster alone

**8× nodes**    **38% × gain** compared to Cluster only  
**34% × gain** compared to Booster only

- 3%-4% overhead per solver for C+B communication (point to point)

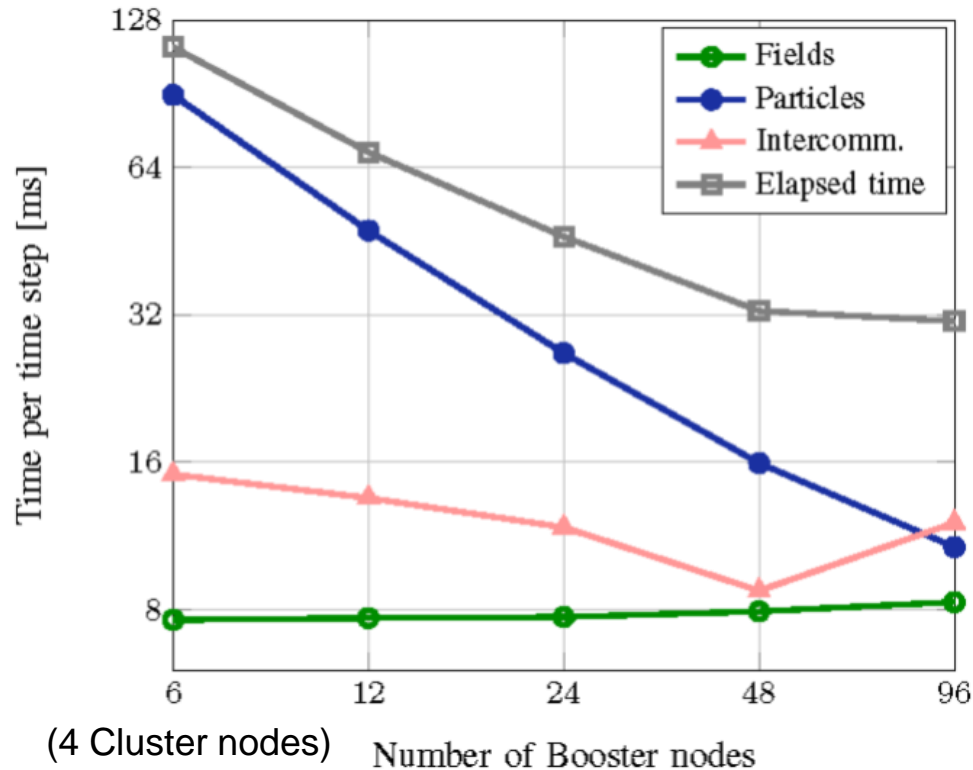


#cells per node	4096
#particles per cell	2048

**A. Kreuzer et al.** "Application Performance on a Cluster-Booster System", 2018 IEEE IPDPS Workshops (IPDPSW), Vancouver, Canada, p 69 - 78 (2018) [[10.1109/IPDPSW.2018.00019](https://doi.org/10.1109/IPDPSW.2018.00019)]

# xPic – Strong Scaling Behaviour

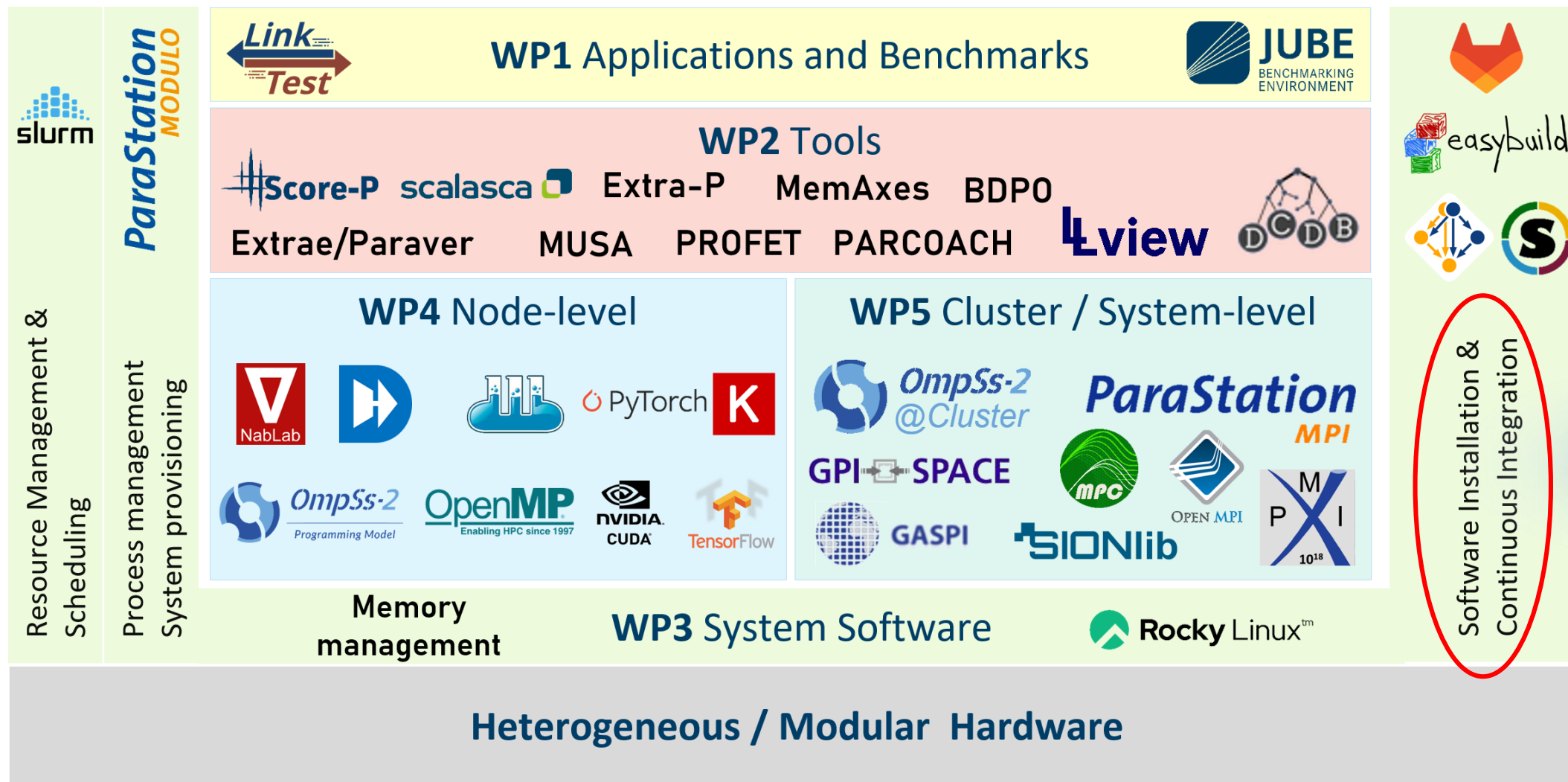
Variable-ratio modular strong scaling



- JSC Jureca system – Intel® Xeon® plus Intel® Xeon Phi™ (KNL)
- Code portions can be scaled-up independently
  - **Particles** scale almost linearly on **Booster**
  - **Fields** kept constant on the **Cluster** (4CNs)
- A configuration is reached where same time is spent on Cluster and Booster
  - Additional 2x time-saving can be reached by co-scheduling “matching” xPic jobs

#cells per node	36864
#particles per cell	1024
#blocks per MPI process	12, 32 or 64

# Integrated Exascale-Ready SW Stack

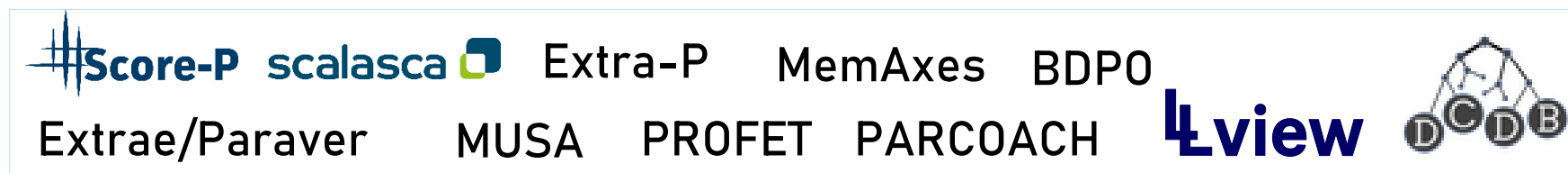


At the heart of the JUPITER system

Public release at <https://gitlab.jsc.fz-juelich.de/deep-sea/wp3/software/easybuild-repository-deep-sea>

# Optimisation Cycles

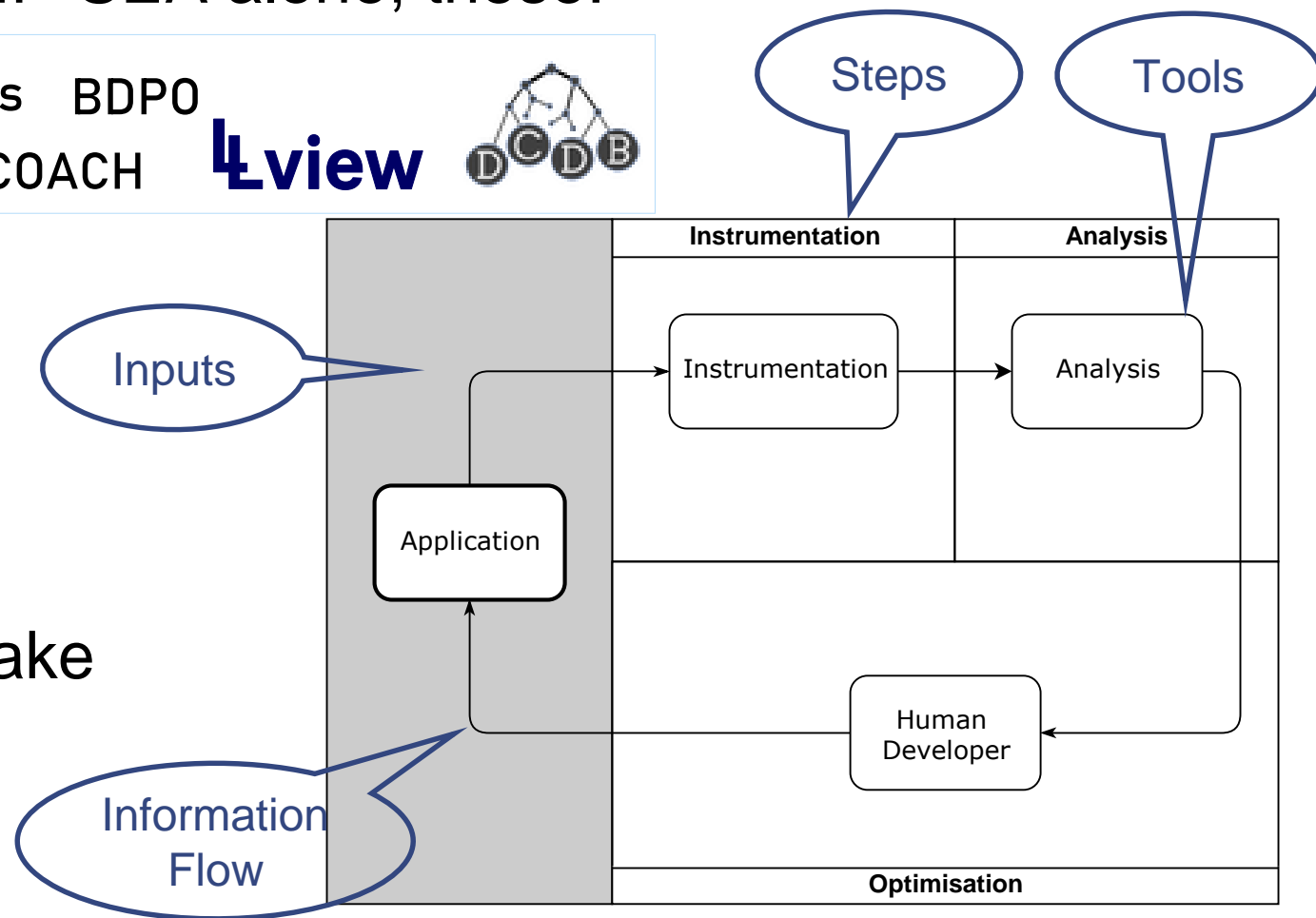
Bewildering variety of SW tools available to HPC SW developers for analysis and optimisation – in DEEP-SEA alone, these:



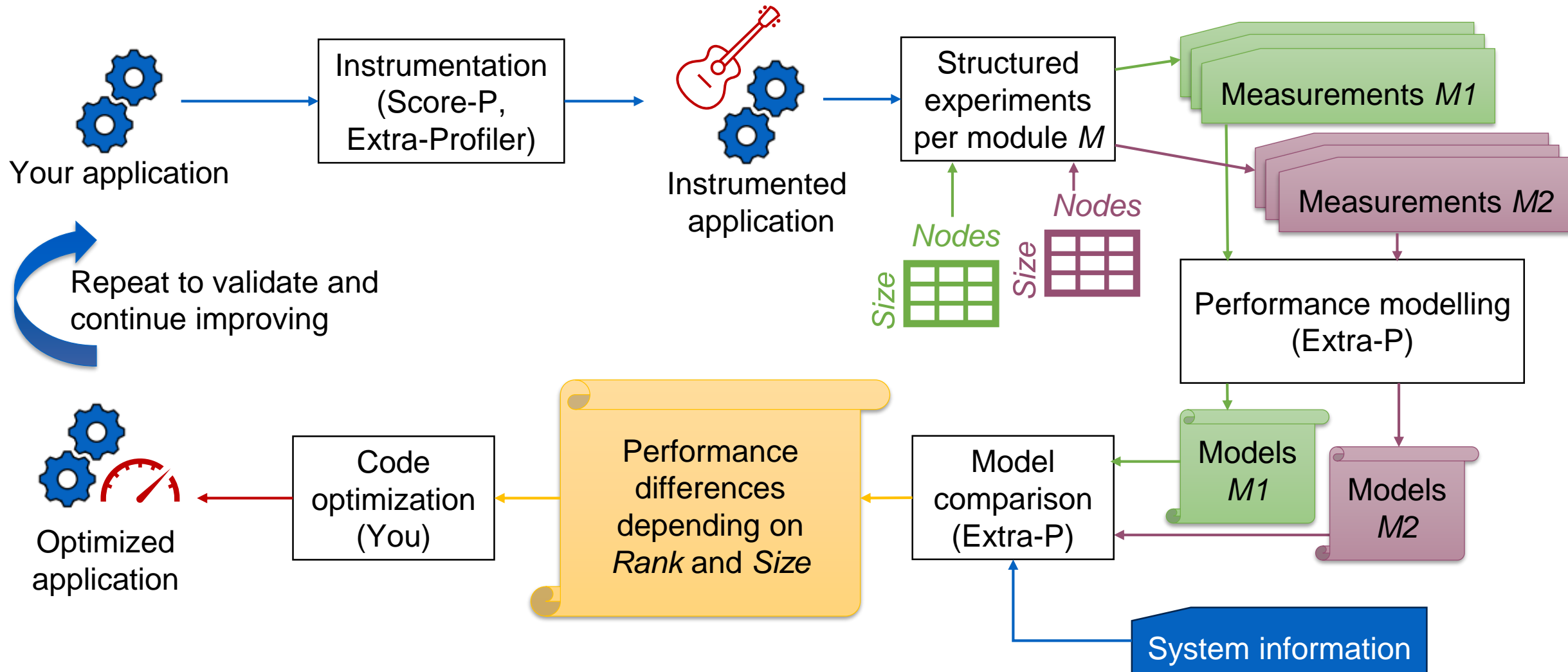
Optimisation cycles encapsulate (complex) tool workflows for *specific purposes*

- Like assessing load balance or optimising energy use

They guide SW developers and make it easier to achieve specific goals



# Application Mapping Optimisation Cycle



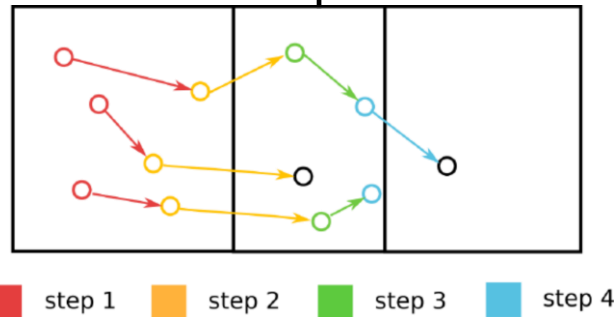


# Use Case: PATMOS

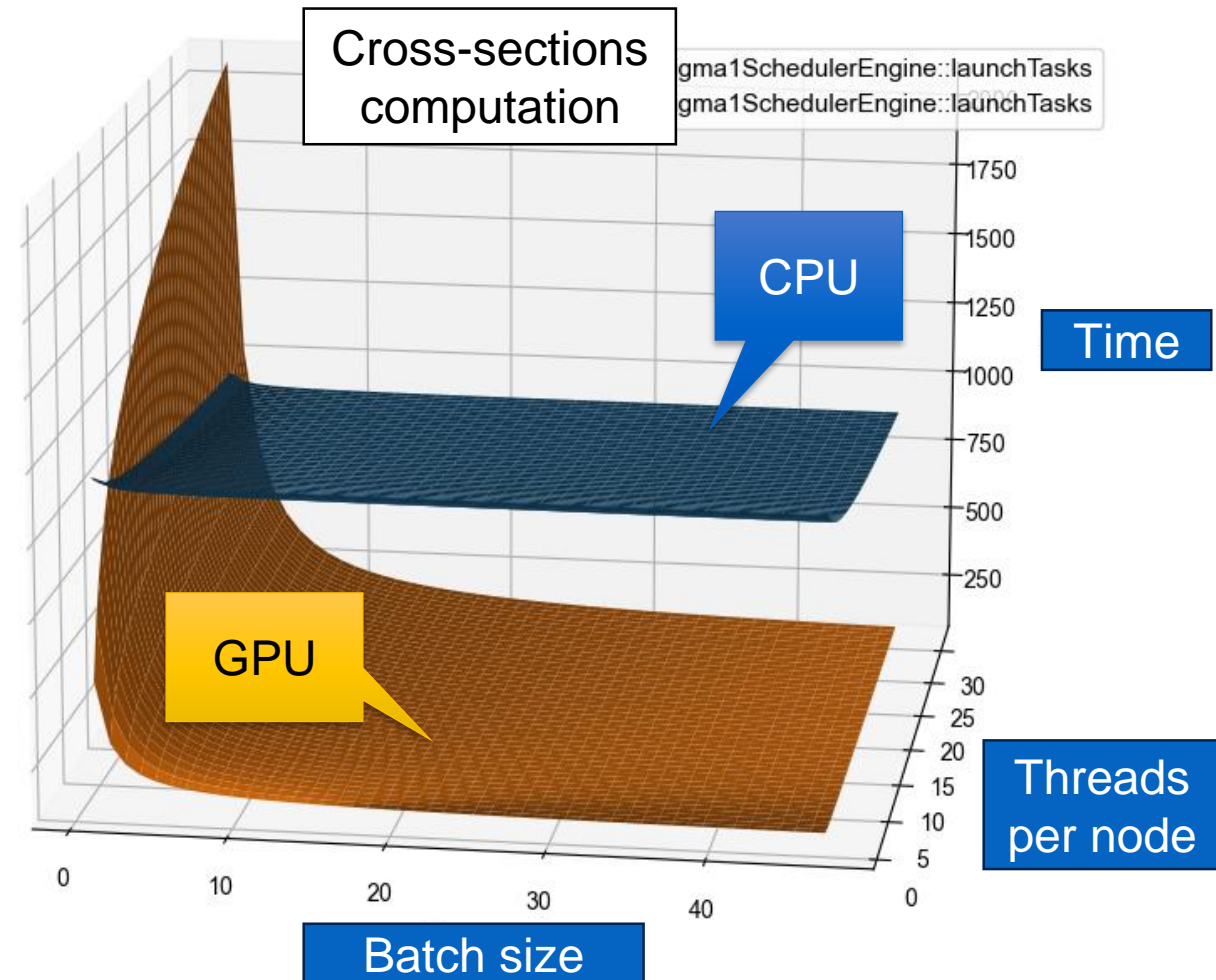
Solves the neutron transport equations to simulate evolution of physical quantities for complex systems

Cross-sections computation represents 60% to 90% of total runtime

- Porting cross section computation to GPU
- Offload batch-size particles at a time



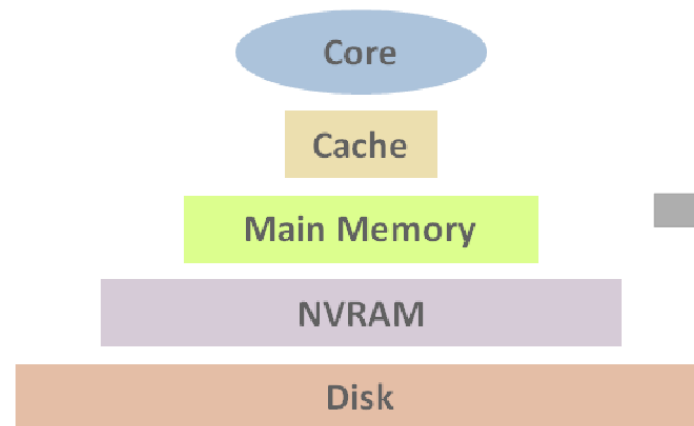
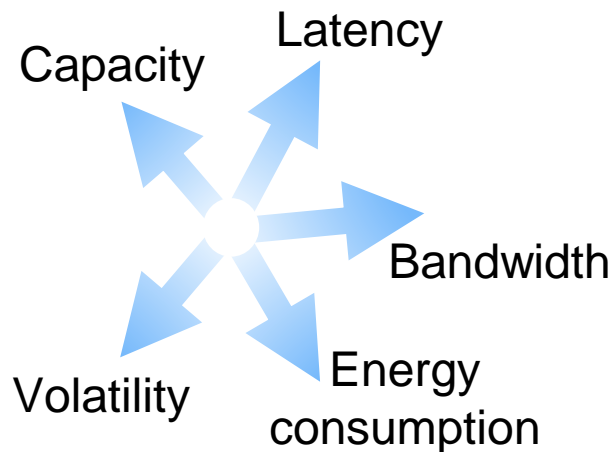
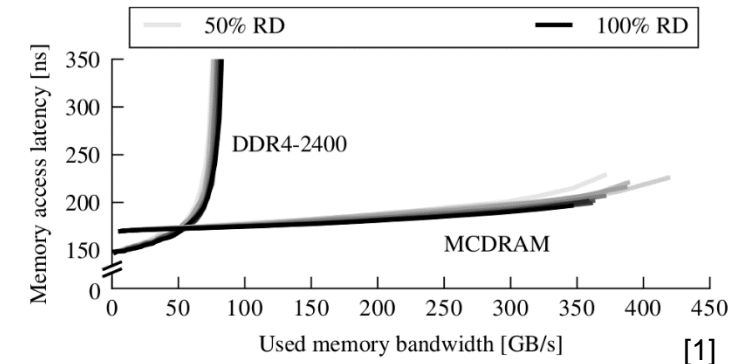
Split of application depends on batch size



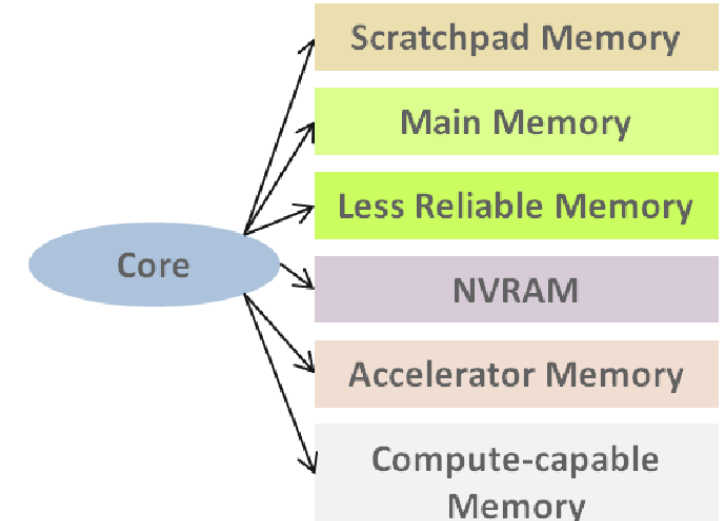
# Heterogeneous/Hierarchical Memory

## Examples...

- DDR DRAM
- Scratchpad (Embedded systems-on-chip, GPUs)
- High bandwidth memory (Intel Xeon Phi, GPUs)
- Byte addressable non-volatile memory (HP's Machine, Intel Optane)
- Compute Express Link (CXL): high-speed interface to accelerators and memory modules



**Memory hierarchy**

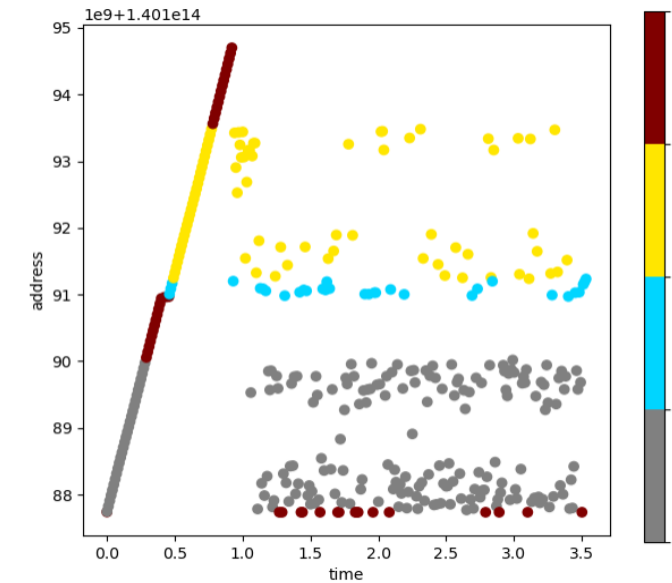
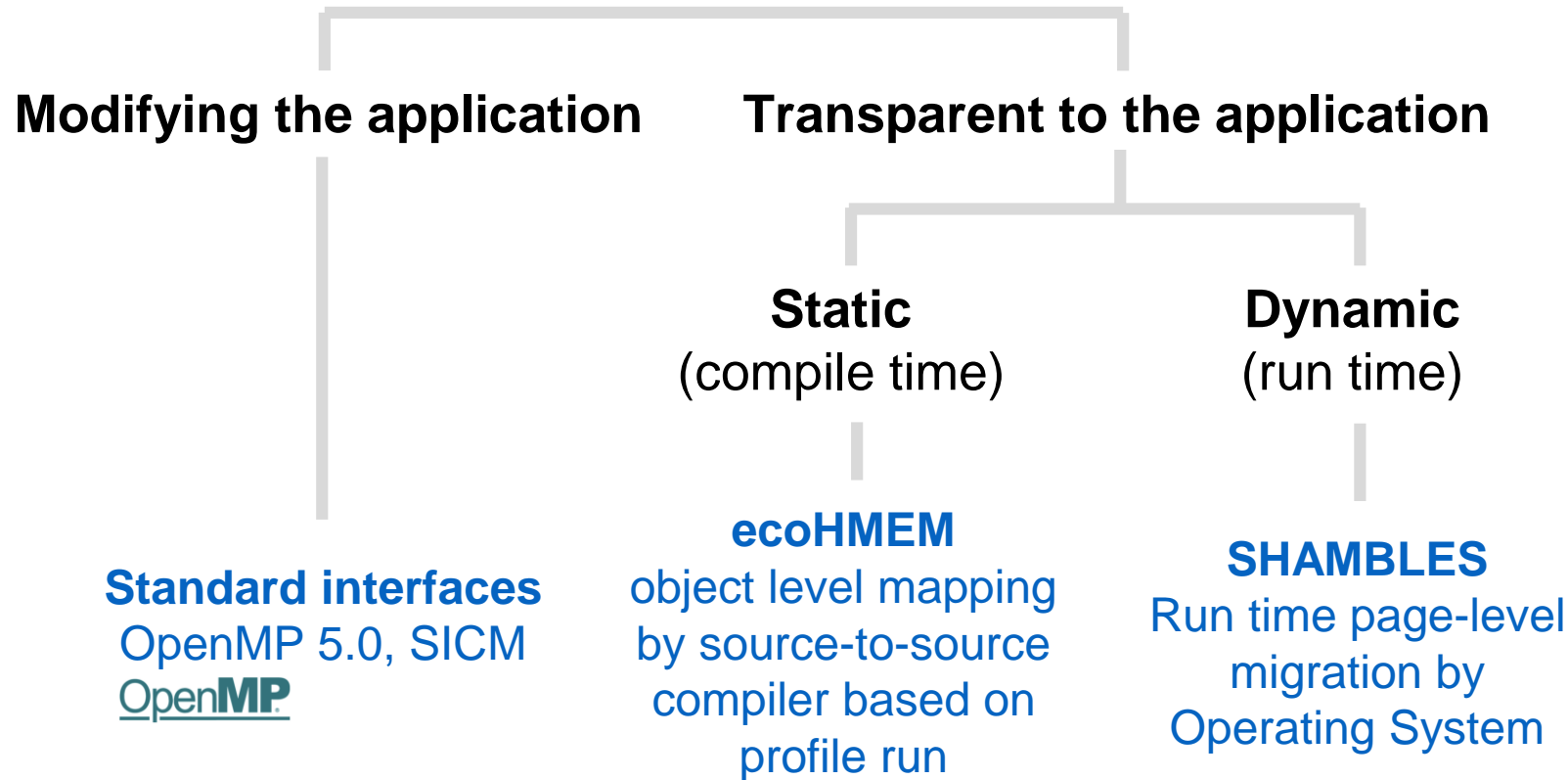


**Explicitly managed**

[1] Milan Radulovic et al. PROFET: Modeling System Performance and Energy Without Simulating the CPU. ACM SIGMETRICS 2019

# Heterogeneous/Hierarchical Memory Tools

- To which degree do the applications need to be modified?
- Which layer manages the memory? When?
- How much can the applications benefit?



**SHAMBLES scatter plot  
example for sparse kernel**

# Malleability

Usual HPC workload resource reservation  
(constant # cores or nodes over time)

Actual use of resources varies over time  
(yellow curve)

Workload is able to use more  
resources in certain phases (arrow)

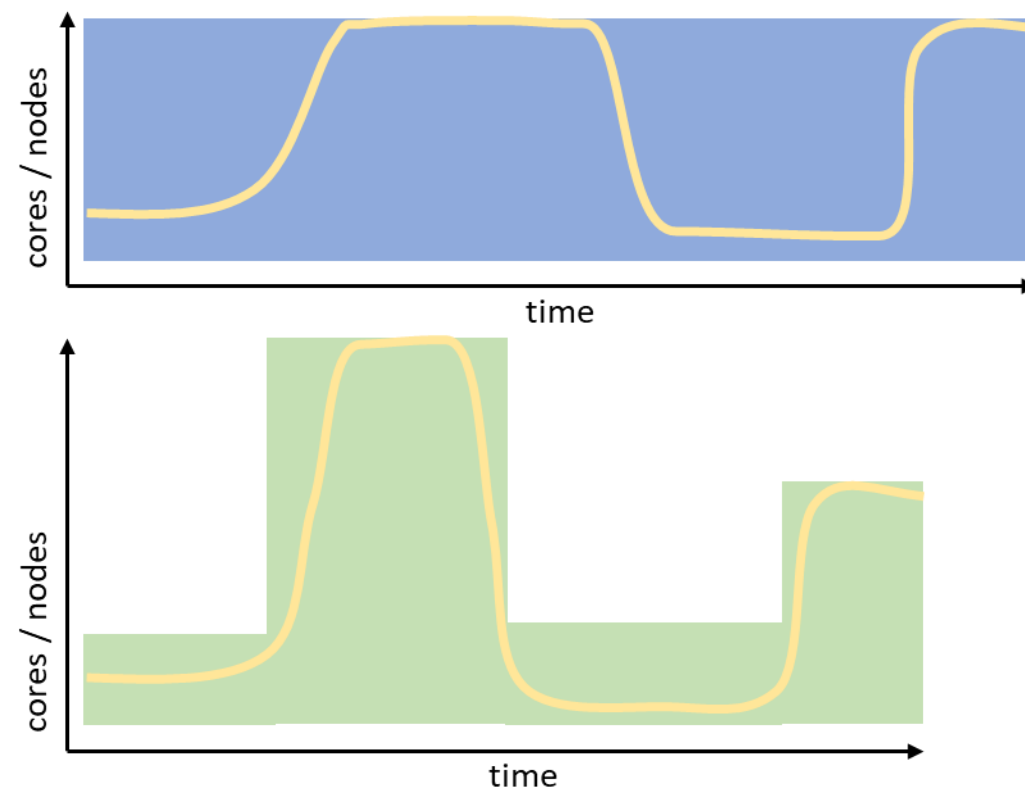
Ideal resource allocation for the workload in  
green

Malleable applications

- Release resources not required
- Acquire more resources if advantageous

Change in # of nodes do require data  
redistribution in the workload

DEEP-SEA provides MPI & Slurm prototypes for  
enabling application-driven (active) malleability

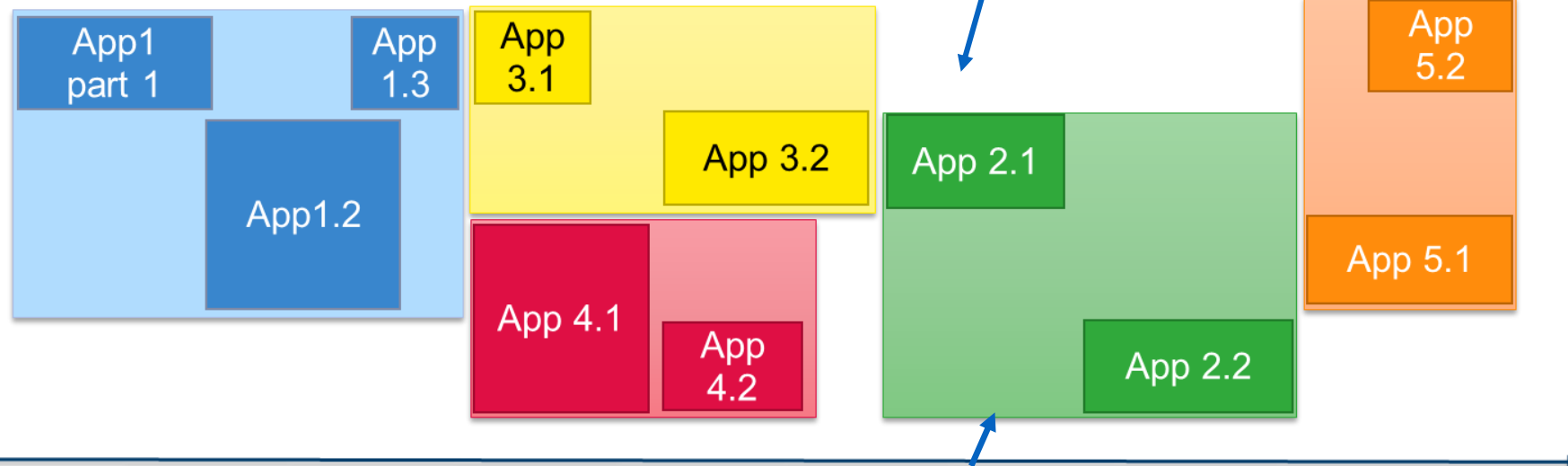
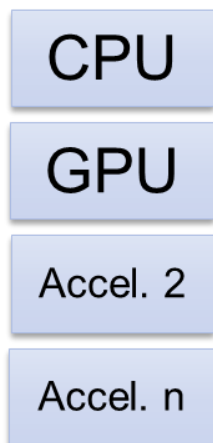


# Scheduling



Current  
behaviour

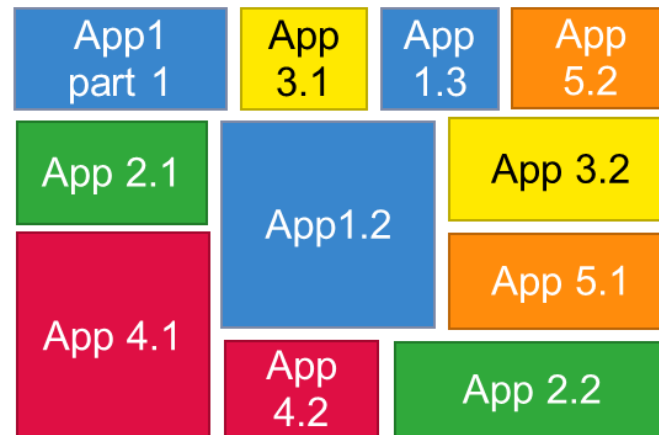
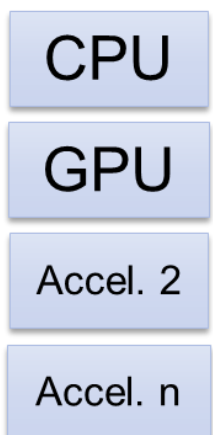
resource



WIP

Ideal  
behaviour

resource



Resource  
reservation window

time

Application part

=

Resource reservation window

time





# Funding Acknowledgement



**EuroHPC**  
Joint Undertaking

SPONSORED BY THE



Federal Ministry  
of Education  
and Research



**ΓΓΕΚ**  
ΓΕΝΙΚΗ ΓΡΑΜΜΑΤΕΙΑ  
ΕΡΕΥΝΑΣ ΚΑΙ ΚΑΙΝΟΤΟΜΙΑΣ



Financiado por  
la Unión Europea  
NextGenerationEU



Swedish  
Research  
Council



The DEEP Projects have received funding from the European Commission's FP7, H2020, and EuroHPC JU Programmes, under Grant Agreements n° 287530, 610476, 754304, and 955606. The DEEP-SEA project receives support from Belgium, France, Germany, Greece, Spain, Sweden, and Switzerland



[www.deep-projects.eu](http://www.deep-projects.eu)



@DEEPprojects