# From Pascaline to "Piz Daint" in the Alps infrastructure: A Modern-Day View of Computing in Science

Thomas C. Schulthess

# "Piz Daint," CSCS' current flagship supercomputer
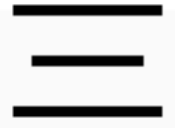


Introduced in 2013 and since 2017 features 5,704 NVIDIA P100 GPU accelerators, dubbed "Pascal"

# World's Most Powerful AI-Capable Supercomputer?

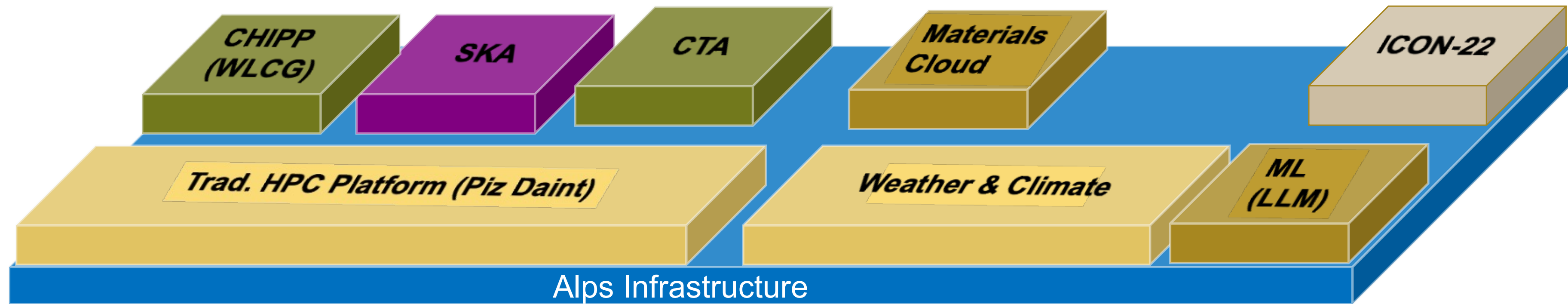## CSCS, Hewlett Packard Enterprise and NVIDIA Announce World's Most...

12.04.2021

"Alps" system to advance research across climate, physics, life sciences with 7x more powerful AI capabilities than...

MORE

MORE SCIENCE

# "Piz Daint" in the "Alps" Infrastructure

To a particular community, a platform will look like a dedicated supercomputer



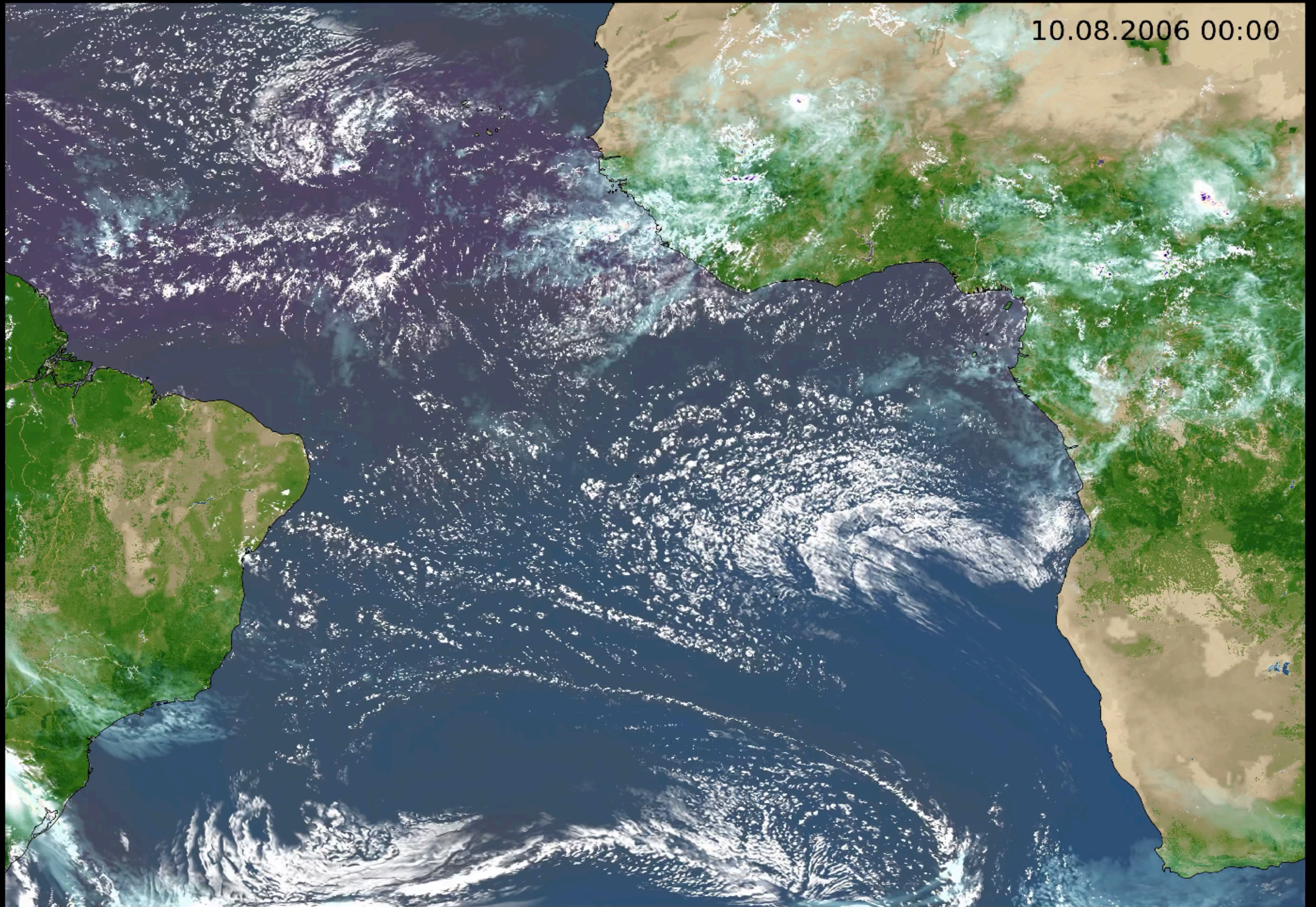Slingshot network
Cray System Management software (μ-service arch.)

"Supercomputers are by definition the fastest and most powerful general-purpose scientific computing systems available at any given time."

–Dongarra et al. in "Numerical Linear Algebra for High-Performance Computers,"
SIAM 1998.

ETH zürich

10.08.2006 00:00

About 300x10⁶ grid points

Christoph Heim, ETH Zürich

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

6

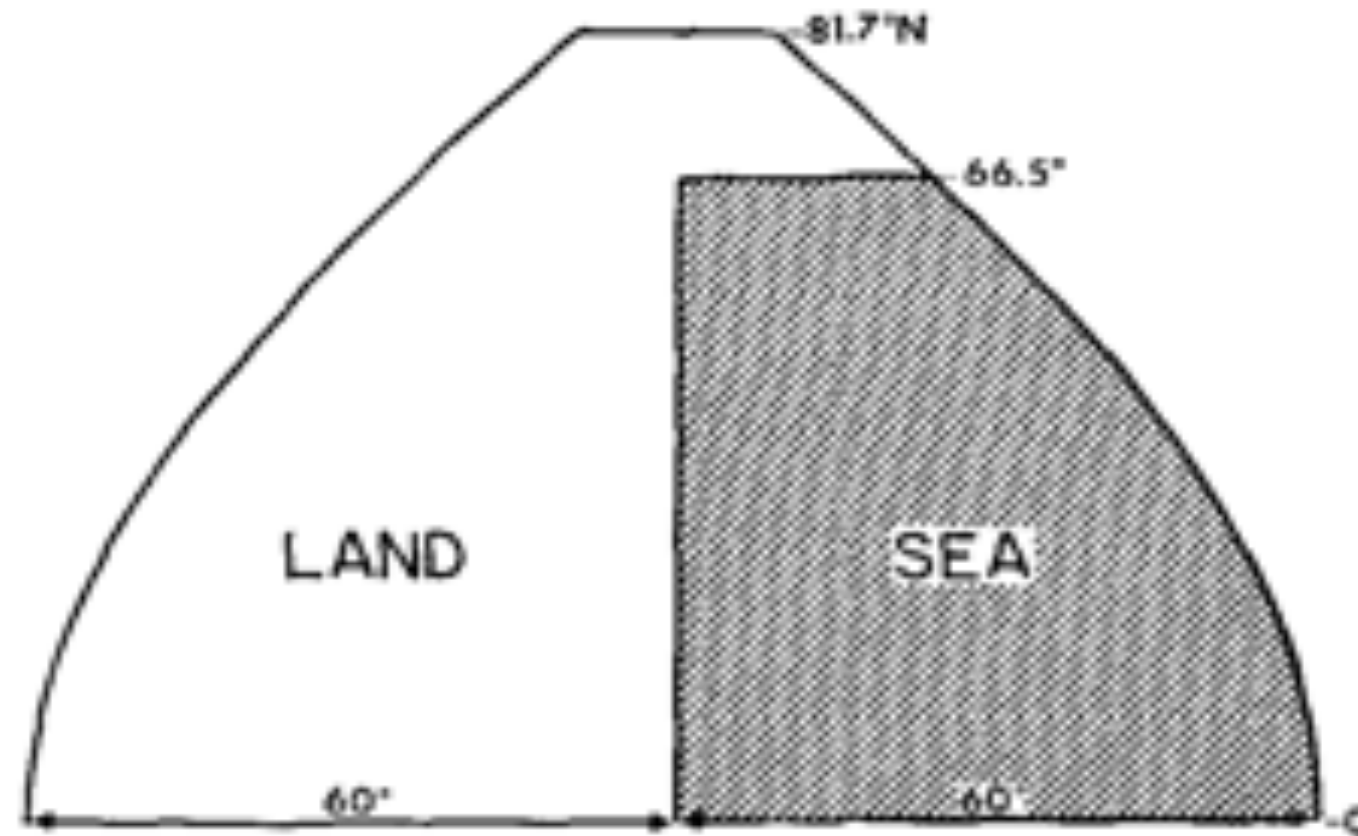# Frist 2 x CO2 general circulation model experiment

*Manabe and Wetherald (1975): The Effects of Doubling the $CO_2$ Concentration on the Climate of a General Circulation Model. J. Atmos. Sci., Vol. 32(1), 3-15*

**Main simplications:**

‣ Atmosphere only
‣ Idealised distribution of land and sea, not global (only 120° longitude & periodic)
‣ Lateral resolution about 450 km with 9 layers (about 20 x 34 x 9 = 5220 grid points)
‣ No seasonal cycle, no diurnal cycle
‣ Prognostic water vapour and snow, bucket model over land, specified clouds
‣ Integrated for only a few 100 days

Nobel Prize 2021
Hasselmann, Manabe & Parisi



**Notable results:**

‣ Equilibrium climate sensitivity of **2.9°C**

| Change of $CO_2$ content (ppm) | R-W model | M-W model | G-C model |
|---|---|---|---|
| 300 → 600 | +1.95 | +2.36 | +2.93 |
| | 1D radiative-convective equilibrium models | | General circulation model |

‣ Polar amplification
‣ Intensification of hydrological cycle
‣ Weakening of extratropical storm tracks

# Scaling of computing and model performance

IBM Stretch 7030 (1961-1982)

About 1 MFLOPs = $10^6$ ops / s

Piz Daint (CSCS, Lugano, 2013-2024)

20 PFLOPs = $20 \times 10^{15}$ ops / s (2017-2024)

Manabe & Weatherald GCM
$\Delta x$ = 450 km
9 layers

Today's runs?
$\Delta x$ = 500 m
150 layers

Scale to today's computers
assuming optimal use of hardware

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# Evolution of computing systems and model capability at ECMWF

**Computational power drives spatial resolution**

Schulthess et al., Comp. Sci. Eng. 21 (1), 31-40 (2018)

# A Bifurcation in Moore's Law?



Nick Trefethen, SIAM News, September 01, 2023

Jim Gray on eScience:
A Transformed Scientific Method

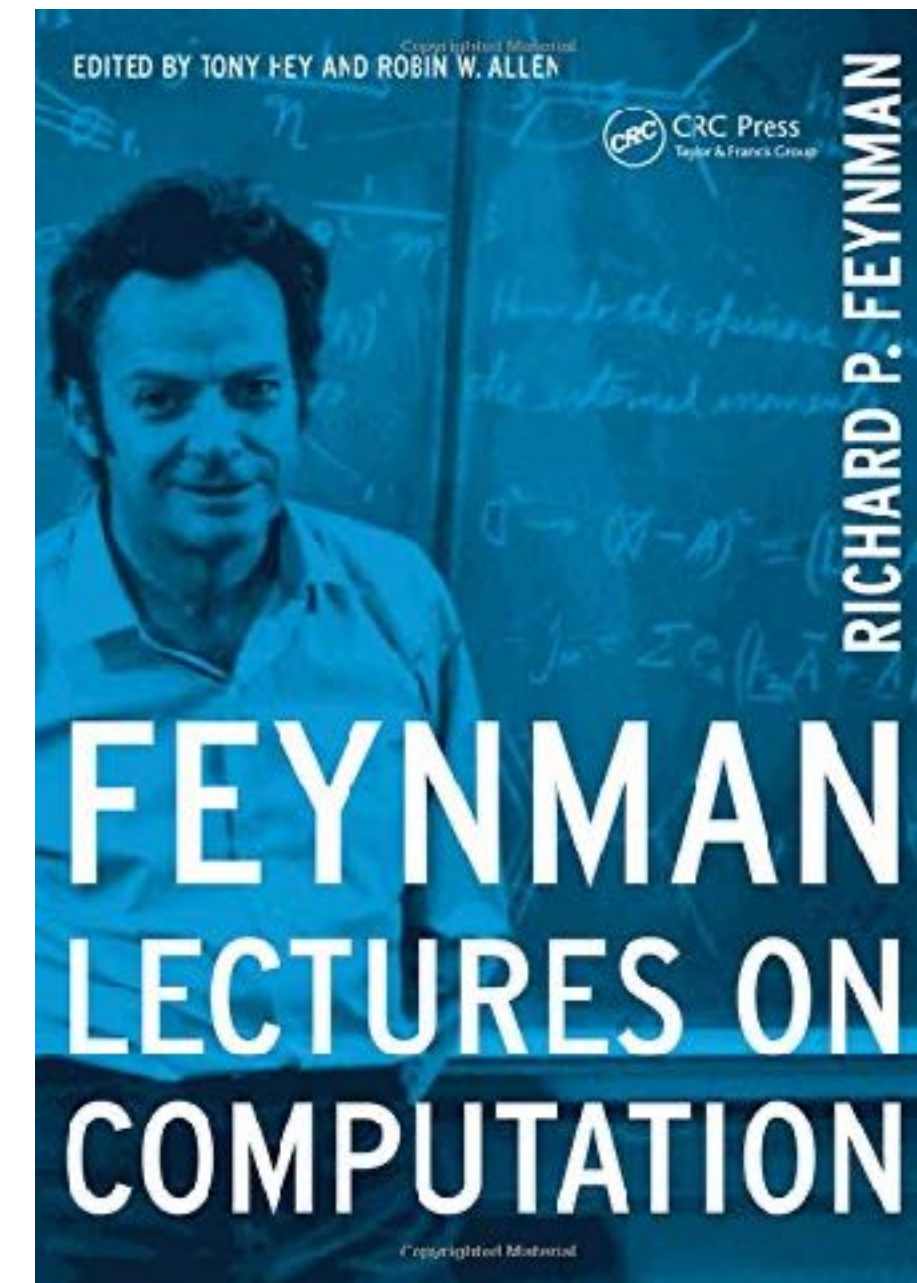National Research Council's Computer Science and
Telecommunications Board, Jan. 11, 2007

# Characteristics of BigData Analytics

Important considerations when dealing with digital data:

$\left\{\begin{array}{l}\end{array}\right.$

1. Velocity
2. Volume
3. Variety
4. Veracity
5. Value

Experiment

Theory

Simulation

BigData Analytics

The FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

"I've studies all available charts of the planets and stars and none of them match the others. There are just as many measurements and methods as there are astronomers and all of them disagree. What's needed is a long term project with the aim of mapping the heavens conducted from a single location over a period of several years."

–Tycho Brahe, 1563

*Experiment ( observation ),Veracity and Variety*

The **FOURTH PARADIGM**

DATA-INTENSIVE SCIENTIFIC DISCOVERY

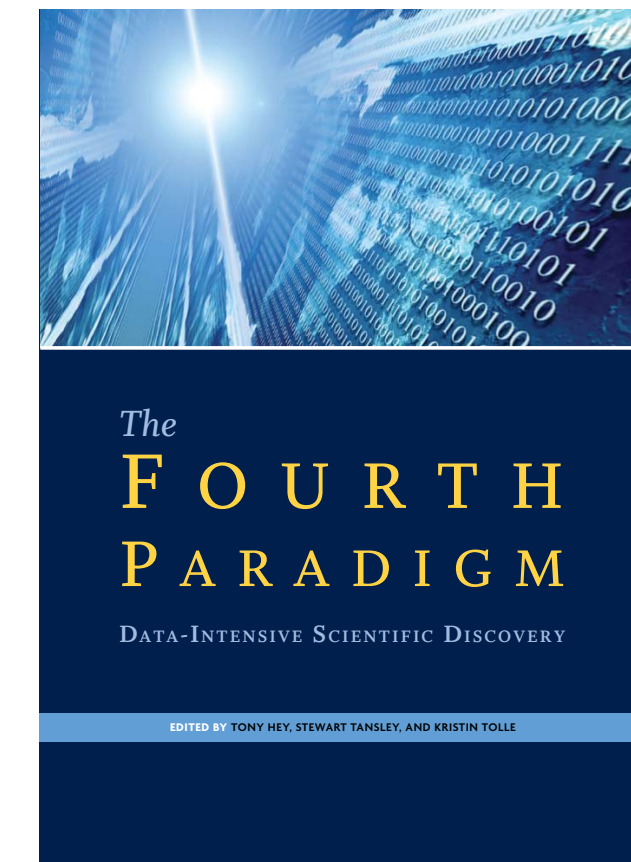EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

Tycho Brahe's Mars Observations

source:

$$M = E - \varepsilon \sin E$$

1. Solve $E(M)$ (Numerics)
2. Solve $\phi(E)$ (Geometry)

source: www.pafko.com/tycho/

# Data Analytics

Tycho Brahe's Mars Observations
The Orbit as Calculated with Modern Methods

source: www.pafko.com/tycho/

# Theory

The
# FOURTH
# PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

Jean Joseph Le Verrier predicts existence and position of Neptune to within 1º, confirmed by Johann Galle on 09/23/1846



Similar predictions around the same time by John Couch Adams

*Simulation?*

# Richardson's forecast factory (1922)



Lewis Fry Richardson:
Weather Prediction by Numerical Process

# *Bulk synchronous parallel (BSP) computing model*

John von Neumann with first "electronic computing instrument" that was built at Princeton's IAS between 1946 and 1952

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# European Center for Medium-Range Weather Forecasts



adapted from Schär, ETH Zürich

**ECMWF**

An independent intergovernmental organization established in 1975

Switzerland was founding member of ECMWF among 18 countries

Today the worldwide leading numerical weather prediction center

Provides input data for the weather predictions of MeteoSwiss

*All of the above + Velocity, Volume and Value*

The
**FOURTH**
**PARADIGM**
DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY: TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# MeteoSwiss' Performance Ambitions in 2013



We need a 40x improvement between 2012 and 2015 at constant cost

# Porting codes to GPUs, Xeon, ARM, etc.

CUDA (C / C++ / Fortran)

OpenCL

```
 8   __global__ void add_pw_ekin_gpu_kernel(int num_gvec__,
 9                                           double alpha__,
10                                           double const* pw_ekin__,
11                                           cuDoubleComplex const* phi__,
12                                           cuDoubleComplex const* vphi__,
13                                           cuDoubleComplex* hphi__)
14   {
15       int ig = blockIdx.x * blockDim.x + threadIdx.x;
16       if (ig < num_gvec__) {
17           cuDoubleComplex z1 = cuCadd(vphi__[ig], make_cuDoubleComplex(alpha__
18                                                                         alpha__
19           hphi__[ig] = cuCadd(hphi__[ig], z1);
20       }
21   }
```

```
13   __kernel void vector_add(const int n, __global float *a, __global float *b, __global float *c) {
14       const int i = get_global_id(0);

         + b[i];
```

OpenACC

OpenMP 4.x

```
76       acc = 0
77       !$acc parallel present(x)
78       !$acc loop reduction(+:acc)
79       do i = 1, N
80           acc = acc + x(i) * x(i)
81       enddo
82       !$acc end parallel
83       call mpi_allreduce(acc, accglobal, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD, err)
```

```
omp target data map(tofrom: x[0:n],y[0:n])

#pragma omp target
#pragma omp for
for (int i = 0; i < n; i++)
    y[i] += a * x[i];
}
```

# COSMO: old and new (refactored) implementation



* two different OpenMP backends

# Where the factor 40 improvement came from

**Investment in software allowed mathematical improvements and change in architecture**



Requirements from MeteoSwiss

6x

Data assimilation

24x

Ensemble with multiple forecasts

10x

Grid 2.2 km → 1.1 km

Constant budget for investments and operations

1.7x from software refactoring (old vs. new implementation on x86)

2.8x algorithmic improvements (resource utilisation, mixed arithmetic precision)

2.3x change in architecture (CPU → GPU)

Bonus: reduction in power!

2.8x Moore's Law & arch. improvements on x86

1.3x additional processors

**There is no silver bullet!**

# Setting a new baseline for atmospheric simulations

The state-of the art implementation of COSMO running at
most weather services on multi-core hardware.

**~10x**

The refactored version of COSMO running at MeteoSwiss
on multi-core or GPU accelerated hardware.

COSMO at $\Delta x = 50$ km
189 x 142 x 60 = 1.6 x 10$^6$ grid points
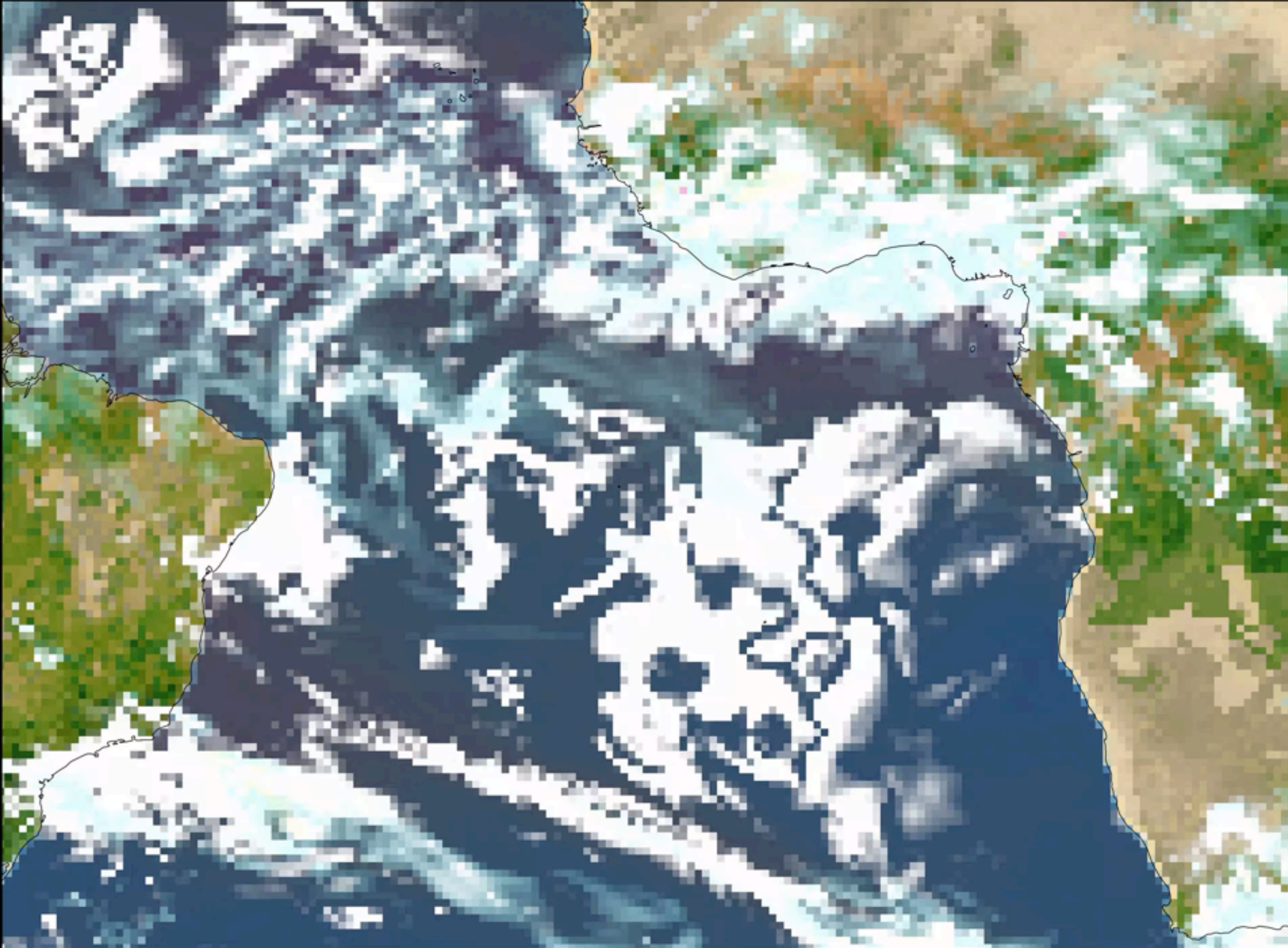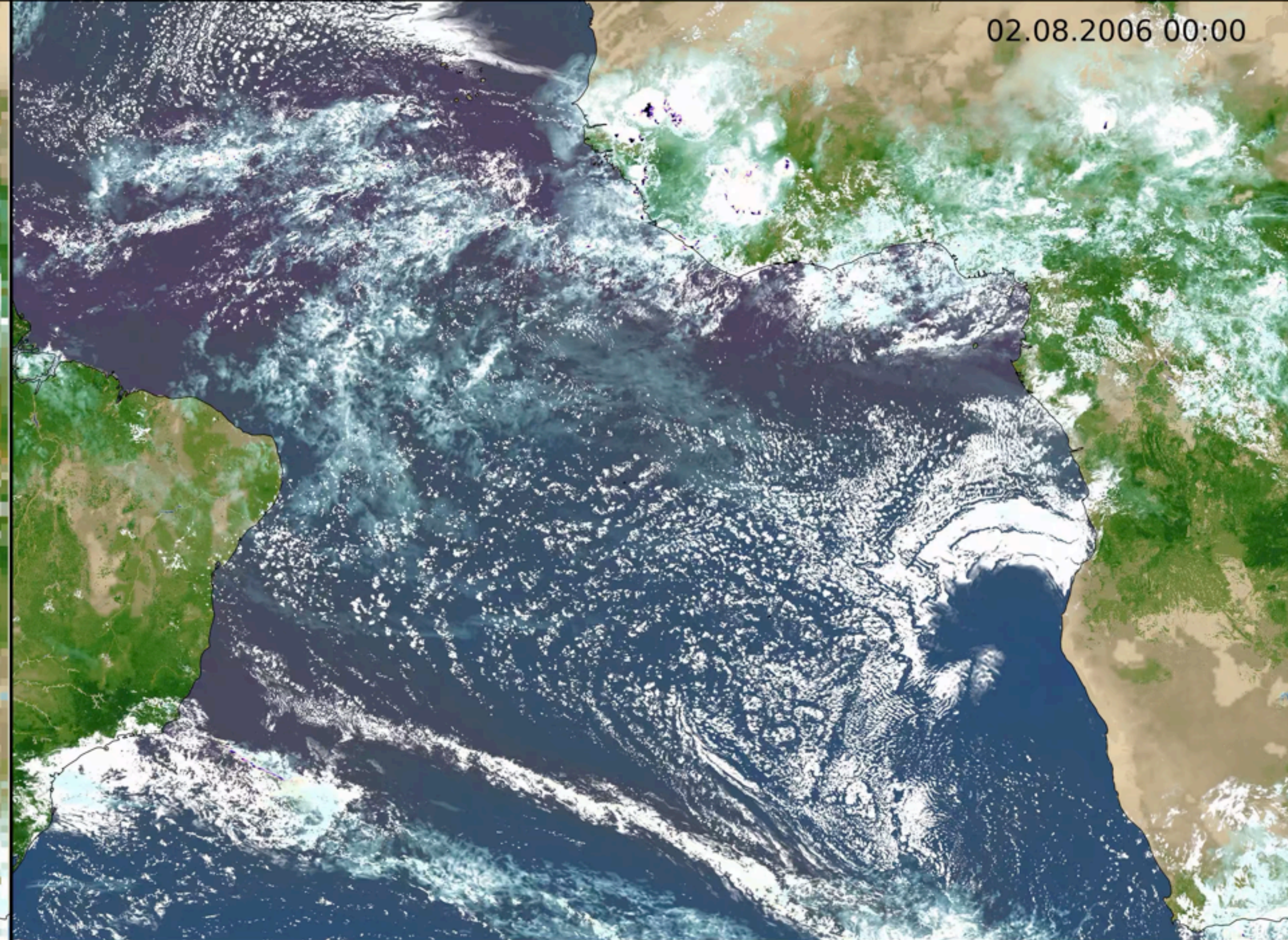
COSMO at $\Delta x = 3$ km
2750 x 2065 x 60 = 340 x 10$^6$ grid points

02.08.2006 00:00

Christoph Heim, ETH Zürich

# "Exascale:" our goal for 2024-2026 climate applications runs

| | |
|---|---|
| Horizontal resolution | 1 km (globally quasi-uniform) |
| Vertical resolution | 180 levels (surface to ~100 km) |
| Time resolution | Less than 1 minute |
| Coupled | Land-surface/ocean/ocean-waves/sea-ice |
| Atmosphere | Non-hydrostatic |
| Precision | Single (32bit) or mixed precision |
| Compute rate | 1 SYPD (simulated year wall-clock day) |

Schulthess, P. Bauer, N. Wedi, O. Fuhrer, Th. Hoefler, Ch. Schär, Comp. Sci. Eng. 21 (1), 31-40 (2018)

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# Baseline in 2018: Running COSMO & IFS ("the European Model") at global scale on "Piz Daint"

Scaling to full system size: ~5300 GPU accelerate nodes available



Running a near-global (±80° covering 97% of Earths surface) COSMO 5.0 simulation & IFS
> Either on the hosts processors: Intel Xeon E5 2690v3 (Haswell 12c).
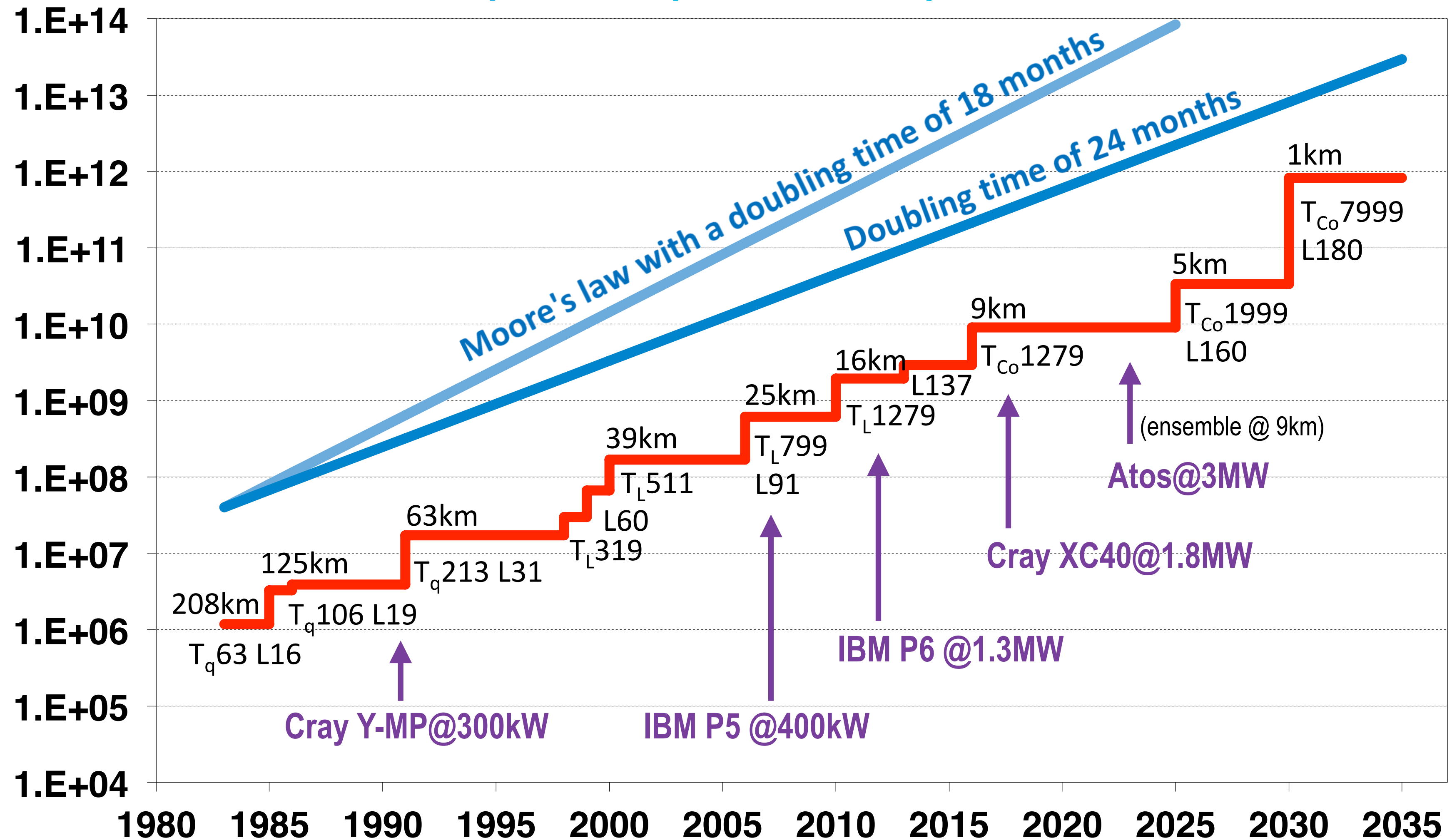> Or on the GPU accelerator: PCIe version of NVIDIA GP100 (Pascal) GPU

# The baseline for COSMO near-global and IFS

| | Near-global COSMO[15] | | Global IFS[16] | |
|---|---|---|---|---|
| | Value | Shortfall | Value | Shortfall |
| Horizontal resolution | 0.93 km (non-uniform) | 0.81× | 1.25 km | 1.56× |
| Vertical reso-lution | 60 levels (surface to 25 km) | 3× | 62 levels (sur-face to 40 km) | 3× |
| Time resolu-tion | 6 s (split-explicit with sub-stepping)* | – | 120 s (semi-implicit) | 4× |
| Coupled | No | 100x (single trajectory) times 50x (ensemble) | | 1.2× |
| Atmosphere | Non-hydrostatic | – | Non-hydro-static | – |
| Precision | Single | – | Single | – |
| Compute rate | 0.043 SY | Goal is to stay within ~ 5MW 3× | 0.088 SYPD | 11× |
| Other (e.g., physics, …) | microphysics | 1.5× | Full physics | – |
| Total short-fall | | 101× | | 247× |

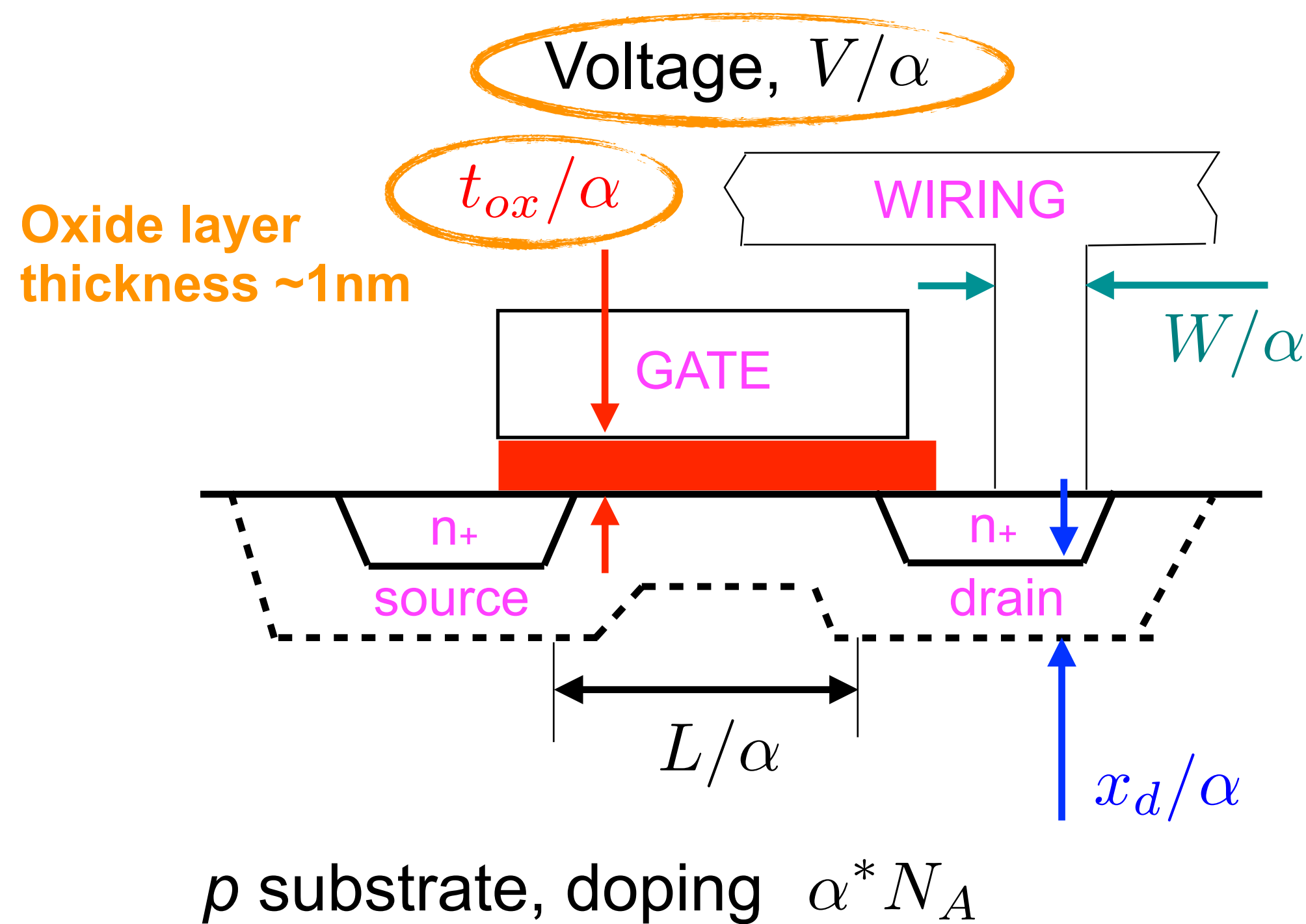Schulthess et al., Comp. Sci. Eng. 21 (1), 31-40 (2018)

# Evolution of computing system and model capability at ECMWF

**Computational power drives spatial resolution**

# The end of Dennard Scaling

Robert H. Dennard (1974)



Voltage, $V/\alpha$

$t_{ox}/\alpha$

**Oxide layer thickness ~1nm**

WIRING

GATE

$W/\alpha$

n+ source

n+ drain

$L/\alpha$

$x_d/\alpha$

*p* substrate, doping $\alpha^* N_A$

SCALING

~~Voltage:~~ $V/\alpha$
~~Oxide:~~ $t_{ox}/\alpha$
Wire width: $W/\alpha$
Gate Width: $L/\alpha$
Diffusion: $x_d/\alpha$
Substrate: $\alpha^* N_A$

CONSEQUENCE:
Higher density: $\sim \alpha^2$
Higher speed: $\sim \alpha$
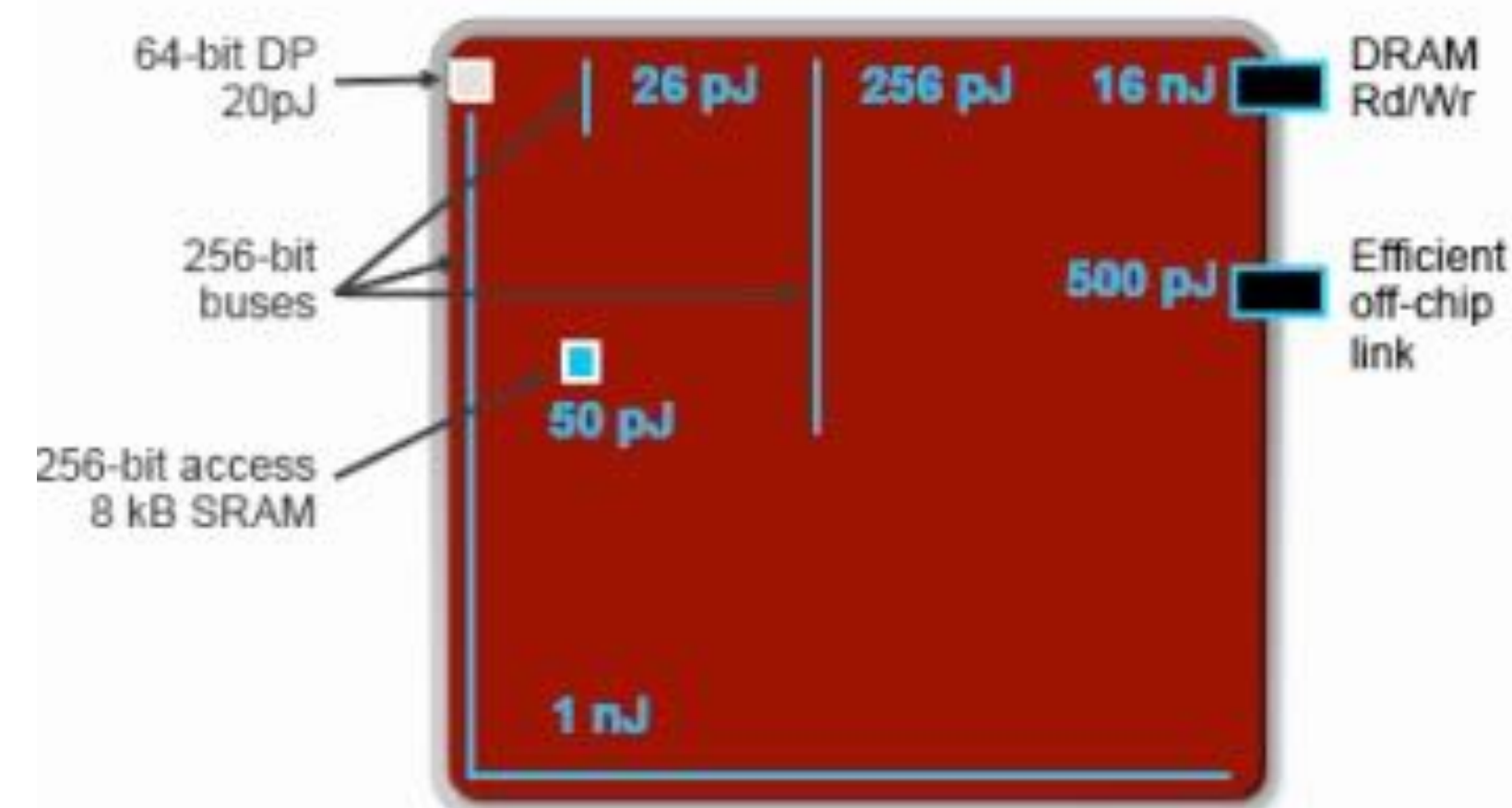Power/ckt: $\sim 1/\alpha^2$
~~Power density: $\sim$ constant~~

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# Who consumes how much energy (28nm)

- 64 bit floating point unit: 20 pJ

- 256-bit access 8kB SRAM:     50 pJ

- 256-bit bus across die:  1,000 pJ

- Read/write to DRAM:     16,000 pJ



Source: Bill Dally, 2011

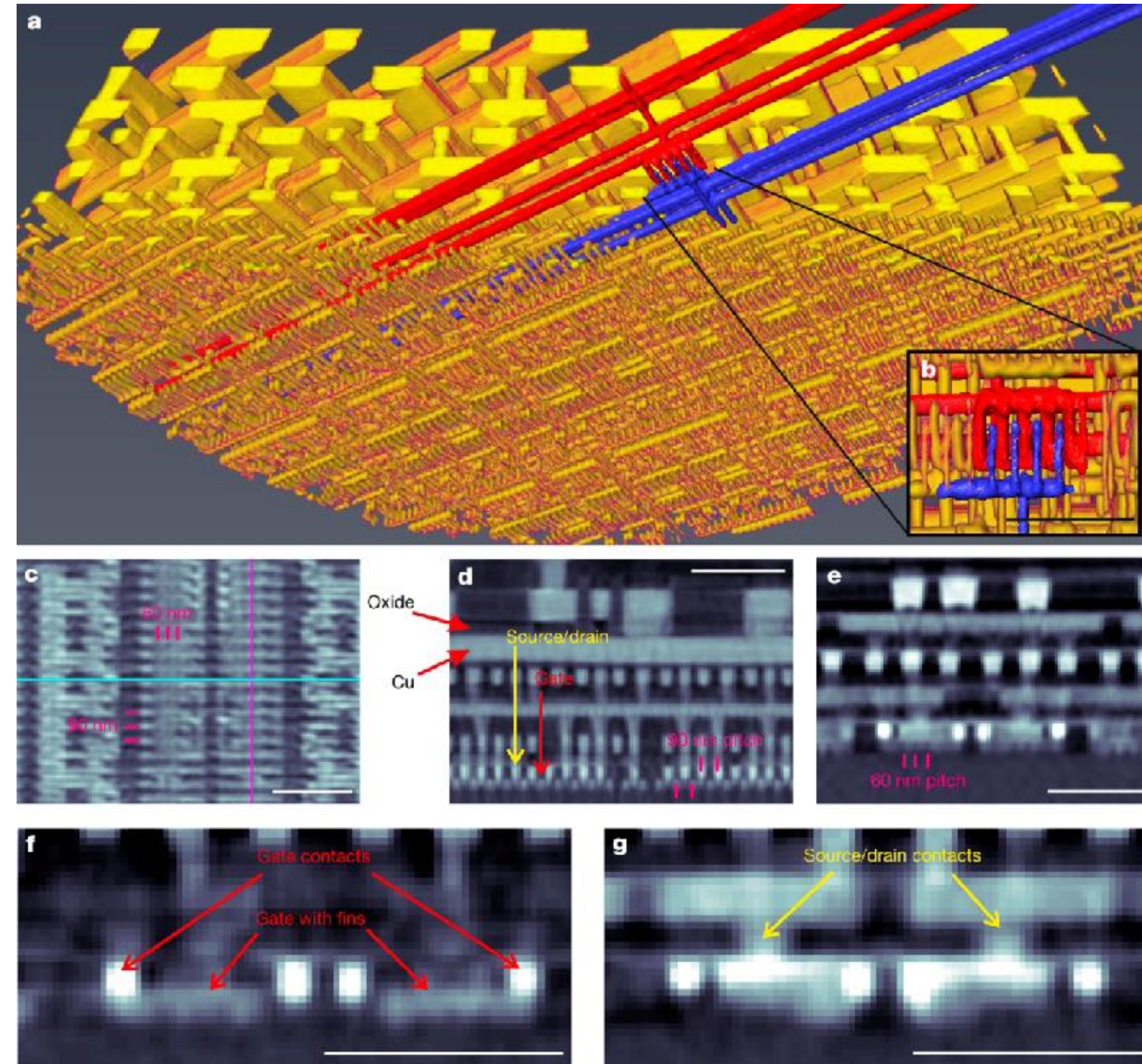**By a wide margin, most energy is spend in moving data on the die and to memory**

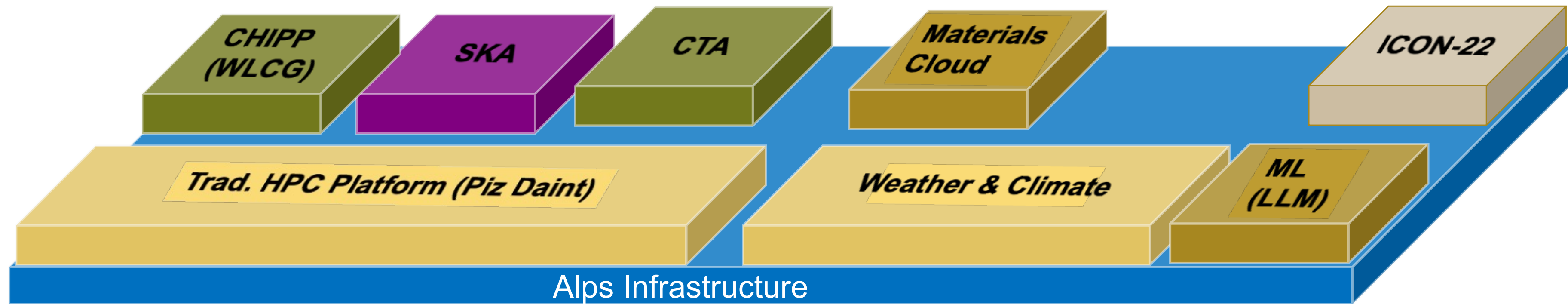**Developing algorithms that maximise data locality should be THE TOP PRIORITY**



20 mm

M Holler *et al. Nature* **543,** 402–406 (2017) doi:10.1038/nature21698

# "Piz Daint" in the "Alps" Infrastructure

To a particular community, a platform will look like a dedicated supercomputer

# Preliminary comparison A100 vs. GH200(*)

(*) A02 engineering samples



ICON benchmark

|        | time (s) | power (W) |
|--------|----------|-----------|
| A100   | 2196     | 388       |
| GH200  | 1518     | 570       |
|        | 1.45     | 1.47      |

B. Cumming and W. Sawyer (CSCS)

# Baseline running on "Piz Daint"

| | Near-global (COSMO) | | ICON global | |
|---|---|---|---|---|
| | Value | Shortfall | Value | Shortfall |
| **Horizontal resolution** | 0.93 km (non uniform) | 0.81x | 5 km (uniform) | 25x |
| **Vertical resolution** | 60 levels | 3x | 90 levels | 2x |
| **Time resolution** | 6s (split-explicit with sub-stepping) | — | 40s (split-explicit with sub stepping) | 5x |
| **Couple** | No | 1.2x | No | 1.2x |
| **Atmosphere** | Non-hydrostatic | — | Non-hydrostatic | — |
| **Precision** | Single | — | ? | — |
| **Simulation rate** | 0.043 SYPD | 23x | 0.4 SYPD | 2.5x |
| **Other (e.g. physics, …)** | Microphysics | 1.5x | Full physics | — |
| **Adjusted to 5300 nodes** | 5300 nodes | — | 1000 nodes | 0.19x |
| | | 101 | | 190 |

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# Conclusions

‣Science can greatly benefit from fully embracing digitalisation

‣Moore's Law is fading: computing and data research infrastructures have their cost

‣A multitude of computer architectures is the consequence

‣Continued investments in algorithms and software development is essential

‣Software engineering has to become a first class citizen

ETH zürich

# Thank you to CSCS, partners such as MeteoSwiss, HPE/ Cray, NVIDIA, as well as many colleagues and collaborators

Tim Palmer (U. of Oxford)

Bjorn Stevens (MPI-M)

Peter Bauer (ECMWF)

Oliver Fuhrer  (MeteoSwiss)

Nils Wedi (ECMWF)

Sadaf Alam (U. of Bristol)

Torsten Hoefler (ETH Zurich)

Christoph Schar (ETH Zurich)

CSCS
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

# Thank you