

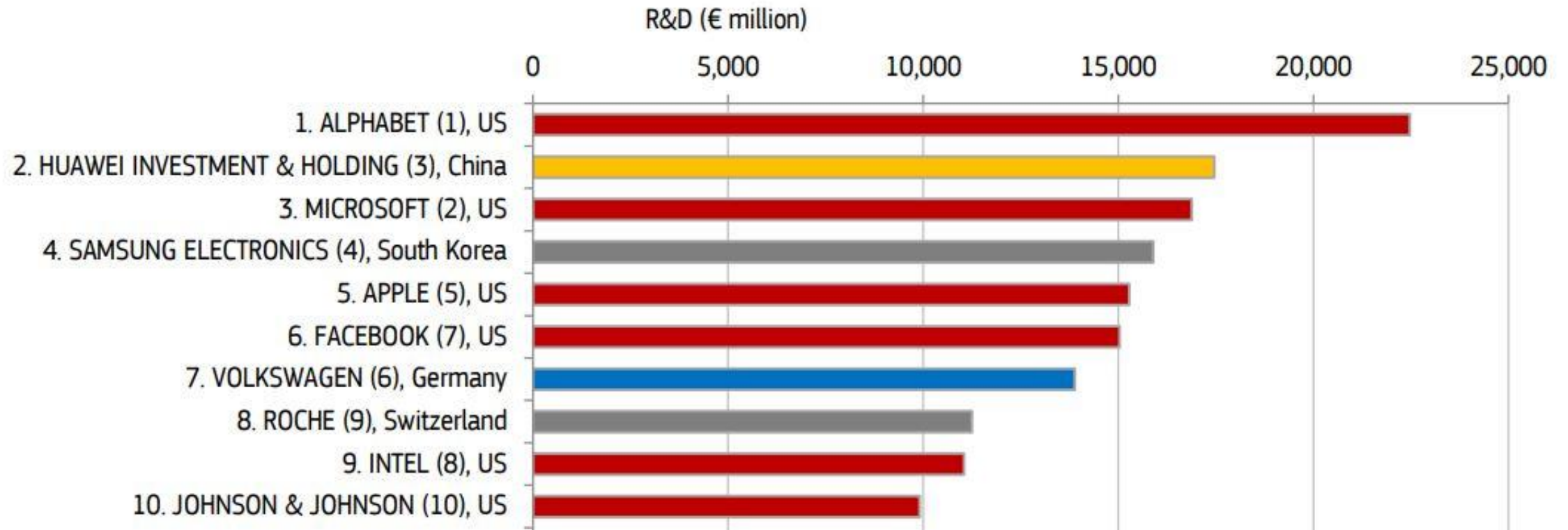
# Future Architectures for Computing and AI

**Bill McColl**  
**Director, Computing Systems Lab**  
**Zurich Research Center**

# Huawei Research

# EU List of World's Top R&D Investors 2021

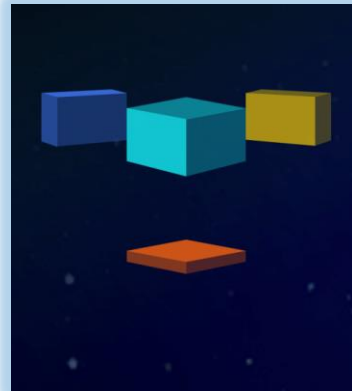
**Huawei is #2, above Microsoft, Samsung, Apple, Meta, Intel**



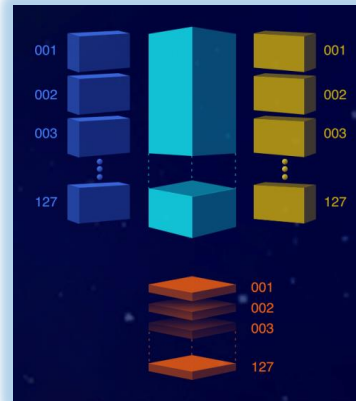
- **Full Stack Tech Company – Driven by Research and Innovation**
  - Hardware: CPUs, GPUs, AI Accelerators
  - AI Supercomputers
  - Communications (5G, 5.5G, 6G, Optical)
  - Servers
  - Storage
  - Huawei Cloud
  - Mobile Devices (Phones)
  - Automotive
  - Open Source Software (Operating Systems, Databases, AI Platforms)

# Huawei Hardware: CPUs, GPUs, AI Accelerators,...

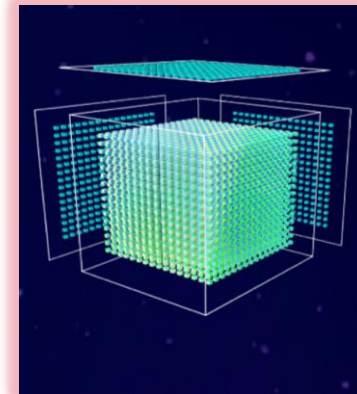
**CPU**  
**Scalar Compute**



**GPU**  
**Vector Compute**



**Ascend AI**  
**Tensor Compute**



# AI Accelerator Clusters



# Huawei Open Source Software

- **OpenEuler**
  - Open source operating system for servers and cloud
- **Hongmeng**
  - Open source operating system for devices – phone, IOT
- **OpenGauss**
  - Open source database
- **Mindspore**
  - Open source AI platform



## Huawei to open new €150 million European cloud services hub in Ireland creating 200 jobs

18th October 2022

Huawei today announced a €150 million investment and the creation of 200 new jobs as it plans to open its first European cloud hub in Dublin. The announcement was made at the Huawei Connect event in Paris, France. The hub will create 60 jobs in the next 2 years and 200 by 2027, including sales, pre-sales, legal, tax, operational, management and research positions.

The new European hub will serve customers across the continent and help ambitious Irish enterprises to expand into new global markets, by providing a secure, compliant, innovative, and sustainable cloud platform for growth and development. Supported by Huawei's industry-leading technology, investment in R&D, and open digital ecosystem, Huawei Cloud will offer a suite of cutting-edge cloud services across Europe, using Ireland as the platform.



# Huawei Research in Europe

# European Research Centers

- Amsterdam
- Cambridge
- Dublin
- Dresden
- Edinburgh
- Gothenburg
- Grenoble
- Helsinki
- Ipswich
- Kyiv
- London
- Lund
- Milan
- Munich
- Nice
- Nuremberg
- Paris
- Pisa
- Stockholm
- Tampere
- Tel Aviv
- Warsaw
- Zurich



# Huawei - Partnering in Research

- **Funding Europe's professors, postdocs and students to drive research and innovation**
  - Projects
  - Gift Funding
  - Joint Labs
  - Research Centers

The image shows a chalkboard with handwritten mathematical equations. The top equation is a square root approximation:  $\sqrt{\frac{p(1-p)}{n}} \approx \sqrt{\frac{p(1-p)}{n}}$ . Below it, there is a complex fraction: 
$$= \frac{\frac{n}{2(\frac{t_p^2}{n} + 1)}}{2w + \frac{t_p^2}{n} + \sqrt{4t_p^2 \frac{w(1-w)}{n}} + \left(\frac{t_p^2}{n}\right)^2} \left( \frac{1}{2} \right)$$
 This is followed by a probability statement: 
$$P\left(|w - \frac{1}{2}| < \frac{t_p}{\sqrt{n}} \sqrt{\frac{p(1-p)}{n}}\right) = 1 - \frac{1}{2}$$
 Then, an arrow points to the limit as  $n \rightarrow \infty$ : 
$$\xrightarrow{n \rightarrow \infty} \frac{1}{2} \left( 2w \pm \sqrt{4t_p^2 \frac{w(1-w)}{n}} \right) = w$$
 At the bottom, there is another equation: 
$$w^2 - p^2 - 2wp < \frac{t_p^2}{n} \left( \frac{p(1-p)}{n} \right)$$

# Huawei Research in Zurich

- **Launched 2020**
  - Center of major growth in future computing research by Huawei in Europe
- **Mission**
  - Carry out open research to drive breakthroughs in fundamental theory, software and architectures for future computing systems
  - Bring together world-leading researchers and top academic research partners to work together on the grand challenges of future computing

- **Full Stack Research - Architectures, Software and Algorithms**
  - CPU Architectures
  - GPU and AI Accelerator Architectures
  - Memory and Storage Architectures
  - Hardware-Software Interfaces and Co-Design
  - Power-Efficiency
  - Hardware Security
  - Cluster, Datacenter and Cloud Architectures
  - Interconnects and Fabrics: New Protocols, Hardware and Software
  - Heterogeneous Systems and Integrated CPU+GPU+AI Devices
  - Parallel Computing and HPC: Programming Models, Architectures, Software, Algorithms
  - New Compiler Technologies and MLIR
  - Dynamic language Runtimes

# **Future Architectures for Computing and AI**

# Computing Today

- **Cloud computing dominant**
  - Majority of all global IT spending by small number of “hyperscalers” (less than 20 companies)
  - Hyperscalers now designing their own chips and architectures
  - Cloud is still mostly business and web apps
- **Open source software dominant**
- **Disaggregation**
  - Flexible dynamic compute pools, memory pools, storage pools
- **New hardware**
  - DPUs and other “infrastructure processors”
  - Converged compute nodes : CPU + GPU + AI



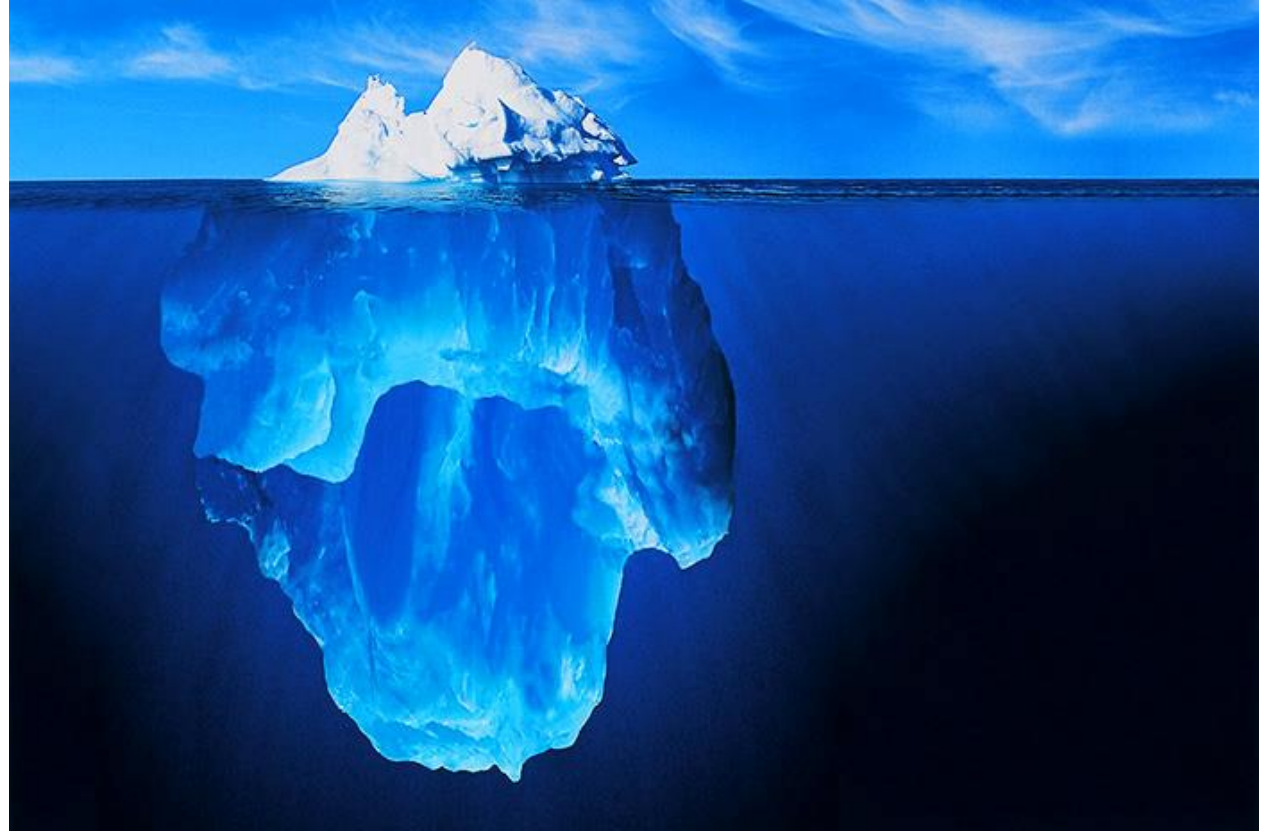
# Future: Moving From Clouds to Superclouds

- **Supercloud Computing**
  - Everything as a commercial cloud service
  - HPC, AI, Physical Simulation, Data Analytics, Knowledge Bases,...
- **New intelligent applications and services combining HPC, AI, Simulation, Analytics, Knowledge**
  - Smart Science, Design, Modelling and Engineering
  - Digital Twins
  - Intelligent Infrastructure for Smart Cities, Transportation and Energy
  - Smart Healthcare

# From Clouds to Superclouds?

Today's architectures for HPC, Simulation, AI, Big Data address only the tip of the iceberg

They are OK for running a single app on a single dedicated machine that functions perfectly with no faults and no long tails

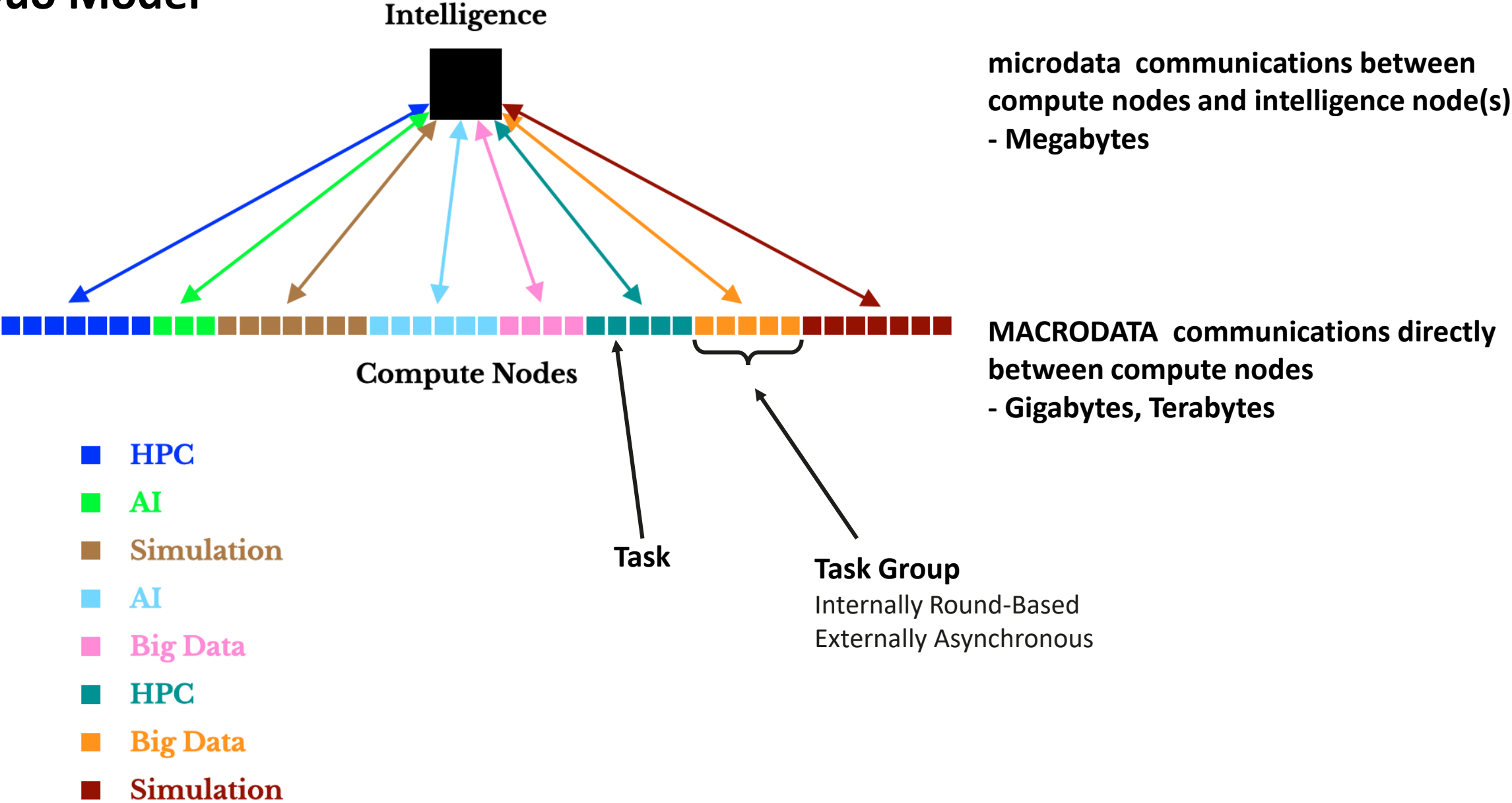


# From Clouds to Superclouds? - Research Challenges

1. Thousands of diverse applications and services running simultaneously
2. Continuous inputs and updates from edge devices
3. Disaggregated hardware architectures at massive scale: flexible compute pools + storage pools
4. Need to guaranteed Quality of Service with constantly changing priorities and deadlines
5. Need continuous nonstop operation with automatic fault tolerance and tail tolerance
6. Need continuously optimized shared infrastructure for maximum performance and efficiency
7. Need to compute on massive shared evolving datasets and knowledge bases

**Need new supercloud model and architecture**

# Duo Model



# Resilience? – Faults and Long Tails

10 October 2022

## World's Fastest Supercomputer Can't Run a Day Without Failure

By Anton Shilov published 1 day ago

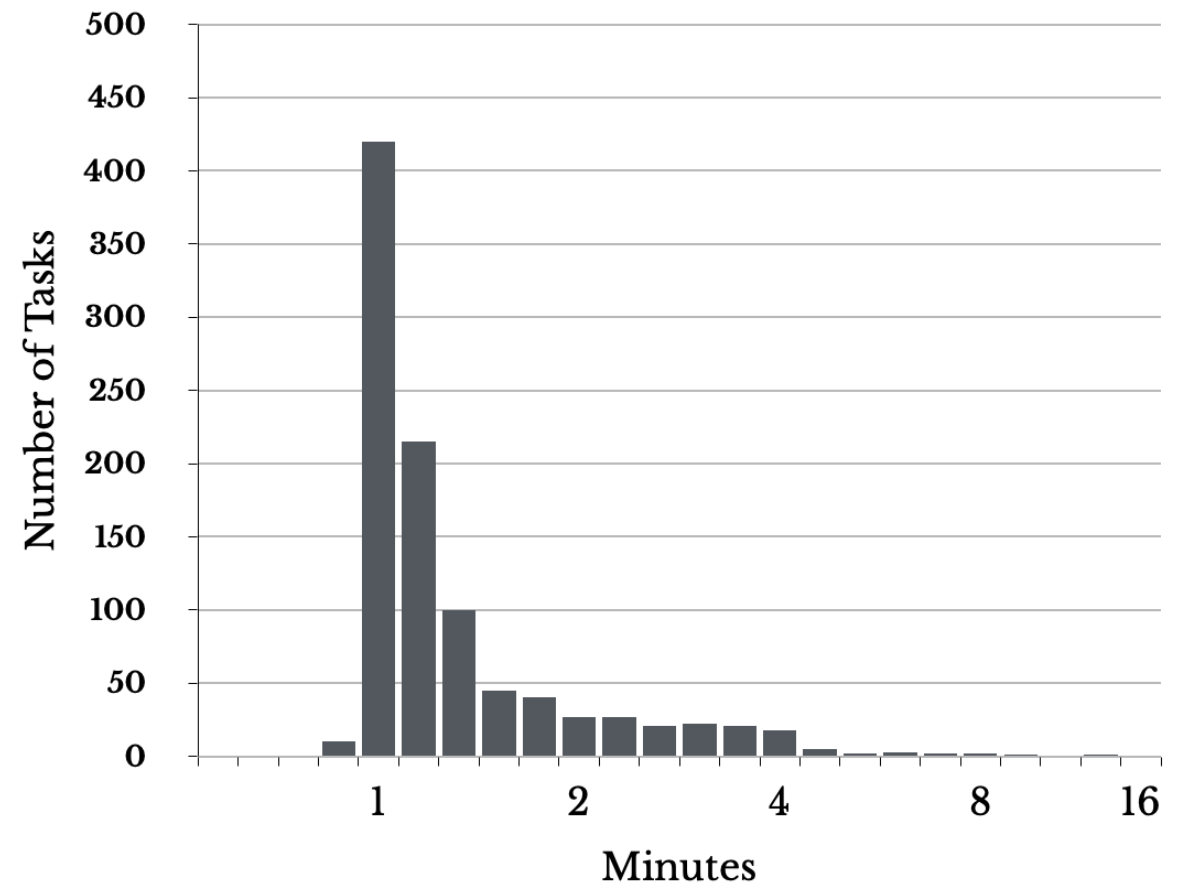
AMD's Instinct GPUs and HPE's Slingshot Interconnects blamed.

[f](#) [t](#) [g](#) [p](#) [F](#) [e](#) [c](#) Comments (24)



(Image credit: OLCF)

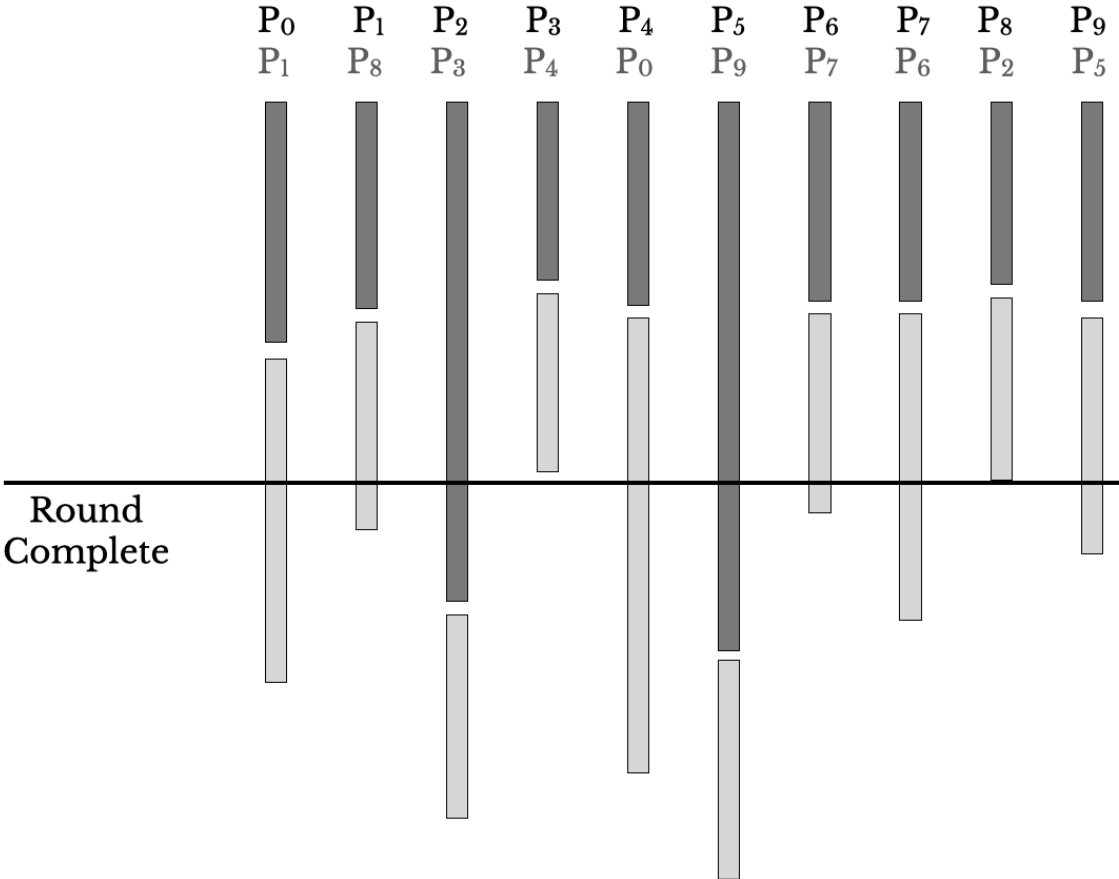
## Tail latency in cloud services



# Duo - Resilience

Intelligent cloning for nonstop operation

Automated optimized high-performance  
fault and tail tolerance for any  
computation at any scale

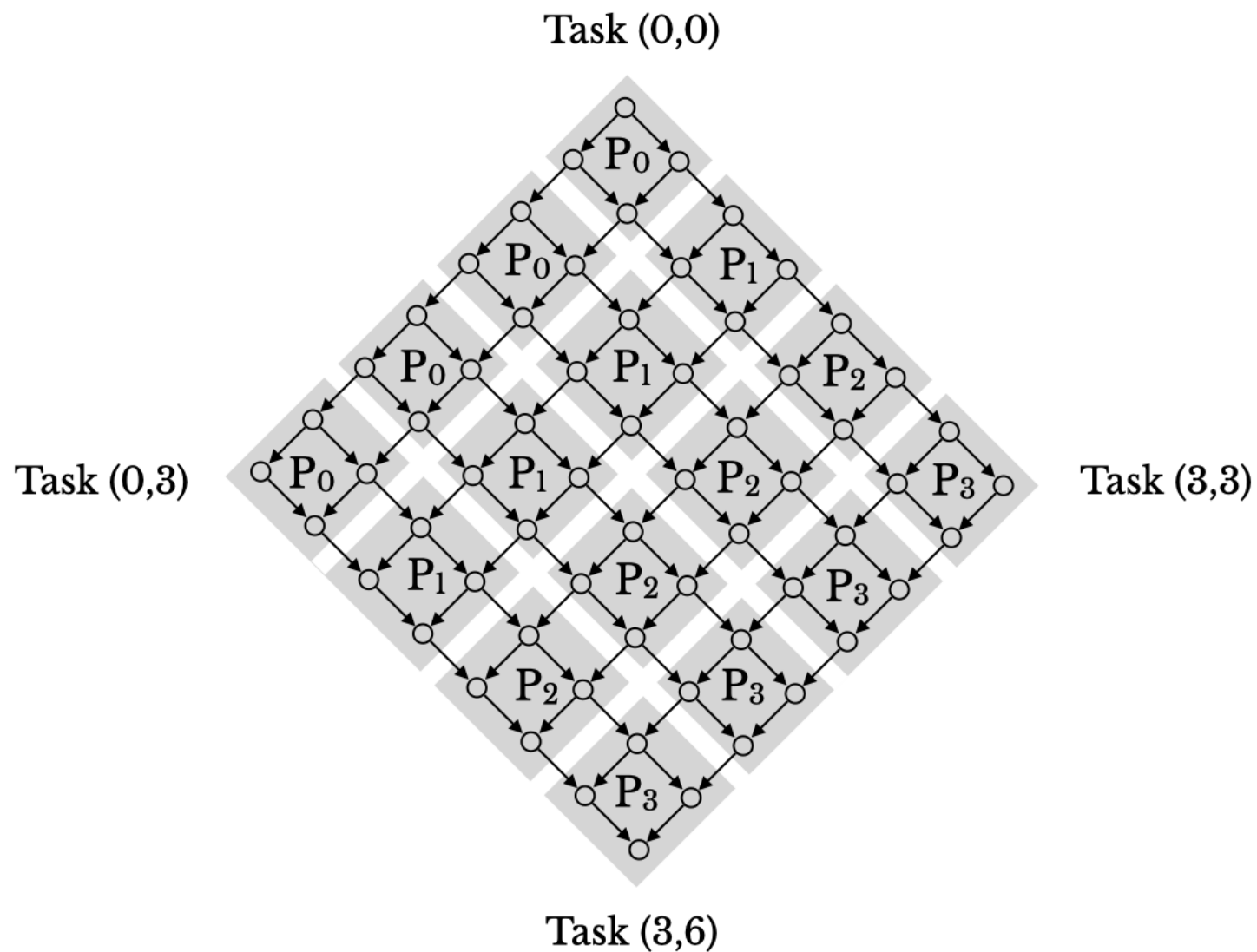


# Duo - Scheduling

Single task group

8x8 2D Grid DAG

Mapped to 4 nodes, 7 rounds

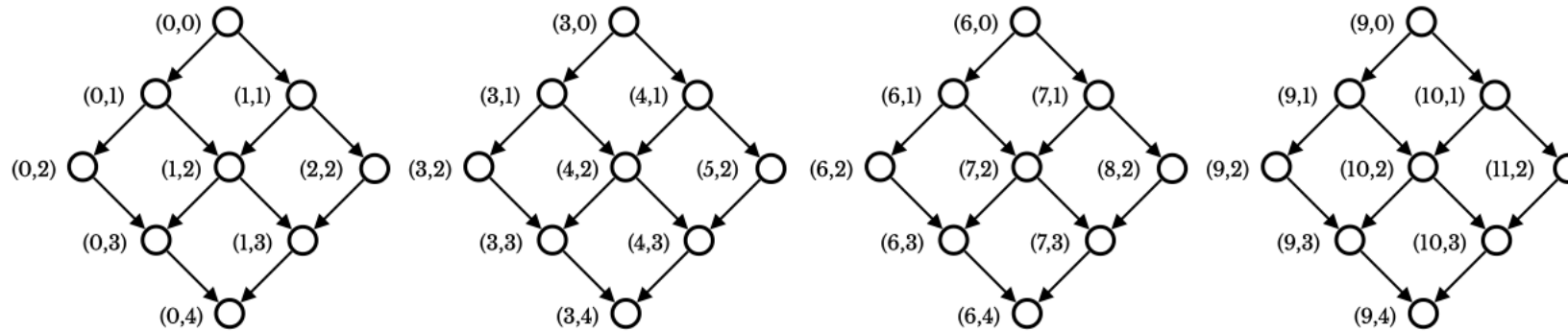




# Duo – Co-Scheduling

## Concurrent scheduling of four 3x3 2D Grid DAGs

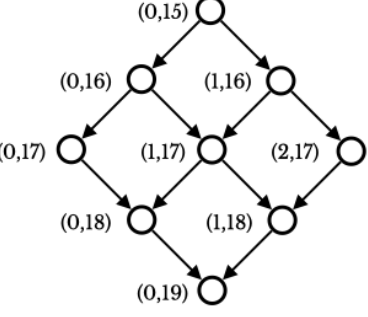
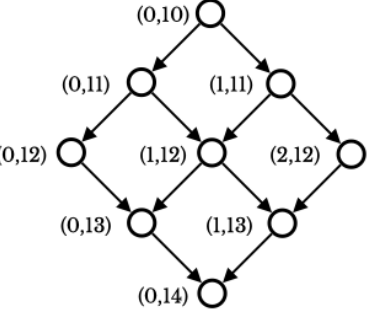
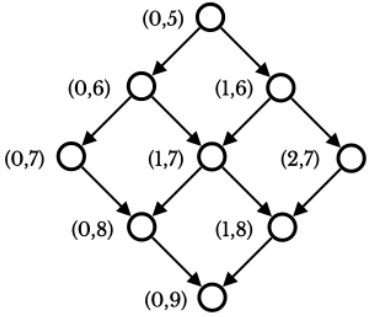
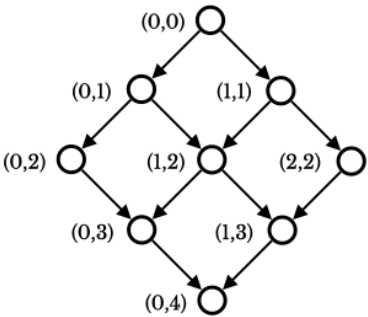
12 nodes, 5 rounds



# Co-Scheduling Many Task Groups

Sequential scheduling of four 3x3 2D Grid DAGs

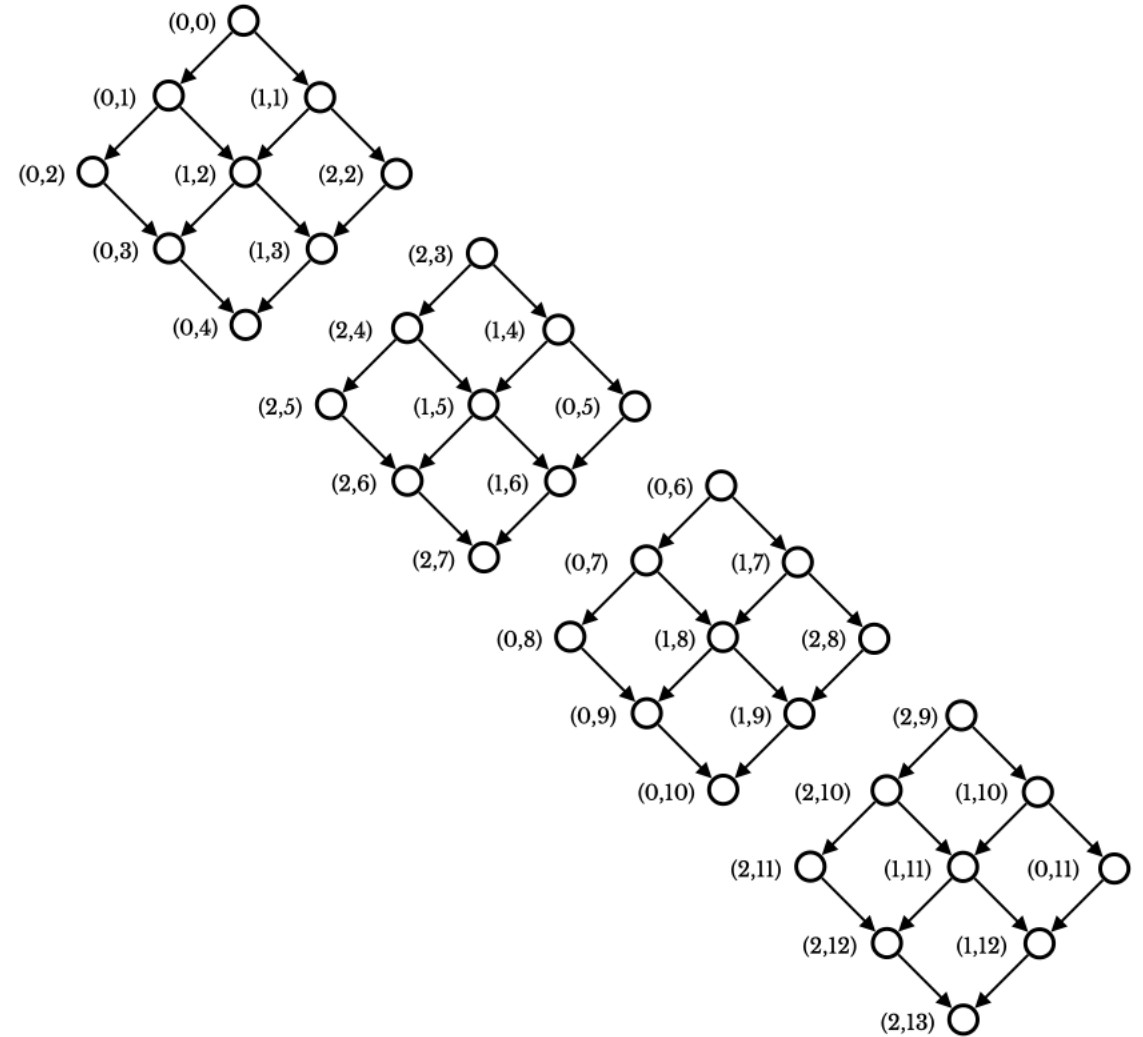
3 nodes, 20 rounds



# Co-Scheduling Many Task Groups

1-pipeline scheduling of four 3x3 2D Grid DAGs

3 nodes, 14 rounds



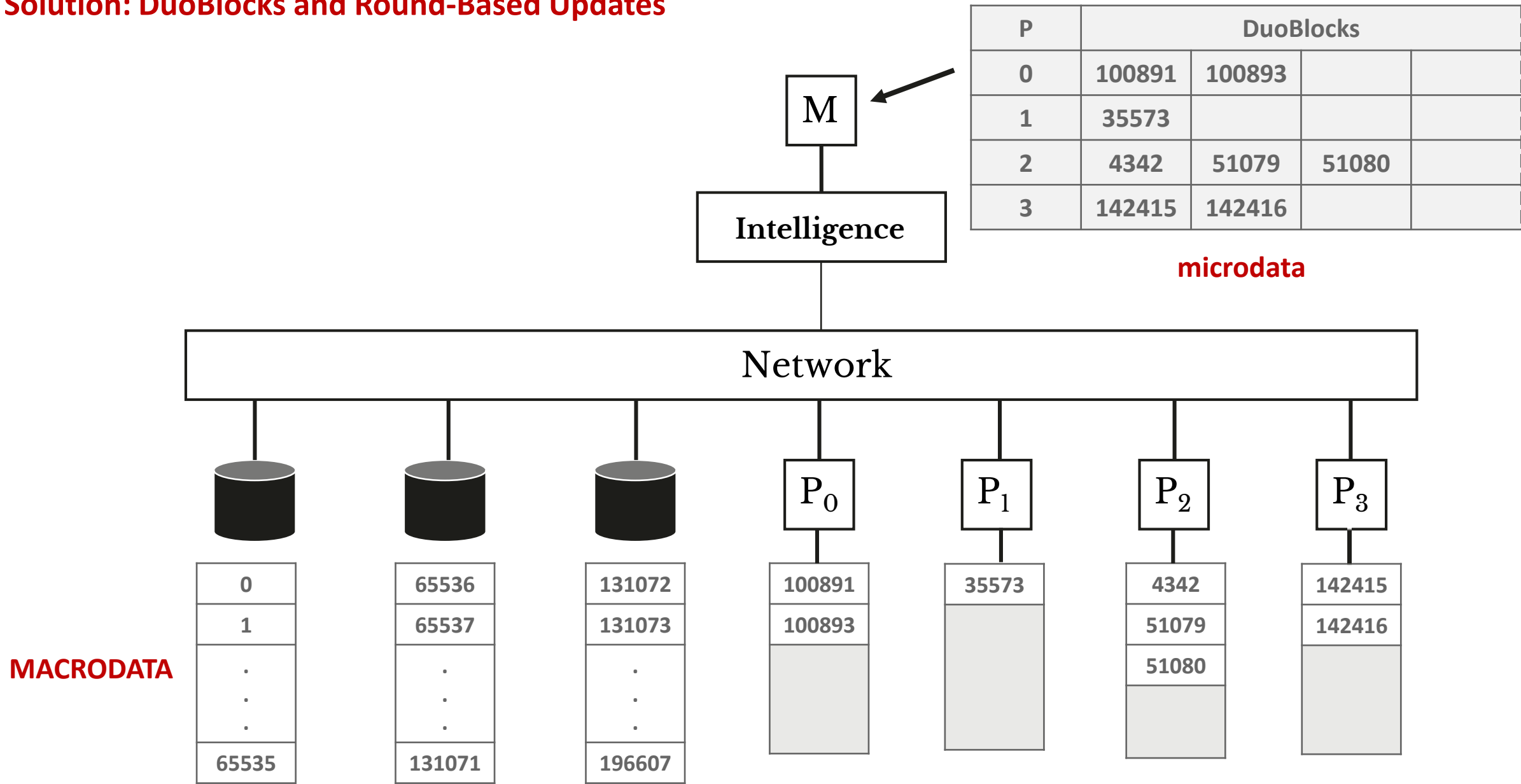
# Scheduling an App and Co-Scheduling Many Apps?

- **Computation load and balance is important**
  - **Communication load and balance is important**
  - **Minimizing local memory size required is important**
  - **Optimizing for priorities, deadlines and cost is important**
  - **There are often unavoidable tradeoffs between these objectives**
- 
- **The new model and architecture enables a single algorithm/program to be developed that is optimal for any scale of architecture, and to be automatically and optimally scheduled for that architecture**

# Shared Data?

- **Massive dynamic (read-write) shared datasets**
  - Key-Value stores
  - Parameter servers in ML
  - Huge geo, bio and other scientific datasets
  - Huge knowledge bases
- **Research Challenge**
  - Parallel computations requiring random access to massive shared dynamic datasets
  - Cannot use shared memory at scale
  - Cannot broadcast and re-broadcast all of dataset after updates
  - Non-oblivious access rules out simple sharding
- **However!**
  - For most apps, each compute node needs only small parts of the massive MACRODATA at any one time

Solution: DuoBlocks and Round-Based Updates





ZURICH RESEARCH CENTER