# EVEREST

Computing Frontiers 2022, BigDaW Session

# The EVEREST SDK

## Mapping big-data applications onto heterogeneous reconfigurable computing systems

**CHRISTOPH HAGLEITNER**

*IBM Research, Zurich; EVEREST Project Coordinator*

hle@zurich.ibm.com

# Introduction

EVEREST     dEsign enVironmEnt foR Extreme-Scale big data analyTics on heterogeneous platforms

SDK        System Development Kit = Tools for
- Application description
- Deployment on Target System
- Compilation
- Runtime environment
- Data management and security

EVEREST

# EVEREST Consortium

**IBM Reseach Lab, Zurich (Switzerland)**
Project Administration, Prototype of the target system

**Università della Svizzera italiana (Switzerland)**
Data security requirements and protection techniques

**Centro Internazionale di Monitoraggio Ambientale (Italy)** Weather prediction models

**Virtual Open Systems (France)**
Virtualization techniques, runtime extensions to manage heterogeneous resources

**Numtech (France)**
Application for monitoring the air quality of industrial sites

**Politecnico di Milano (Italy)**
Project Administration, High-Level System, Flexbile Memory Manager, Autotuning

**TU Dresden (Germany)**
Domain-specific extensions, code optimizations and variants

**IT4Innovations (Czech Republic)**
Exploitation leaders, Large HPC infrastructure, Workflow libraries

**Duferco Energia (Italy)**
Application for prediction of renewable energies

**Sygic A/S (Slovakia)**
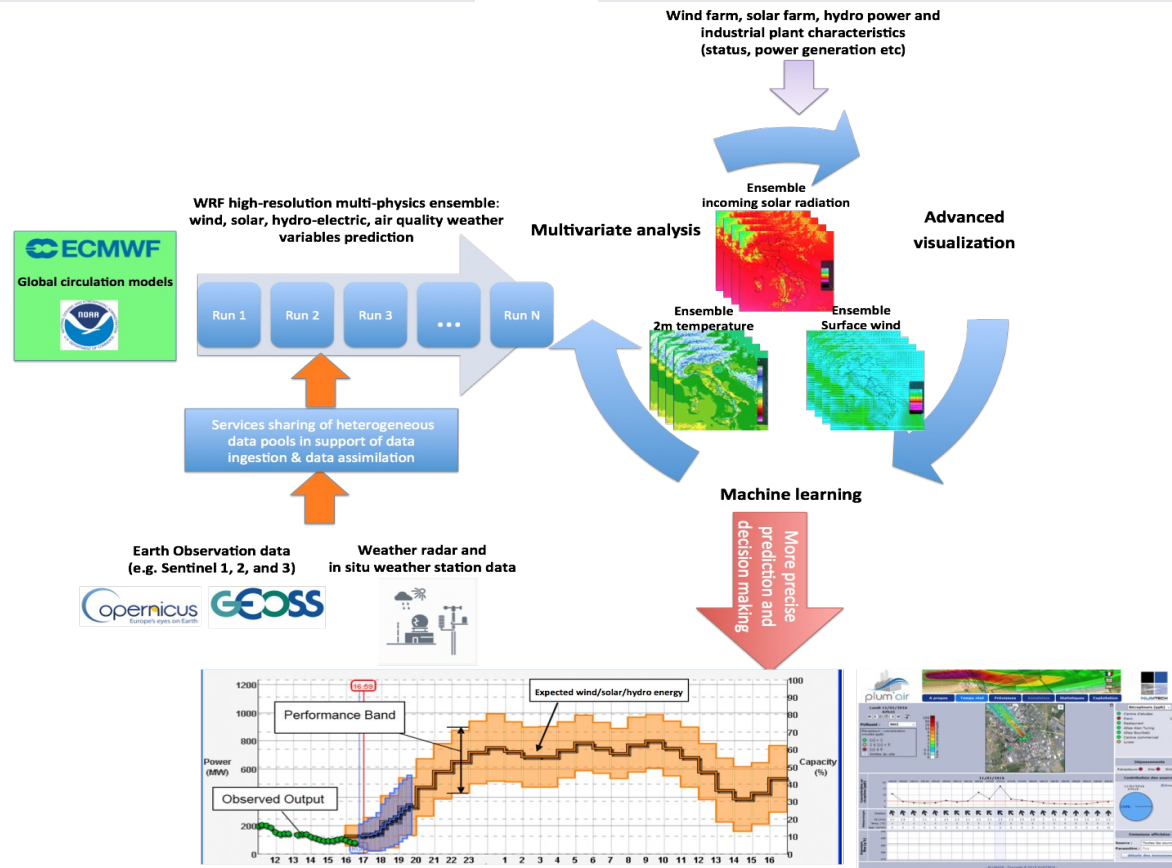Application for intelligent transportation

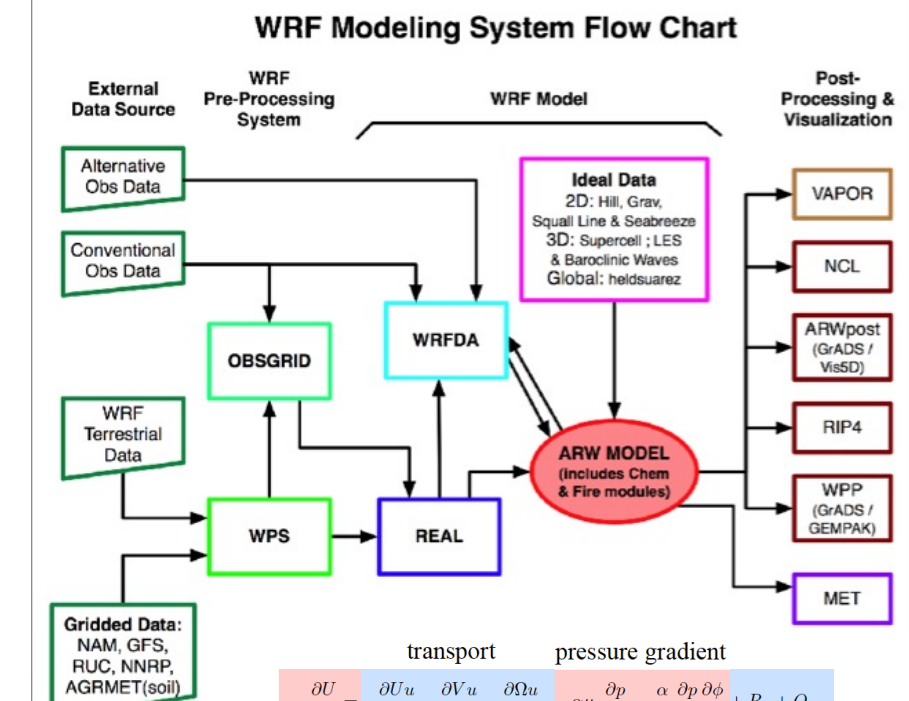EVEREST

# EVEREST Use Cases

EVEREST

# The WRF Model

## First step of two use cases…



WRF is an open-source model supported primarily by the US National Center for Atmospheric Research (NCAR), the US National Oceanic and Atmospheric Administration and the US National Center for Environmental Prediction – NCEP
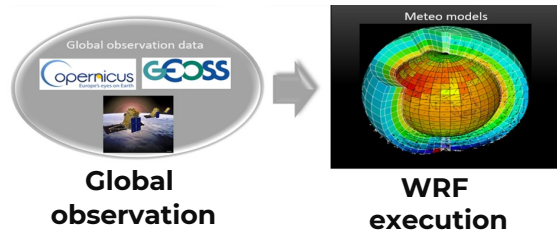


WRF Modeling System Flow Chart

# Air-quality use case: Workflows and Challenges

## 1. WRF Deterministic weather forecast



Global observation

WRF execution

**Improve speed to produce forecast**

## 2. Ensemble prediction

### N x deterministic weather forecast



GFS
* 4 cycles

METEO France
* 4 cycles

NUMTECH
* 2 cycles

EVEREST

**+**

Local weather observation on-site



**ML**

One aggregated weather forecast forced by observation

**Improve quality of local weather forecast**

## 3. Air-quality dispersion forecast



Landuse, Topography

Industrial site data

Emission forecast

Local weather forecast



Air-quality forecast

**Improve speed to produce air-quality forecast and its quality**

**EVEREST**

# Renewable Energy use case: Context and Challenge

Different challenges due to intermittency of the wind power generation:

- Transmission System Operator (in Italy TERNA) to ensure the balance of grid (very short term horizon: 1s to 1h)
- Traders to forecast the power to sell on energy market, intraday or day ahead (short term horizon: 1h to 24h)
- Wind farm owners to schedule their maintenance programs (long term horizon)
  → **great value of improved wind power forecast accuracy**



**Wind Power Global Capacity and Annual Additions, 2008-2018**

EVEREST

# Advanced Traffic Modeling

- Mobility platform supporting cities with advanced traffic modelling

- **Data sources**
  - Historical and real-time **Floating Car Data** (FCD)
  - **Origin-destination matrix** (ODM) defining city
  - **Road network graph** including road restrictions;
  - Historical **weather data** (temperature, precipitation)

- **Traffic services**
  - **What-if analysis** for given scenarios, e.g. road closure;
  - **Intelligent routing** for large amount of vehicles
  - **Traffic prediction** for major road elements of cities

EVEREST

# EVEREST Target System

# EVEREST Target System: Brief Overview

## Network-attached and PCIe-attached FPGA nodes

- Off-the-shelf FPGA devices
- User logic can be easily designed and customized with HLS tools

## DC infrastructure and Supercomputers

- workflow orchestration
- reference implementation



**FPGA as a Co-Processor**

**FPGA as a Peer-Processor**

# cloudFPGA

## FPGAs as 1st-class citizens within a DC ...

- disaggregated from the server nodes

- connected directly to the DC network for its access and to communicate with CPUs and other FPGAs

- densely packed into DC chassis and racks and distributed across the DC



Xilinx Kintex UltraScale XCKU060 FPGA with 2x8GB of DDR4 memory

*Figure 1: (a) The disaggregated FPGA and (b) the carrier board.*

EVEREST

# cloudFPGA Development Kit (cFDK)

- **network-attached solution composed of:**
  - Interface logic already designed (**cF Shell**) to support system integration
    - TPC/UDP communication is managed transparently to the user logic
  - User logic (**ROLE**) that can be easily designed and customized with traditional HLS tools

- **application code running on host**
  - FPGA accessible through the network
  - Low-level libraries for host-FPGA communication

- **create clusters of FPGAs**
- **IDE incl allocation and mgmt of resources**

cFDK released at
https://github.com/cloudFPGA

# Conceptual System Overview

- Envisioned for demonstration purposes
- Multi-node demonstrator based on the technology and the components available during the project's timeline



**EVEREST Heterogeneous Node**

- 4 Nodes equipped with:
  - **2x AMD EPYC Milan 7643P**
    48 Cores, 96 Threads, 256GB
    DDR4, 2x960GB SATA SSD
    Dual ConnectX®-5 EN 100GbE
  - **2x XILINX Alveo U55C, single slot**
    PCIe® Gen3 x16, 2xGen4 x8 ,
    CCIX, 16GB HBM2, 2x QSFP28
  - **1x XILINX Alveo U280, dual slot**
    PCIe® Gen3 x16, 2xGen4 x8, CCIX,
    8GB HBM2, 2x QSFP28 (100 GbE)
  - **4x Spare slots** PCIe® 4.0 x16

29" Depth

24x 2.5" Hot-swap Drive Bays

**(Wistron Mihawk)**

Dual-port Mellanox
ConnectX-5 100G

**EVEREST Node**
Network-attached FPGA

10
Gb/s

up to 64x per 2U node
**(cloudFPGA)**

TCP/UDP          TCP/UDP          TCP/UDP

**DC Network**

EVEREST

# FPGA-accelerated HPC System Overview

- Multi-node demonstrator based on EVEREST SDK

- 4 Nodes equipped with:
  - **2x AMD EPYC Milan 7643P**
    48 Cores, 96 Threads, 256GB
    DDR4, 2x960GB SATA SSD
    Dual ConnectX®-5 EN 100GbE
  - **2x XILINX Alveo U55C, single slot**
    PCIe® Gen3 x16, 2xGen4 x8 ,
    CCIX, 16GB HBM2, 2x QSFP28
  - **1x XILINX Alveo U280, dual slot**
    PCIe® Gen3 x16, 2xGen4 x8, CCIX,
    8GB HBM2, 2x QSFP28 (100 GbE)
  - **4x Spare slots** PCIe® 4.0 x16



29" Depth

24x 2.5" Hot-swap Drive Bays

EVEREST

# EVEREST Target System

**Computing continuum** to enable **cloud-to-edge integration**



We aim at outperforming centralized/homogeneous solutions

# EVEREST Compilation Framework

Multi-location use cases



HPC ⟷ ML

Secure connections

Orchestration

**Dataflow/task graphs (routing)**

**HPC kernels (weather simulation)**

**ML (predictive models and decision making)**

Data management techniques

Virtualized runtime environment

EVEREST

# EVEREST Compilation Framework

Multi-location use cases



HPC ↔ ML

Secure connections

Orchestration

| Dataflow/task graphs (routing) | HPC kernels (weather simulation) | ML (predictive models and decision making) | |
|---|---|---|---|
| **Language/framework support** | | | WP3: Data management techniques |
| - Sequential syntax, implicit parallelism<br>- Deterministic execution<br>- Support for shared state | - DSLs for kernels (numerics, WRF model)<br>- Integration in HPC distributed infrastructure | - Interoperability (pytorch, TensorFlow, TVM)<br>- Support for exchange formats (ONNX) | |
| **High-level source-to-source compilation** | | | |
| - Dataflow IR: I/O, batching, pipelining | - MLIR dialect and polyhedral analysis | - Model partitioning based on TVM/RelayIR | |
| **High-level synthesis and memory design flow** | | | |
| - HLS support for OpenMP and irregular accesses | - Generation of memory controllers | - Re-use existing HLS tools and libraries | |
| **Low-level compilation, bitstream generation, and system integration** | | | |
| WP5: Virtualized runtime environment | | | |

EVEREST

# EVEREST Compilation Framework



| Dataflow/task graphs | HPC kernels | Machine learning |
|---|---|---|

**Frontend and IR**
- **Ohua** (implicit dataflow): Rust and Python syntax
- **Cfdlang** (tensors), Fortran integration (for WRF)
- TVM integration (multiple input languages, ONNX)

**MLIR-based representations**
- Dataflow IR (internal), task graphs
- Dataflow-IR (planned)
- Tensors, stencils (**teil**), linear algebra, number representations (**base2**)
- ONNX-MLIR (external)
- TVM Relay IR (external)

**Middle-end optimization**
- Batching, synch. minimization, irregular data parallelism, pipelining
- Algebraic and polyhedral transforms, memory partitioning, buffering, pipelining, HLS-pragma insertion
- TVM optimizations (external) and trade-off streaming vs batching

**HW and HLS**

Multi-variant (high-level) code generation: Stand-alone host code (Rust,C/C++), interfacing for HLS via MLIR, LLVM IR, Relay IR or C/C++ with HLS pragmas, generation of kernel code and wrapper code

- Python/Rust FPGA offloading via LLVM-IR
- **Mnemosyne**: Memory subsystem-optimization
- **Bambu**: Number representations, polyhedral optimization
- **Dosa:** Select implementations of ML operators

Kernel/function-level HLS: Bambu (internal) or Vivado/Vitis (external)

**Code-gen and system integration**
- **Olympus**: Kernel spatial replication, HW wrappers and adaptors, streaming interfaces, MLIR-based graph transformation
- **Dosa:** Legalizing, system-generation

Bitstream generation, host-code compilation with interfaces for the runtime (auto-tuning, resource management), and application integration

**WP3: Data management techniques**

**WP5: Virtualized runtime environment**
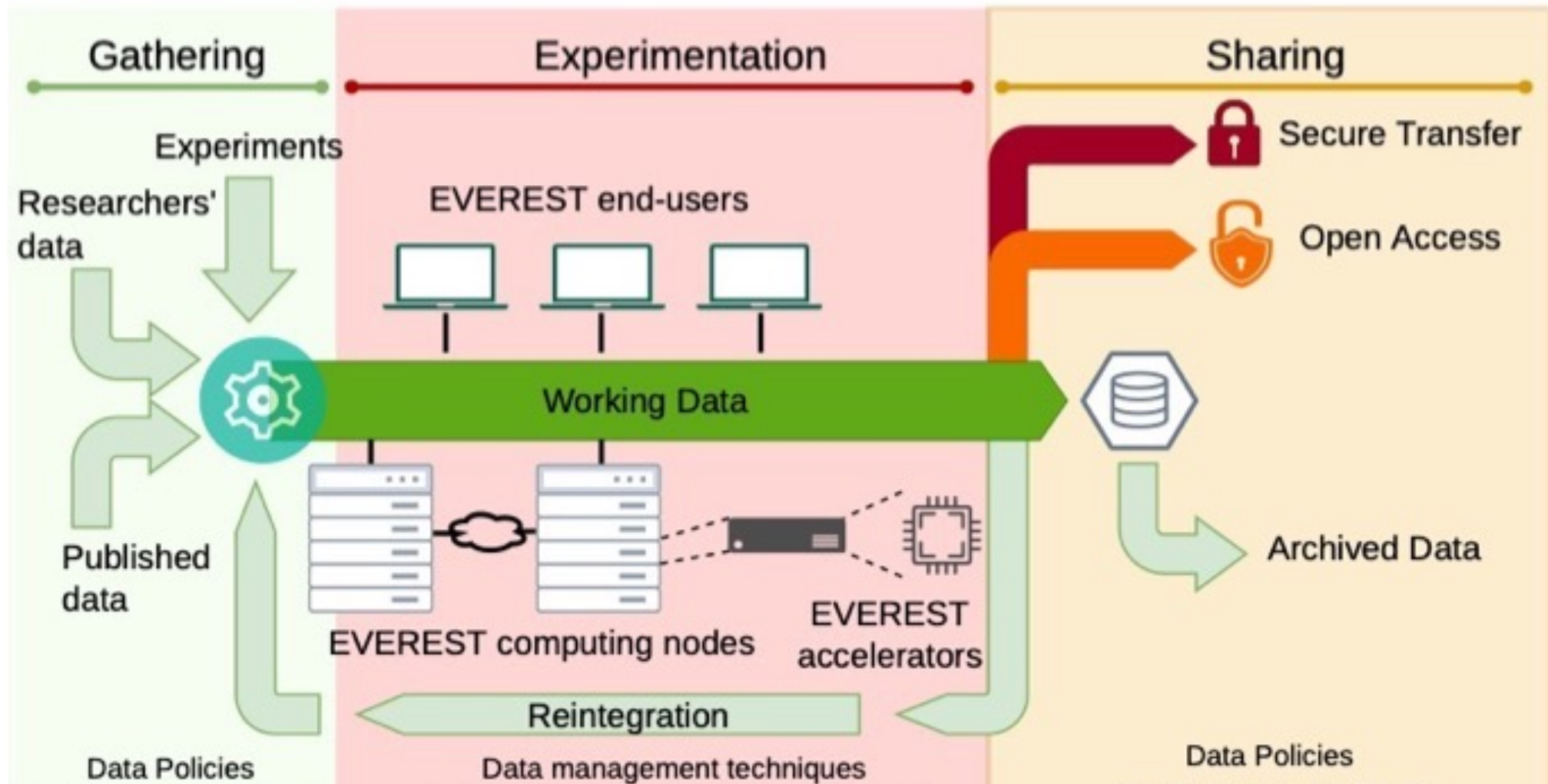
EVEREST

# DNN on FPGAs: no yet-another-narrow-framework

- DNN on FPGAs is a highly active area of research
  - →chances are that for a particular problem (i.e. some-convolution-to-put-on-FPGA with certain performance) someone has already developed and published a good implementation
    (e.g. haddoc2, FINN, hls4ml, VTA, VitisAI, and other open source frameworks...)
  - → Why to re-invent the wheel and not reuse it?
- BUT: Who knows what is the best available implementation for the current problem (I.e. the ONNX input by the user)?
- Standardized way to include all available 3rd-party libraries (including Everest flows) within architecture generation
  - Automatic DSE of best available framework (depending on: operation, precision, target device)
- Frontend currently based TVM, but plan to integrate also MLIR interfaces/modules
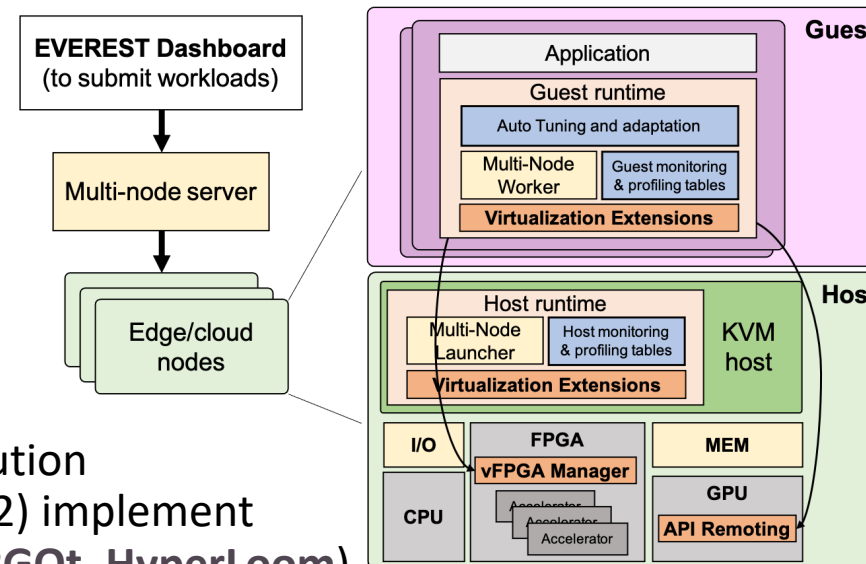
EVEREST

# EVEREST Data Management & Protection

# EVEREST Runtime Environment

... implements the *selection of "variants"* and the *hardware configuration* based on the *system status*

- **Dynamic adaptation and autotuning** (**mARGOt**)
- **Two-level runtime** for (1) virtualization of hardware resources regardless their distribution and the low-level details of the platforms; (2) implement functional decisions (**VOSYS solutions**, **mARGOt**, **HyperLoom**)



**How to collect system status and expose it to the runtime?**
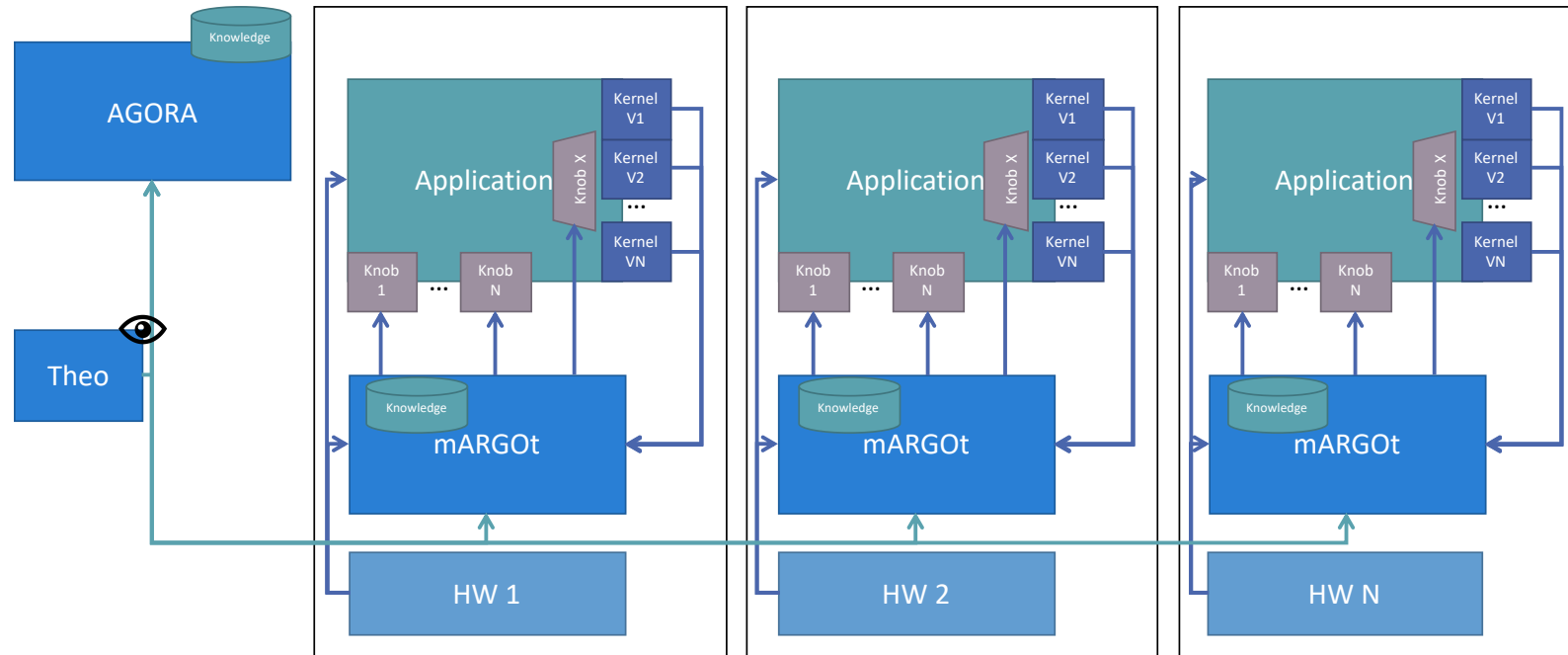
**Runtime API**

**Autotuning API**

**Hiding communication latency (e.g., prefetching)**

**Seamless execution when varying the system configuration (resources, nodes, data, etc.)**

# EVEREST Runtime Environment

The EVEREST FPGA systems include a **monitoring and decision infrastructure** for **dynamic autotuning** based on workload conditions



- **Application variants** (either software or hardware) are generated at design time (compilation and hardware synthesis), and selected at run time based on the actual available hardware resources

# Conclusion

- EVEREST is a dEsign enVironmEnt foR Extreme-Scale big data analyTics on heterogeneous platforms
  - built on the assumption that the future of computing is heterogeneous but the current tools do not support it
  - focus on building support for FPGAs
- The work towards an universal IR facilitates the re-use of innovations across the full stack including
  - extensions to new application domains / languages
  - extensions to different accelerator architectures
  - integration with different workflow engines / runtime environments
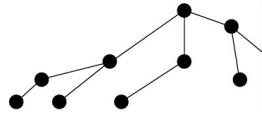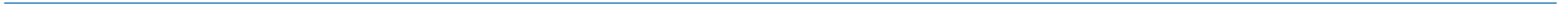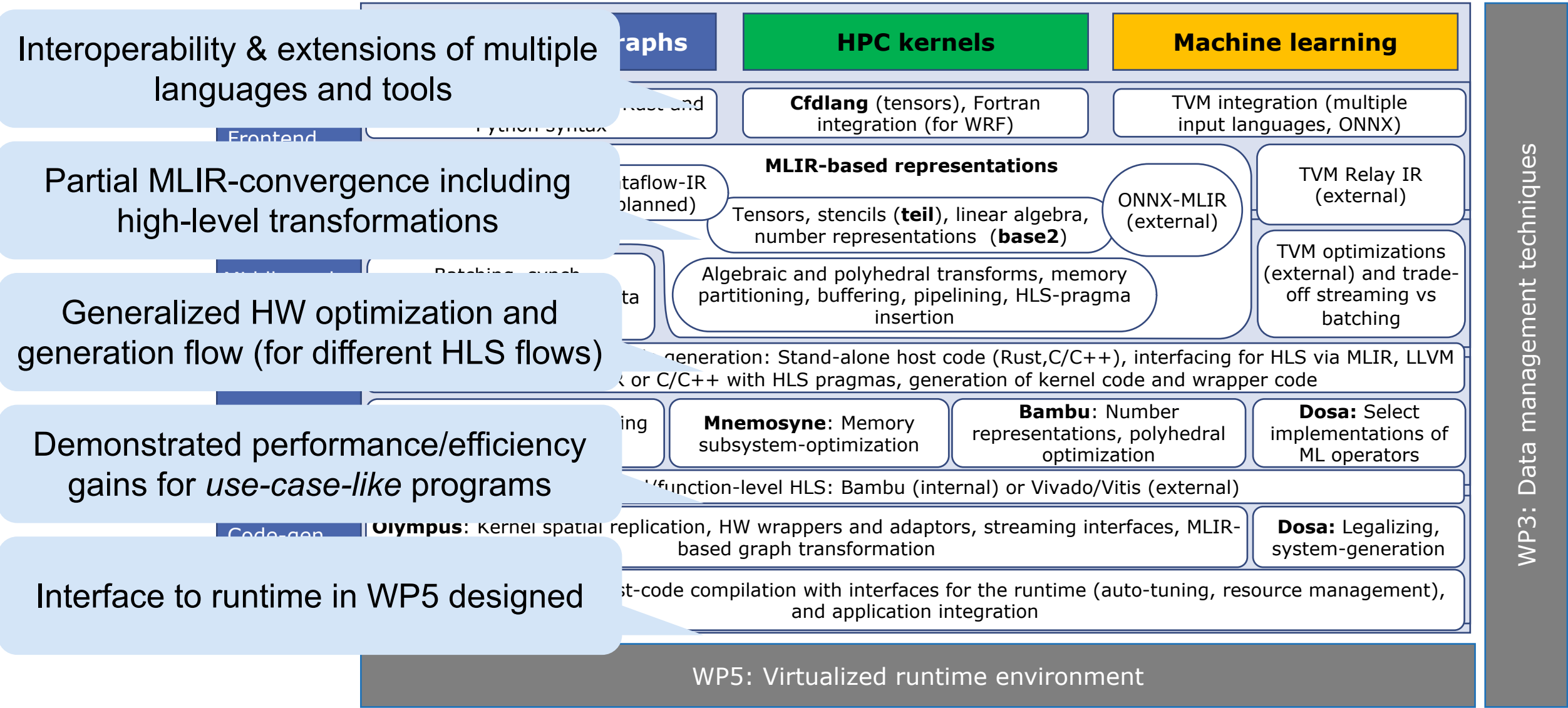- Stay tuned ;-) ... https://everest-h2020.eu/

EVEREST

# Thank You!

EVEREST

EVEREST

DESIGN ENVIRONMENT
FOR EXTREME-SCALE BIG DATA ANALYTICS
ON HETEROGENEOUS PLATFORMS

# EVEREST Compilation Framework

Interoperability & extensions of multiple languages and tools

Partial MLIR-convergence including high-level transformations

Generalized HW optimization and generation flow (for different HLS flows)

Demonstrated performance/efficiency gains for *use-case-like* programs

Interface to runtime in WP5 designed

**raphs**

**HPC kernels**

**Machine learning**

Frontend

…Rust and Python syntax

**Cfdlang** (tensors), Fortran integration (for WRF)

TVM integration (multiple input languages, ONNX)

**MLIR-based representations**

…taflow-IR (planned)

Tensors, stencils (**teil**), linear algebra, number representations (**base2**)

ONNX-MLIR (external)

TVM Relay IR (external)

…ta

Algebraic and polyhedral transforms, memory partitioning, buffering, pipelining, HLS-pragma insertion

TVM optimizations (external) and trade-off streaming vs batching

Batching, synch…

…generation: Stand-alone host code (Rust,C/C++), interfacing for HLS via MLIR, LLVM …R or C/C++ with HLS pragmas, generation of kernel code and wrapper code

…ing

**Mnemosyne**: Memory subsystem-optimization

**Bambu**: Number representations, polyhedral optimization

**Dosa:** Select implementations of ML operators

…/function-level HLS: Bambu (internal) or Vivado/Vitis (external)

Code-gen

**Olympus**: Kernel spatial replication, HW wrappers and adaptors, streaming interfaces, MLIR-based graph transformation

**Dosa:** Legalizing, system-generation

…st-code compilation with interfaces for the runtime (auto-tuning, resource management), and application integration

WP3: Data management techniques

WP5: Virtualized runtime environment

EVEREST

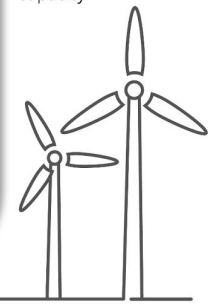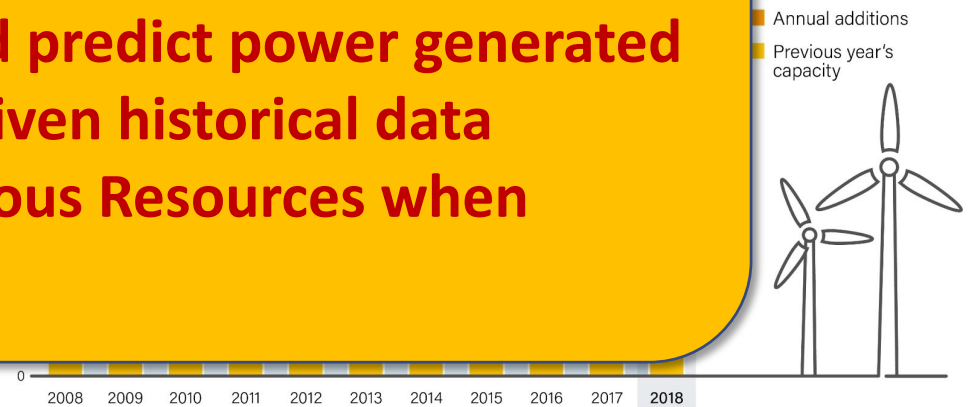# Renewable Energy use case: Context and Challenge

- Different challenges due to <u>intermittency of the wind power generation</u>:
  - <u>Transmissio...</u> ... e of grid (very short t...
  - Traders to f... ...y ahead (short term ...
  - Wind farm o... ... horizon)
  - → **great eff...**

**CHALLENGES:**

1. **Integration of the the data coming from sensors available on wind turbine for improving weather forecast model**
2. **Extend the wind speed estimation at 80m instead of the 10m default in WRF to improve prediction capabilities**
3. **Adopt AI models to learn and predict power generated by the specific site/turbine given historical data**
4. **... making use of Heterogeneous Resources when needed**

Annual additions
Previous year's capacity

0
2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 **2018**

EVEREST

# Advance Traffic Modeling for Smart Cities use case

- Mobility platform for supporting cities with advanced traffic modelling

- **Data source**
  - Historical a
    - e.g. GPS
  - **Origin-des**
  - **Road netw**
  - Historical w

- **Traffic servi**
  - **What-if ana**
  - **Intelligent routing** for large amount of vehicles
  - **Traffic prediction** for major road elements of cities

**CHALLENGES:**

1. **Boost the FCD data collection and processing to compute 3D traffic models**
2. **Develop and efficiently deploy a traffic simulator on a multinode architecture to predict traffic behaviour**
3. **Adopt AI techniques to learn the traffic patterns resulting into a traffic prediction service.**
4. **... making use of Heterogeneous Resources**

# The WRF Model

- First step of two use cases…

## WRF Modeling System Flow Chart



**CHALLENGES**
- **Automatic deployment of several instances of the model in HPC resources with multiple boundaries conditions**
- **Efficient High Resolution output evaluating the usage of FPGA resources**

- WRF is an collab... supported principally by the US National Center for Atmospheric Research (NCAR), the US National Oceanic and Atmospheric Administration and the US National Center for Environmental Prediction – NCEP)

EVEREST

# EVEREST Consortium



[1]IBM Research Europe, Switzerland, [2]Politecnico di Milano, Italy, [3]Università della Svizzera Italiana, Switzerland, [4]Technische Universitat Dresden, Germany, [5]Centro Internazionale di Monitoraggio Ambientale, Italy, [6]IT4Innovations, VSB – Technical University of Ostrava, Czech Republic, [7]Virtual Open System, France, [8]Duferco Energia, Italy, [9]NUMTECH, France, [10]Sygic, Slovakia