

Towards Integrated Hardware/Software Ecosystems for the Edge-Cloud-HPC Continuum: the TransContinuum Initiative



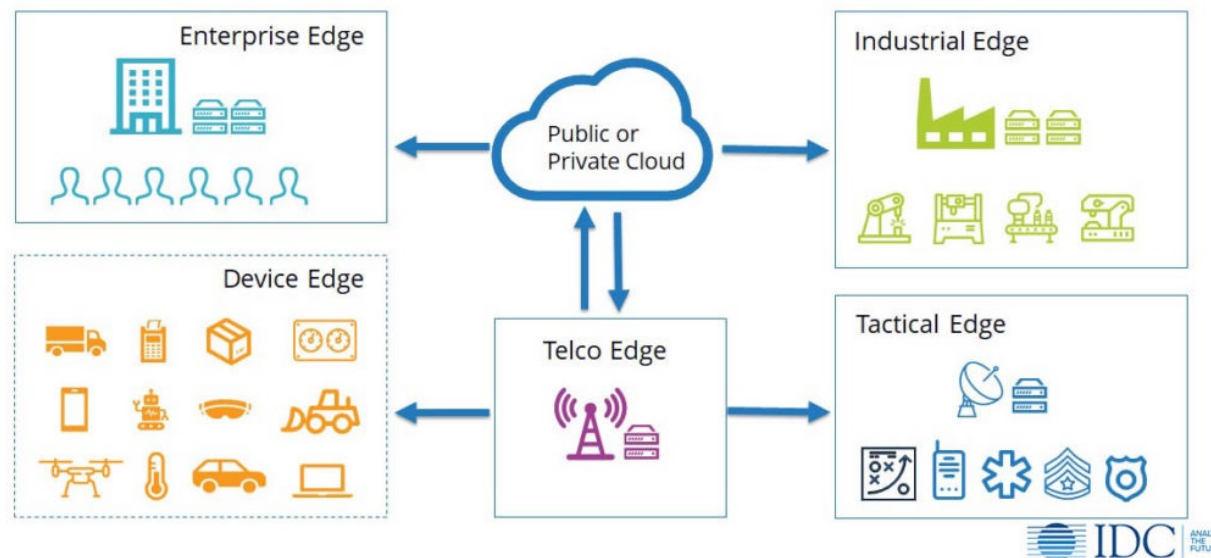
Gabriel Antoniu, Inria

With contributions from Patrick Valduriez, Inria, Hans-Christian Hoppe, Jens Krüger - scapos AG
and from the KerData team at Inria, Rennes

5 March 2023

HPC workshop, Paris

Context: IT Investments Shift to the Edge



IDC predictions:

- In 2022, Enterprise and service provider spending on Edge computing will reach \$40 billion in Europe, and increase with a five year annual growth rate of 16.4%
- By 2023, over 50% of new Enterprise IT infrastructure deployed will be at the Edge rather than corporate datacenters
- By 2024, the number of apps at the Edge will increase 800% (compared to 2020)

Sources: IDC Press release on "IDC Forecasts Double-Digit Growth for European Edge Investments", January 18, 2022
IDC FutureScape: Worldwide IT Industry 2020 Predictions
IDC blog Edge Computing: Not All Edges are Created Equal (<https://blogs.idc.com/2020/06/01/edge-computing-not-all-edges-are-created-equal/>)

Is This a Directional Evolution?



Supercomputer

Cloud

Edge

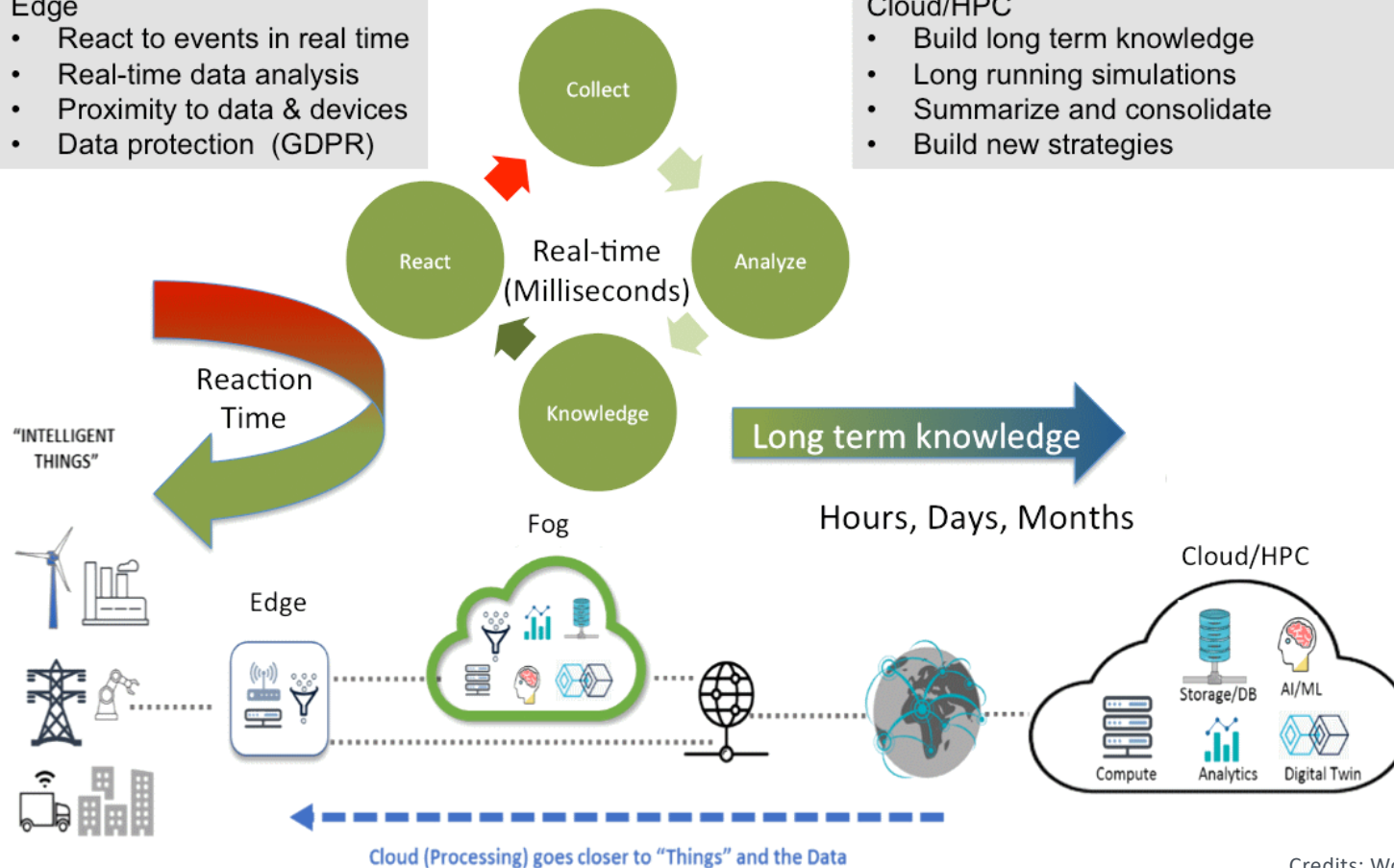
The Computing Continuum: a Rather Circular View for Dynamic, Continuous Workflows

Edge

- React to events in real time
- Real-time data analysis
- Proximity to data & devices
- Data protection (GDPR)

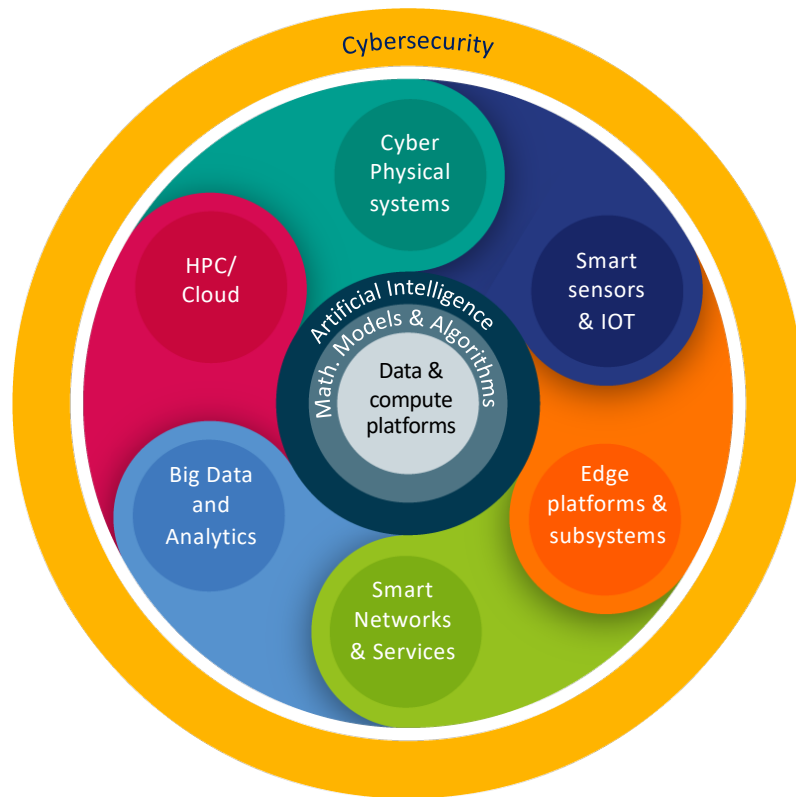
Cloud/HPC

- Build long term knowledge
- Long running simulations
- Summarize and consolidate
- Build new strategies



Credits: Woodside Capital

The Digital Continuum: Cross-area Challenges



Original version courtesy HiPEAC

A continuous dynamic workflow

Between
Smart Sensors and IOT devices
and
HPC / cloud centers
passing through
Edge platforms & subsystems
as well as
Smart Networks and Services
executing
Simulation & Modelling, Big Data Analytics and ML*
based on
Math. Methods & Algorithms incl. MSODE**
pervasively augmented by
Artificial Intelligence
protected and secured by
Cybersecurity
back to
Cyber-Physical Systems,
all based on
Data and compute platforms (hw and sw)

* ML: Machine Learning
** MSODE: Modelling, Simulation and Optimization in Data-rich Environment

TransContinuum Initiative (TCI): our vision

Introduction

This document outlines a vision for a horizontal collaboration between European associations and projects involved in IT technology, application and services provisioning for the Digital Continuum.

The term TransContinuum describes the defining characteristic of the infrastructure required for the convergence of data and compute capabilities in many leading edge industrial and scientific use scenarios. A paradigm change is needed: we will have to design systems encompassing millions of compute devices distributed over scientific instruments, IoT, supercomputers and Cloud systems through LAN, WLAN and 5G networks.

The Digital Continuum



A continuous dynamic workflow

Between
Smart Sensors
and IOT devices at the edge
and
HPC / cloud centers
over
Smart Networks and Services
executing
Simulation & Modelling,
Big Data Analytics, ML*
based on
Math. Methods & Algorithms incl.
MSODE**
pervasively augmented by
Artificial Intelligence
protected and secured by
Cybersecurity
back to
Cyber-Physical Systems

* ML: Machine Learning

** MSODE: Modelling, Simulation and Optimization in Data-rich Environment

Original version courtesy of HiPEAC

Participating organisations



Jean-Pierre Panziera, ETP4HPC chairman



Luigi Rebuffi, ECSO Secretary General



Thomas Hahn, president of BDVA



Colin Willcock, 5GIA chairman



Wil Schilders, president of EU-MATHS-IN



Philipp Slusallek, CLAIRE Co-Founder and Member of the Board



Koen De Bosschere, HiPEAC coordinator



Jürgen Sturm, Chairman of the Management Board



Destination Earth Initiative



Goal: create and operate high precision digital models (digital twins) of the Earth to monitor and predict environmental change and human impact

Approach:

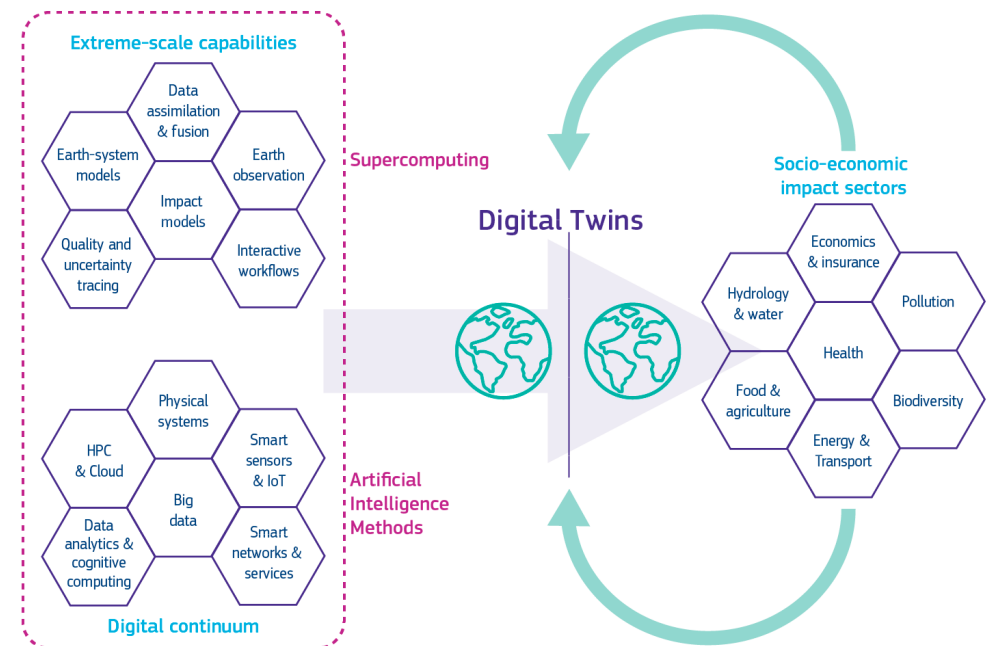
- Simulate atmosphere, oceans/ rivers and cryosphere at very high resolution
- Feed in data from a huge number & variety of sensors
- Monitor & model bio- and anthroposphere

Key aspects:

- Anticipate environmental disasters and resultant socio-economic crises to save lives and avoid large economic downturns
- Enable the development and testing of scenarios for ever more sustainable development
- Near-real time analysis of climate and plant observation to correct simulation scenarios
- Integrate simulations with ML/AI-based analysis/forecasting

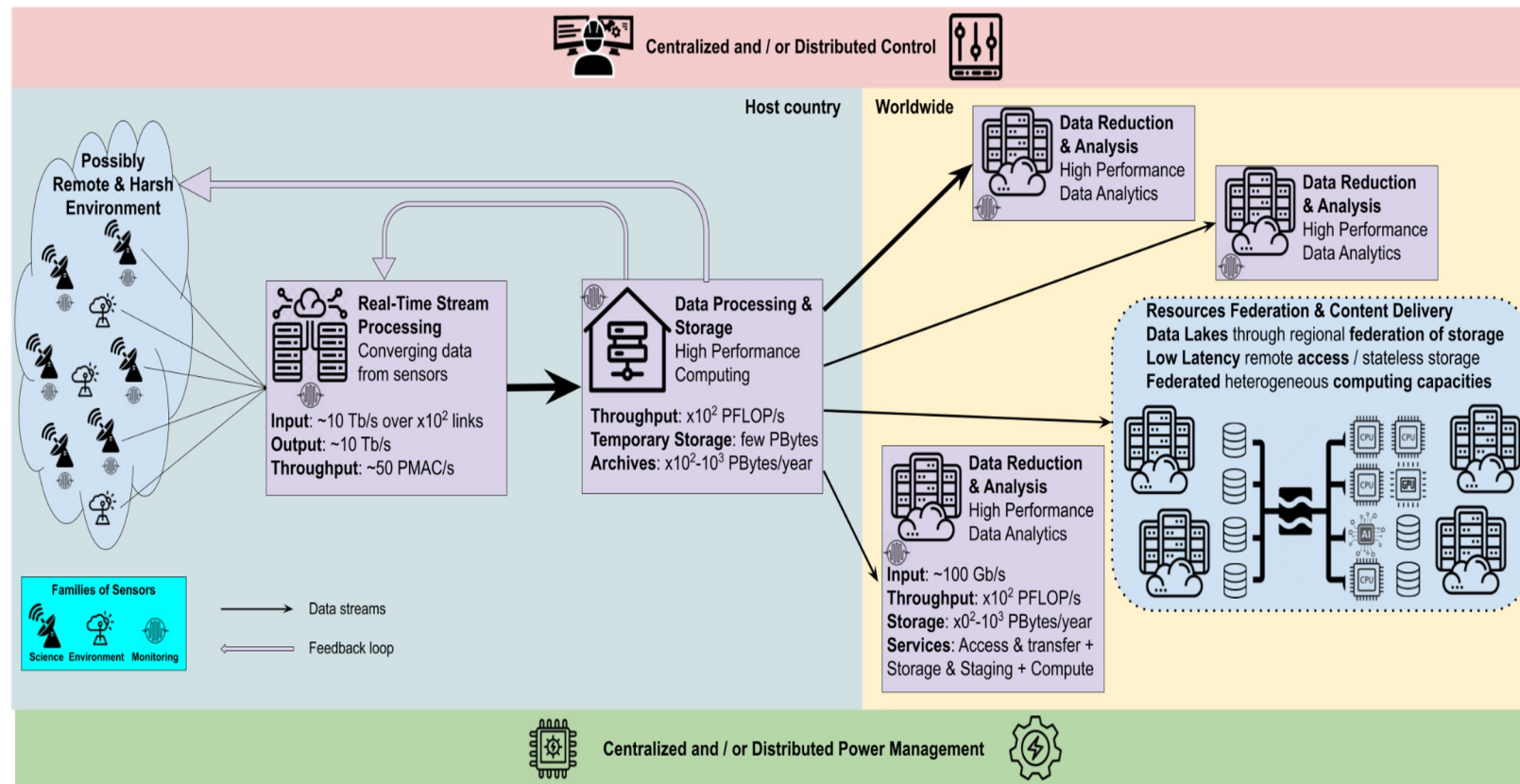
Need:

- Huge requirements for HPC compute and data capability
- Assimilation of diverse sensors and data streams (satellites through mobile phones)
- Disaster avoidance/recovery needs Urgent Computing capabilities



Source: <https://digital-strategy.ec.europa.eu/en/library/destination-earth>

The SKA Data Workflow from Sensors to HPC Centers



A Use Case in Smart Buildings & Cities

Goal: Optimize the energy consumption of buildings

Approach:

- Combine historical operation data with real-time sensor data using ML to predict high-resolution local weather forecasts
- Control building lighting and heating/cooling/ventilation accordingly
- Perform predictive maintenance actions

Key aspects:

- Real-time Edge pre-processing of local sensor data
- Cloud-based analytics
- Further computation-bound simulations may need HPC facilities
- The input data can vary in frequency, relevance, amount

Need:

- Deployment and the execution of a distributed, cross platform workflow from the Edge to Cloud/HPC
- Dynamic reallocations of tasks and resources across the continuum



Where Are We Now?

An archipelago of disconnected solutions

- Separate software stacks (HPC, Big Data on Cloud, Edge Analytics, AI) optimised for different goals, with different infrastructure requirements

What is Difficult?

- Deploy and orchestrate a combination of consistent interoperable components across the full continuum
- The different compute, storage and communication systems of a complex CC installation belong to different owners
 - Authentication
 - Interoperability
 - Heterogeneity
- How to ensure security across the continuum?
- Flexible and efficient operation of CC infrastructures



Building Integrated Software Ecosystems for the Continuum: Challenges

Application-level challenges

- Traditional physics-based simulations (HPC) need to smoothly cooperate with data-driven, learning-based analytics and prediction engines (Cloud)
- Hybrid workflows: programming models, composability

Storing and processing data across the continuum: **how to deal with the 3 Vs of Big Data?**

- **Extreme Volume across the continuum:** support the access and processing of “cold”, historical data and “hot”, real-time data + (virtually infinite) simulated data
- **Extreme Velocity across the continuum:** unified data processing (in situ/in transit, stream-based) in a common software ecosystem
- **Extreme Variety across the continuum:** unified data storage abstractions to enable distributed processing and analytics across the continuum
 - Interoperable data formats
 - "Semantic interoperability" through shared ontologies
 - Storage interfaces should match the needs

Building Integrated Software Ecosystems for the Continuum: Challenges

Managing computation across the continuum

- **Heterogeneity**: wide variety of processors, accelerators, storage devices and systems, and communication systems
- Dynamic **scheduling and orchestration of workflows** which evolve at runtime, to optimize performance and energy
- **Support seamless deployment and migration** of workflow components despite heterogeneity
- Definition and automatic derivation of **performance models**

Managing **dynamic** workflows with ad-hoc load variation

- React to certain events, **depending on data contents** or on **interactive requests**
- **Dynamically adapt** the mapping of the workflow onto the infrastructure
- In some applications (e.g., disasters) parts of the infrastructure suddenly become unavailable
- Requires efficient coupling between Cloud-oriented dynamic orchestrators and traditional batch-based resource management systems, as a step towards more integrated software approaches to dynamic resource management across the continuum

Building Integrated Software Ecosystems for the Continuum: Challenges

AI-related challenges

- New heterogeneity of use cases and hardware
- The deep learning software stacks must be supported (python dependencies handling, containerisation)
- **Ad-hoc training and inference runs with tight timing constraints** must be supported (urgent and interactive computing)

Cybersecurity challenges

- **Federated authentication**, authorisation and accounting, monitoring, resource allocations, encryption, user insulation, container certification, etc.
- GDPR-related constraints: HPC centres must provide all tools necessary to address regulatory requirements.

Cooperation challenges

- Interaction of multiple expert communities (HPC, Big Data, AI, cybersecurity, IoT, 5G, etc.).
- **Establishing commonly agreed, shared goals and priorities**
- Need a common vocabulary and common roadmaps
- This is precisely the core motivation underlying the TransContinuum Initiative (TCI)!

Summary

- Edge computing is impacting the HPC Research agenda!
- Use cases combining HPC simulation, analytics and AI are emerging
 - They require Edge, Cloud and HPC
 - Adopt and evolve the “Digital Twin” approach
 - Targeting scientific, societal and business benefits
- An integrated SW ecosystem spanning across Edge, Cloud and HPC systems is key for sustained success – and it is evolving
 - Programming environment(s) to develop applications combining simulation, AI/ML and data analytics
 - Data storage, transfer, processing, assimilation across the continuum
 - Manage complex computations on large distributed, heterogeneous infrastructures
 - Efficient handling of dynamic, distributed workflows
 - Cybersecurity mechanisms to protect infrastructure & data
 - Requires cooperation across multiple areas

An ETP4HPC White Paper on Edge-Cloud-HPC Continuum Challenges

ETP 4
HPC

EUROPEAN TECHNOLOGY
PLATFORM FOR HIGH
PERFORMANCE COMPUTING

2022

ETP4HPC's SRA 5

STRATEGIC RESEARCH
AGENDA FOR
HIGH-PERFORMANCE
COMPUTING IN EUROPE

SEPTEMBER 2022
**EUROPEAN
HPC RESEARCH
PRIORITIES
2023 - 2027**

Quantum for HPC

Sustainability
the Next Big Thing

HPC in the Digital
Continuum

Industrial Use of HPC

On the path to
Exascale

Heterogeneous
High-Performance
Computing

Unconventional HPC
Architectures

Centre to edge
framework

HPC for Urgent
Decision Making

Federated HPC,
Cloud and Data
Infrastructures

Application
co-design

Towards Integrated Hardware/Software Ecosystems for the Edge-Cloud-HPC Continuum

Supporting integrated applications across the Edge-Cloud-Supercomputer layers to address critical scientific, engineering and societal problems

White Paper

Gabriel Antoniu (Inria), Patrick Valduriez (Inria),
Hans-Christian Hoppe (SCAPOS), Jens Krüger
(Fraunhofer ITWM)

28/09/2021

etp4hpc.eu

@etp4hpc

ETP 4
HPC

EUROPEAN
TECHNOLOGY
PLATFORM
FOR HIGH
PERFORMANCE
COMPUTING

<https://www.etp4hpc.eu/white-papers.html#continuum>



What Are We Doing at Inria About This?



- Reproducible Data Processing, Analytics and AI on the Edge-to-Cloud Continuum
 - E2Clab: Explore the continuum through reproducible experiments
 - Efficient deep learning training on the Computing Continuum
 - Supporting online learning and inference in parallel across the continuum
 - Projects:
 - Inria-DFKI ENGAGE: NExt GeNeration ComputinG Environments for Artificial IntelliGence
 - Joint Lab for Extreme-Scale Computing, UNIFY Associate Team
 - **STEEL project within the CLOUD PEPR National project (2023-2030)**
- Convergence of Extreme-Scale Computing and Big Data Infrastructures through Data orchestration
 - HPC/Cloud convergence: dynamic allocation of storage resources, cross-architecture I/O orchestration
 - A motivating use case: SKA (Square Kilometre Array)
 - Framework: **ExaDoST project within NumPEX (2023-2027)**

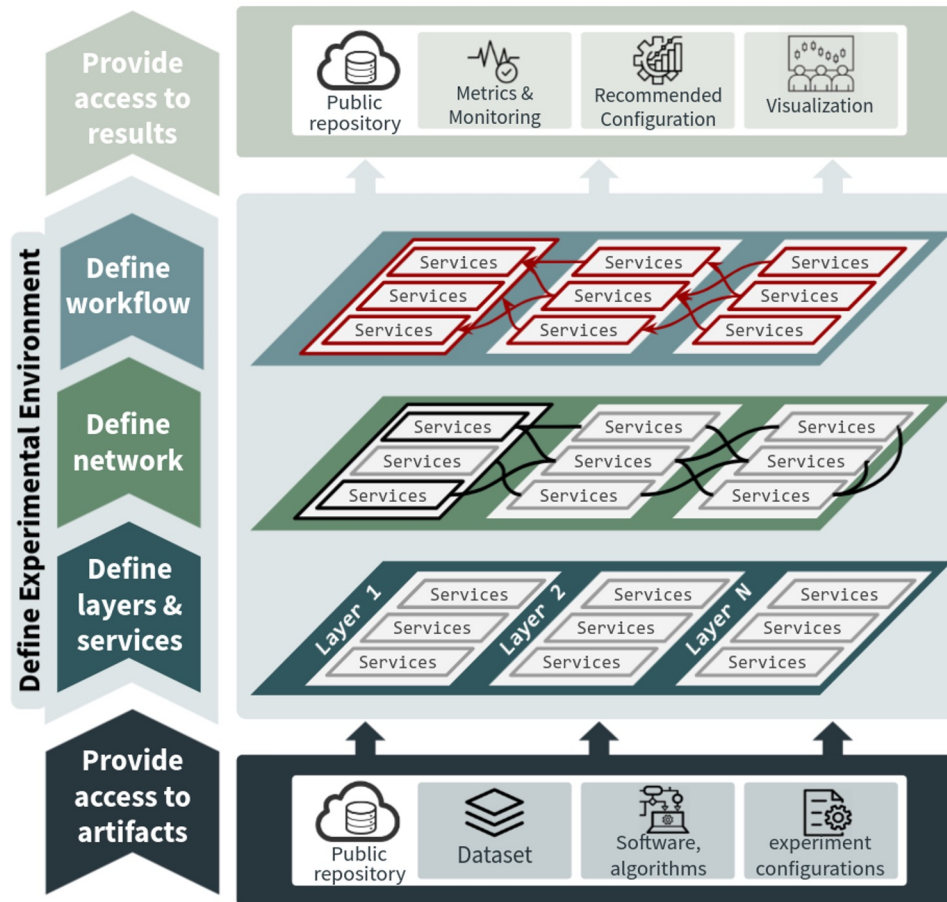
Alexandru
Costan



François
Tessier



E2Clab: Reproducible Performance Optimization of Complex Applications on the Edge-to-Cloud Continuum



Methodology

- ★ **Reproducible Experiments**
Repeatability, Replicability & Reproducibility
- ★ **Mapping**
Application parts & physical testbed
- ★ **Variation & Scaling**
Experiment variation and transparent scaling
- ★ **Network Emulation**
Edge-to-Cloud communication constraints
- ★ **Experiment Management**
Deployment, Execution & Monitoring



IEEE Cluster 2020

<https://hal.archives-ouvertes.fr/hal-02916032>

E2Clab Optimization Methodology



IEEE Cluster 2021

<https://hal.archives-ouvertes.fr/hal-03310540>

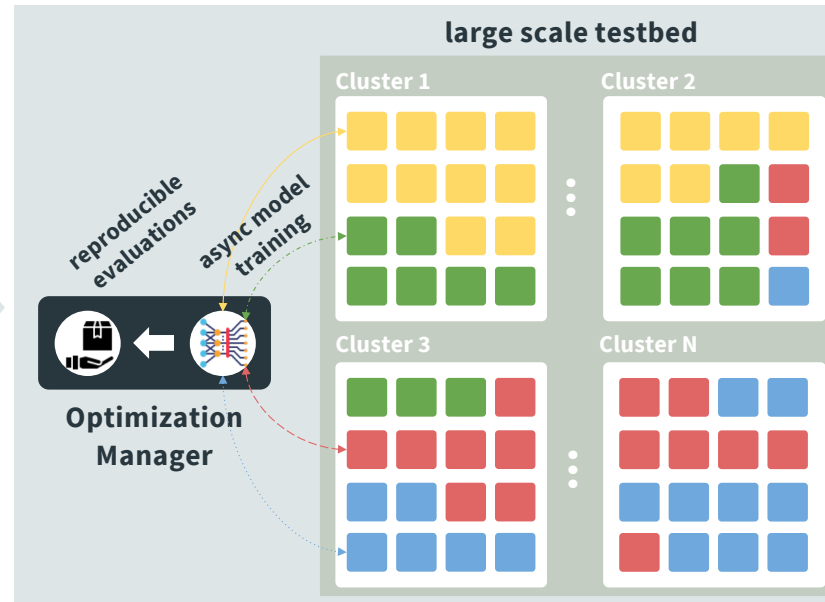
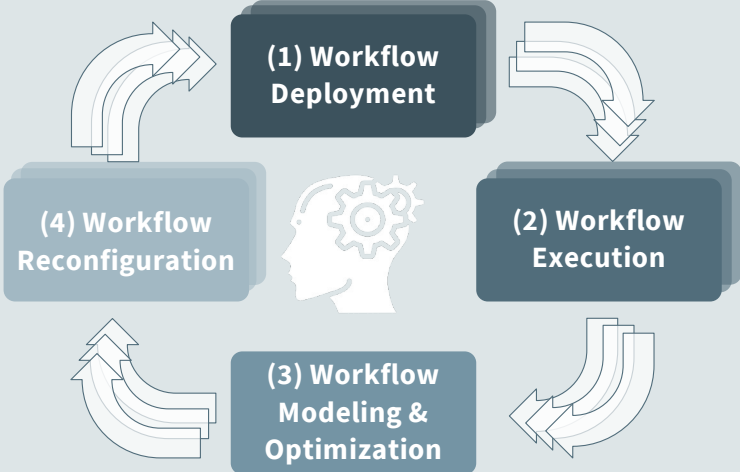
Phase I

Define Optimization Problem

$$\begin{aligned} \min_x & f_m(x), & m = 1, 2, \dots, M \\ \text{subject to} & g_j(x) \leq 0, & j = 1, 2, \dots, J & \text{Inequality constraints.} \\ & h_k(x) = 0, & k = 1, 2, \dots, K & \text{Equality constraints.} \\ & x_i^L \leq x_i \leq x_i^U, & i = 1, 2, \dots, I & \text{Bounds on variables.} \end{aligned}$$

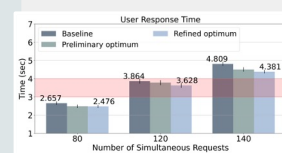
Phase II

Parallel + Scalable + Reproducible
application optimization on large scale testbeds



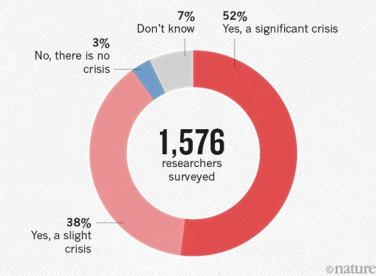
Phase III

Summary of Computations



Reproducible Research

IS THERE A REPRODUCIBILITY CRISIS?

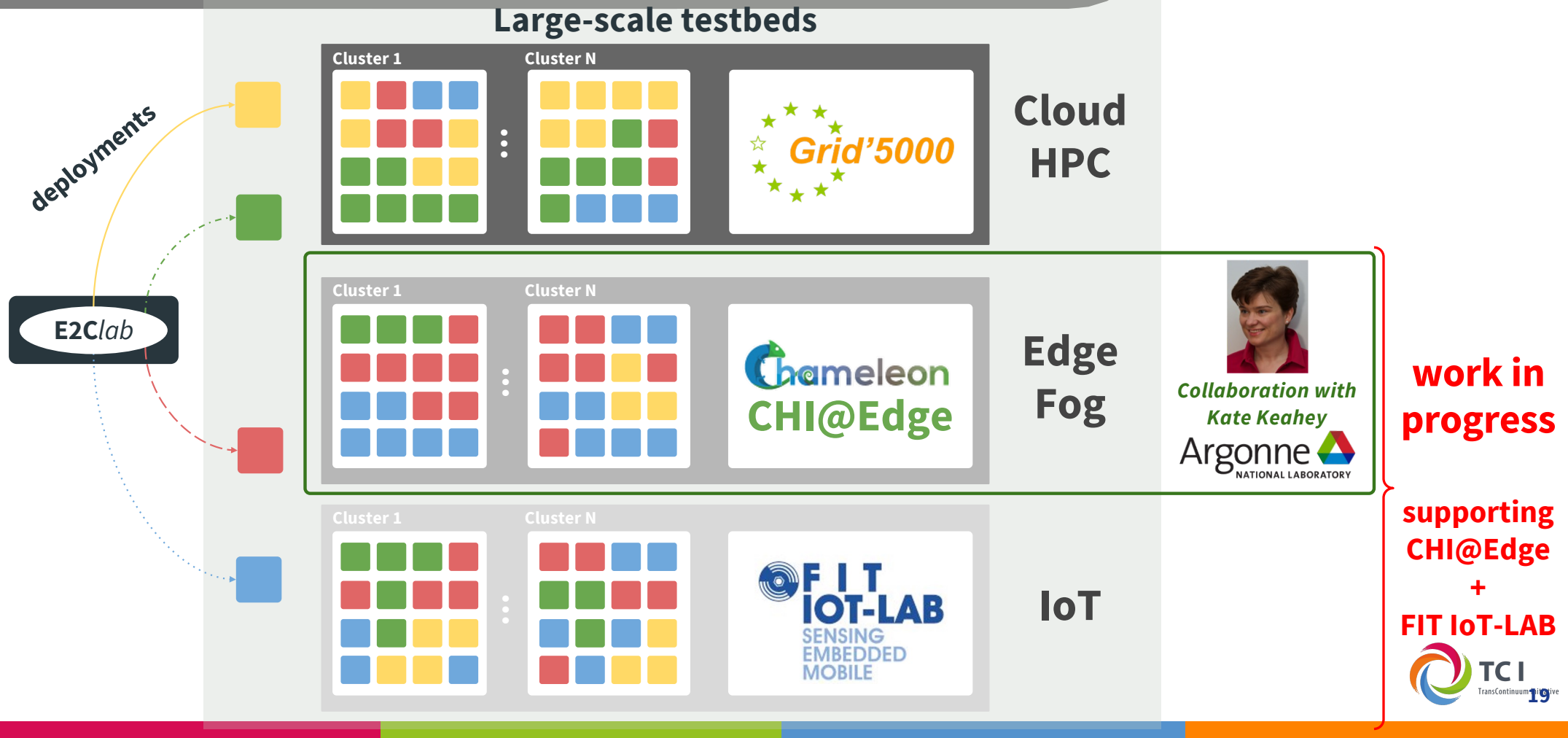


source: <https://www.nature.com/news/1.19970>

“+70% failed to reproduce another scientist's experiments”

“+50% failed to reproduce their own experiments”

Ongoing work: From IoT/Edge to Cloud/HPC environments





Thank you !
