# The International Post-Exascale (InPEx) Project

SC23, Birds of a Feather
November 15, 2023 – rooms 301-303

Prof. Jack Dongarra, UTK

Prof. Pete Beckman, ANL

Dr. Jean-Yves Berthou, Inria

Prof. Satoshi Matsuoka, Riken R-CCS

Dr. Sergi Girona, BSC-CNS

Prof. Bernd Mohr, JSC

**Dr. Emmanuel Jeannot, Inria**

# The International Post-Exascale (InPEx) Project

**InPEx expected outcomes**

- Formation of a solid network of exascale computing leaders, all around the globe
- Landmark documents largely exploited, worldwide, for supporting future post-exascale science
- Contribute to the implementation of an international, shared, high-quality computing environment based on the principles and practices of co-design

**Actions**:

- International Post-Exascale (InPEx) workshop series
- Dedicated international working groups

**Participants:**

Researchers, engineers (comp. science, math, application domains), HW&SW, industry, funding bodies

**How to contribute?** Send 2 pages white paper to: inpex@inpex.science

# Feedback from the subgroups of the pre-meeting

- Meeting in Reims (France), October 19/20, 2023 :
  https://numpex.irisa.fr/international-collaborations-and-inpex-workshops/

- 6 subgroups :

  1. **Software production and management:** packaging, documentation, builds, results, catalogs, continuous integration, containerization, LLVM, parallel tools and sustainability. – Bernd Mohr (JCS), Bruno Raffin (Inria)

  2. **HPC/AI convergence:** ML, LLM for science, open models and datasets for AI training – Pete Beckman (NAISE), Jérôme Bobin (CEA)

  3. **Energy and environmental impact and sustainability** – Michèle Weiland (EPCC – University of Edinburgh), Georges Da Costa (IRIT)

  4. **Future and disruptive SW & HW technologies and usages (including accelerators)**: roadmaps, adoption... – Jack Dongarra (Univ. Tennessee), Jean-Yves Berthou (Inria)

  5. **Co-design, benchmarks/mini-Apps/Proxy and evaluation (HW & SW & Applications)** – Jean-Pierre Vilotte (CNRS), Masaaki Kondo (Riken CCS)

  6. **Digital Continuum and Data management** – Francesc Lordan (BSC), François Bodin (Univ. Rennes)

We neet to focus on the « HOW ». We want concrete outcome: science and technology (not only problems)

# 1. Software production and management

**Context:**
- Machines are getting more heterogeneous (CPUs, GPUs, mem. hierarchy, storage, network), and thus more difficult to program
- HPC software stack is getting more complex, build from an assembly of different components (HPC+HPDA+AI)
- HPC software is expected to interconnect with non-HPC components for building workflows in the Compute Continuum

**Problematic(s):**
- Development productivity
- Compilation and deployment
- Performance testing/portability
- Reproducibility
- Common software stack for HPC+HPDA+AI ??

# 1. Software production and management

**Conclusions:**

- Make it easier for exisiting users to share software products

- Can me make spack / module / easybuild / guix-hpc interoperate?

- Collaborate on container sharing (if software stack somewhat converge) and test within CI/CD pipelines

- Survey AI stacks  and ensure we have a compatible software environment

- Increase productivity for users on our systems, situation is different than for AI as we have more variety

**Main Actions:**

- Survey on software tools/libs used in the different groups

- Create a mailing list (using emails of people in this working group room).

- Share working group spreadsheets (US/RIKEN working group list).

- Share with European colleagues and identify additional people to join collaboration, focus on action items above).

- Collaboration within the high performance software foundation (https://hpsfoundation.github.io/)?

# 2. AI for science / Science for AI

**AI for science (HPC)**

- As an "accelerator" (solvers, surrogate models, in-situ data analysis, code coupling etc.)
- A game-changer in applications (e.g. very fine-grain simulations, digital twins, inverse design, HPDA)
- AI-centric HW for HPC
- Dev. Productivity (e.g. LLM for code generation/conversion)
- Other impact ?

**Science (HPC) for AI**

- HPC data/compute workflows
- HPC tools for large model training/meta-learning (e.g. solvers, IO, task scheduling, etc)
- Energy/performance measurement/profiling
- Other impact ?

**Problematics**

1. **AI models for HPC**
    1. Common models for HPC applications ?
    2. Model validation/robustness/trust
    3. Flexible/extendable models
2. **How HPC can benefit from AI-centric dev ?**
    1. HW convergence, WS stack convergence ?
    2. Data everywhere in current/future systems/applications, can we benefit/reuse parts of the AI-centric tools ?
    3. Composability: e.g. interfacing with AI frameworks (e.g. pytorch, tensorflow, etc.)

# 2. AI for science / Science for AI

**Action plan:**

- The hard work is going from raw data to usable AI data and model
  - **Each InPeX community provides: Fully explained: open code, published data, from start to finish, presents their work at future workshops**
- **"The Pile" for Science**: large training data set available for foundation mode
  - Work together to share large, open scientific datasets for training and testing
  - **Each InPeX community ADDS their dataset to a unified, large set useful for building LLMs or multimode GPT**
- We must publish data in the form suitable for the AI community
  - We can identify and publish challenge problems
- Large foundation models for science
  - TPC (https://tpc.dev)
  - Many groups are working on building their own small models, specialized for their community
  - FugakuGPT and AuroraGPT. The EU strategy on AI4Science must be clarified
- Can we make strong statements about *Scientific AI* (reproducibility, trustworthy, explainable)?
  - **Development of shared concepts and language for discussing and comparing**
- Coordinated, shared path to connect with broader AI community
  - Clear explanation that international AI infrastructure is available in national investments
  - Better organize linking of AI and Science communities
  - HuggingFace, Allen Inst, etc.

# 3. Energy and environmental impact and sustainability

**Context:**

- Several levels: Hardware (including datacentre); SW Stack; Applications.
- Available leverages:
  - Conventional: Improve hardware, software, compilers, numerical libraries, scheduling, dynamism (tasks-based applications)
  - Unconventional: reconfiguration of applications, of hardware; power capping, approximate computing; cross-system scheduling
- New constraints: curtailment; $CO^2$ and energy reporting and budget.

**Problematics:**

- Power and Energy reduction
- $CO^2$ impact
- Cost for users
  - At the cost of performance?
  - At the cost of usability / portability / code sustainability?
- Education: even fundamental understanding is lacking

# 3. Energy and environmental impact and sustainability

**Conclusion:**

- Priority is science, and the goal is to optimize $CO^2$ per « *Nobel Prize* »

**Actions:**

- [Lobbying] Clarify the political expectation and their societal impact
- [Tool] Providing feedback to users: Eq. $CO^2$, Wh, up to abnormal behavior for certain libraries
- [Workshop] Discussion on metrics: how to improve Green500, etc.
- [Workshop] Session on success stories and actual failures of pre-exascale operators
- [Workshop] Challenge for students/researchers: use actual application on a datacenter to reduce power consumption while keeping the same/acceptable performance
- [Workshop] Processor technologies - e.g. reducing standby power consumption
- [MOOC] User education

# 4. Future and disruptive SW & HW technologies and usages

**Problematics**

- What technologies **available today** are disruptive?
- What **future** disruptive technologies ?
- **Other impact** of disruption?
- For what **uses, for what impact?**

**Opportunities/game changer and threats**

- IA and HPC (IA-code generator, IA-solver, …)
- Disruptive Hardware Technologies  (Chiplet, Neuromorphic Computing, Optical Computing, DNA Storage and Computing, Graphene-Based Processors, Silicon Photonics , Quantum Computing )
- Disruptive Software Technologies (Containers, Quantum Computing Software ?)
- Disruptive trends for data management (Emerging storage technologies, disaggregated memory (CXL), compute (FPGA) and storage resources, Leverage fine-grain I/O monitoring information, energy limit for the power consumption, Risk: some storage technologies could be stopped (e.g., Intel Optane))

# 4. Future and disruptive SW & HW technologies and usages

**Action Plan**

**Action 1**, InPEx workshops, applications + AI:

- Invite different communities to exchange on opportunities and results obtained using IA
- Use ambassadors who are already convinced to spread the message

**Action 2** (longer term): organize domain scientific challenges (climate, astro, bio, ….) where AI might be a game changer

**Action 3:** Quantum computing: share access to our infrastructures in EU/Japan/US, testing different implementations of quantum facilities, experimenting different programming models, mathematical libraries, training, …

**Action 4.** Identifying game changer among possible disruptions (disruptions that change practices and impact)

    **Action 4.1** InpEx workshop, organize a session dedicated to

- Chiplets
- DNA storing is coming (may replace tape rapidly)
- Photonics
- Quantum algorithms (coupled infra)
- Disagregation memory
- Goal: identify opportunities it opens, new practices, impact

    **Action 4.2** Funding, support exploratory projects with funding

# 5. Co-design, benchmarks/mini-Apps/Proxy and evaluation

**Context:**
- Exascale Computing Project (ECP)
- Fugaku & Fugaku NEXT co-design projects
- Euro-HPC JU initiatives
- ETP4HPC

**Problematics:**
- Comp. sc. eng. applications development methodologies, accuracy & performance portability
- Co-designed software-Stack/Applications
- Proxy-apps / mini-apps suites
- HW & SW Integration, Testing & Profiling tools, Benchmarking specifications

# 5. Co-design, benchmarks/mini-Apps/Proxy and evaluation

**Actions**

3 identified sessions (software stack developers, application developers) for the next InPEx workshop

1. Efficient application developpment at exascale and beyond
   - leveraging existing tools (e.g., MFEM, Lib-Paranuma, Lib-CEED, MAGMA, PetSc, OCCA-API, Kokkos, Raja, …)
   - address gaps and missing functionalities
2. Develop, coordinate and shared application-driven proxy-apps and mini-apps suites
   - identify and share proxy-apps and mini-apps with standardised specifications, performance analysis methodologies, metrics and shared results
   - build shared distributed information system (repositories, GitHub, gitlab ….)
3. Performance portable programming models and abstraction level:
   - International coordination, collaboration in the development of core components in a co-design way
   - Increase awareness and use of performance programming models in CSE applications development

# 6. Digital Continuum and Data management

**Context:**
- Edge low-adoption technology but there is interest
- Cyber-physical systems are being used for data collection
- Digital Continuum is currently being driven by large Cloud Providers

**Problematics**
- HPC centers compete with Cloud Providers
  - What has HPC to offer compared to Cloud Providers?
- Real-time data with real-time processing requirements
- Digital Continuum is a multi-tenant environment.
  - Collected data used with multiple purposes
  - Computing Infrastructure is also shared

# 6. Digital Continuum and Data management

**Actions:**

- Build a continuum of trusted entities:
  - How to trust data exchanges in the continuum?
    - How to ensure sensor data veracity (proof of provenance)
    - How to ensure connections  source / destination
  - How to cloudify HPC/IA services?
- Building a continuum related PoC using multiple infrastructures
  - Find a good candidate application
- Participation of the "architecture" EU definition of the continuum
  - https://eucloudedgeiot.eu/task-forces/architecture-tf3/

# The International Post-Exascale (InPEx) workshop series

**Organization, agenda and funding of Inpex workshops**

- Host country covers all accommodation and food costs, each participant cover their own travel costs.
- Two or three days, keynotes and breakout sessions on specific subjects
- Participants : researchers, engineers (comp. science, math, application domains), HW&SW,  industry, funding bodies
- Reducing $CO^2$ impact: enabling remote participation

**Pre-workshop InPEx, October 2023, Reims, Fr**
https://numpex.irisa.fr/international-collaborations-and-inpex-workshops/

| Date | (10/2023) | 11/2023 | 06/2024 | 06/2025 |
|------|-----------|---------|---------|---------|
| Location | Preparatory phase EU (France) | SC'23 - BOF | Workshop1 EU | Workshop2 Japan |
| Date | 03/2026 | 09/2026 | 06/2027 | 09/2027 |
| Location | Workshop3 US | Workshop4 EU | Workshop5 Japan | Workshop6 US |

**How to contribute?** Send 2 pages white paper to :  inpex@inpex.science