# A journey of cloud FPGAs on converged data centres

Dr. Dionysios Diamantopoulos
Research Staff Member
Cloud FPGAs & Tape Group
Cloud & AI Systems Research Department
IBM Research Europe

did@zurich.ibm.com

*Hipeac Workshop 20-22 June 2022, Budapest, Hungary*

*EVEREST: Design and Programming High-performance, distributed,*

*reconfigurable and heterogeneous platforms for extreme-scale analytics*

*https://www.hipeac.net/2022/budapest/*

IBM

# IBM Legal Disclaimer

This content was provided for informational purposes only. The opinions and insights discussed are those of the presenter and guests and do not necessarily represent those of the IBM Corporation.

Nothing contained in these materials or the products discussed is intended to, nor shall have the effect of, creating any warranties or representations from IBM or its suppliers, or altering the terms and conditions of any agreement you have with IBM.

The information presented is not intended to imply that any actions taken by you will result in any specific result or benefit and should not be relied on in making a purchasing decision. IBM does not warrant that any systems, products or services are immune from, or will make your enterprise immune from, the malicious or illegal contact of any party.

All product plans, directions and intent are subject to change or withdrawal without notice. References to IBM products, programs or services do not imply that they will be available in all countries in which IBM operates. IBM, the IBM logo, and other IBM products and services are trademarks of the International Business Machines Corporation, in the United States, other countries or both. Other company, product, or services names may be trademarks or services marks of others.

For copyright and trademark information go to: http://www.ibm.com/legal/us/en/copytrade.shtml
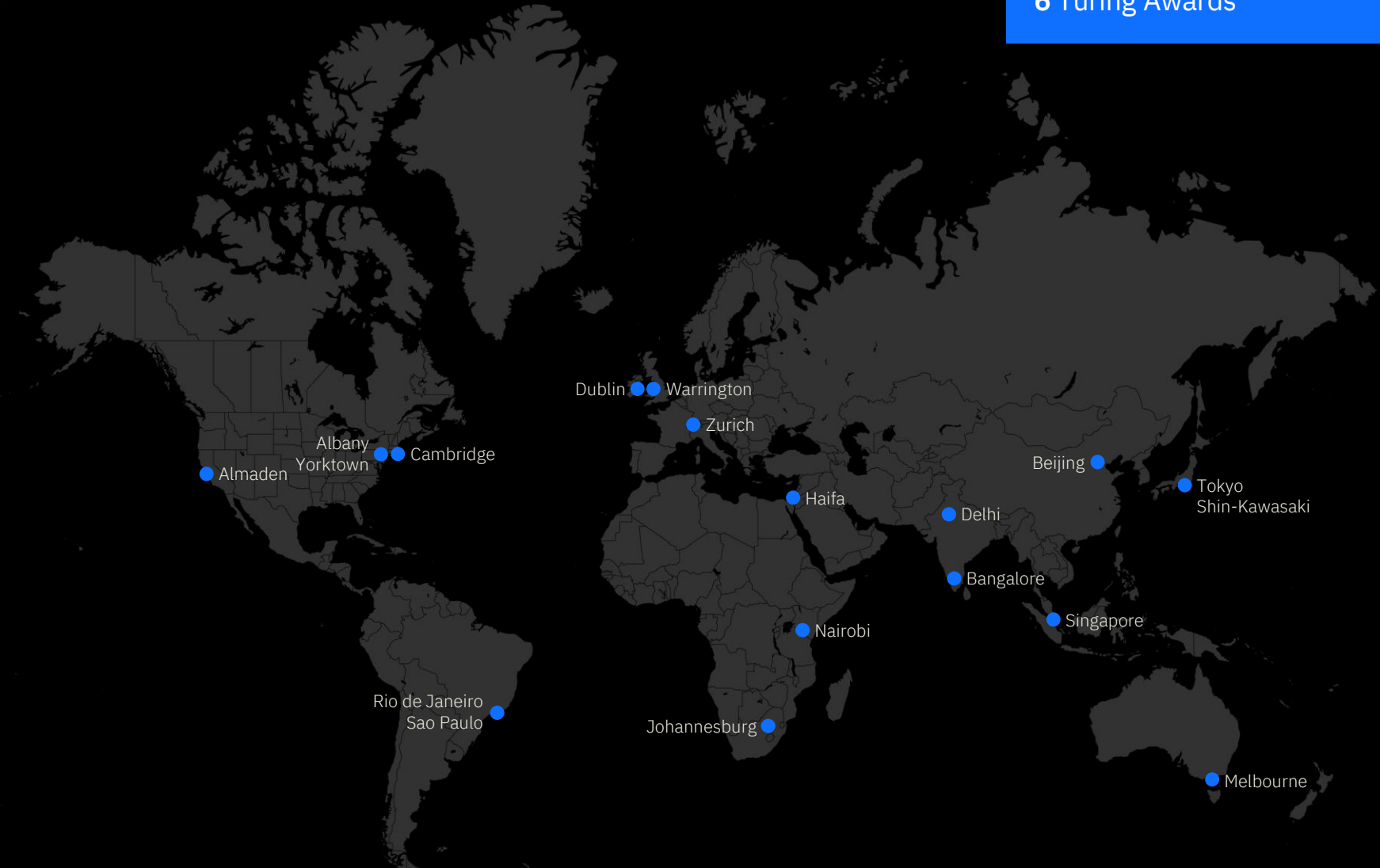
# 3000
Researchers

# 19
Locations

# 6
Continents

Dublin  Warrington
Zurich
Albany
Yorktown  Cambridge
Almaden
Haifa
Beijing
Tokyo
Shin-Kawasaki
Delhi
Bangalore
Nairobi
Singapore
Rio de Janeiro
Sao Paulo
Johannesburg
Melbourne

# IBM Research – Zurich

**Established in 1956**
45+ different nationalities
Open Collaboration:
- Horizon2020: 50+ funded projects and 500+ partners

Two Nobel Prizes:
- 1986: Nobel Prize in Physics for the invention of the scanning tunneling microscope by Heinrich Rohrer and Gerd K. Binnig
- 1987: Nobel Prize in Physics for the discovery of high-temperature superconductivity by K. Alex Müller and J. Georg Bednorz

European Physical Society Historic Site

Binnig and Rohrer Nanotechnology Centre (Public Private Partnership with ETH Zürich and EMPA)

7 European Research Council Grants

We're Inventing What's Next in:

# Hybrid Cloud
# AI
# Quantum
# Science



2020 U.S. Patents: Select Technology Companies

| Company | Patents |
|---|---|
| IBM | 9,130 |
| Microsoft | 2,905 |
| Apple | 2,792 |
| Amazon | 2,244 |
| Google | 1,817 |

Source: 2020 patent data is sourced from IFI Claims Patent Service

**IBM's innovation: Topping the US patent list for 28 years running**

https://www.ibm.com/blogs/research/2021/01/ibm-patent-leadership-2020/

From automated teller machine (ATM), speech recognition technology, DRAM to a novel way to search multilingual documents using NLP, 2300 AI patents !

# Agenda

**Context & Competitive Landscape**

**EVEREST use cases**

**cloudFPGA**

**Converged Infra**

**Example on Traffic Simulation**

# Context & Competitive Landscape for EVEREST

# EVEREST Use Cases



**Renewable energy production prediction**

★ Improve **quality of the predictions**

**Weather prediction modelling (WRF)**

**Air-quality monitoring of industrial sites**

★ Improve the **response time of predictions**

★ Accelerate kernels to execute more tests

**Traffic modeling for intelligent transportation**

★ Improve the **overall performance of traffic simulation**

**Accelerated** computationally-intensive kernels **+** **Machine-learning** kernels

EVEREST

# A Converged computing journey

## Massive CPU pools

## Pools of Accelerators

**Private Clouds**

**Public Clouds and SaaS**

**Dedicated Clouds**

**Super-boosting Traditional IT Silos**

**Accelerators**

*Super-boosting is not enough, breaking silos is imperative ...*
*Converged infra to the rescue*

## An optimal strategy will address key challenges:

**50%** Reduce by 50% the development costs

**10x** Reduce programming efforts by one order of magnitude

## Unleashes the full potential

**30%** energy cost

**10x** the performance of simulations for renewable energies prediction

**2x** the response time of the air-quality predictions

**3x** the overall performance of traffic model framework

# Computing continuum to enable cloud-to-edge integration

Cloud & bus-attached FPGAs at H2020 EVEREST Project

There's no AI without IA

# cloudFPGA concept



Highlights

- dense
  → chassis w/ 64 compute units
  → ~1000 FPGAs / rack

- integration of 1st level switch
  → full cross-sectional BW
  → low cost (cables / rack space)

- energy efficient
  → no SW/FW overhead
  → no CPU overhead
  → (hot) water cooling

- self-hosted / network-attached
  → bare-metal support
  → scalabl

IP Address: 10.10.1.9
DRAM: 8GB, BRAM: 38MB
CLBs:660.000,
DSPs: 2760

The FPGA becomes the node !

IP Address: 10.10.1.50
DRAM: 32GB, Cores: 4

Goal → Deploy FPGAs at large scale in data centers
⬚ 1-10s of thousands per DC

# Standalone network-attached FPGA



1) Integrated NIC (iNIC)

2) PSoC ARM Cortex-M3

iNIC ~15%

GTH

IBM
IBM Research GmbH

Connector

Backplane

3) 10GBASE-KR

USB

1. Replace PCIe I/F with integrated NIC (iNIC).

2. Turn FPGA card into a self-contained appliance.

3. Replace transceivers w/ backplane connectivity.

# One carrier sled = 32 FPGA modules

— 

1. Our first FPGA module uses a Xilinx Kintex Ultrascale KU060
   o A mid-range FPGA with high performance/price and low wattage

# One carrier sled = 32 FPGA modules

—

1. Our first FPGA module uses a Xilinx Kintex Ultrascale KU060
   o A mid-range FPGA with high performance/price and low wattage



×16   ×16

40GE x8

640 Gb/s Ethernet switch

# Two carrier sleds per chassis = 64 FPGAs

—



Figurative picture

Legend (per slice):
```
[==]  x8 40GbE up links              (320 Gb/s)  } Balanced (i.e. no over-subscription)
[--]  x32 10GbE FPGA-to-Switch links (320 Gb/s)  } between the north and south links of
[--]  x32 10GbE redundant links                    the Ethernet switch
[--]  x32 10GbE FPGA-to-FPGA links
[■]  x16 PCIe x8 Gen3
 SP   x1 Service Processor
```

# Sixteen chassis per rack = 1024 FPGAs

—



Spare (5U)

1024 FPGAs (32U)

CDU (5U)

1024 FPGAs → 2.8M DSPs,
$2\times10^{15}$ Fixed-Point Multiply-Accumulates/s
10 Tb/s bi-sec. Bw – 16 TB DDR4 – 40 kW max.

# cloudFPGA

**The FPGA Platform**

**Standalone Network-attached FPGAs over TCP/IP/Ethernet**

IBM

**10 Tb/s full-duplex**

Chassis-1 · · · · Chassis-16

Rack-1

50 kW   5U CDU

64/chassis          1024/rack          Plentiful/DC

# EVEREST Data lifetime

# Key Components

# Building the stack…

# Converged FPGAs

# Use case: Traffic simulation on converged EVEREST platform

EVEREST users



Traffic simulation workflow

https://portal.lexis.tech/

IT4I Datacenter @ Ostrava, Czechia
- x86 nodes
- Lexis portal

IBM Datacenter @ Zurich, Switzerland
- x86 nodes
- cloudFPGA
- Alveo-FPGA cluster (in-progress)

# Use case: Traffic simulation on converged EVEREST platform

## EVEREST project within LEXIS portal



IT4I Datacenter @ Ostrava, Czechia
- x86 nodes
- LEXIS portal



IBM Datacenter @ Zurich, Switzerland
- x86 nodes
- cloudFPGA
- Alveo-FPGA cluster (in-progress)

# Use case: Traffic simulation on converged EVEREST platform

EVEREST HEAppE API

# HPC requires more than just Cloud infrastructure

Without scheduling, workloads are dispatched in a haphazard fashion leading to sub-optimal resource use and longer execution times
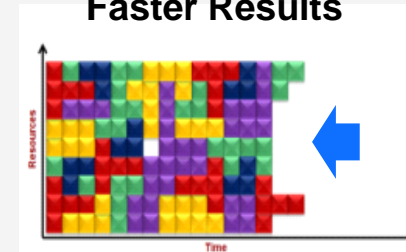


**Without Scheduling**

With sophisticated scheduling, we use resources more efficiently and can optimize for high-utilization, better performance, or both
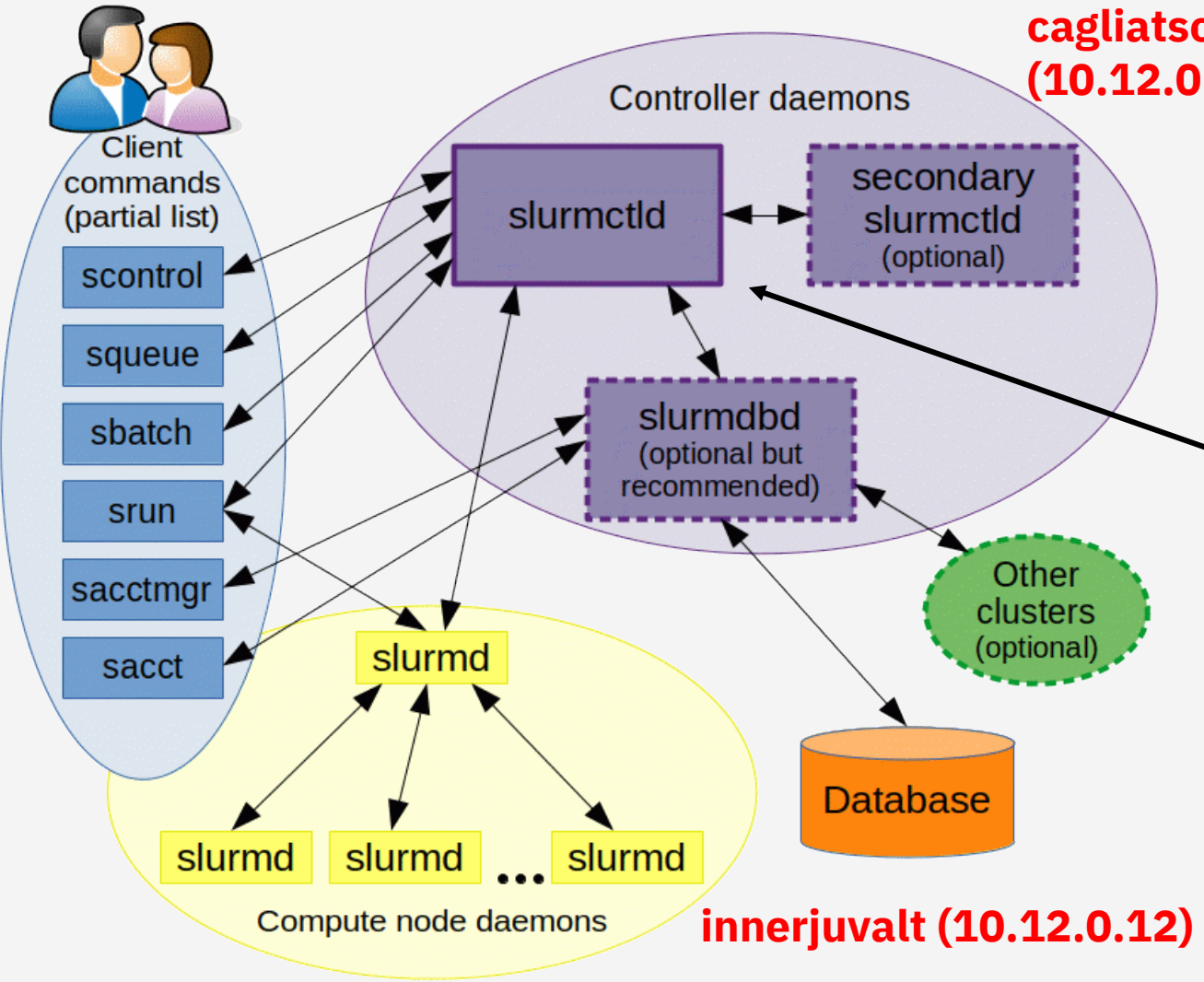


**Fewer Resources**



**Faster Results**

# Use case: Traffic simulation on converged EVEREST platform
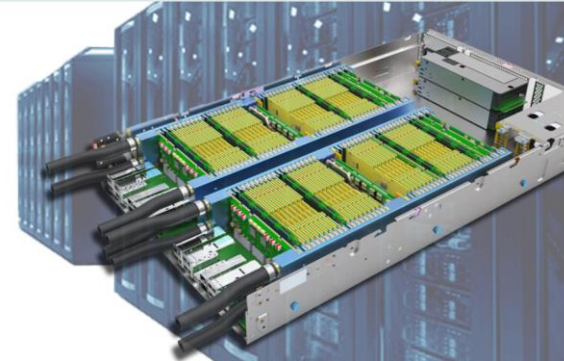
**EVEREST user**

**SLURM setup at IBM:**
- SLURM controller & daemon up and running
- Connection to cloudFPGA
- NFS shared storage
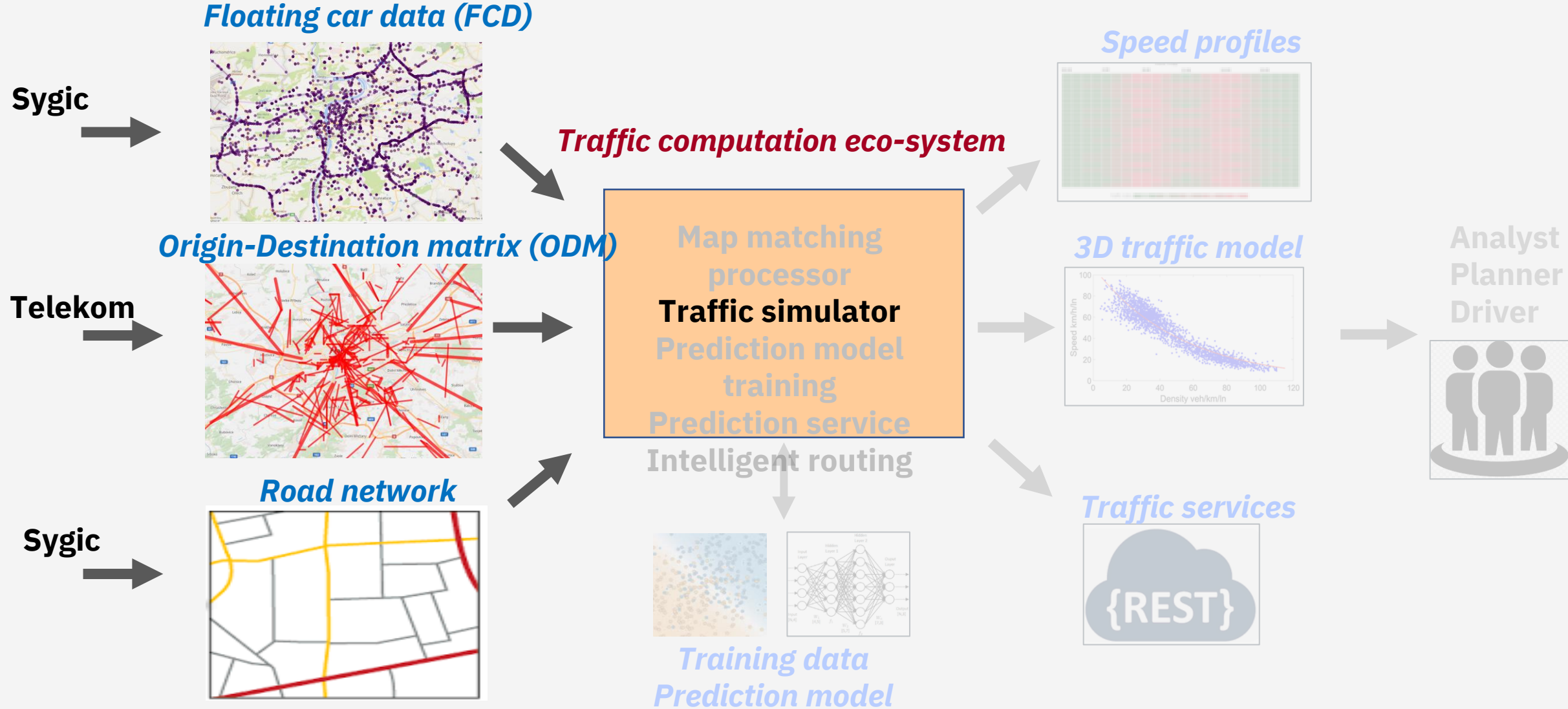- Exposure of a task to workflows out of IBM VPN (HEAppE, SOCKS5)

**cagliatscha (10.12.0.11)**

**cloudFPGA**

Client commands (partial list)
- scontrol
- squeue
- sbatch
- srun
- sacctmgr
- sacct

Controller daemons
- slurmctld
- secondary slurmctld (optional)
- slurmdbd (optional but recommended)

Other clusters (optional)

slurmd

slurmd  slurmd  ...  slurmd

Compute node daemons

**innerjuvalt (10.12.0.12)**

Database

**EVEREST**

# Use case: Traffic simulation on converged EVEREST platform

**Floating car data (FCD)**

**Sygic**

**Speed profiles**

**Traffic computation eco-system**

**Origin-Destination matrix (ODM)**

**Telekom**

Map matching processor
**Traffic simulator**
Prediction model training
Prediction service
Intelligent routing

**3D traffic model**

**Analyst Planner Driver**

**Road network**

**Sygic**

**Training data Prediction model**

**Traffic services**

{REST}

**EVEREST**

# Use case: Traffic simulation on converged EVEREST platform

Create an EVEREST workflow for traffic simulation inside LEXIS

# Use case: Traffic simulation on converged EVEREST platform

Traffic simulation initiated by LEXIS workflow is being executed at IBM SLURM computing node.

# Conclusion for EVEREST Converged infra:


Simplified cluster management


Design your
own cluster with cloud FPGAs


Security
& isolation


Onboard apps within
HPC workflows


Native open-source
experience


Integrated
operational tools

# *Köszönöm ! - Thank you!*

**IBM**