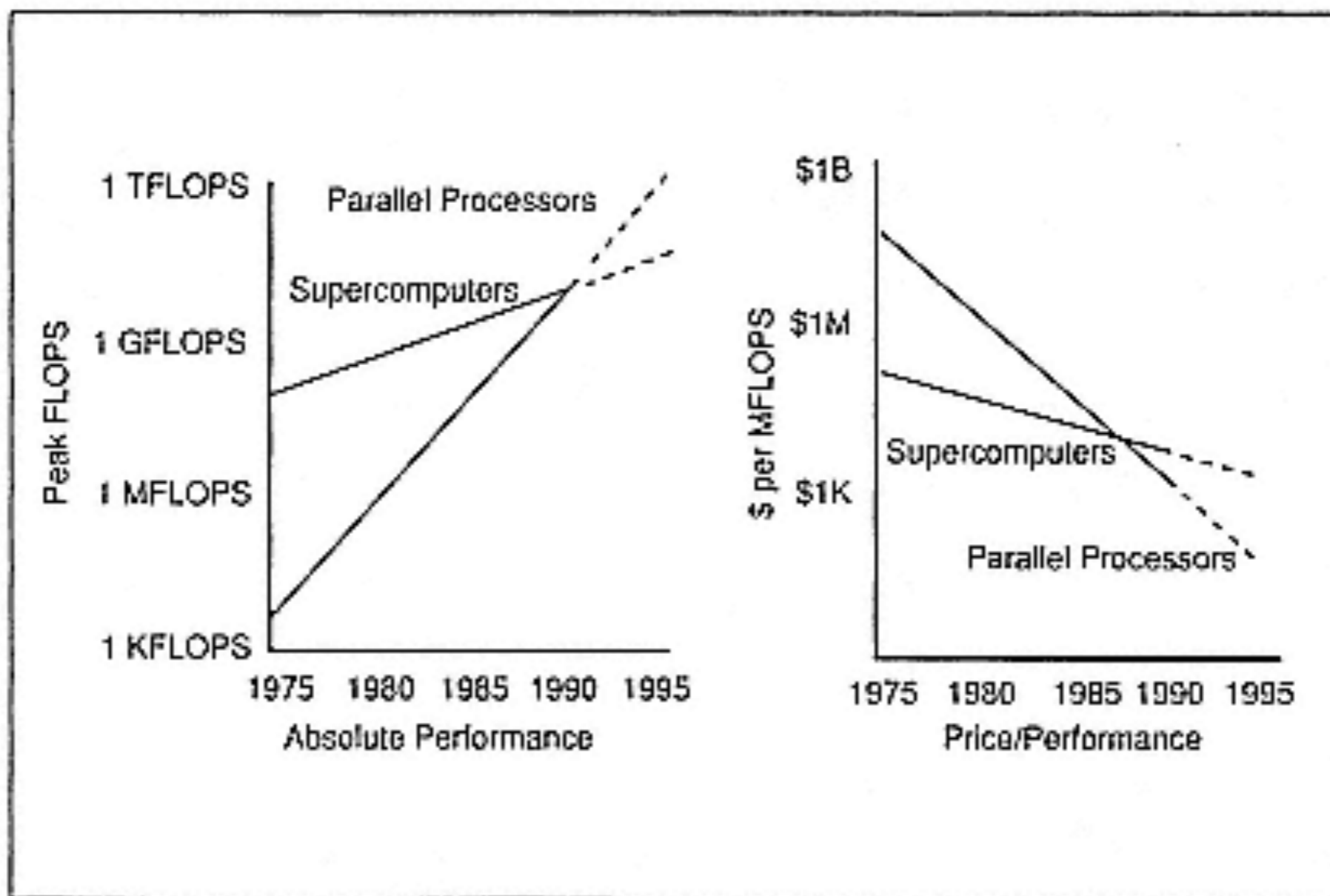




Dusk of Moore's Law: opportunities for weather and climate modelling?

Thomas C. Schulthess

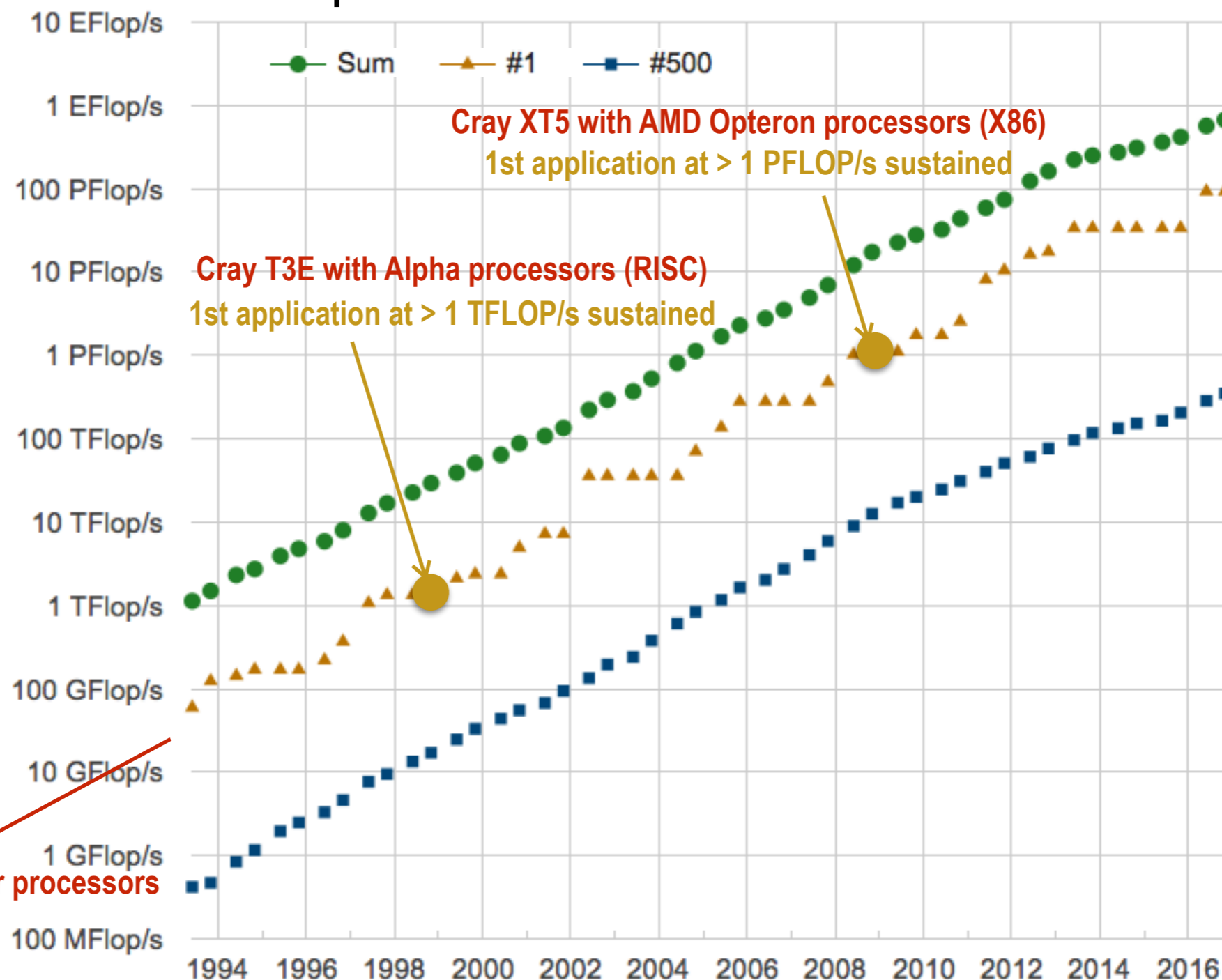
Beginning of change: “Attack of the Killer Micros”



Eugene Brooks (LLNL) @ SC90

The good old days of tera- and petascale computing

Linpack benchmark solves: $Ax = b$



Cray 1, .. X/YMP: vector processors

for the historic development of supercomputing performance, see www.top500.org



\$500,000,000

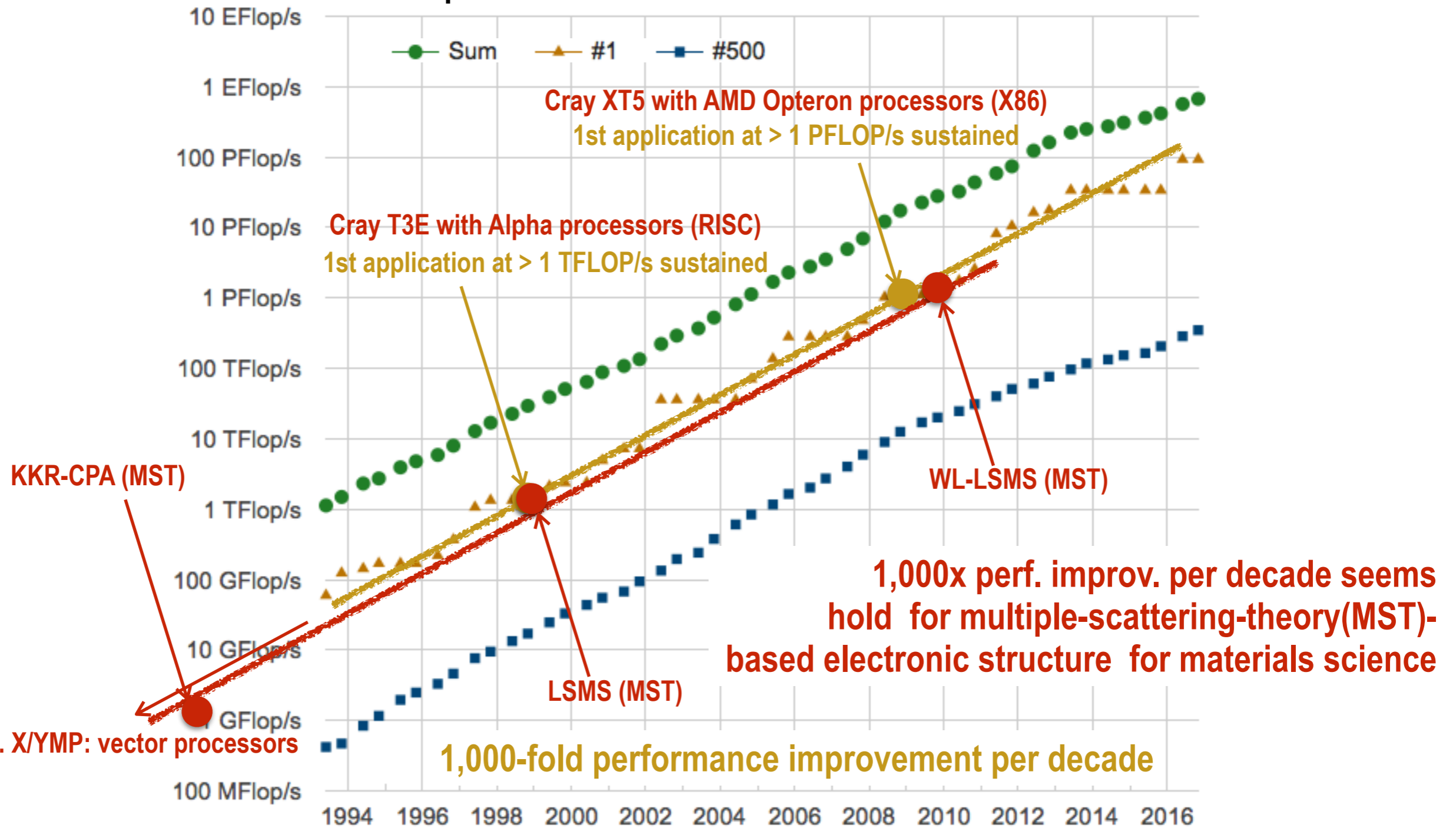
\$2,000,000,000

\$13,000

Source: Andy Keane @ ISC'10

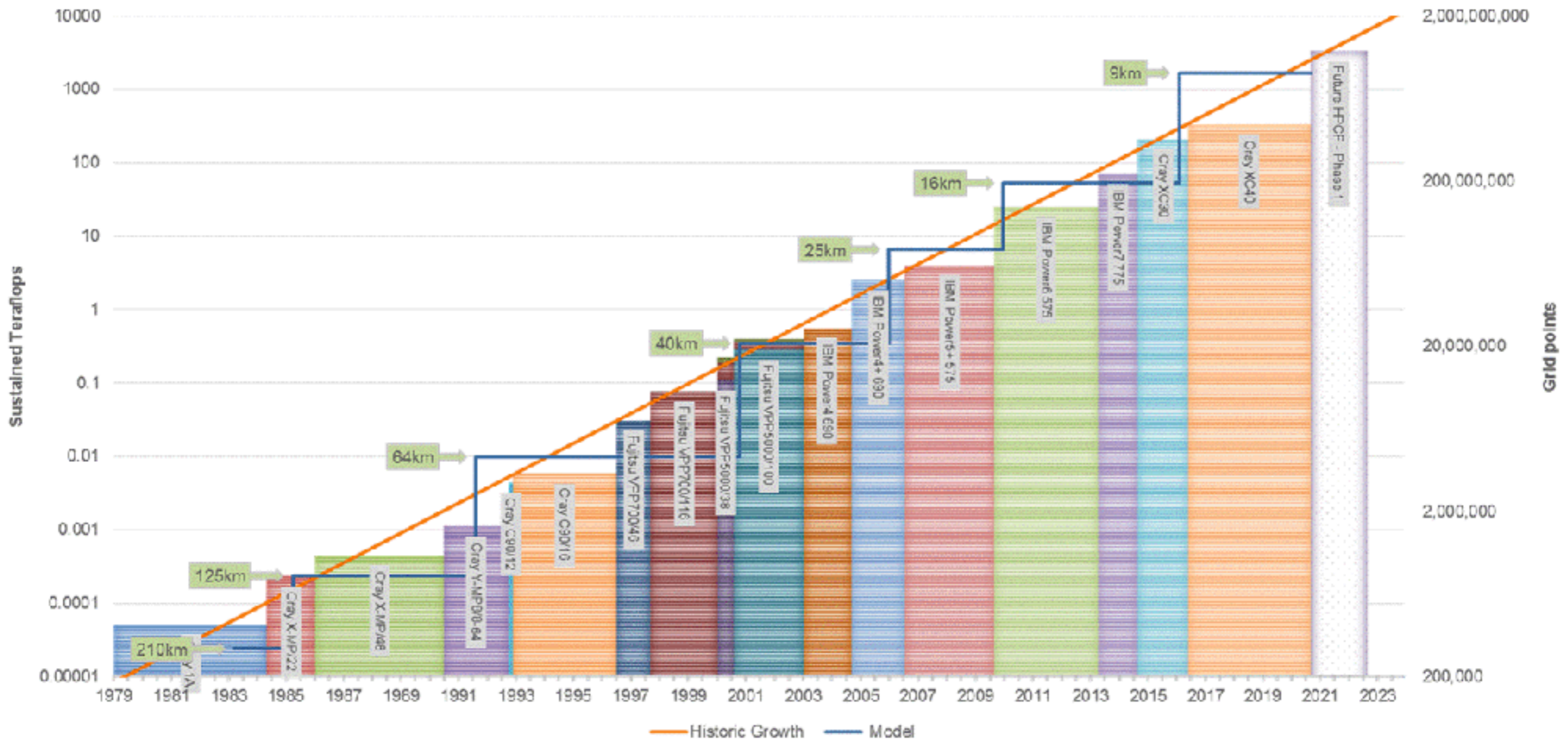
The good old days of tera- and petascale computing

Linpack benchmark solves: $Ax = b$



for the historic development of supercomputing performance, see www.top500.org

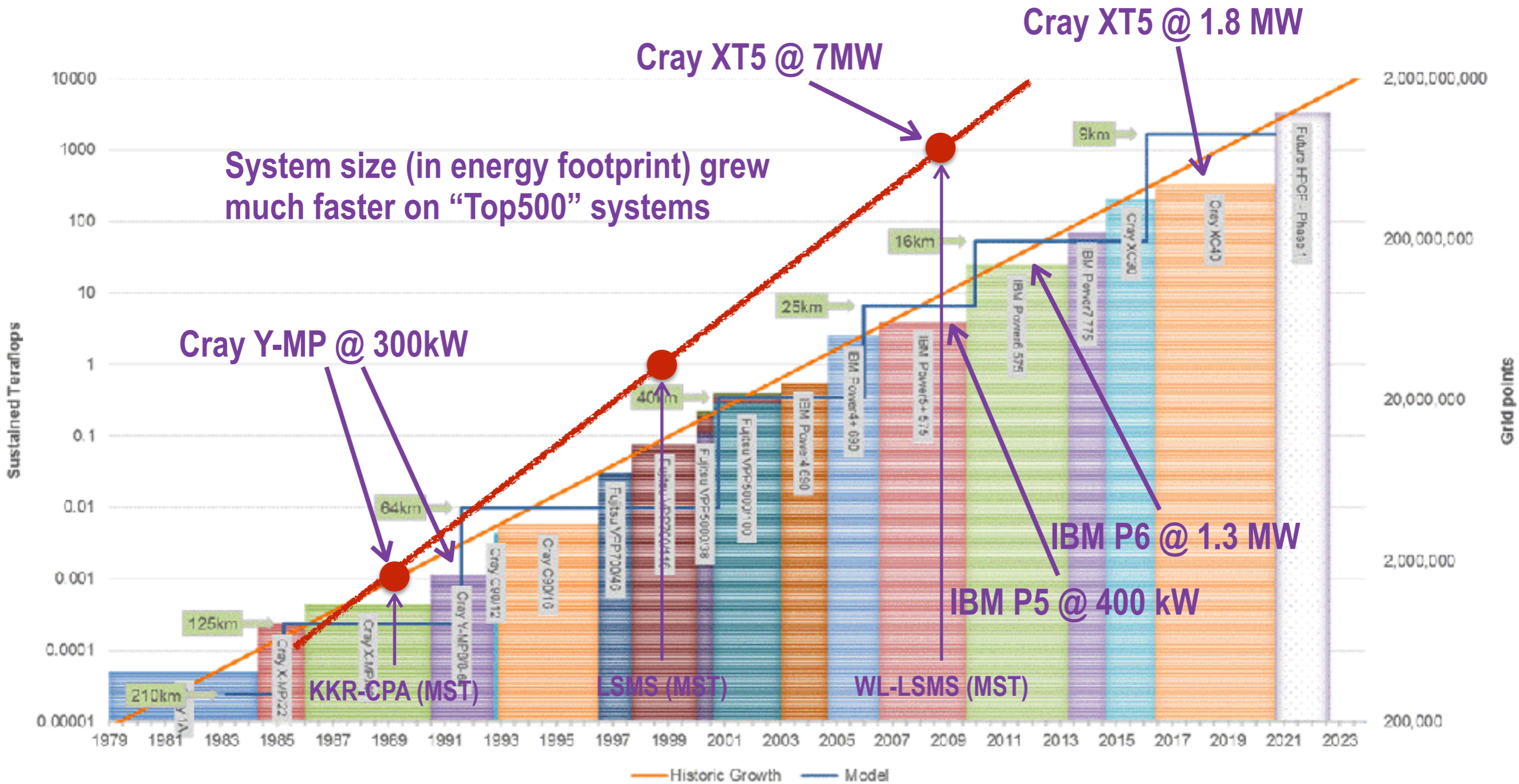
“Only” 100-fold performance improvement in climate codes



Source: Peter Bauer, ECMWF

Has the efficiency of weather & climate codes dropped 10-fold every decade?

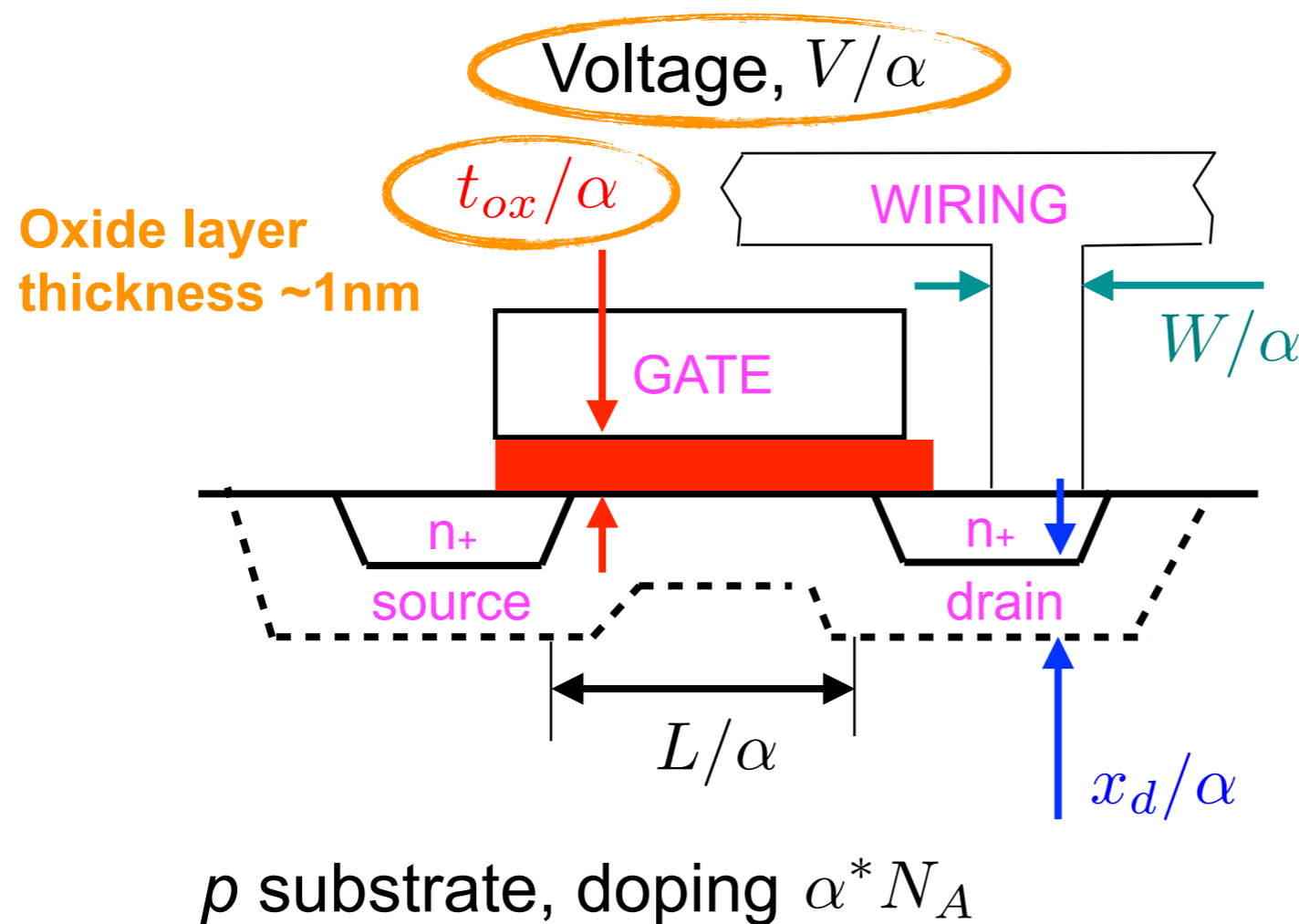
Floating points efficiency dropped from 50% on Cray Y-MP to 5% on today's Cray XC (10x in 2 decades)



Source: Peter Bauer, ECMWF

The end of Dennard Scaling

Robert H. Dennard (1974)



SCALING

~~Voltage: V/α~~
~~Oxide: t_{ox}/α~~

Wire width: W/α
 Gate Width: L/α
 Diffusion: x_d/α
 Substrate: $\alpha^* N_A$

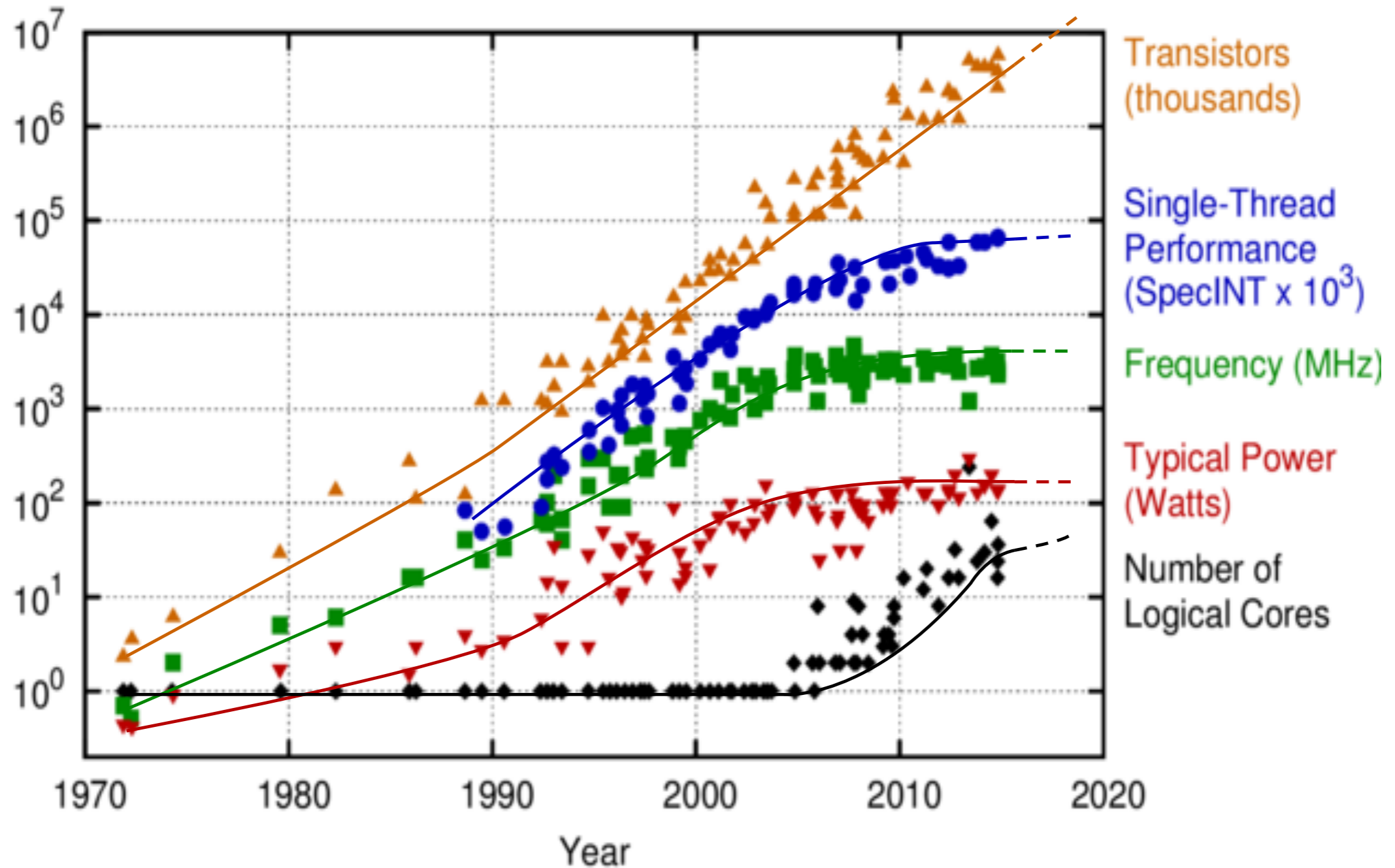
CONSEQUENCE:

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

~~Power density: $\sim \text{constant}$~~

Source: Ronald Luitjen, IBM-ZRL

40 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
 New plot and data collected for 2010-2015 by K. Rupp

Moore's Law 2008-2020

Semiconductor Device Scaling Factors

Technology (High Volume)	45nm (2008)	32nm (2010)	22nm (2012)	16nm (2014)	11nm (2016)	8nm (2018)	5nm (2020)
Transistor density	1.75	1.75	1.75	1.75	1.75	1.75	1.75
Frequency scaling	15%	10%	8%	5%	4%	3%	2%
Voltage (V _{dd}) scaling	-10%	-7.5%	-5%	-2.5%	-1.5%	-1%	-0.5%
Dimension & Capacitance	0.75	0.75	0.75	0.75	0.75	0.75	0.75
SD Leakage scaling/micron	1X Optimistic to 1.43X Pessimistic						

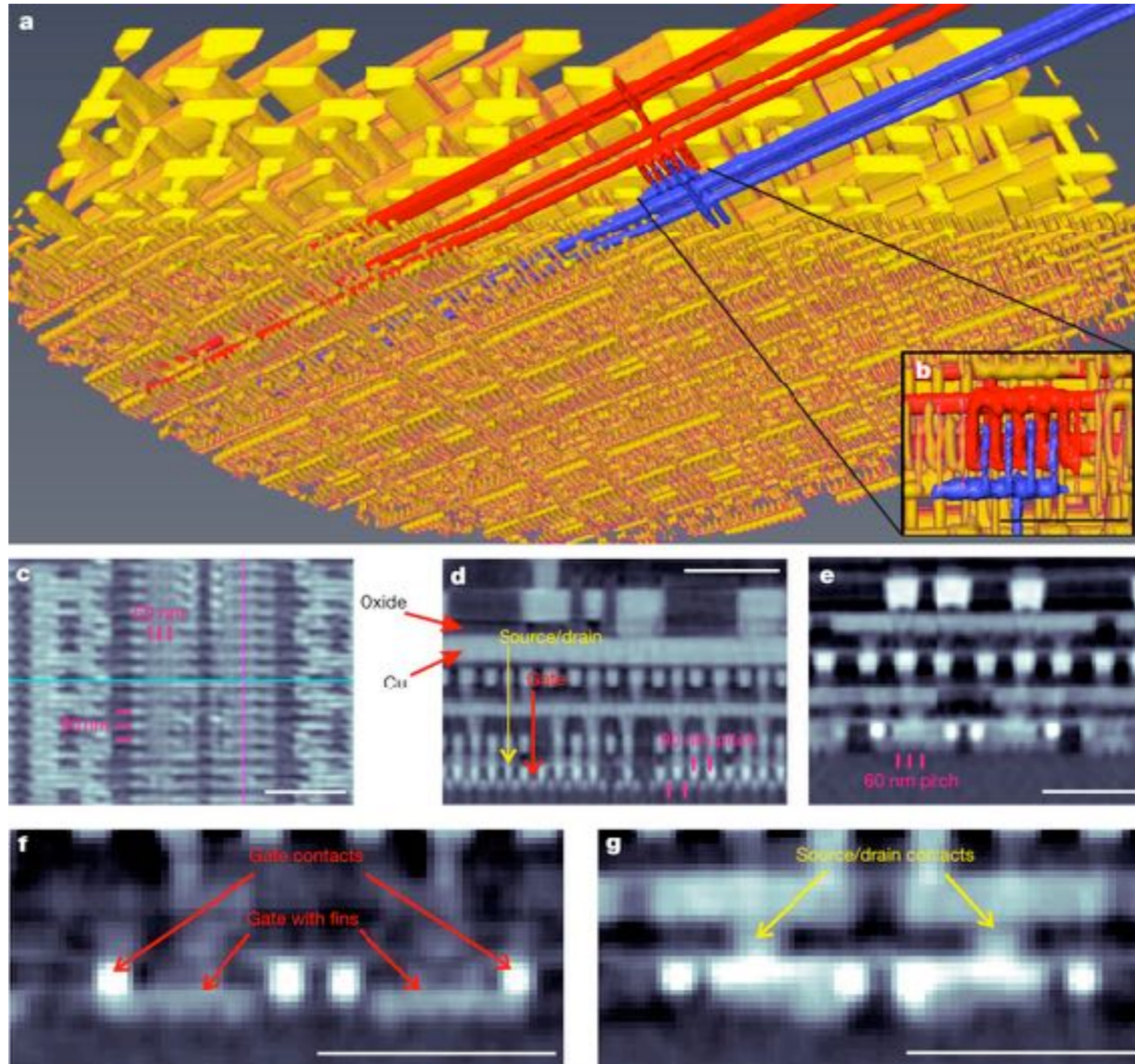
Sources: International Technology Roadmap for Semiconductors and Intel

Moore's Law Takes Miracles ... But It Isn't The Miracle That Will Carry The Day

8

Source: Rajeeb Hazra's (HPC@Intel) talk at SOS14, March 2010

PXCT imaging of Intel processor



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



NVIDIA DGX-1

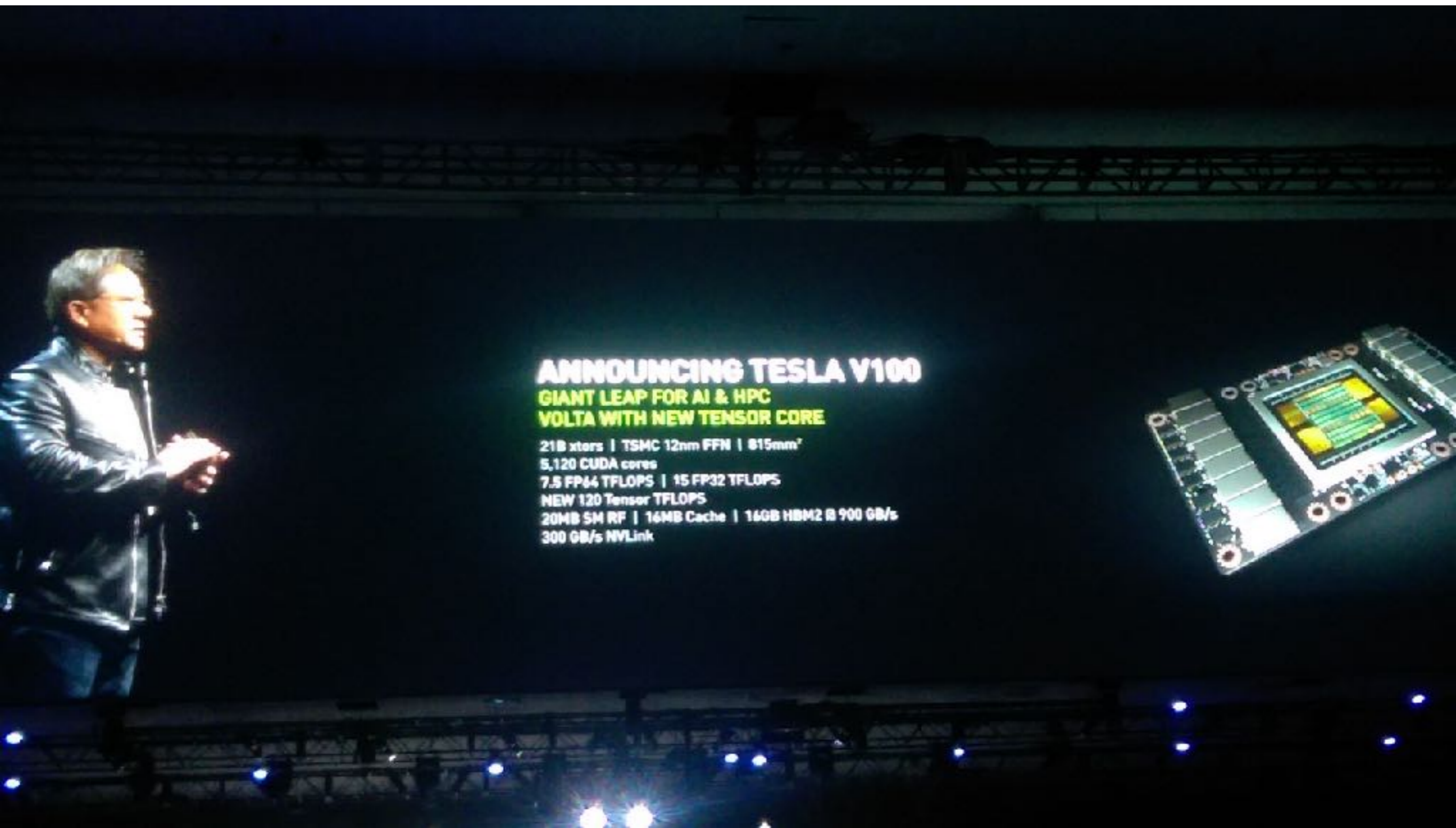
WORLD'S FIRST
DEEP LEARNING SUPERCOMPUTER

170TF | "250 servers in-a-box" | nvidia.com/dgx1

\$129,000 for 8 GPUs, or \$16k a piece



Source: Andy Keane @ ISC'10



ANNOUNCING TESLA V100

**GIANT LEAP FOR AI & HPC
VOLTA WITH NEW TENSOR CORE**

21B xtors | TSMC 12nm FFN | 815mm²
5,120 CUDA cores
7.5 FP64 TFLOPS | 15 FP32 TFLOPS
NEW 120 Tensor TFLOPS
20MB SM RF | 16MB Cache | 16GB HBM2 @ 900 GB/s
300 GB/s NVLink



Porting codes to GPUs, Xeon (Phi), ARM, etc.

CUDA (C / C++ / Fortran)

```

8  __global__ void add_pw_phi_gpu_kernel(int num_gvec__,
9                                     double alpha__,
10                                    double const* pw_ckin__,
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 }

```

OpenACC

```

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)

```

OpenCL

```

13  __kernel void vector_add(const int n, __global float *a, __global float *b, __global float *c) {

```



OpenMP 4.x

```

data map(tofrom: x[0:n],y[0:n])
target
#pragma omp for
for (int i = 0; i < n; i++)
    y[i] += a * x[i];
}

```

**Architectural diversity is here to stay, because it is
a consequence of the dusk of CMOS scaling
(Moore's Law)**

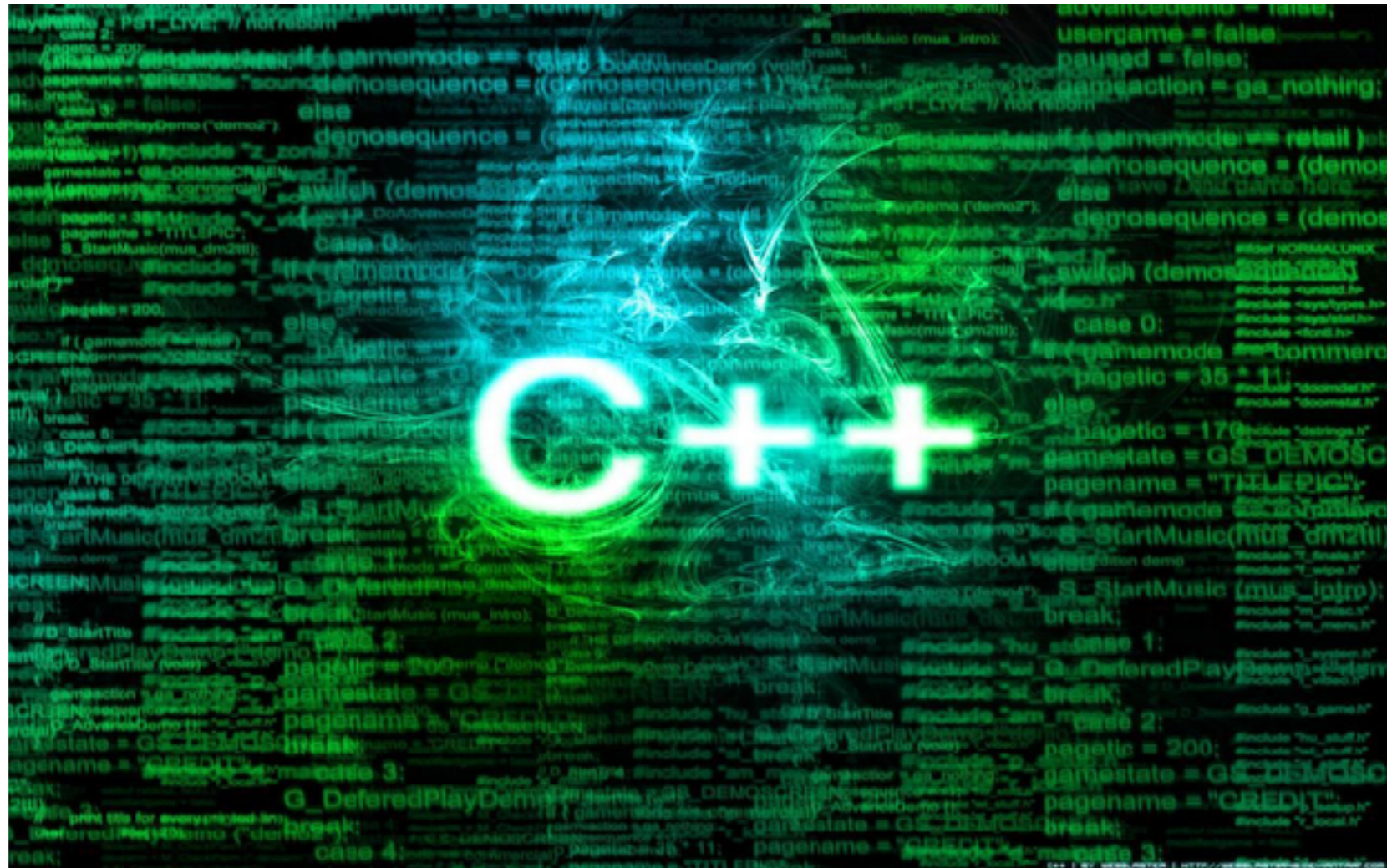
What are the implications?

Complexity in software is one,
but we don't understand all implications

Physics of the computer matters more than ever

The good news

C++ standard is evolving quickly and implementations follow!



C++ 11, 14, (HPX-3/Kokkos), ... 17, 20, ...

Who will pay for the implementation of Fortran, OpenACC, OpenMP, ...?

The top ranking programming languages in 2017

spectrum.ieee.org

Language Types (click to hide)



Web



Mobile



Enterprise

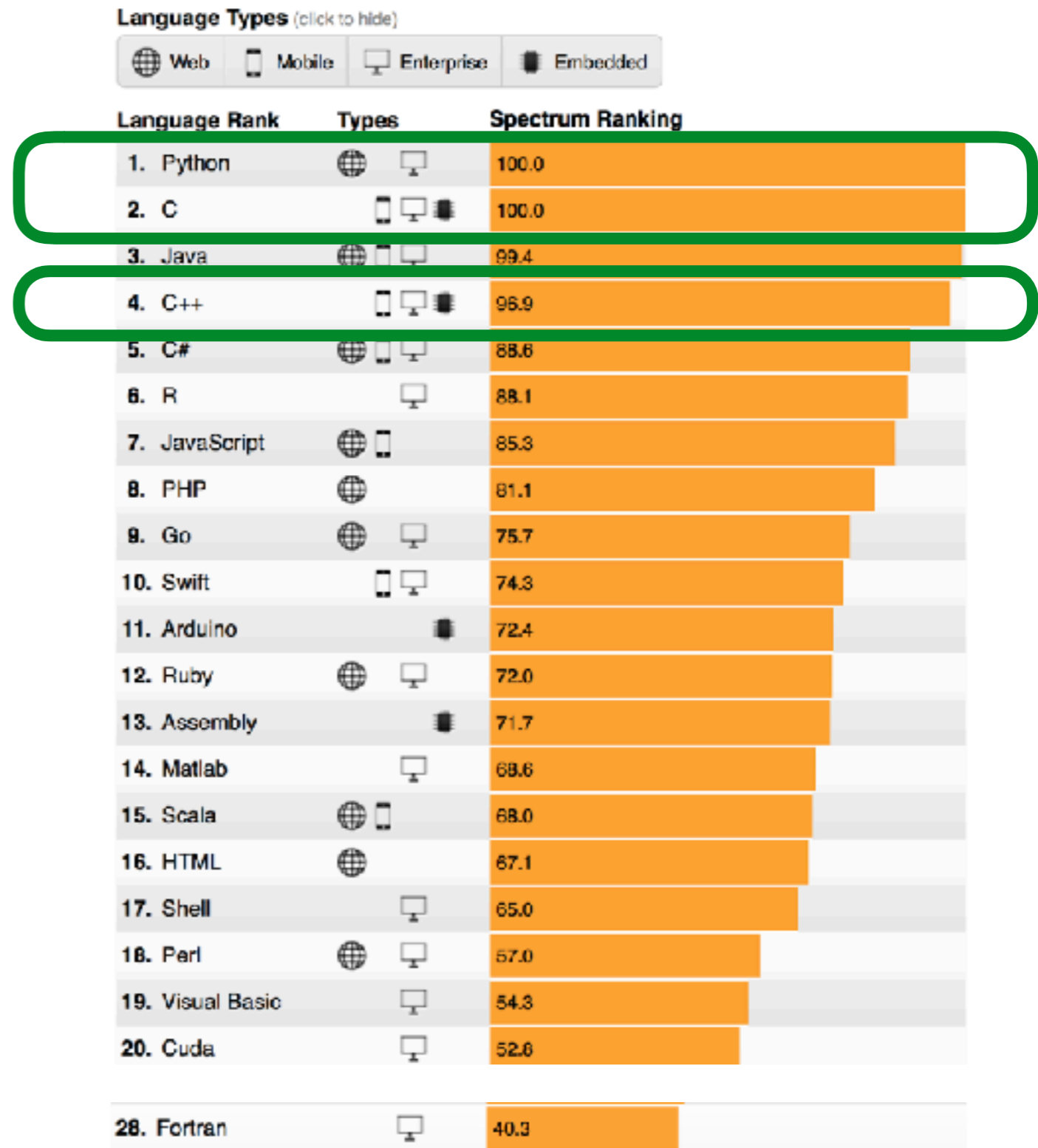


Embedded

Language Rank	Types	Spectrum Ranking
1. Python		100.0
2. C		99.7
3. Java		99.5
4. C++		97.1
5. C#		87.7
6. R		87.7
7. JavaScript		85.6
8. PHP		81.2
9. Go		75.1
10. Swift		73.7

The top ranking programming languages in 2017

spectrum.ieee.org



$$\text{Wind } \rho \dot{\mathbf{v}} = -\nabla p + \rho \mathbf{g} - 2\boldsymbol{\Omega} \times (\rho \mathbf{v}) + \mathbf{F}$$

$$\text{Pressure } \dot{p} = - (c_{pd}/c_{vd}) p \nabla \cdot \mathbf{v} + (c_{pd}/c_{vd} - 1) Q_h$$

$$\text{Temperature } \rho c_{pd} \dot{T} = \dot{p} + Q_h$$

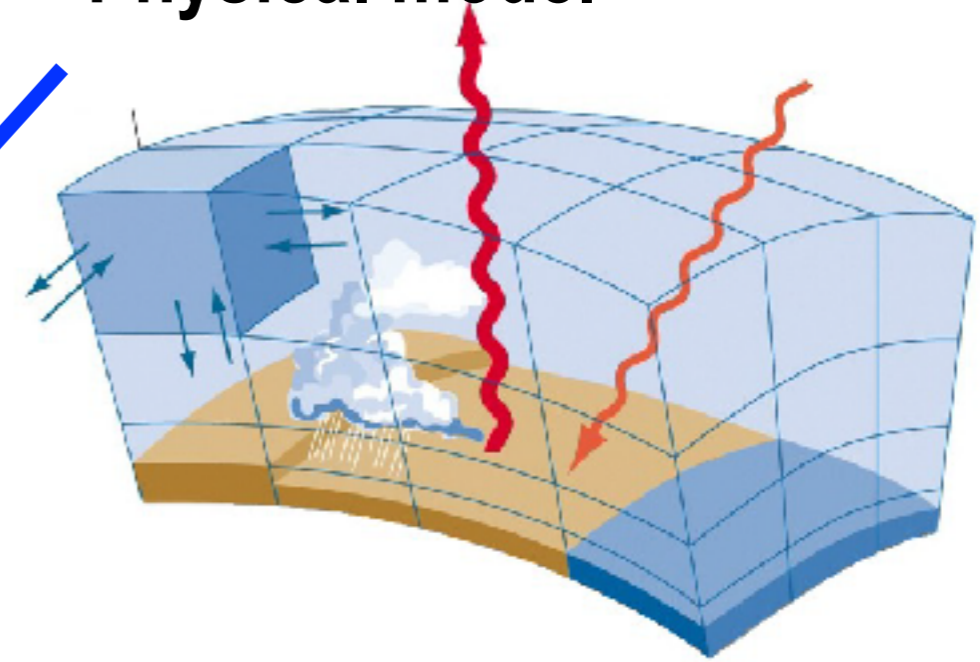
$$\text{Water } \rho \dot{q}^v = -\nabla \cdot \mathbf{F}^v - (I^l + I^f)$$

$$\rho \dot{q}^{l,f} = \nabla \cdot (\mathbf{P}^{l,f} + \mathbf{F}^{l,f}) + I^{l,f}$$

$$\text{Density } \rho = p [R_d (1 + (R_v/R_d - 1) q^v - q^l - q^f) T]^{-1}$$

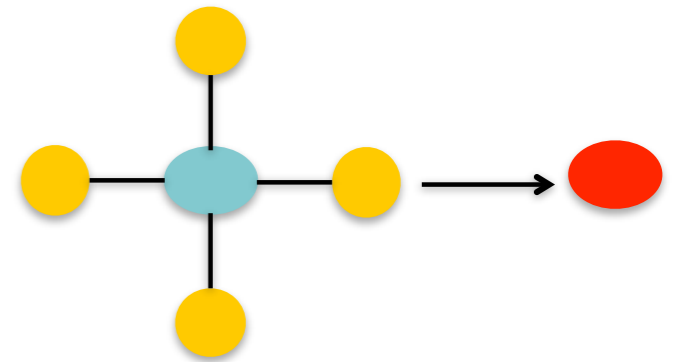
Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description



```
lap(i,j,k) = -4.0 * data(i,j,k) +
data(i+1,j,k) + data(i-1,j,k) +
data(i,j+1,k) + data(i,j-1,k);
```

Imperative code

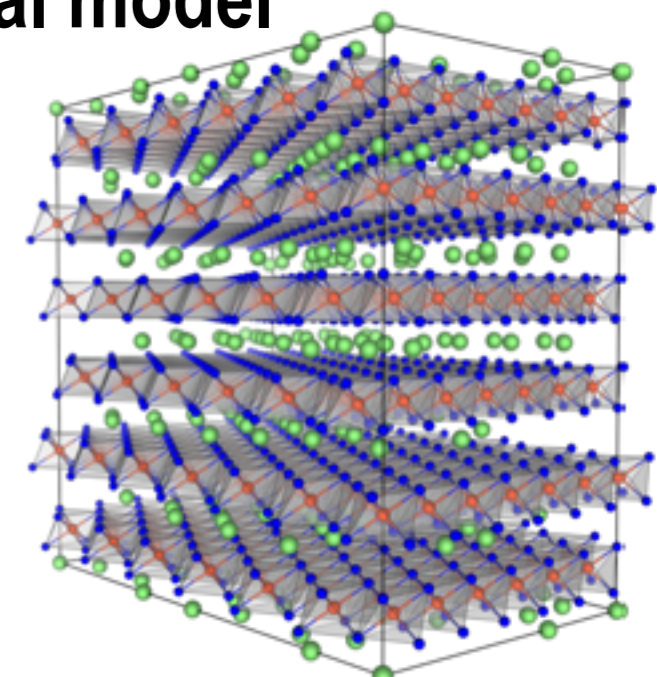
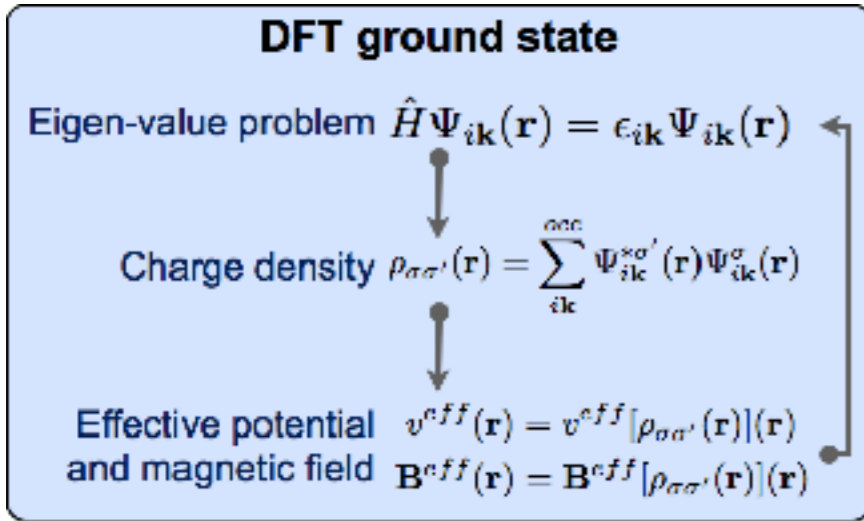
Compilation



Computer

Computer engineering

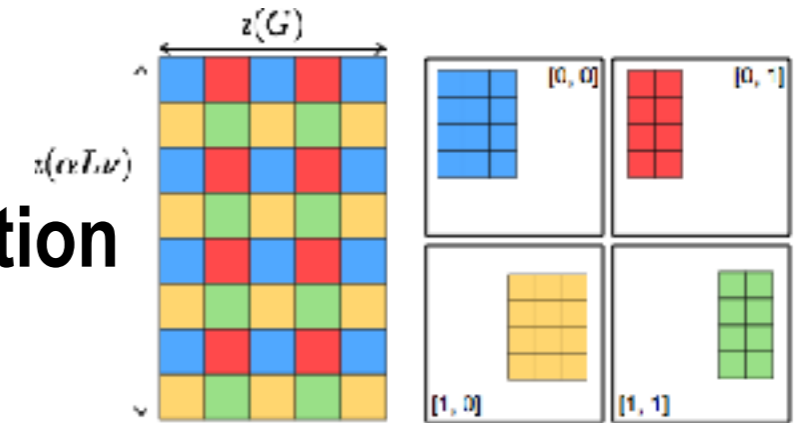
Schulthess, Nature Physics, vol 11, 369-373 (2015)



Mathematical description

Domain science & applied mathematics

Algorithmic description



Imperative code

Compilation

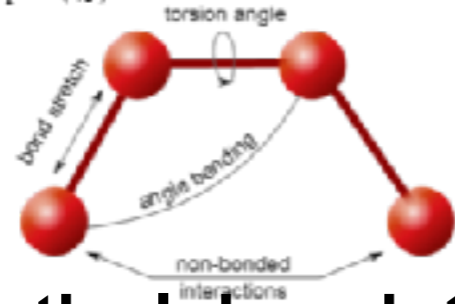


Computer

Computer engineering

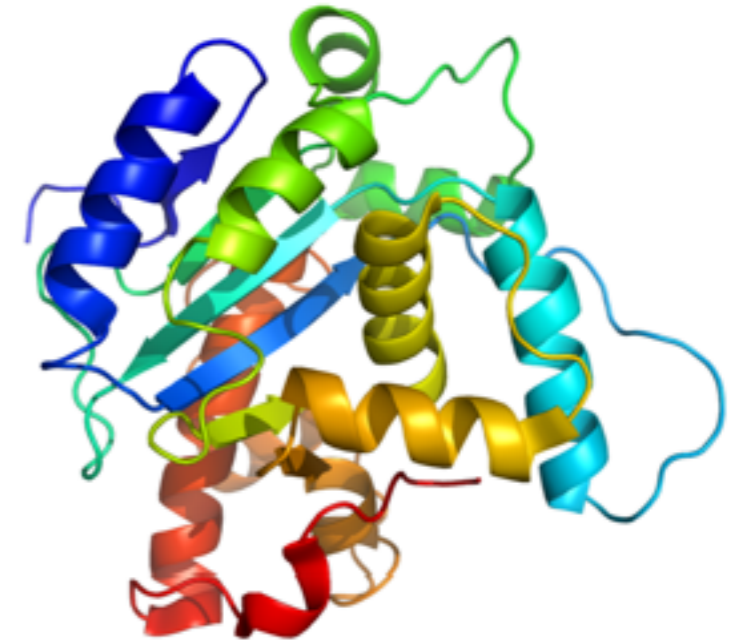
Schulthess, Nature Physics, vol 11, 369-373 (2015)

$$\begin{aligned}
 V(r) = & \sum_{\text{bonds}} k_b(b - b_0)^2 + \sum_{\text{angles}} k_\theta(\theta - \theta_0)^2 \\
 & + \sum_{\text{dihedrals}} k_\phi(1 + \cos(n\phi - \phi_0)) + \sum_{\text{impropers}} k_\psi(\psi - \psi_0)^2 \\
 & + \sum_{\text{non-bonded pairs}(i,j)} 4\epsilon_{ij} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] + \sum_{\text{non-bonded pairs}(i,j)} \frac{q_i q_j}{\epsilon_{ij} r_{ij}}
 \end{aligned}$$



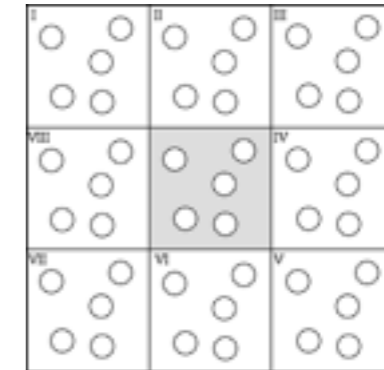
Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description



```

Cray PAT API *11 loop init mtran* in UCCA CDAGWETHMUCI
call pat_region_begin(11,'11 loop init mtran',pat_stat)
!
!!! DO IMODE=1,IMODES
!!! MDOLD(:,IMODE)=MD(:,IMODE)
!!! DO ICP=1,ICPS
!!! MDOLD(:,IMODE,ICP)=MD(:,IMODE,ICP)
!!! DO JMODE=1,JMODES
!!! MTRAN(:,IMODE,JMODE,ICP)=0.0
!!! ENDDO
!!! ENDDO
!!! ENDDO
!
! replace triple loops above with F90 array syntax (let
! compiler decide)
MDOLD=MD
MDOLD=MD
MTRAN=0.0
call pat_region_end(11,pat_stat)

```

Imperative code

Compilation

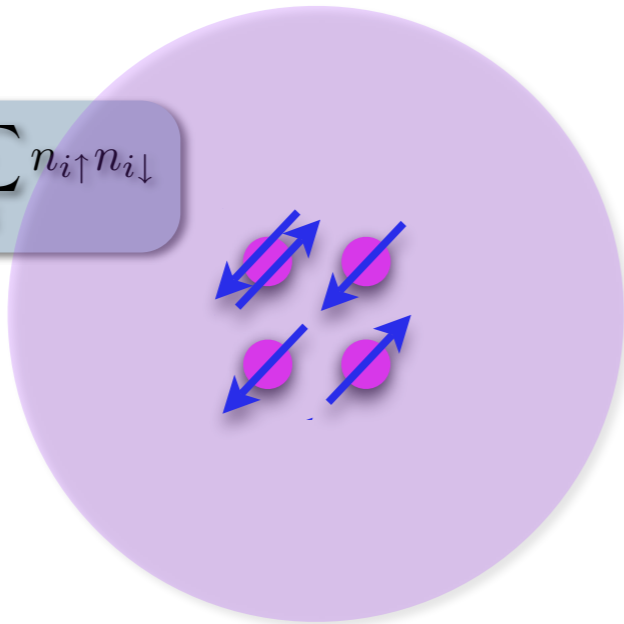


Computer

Computer engineering

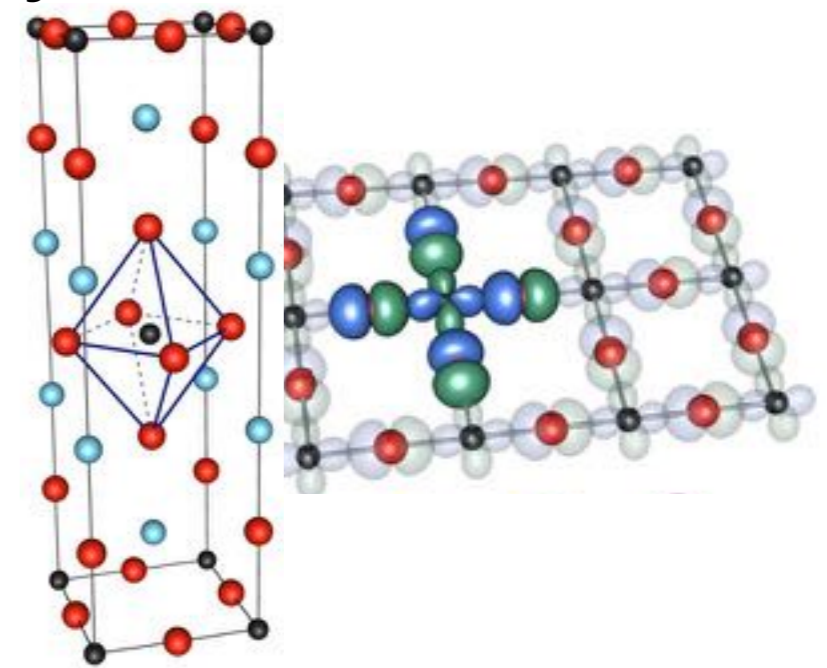
Schulthess, Nature Physics, vol 11, 369-373 (2015)

$$\mathcal{H} = -t \sum_{\langle ij \rangle, \sigma} c_{i\sigma}^\dagger c_{j\sigma} + U \sum_i n_{i\uparrow} n_{i\downarrow}$$



Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description

$$\mathbf{G}_c(\{s_i, l\}_{k+1}) = \mathbf{G}_c(\{s_i, l\}_k) + \mathbf{a}_k \times \mathbf{b}_k^t$$
$$\mathbf{G}_c(\{s_i, l\}_{k+1}) = \mathbf{G}_c(\{s_i, l\}_0) + [\mathbf{a}_0 | \mathbf{a}_1 | \dots | \mathbf{a}_k] \times [\mathbf{b}_0 | \mathbf{b}_1 | \dots | \mathbf{b}_k]^t$$

Imperative code

Compilation



Computer

Computer engineering

Schulthess, Nature Physics, vol 11, 369-373 (2015)

Wind $\rho \dot{\mathbf{v}} = -\nabla p + \rho \mathbf{g} - 2\Omega \times (\rho \mathbf{v}) + \mathbf{F}$

Pressure $\dot{p} = -(c_{pd}/c_{vd}) p \nabla \cdot \mathbf{v} + (c_{pd}/c_{vd} - 1) Q_h$

Temperature $\rho c_{pd} \dot{T} = \dot{p} + Q_h$

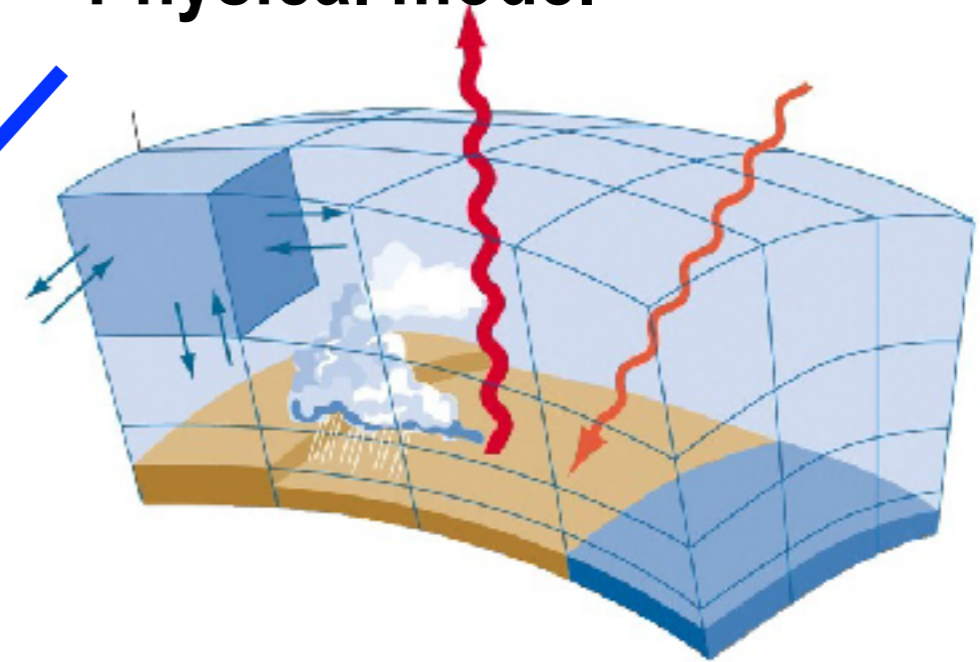
Water $\rho \dot{q}^v = -\nabla \cdot \mathbf{F}^v - (I^l + I^f)$

$\rho \dot{q}^{l,f} = \nabla \cdot (\mathbf{P}^{l,f} + \mathbf{F}^{l,f}) + I^{l,f}$

Density $\rho = p [R_d (1 + (R_v/R_d - 1) q^v - q^l - q^f) T]^{-1}$

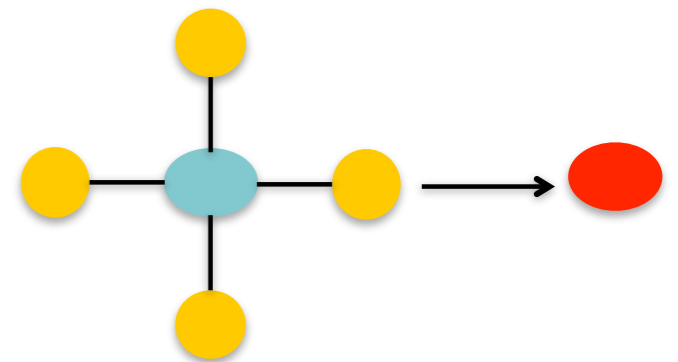
Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description



```
lap(i,j,k) = -4.0 * data(i,j,k) +
            data(i+1,j,k) + data(i-1,j,k) +
            data(i,j+1,k) + data(i,j-1,k);
```

Imperative code

Compilation



Computer

Computer engineering

Schulthess, Nature Physics, vol 11, 369-373 (2015)

Wind $\rho \dot{\mathbf{v}} = -\nabla p + \rho \mathbf{g} - 2\boldsymbol{\Omega} \times (\rho \mathbf{v}) + \mathbf{F}$

Pressure $\dot{p} = -(c_{pd}/c_{vd}) p \nabla \cdot \mathbf{v} + (c_{pd}/c_{vd} - 1) Q_h$

Temperature $\rho c_{pd} \dot{T} = \dot{p} + Q_h$

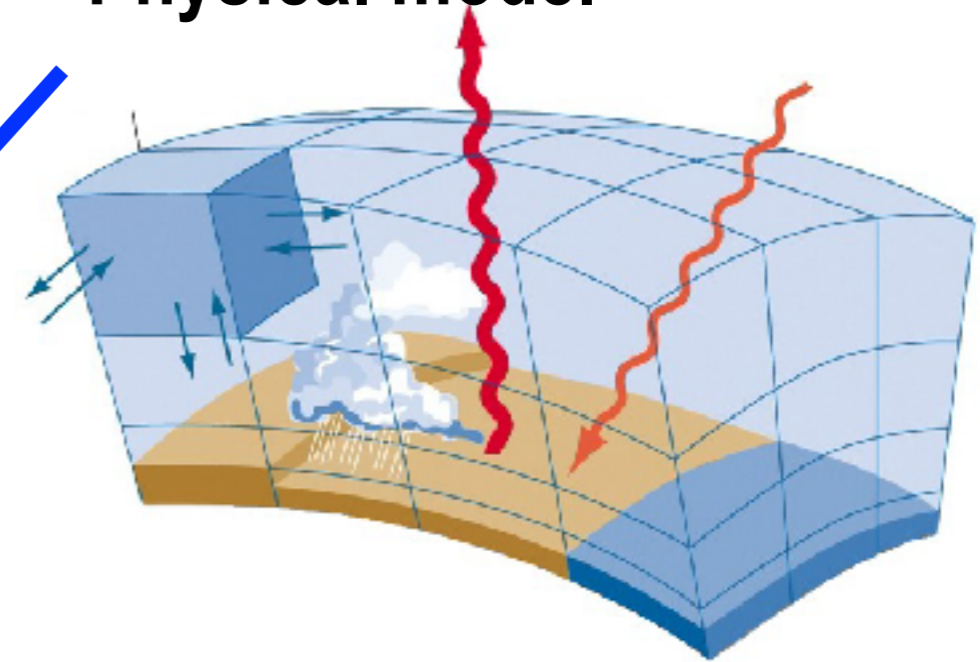
Water $\rho \dot{q}^v = -\nabla \cdot \mathbf{F}^v - (I^l + I^f)$

$\rho \dot{q}^{l,f} = \nabla \cdot (\mathbf{P}^{l,f} + \mathbf{F}^{l,f}) + I^{l,f}$

Density $\rho = p [R_d (1 + (R_v/R_d - 1) q^v - q^l - q^f) T]^{-1}$

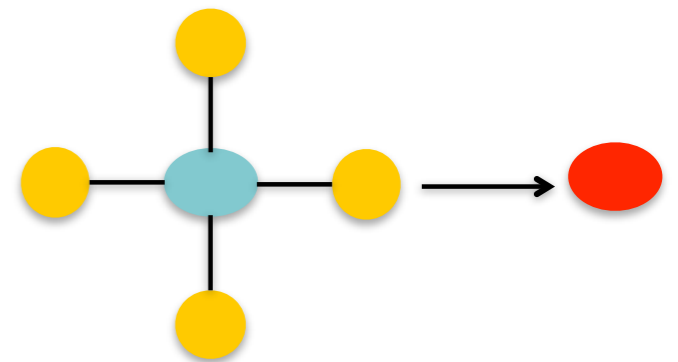
Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description



```
lap(i,j,k) = -4.0 * data(i,j,k) +
            data(i+1,j,k) + data(i-1,j,k) +
            data(i,j+1,k) + data(i,j-1,k);
```

Imperative code

Compilation



Computer engineering

Schulthess, Nature Physics, vol 11, 369-373 (2015)

$$\text{Wind } \rho \dot{\mathbf{v}} = -\nabla p + \rho \mathbf{g} - 2\boldsymbol{\Omega} \times (\rho \mathbf{v}) + \mathbf{F}$$

$$\text{Pressure } \dot{p} = - (c_{pd}/c_{vd}) p \nabla \cdot \mathbf{v} + (c_{pd}/c_{vd} - 1) Q_h$$

$$\text{Temperature } \rho c_{pd} \dot{T} = \dot{p} + Q_h$$

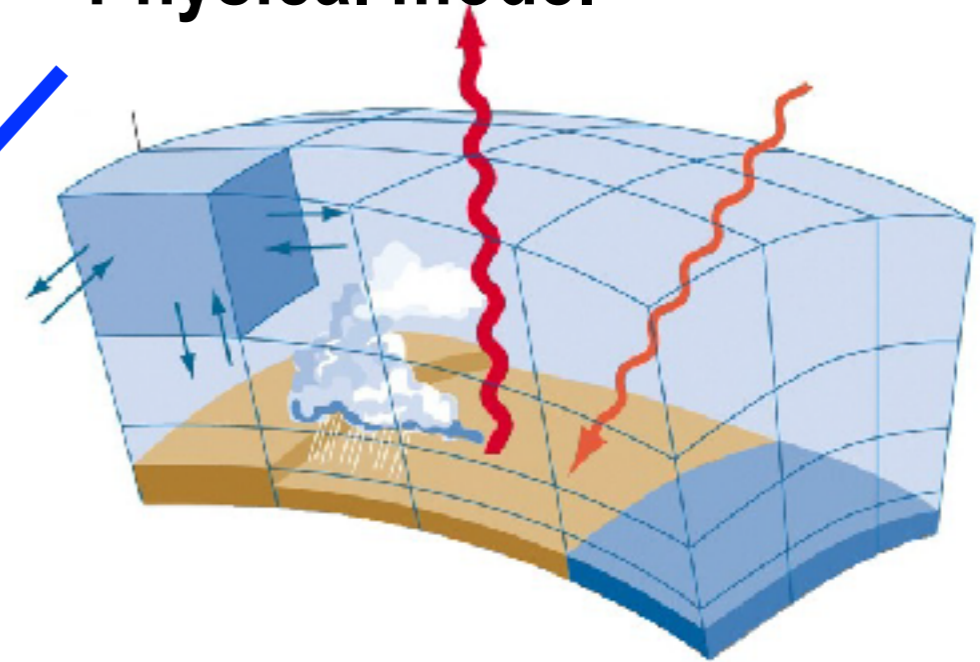
$$\text{Water } \rho \dot{q}^v = -\nabla \cdot \mathbf{F}^v - (I^l + I^f)$$

$$\rho \dot{q}^{l,f} = \nabla \cdot (\mathbf{P}^{l,f} + \mathbf{F}^{l,f}) + I^{l,f}$$

$$\text{Density } \rho = p [R_d (1 + (R_v/R_d - 1) q^v - q^l - q^f) T]^{-1}$$

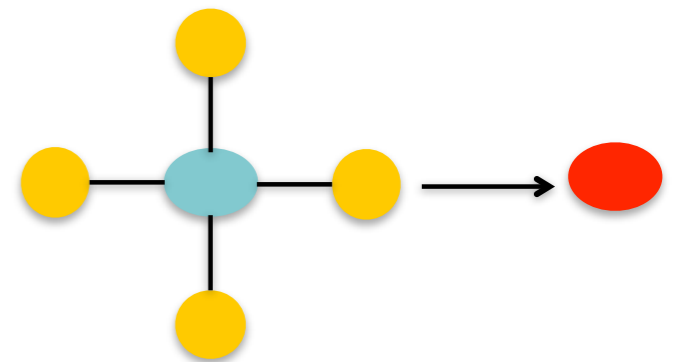
Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description



```
lap(i,j,k) = -4.0 * data(i,j,k) +
data(i+1,j,k) + data(i-1,j,k) +
data(i,j+1,k) + data(i,j-1,k);
```

Imperative code

Compilation



Computer engineering

Wind $\rho \dot{\mathbf{v}} = -\nabla p + \rho \mathbf{g} - 2\boldsymbol{\Omega} \times (\rho \mathbf{v}) + \mathbf{F}$

Pressure $\dot{p} = -(c_{pd}/c_{vd}) p \nabla \cdot \mathbf{v} + (c_{pd}/c_{vd} - 1) Q_h$

Temperature $\rho c_{pd} \dot{T} = \dot{p} + Q_h$

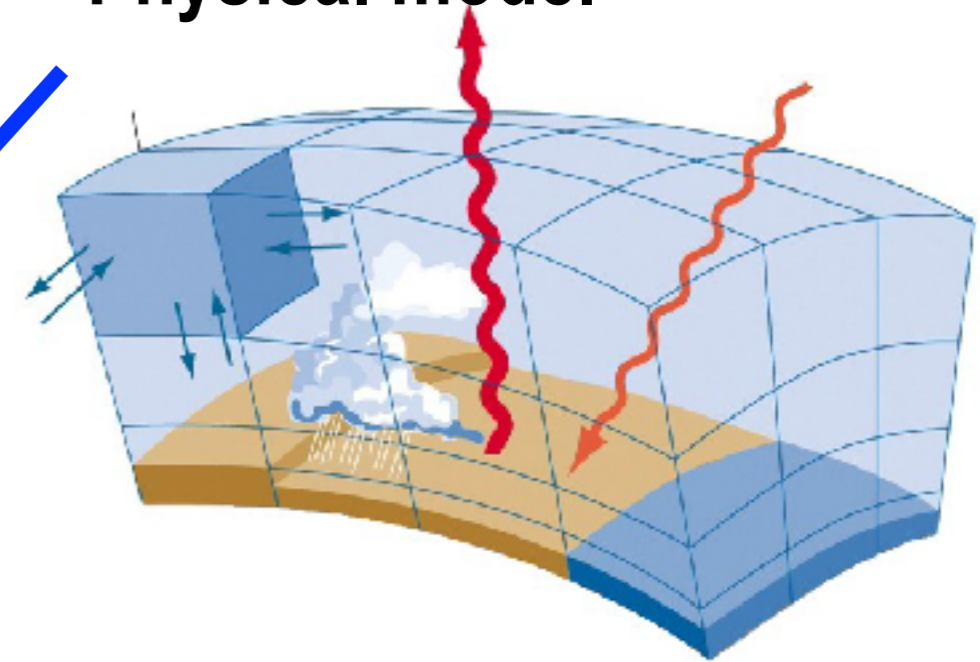
Water $\rho \dot{q}^v = -\nabla \cdot \mathbf{F}^v - (I^l + I^f)$

$\rho \dot{q}^{l,f} = \nabla \cdot (\mathbf{P}^{l,f} + \mathbf{F}^{l,f}) + I^{l,f}$

Density $\rho = p [R_d (1 + (R_v/R_d - 1) q^v - q^l - q^f) T]^{-1}$

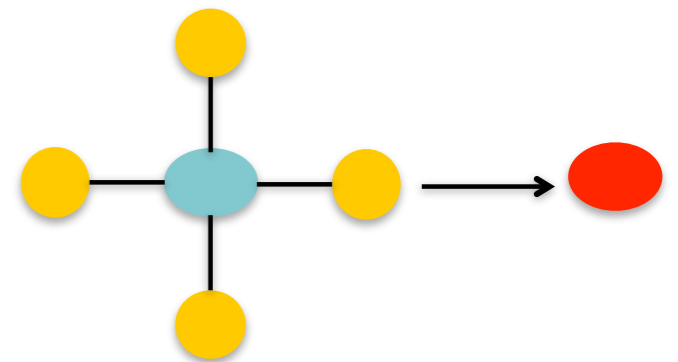
Mathematical description

Physical model



Domain science & applied mathematics

Algorithmic description



```
lap(i,j,k) = -4.0 * data(i,j,k) +
            data(i+1,j,k) + data(i-1,j,k) +
            data(i,j+1,k) + data(i,j-1,k);
```

Imperative code

Compilation



Computer engineering

Science applications using a descriptive and dynamic developer environment

Mathematical description

Physical model

Algorithmic description

Imperative code

Compiler frontend

Optimisation / low-level libraries / runtime

Architecture specific backends

Multi-disciplinary co-design of tools, libraries, programming environment

tools for high-performance scientific computing

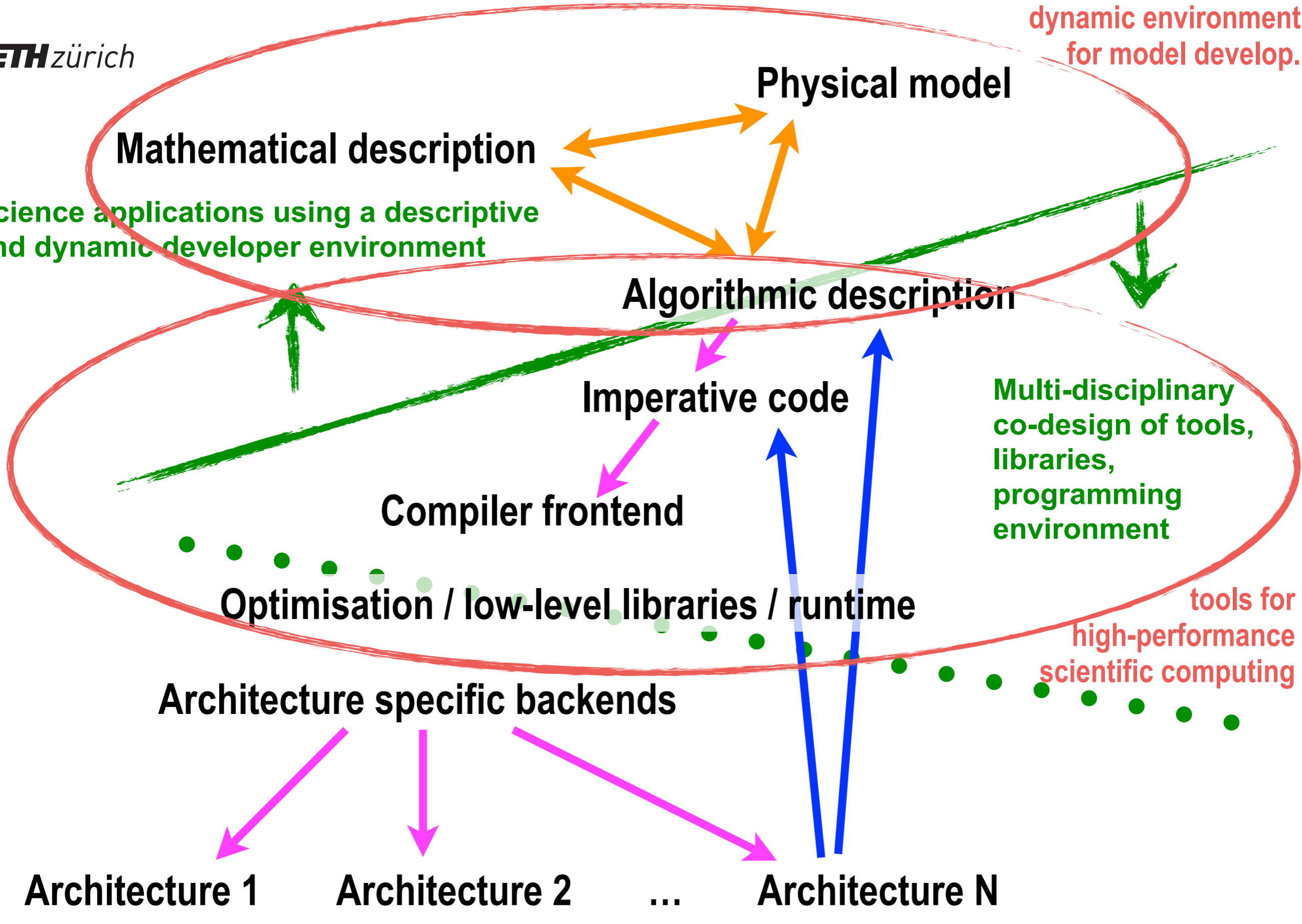
Architecture 1

Architecture 2

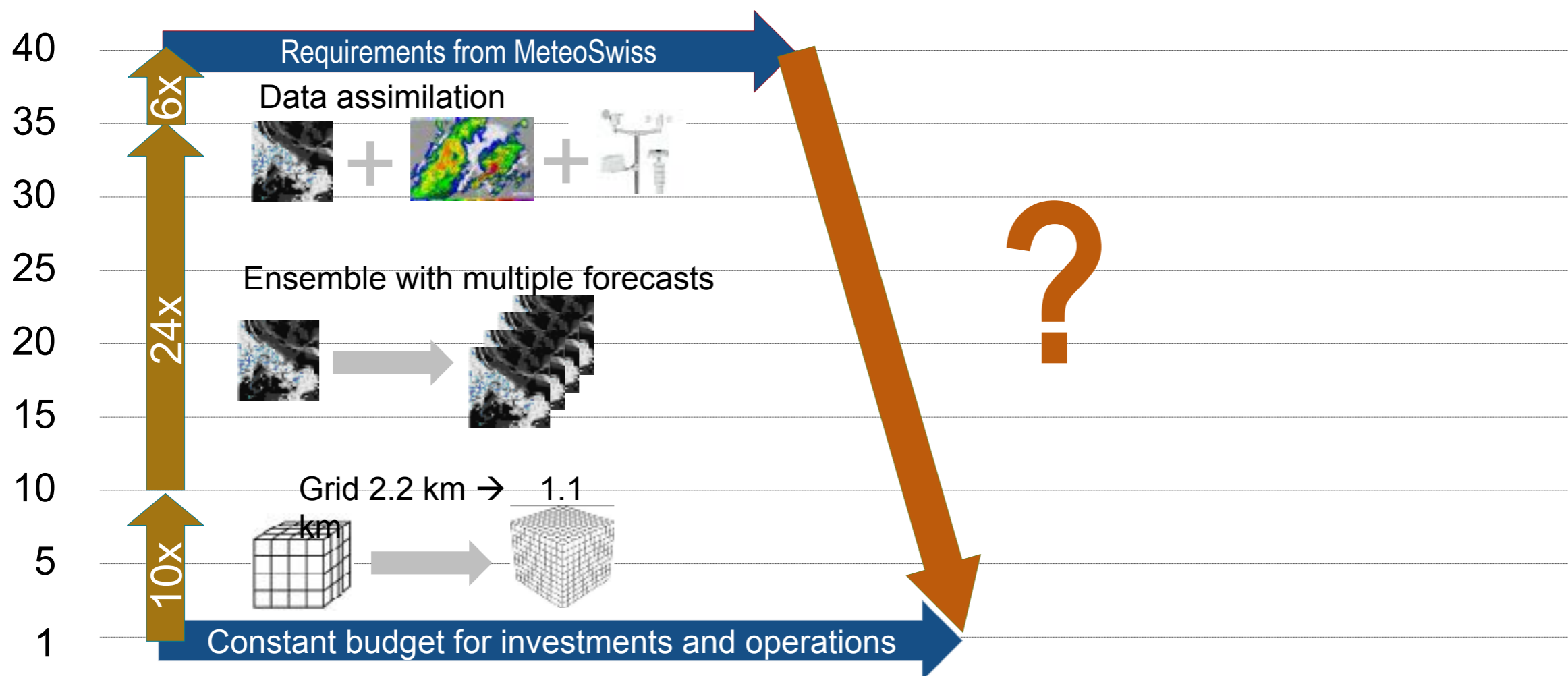
...

Architecture N

Schulthess, Nature Physics, vol 11, 369-373 (2015)

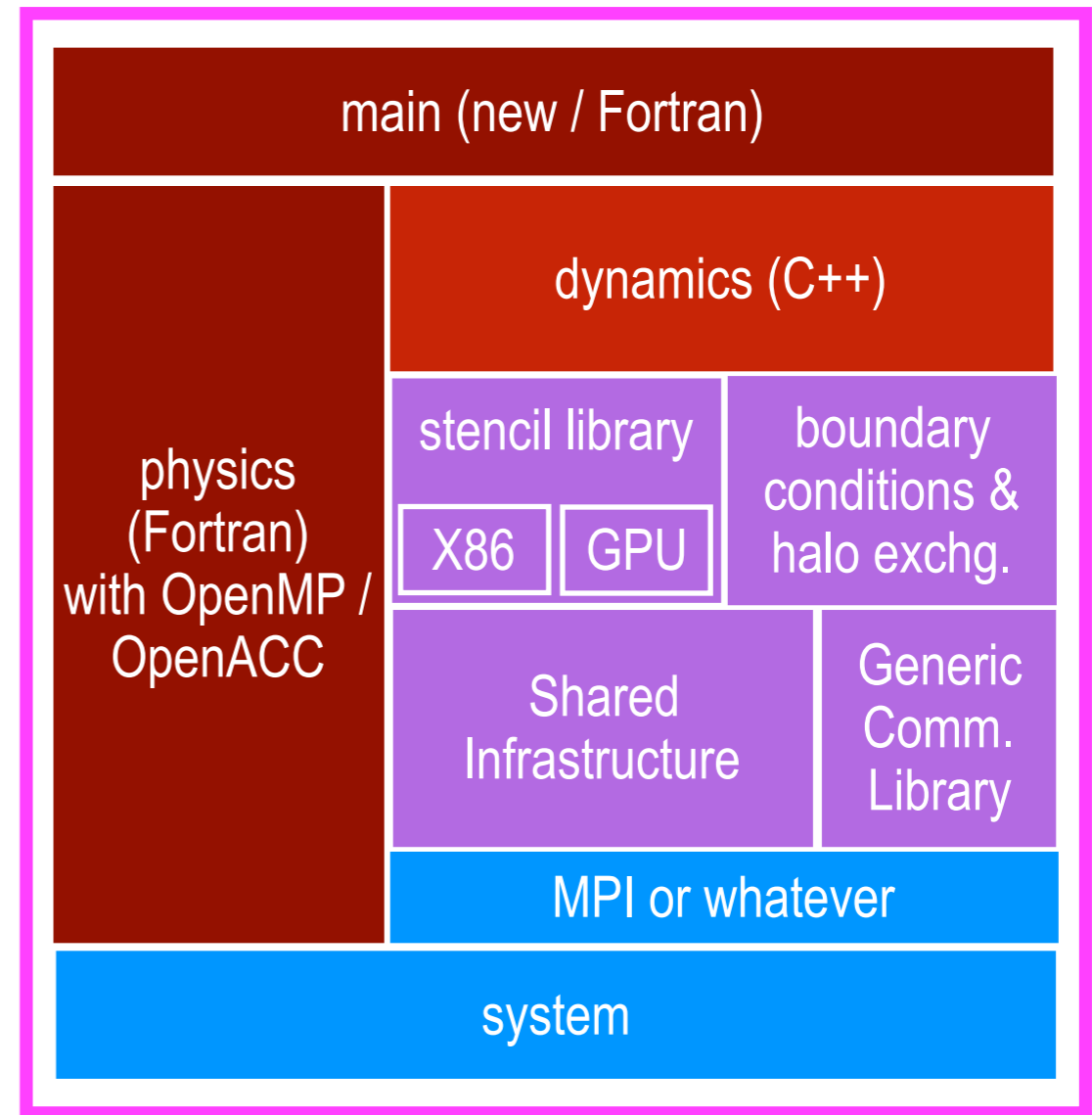
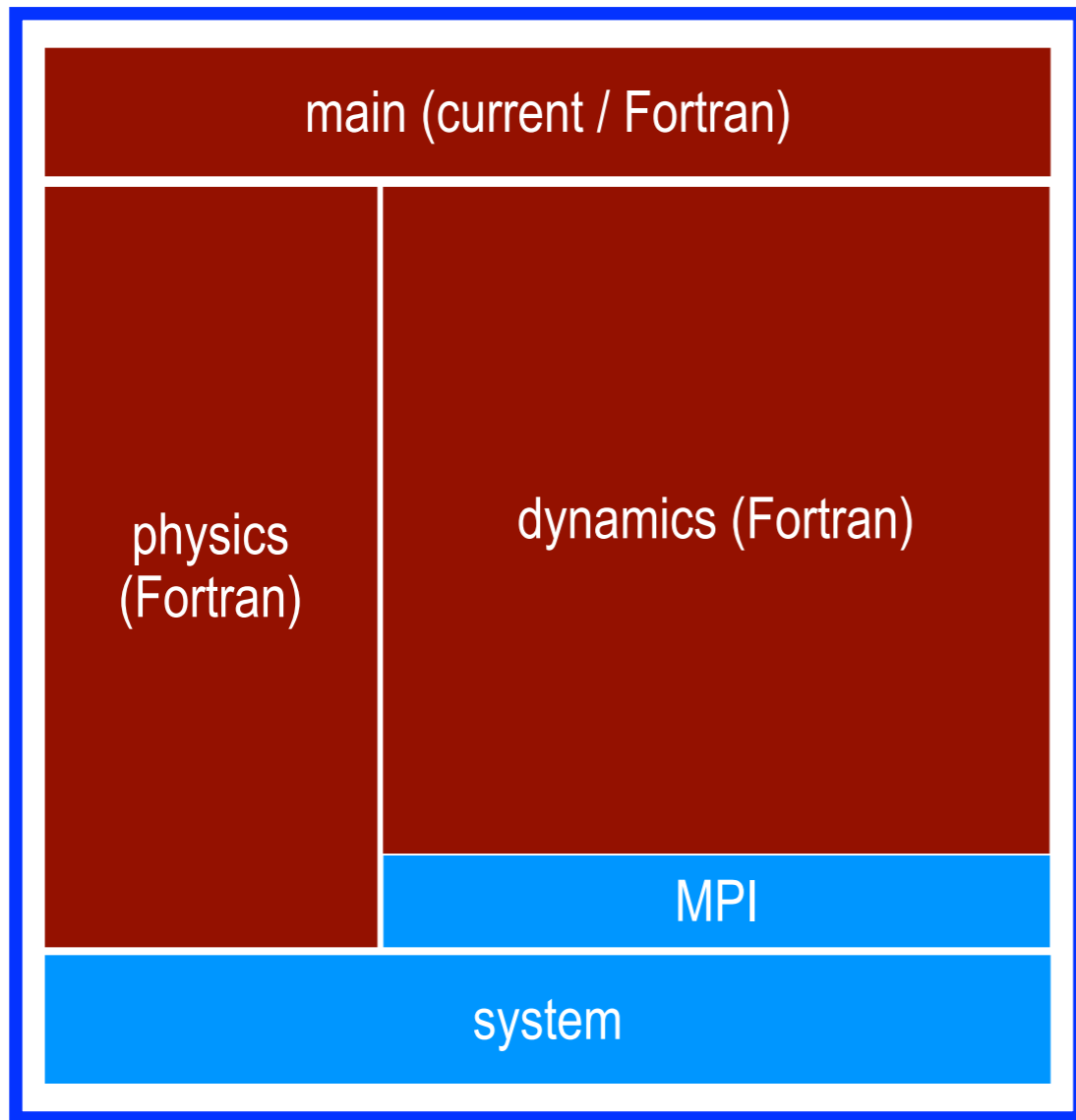


MeteoSwiss' performance ambitions in 2013



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

COSMO: **old** and **new** (refactored) code



MeteoSwiss New Weather Supercomputer

World's First GPU-Accelerated Weather Forecasting System



2x Racks

48 CPUs

192 Tesla K80 GPUs

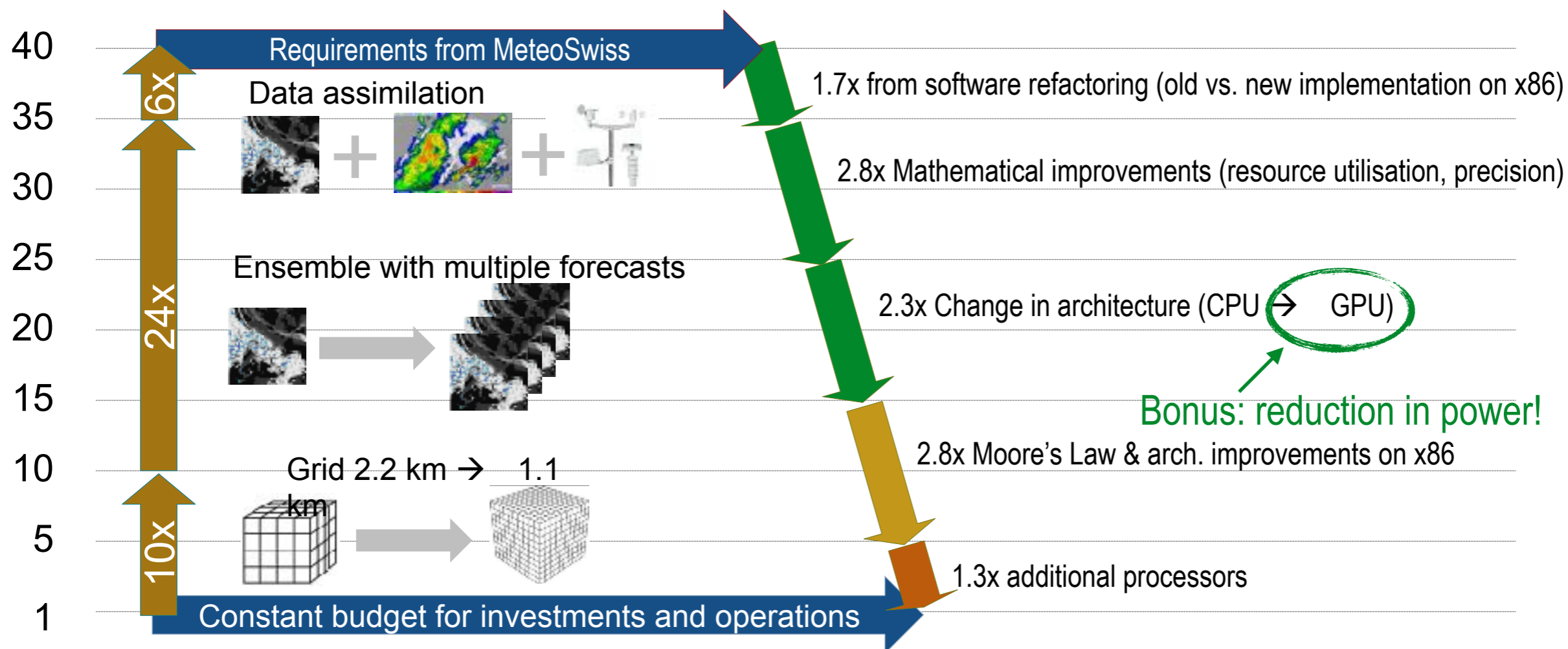
> 90% of FLOPS from GPUs

Operational in 2016



Where the factor 40 improvement came from

Investment in software allowed mathematical improvements and change in architecture



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.



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

Physical model

Mathematical description

Science applications using a descriptive and dynamic developer environment

main (new / Fortran)

dynamics (C++)

stencil library
X86 GPU

boundary conditions & halo exchg.

Shared Infrastructure

Generic Comm. Library

MPI or whatever

system

Multi-disciplinary co-design of tools, libraries, programming environment

tools for high-performance scientific computing

Compiler front

Optimisation / low-level

Architecture specific back

Architecture 1

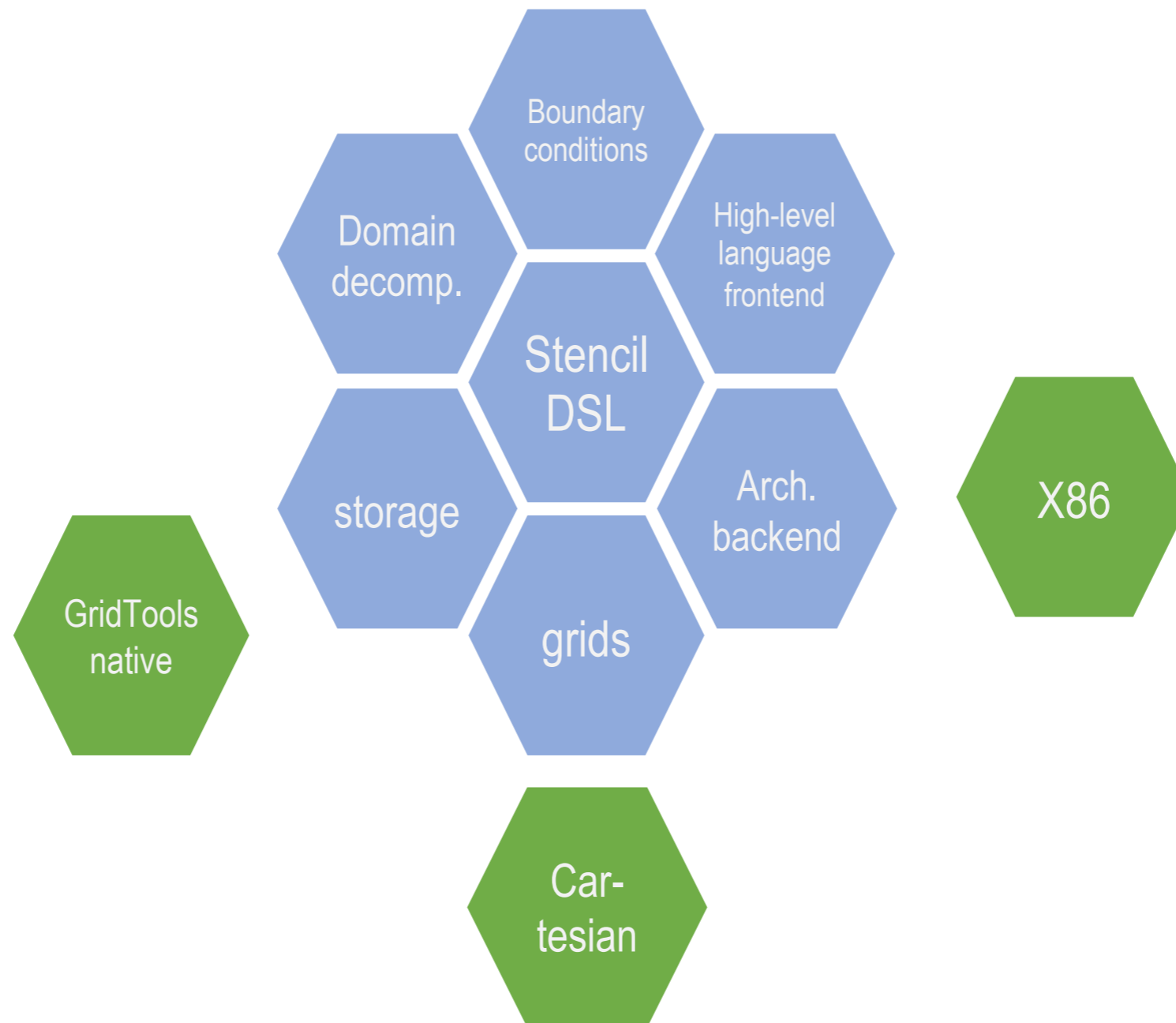
Architecture 2

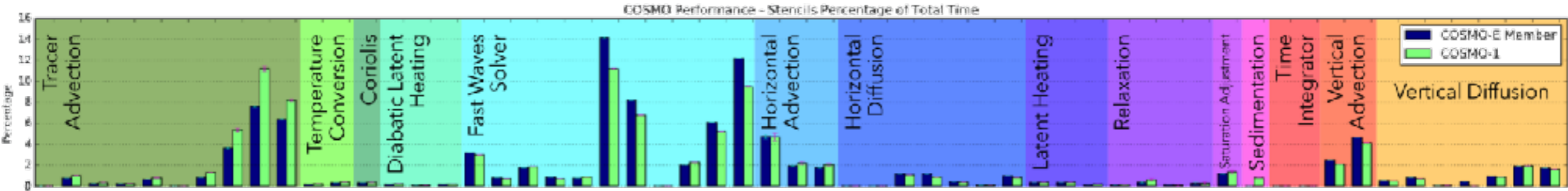
...

Architecture N

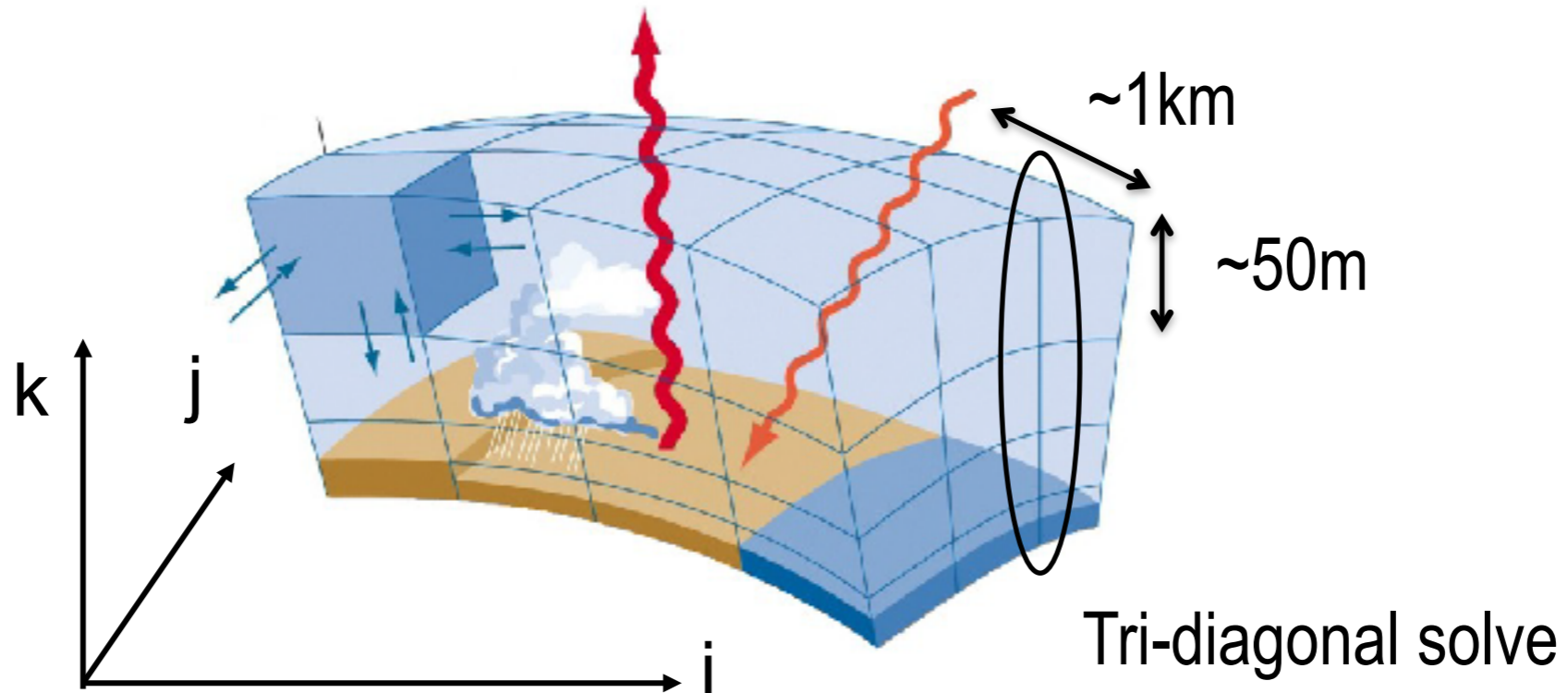
GridTools Framework

COSMO @ MeteoSwiss





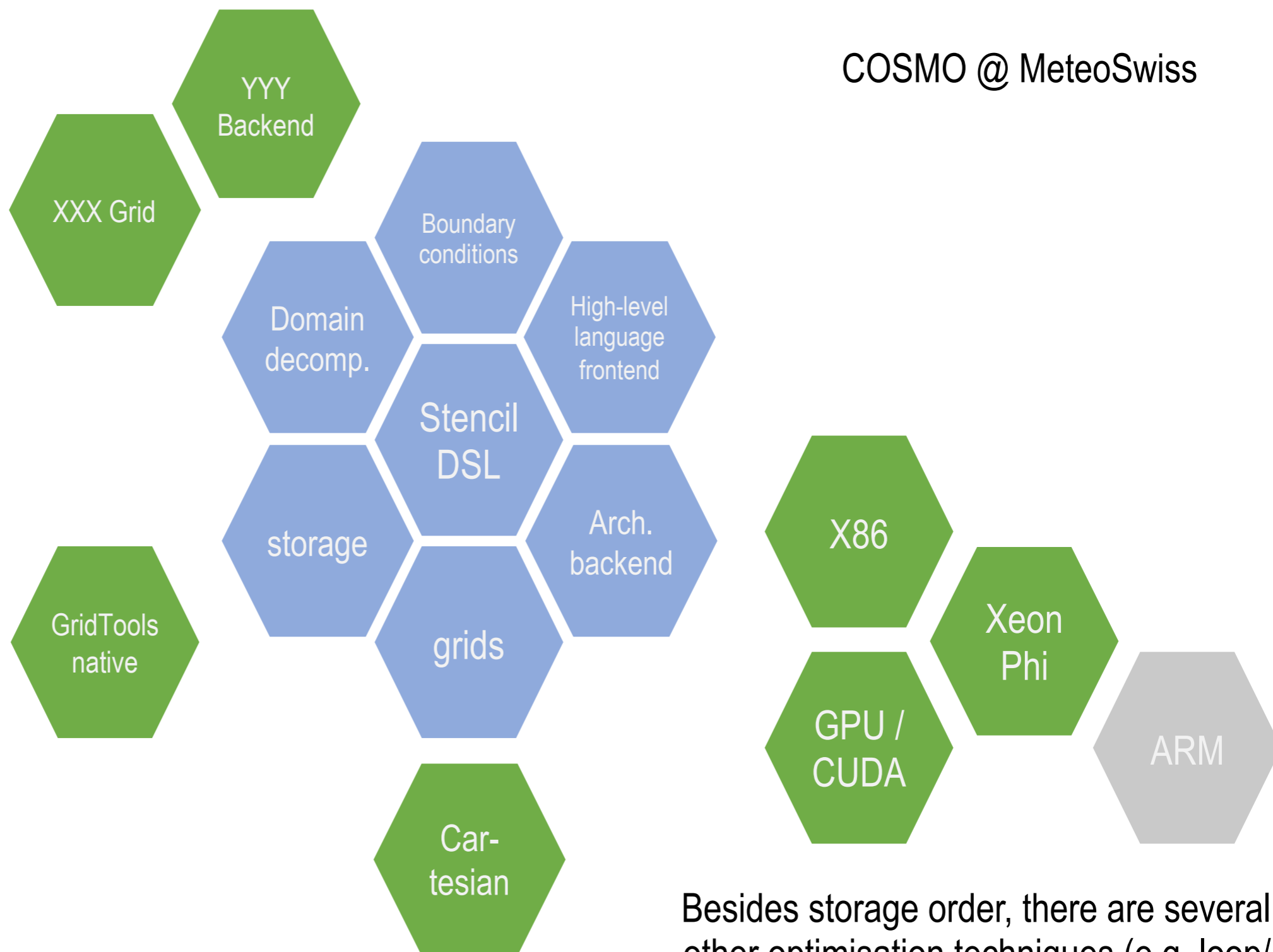
No performance hotspots



Store fields in (i,j,k) or (k,i,j) order?

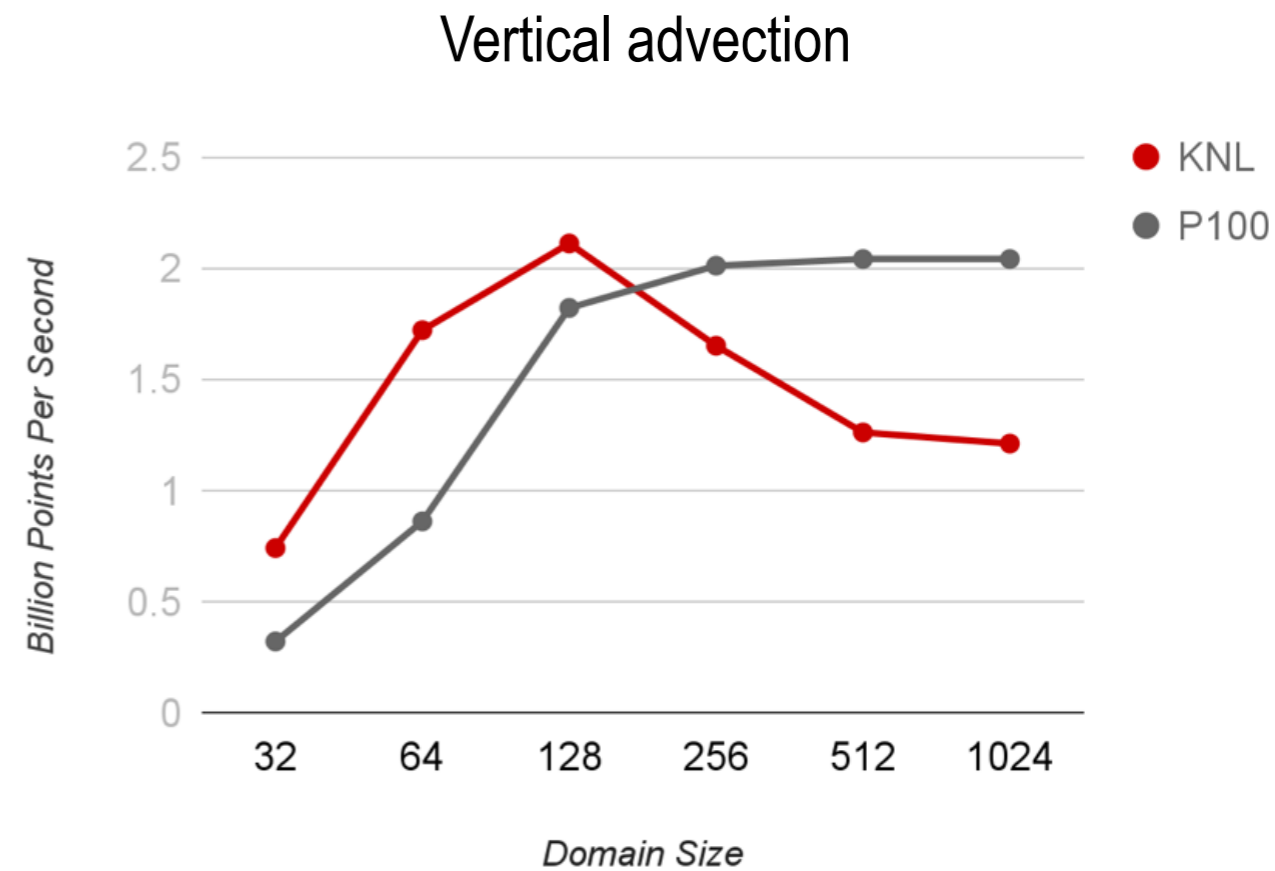
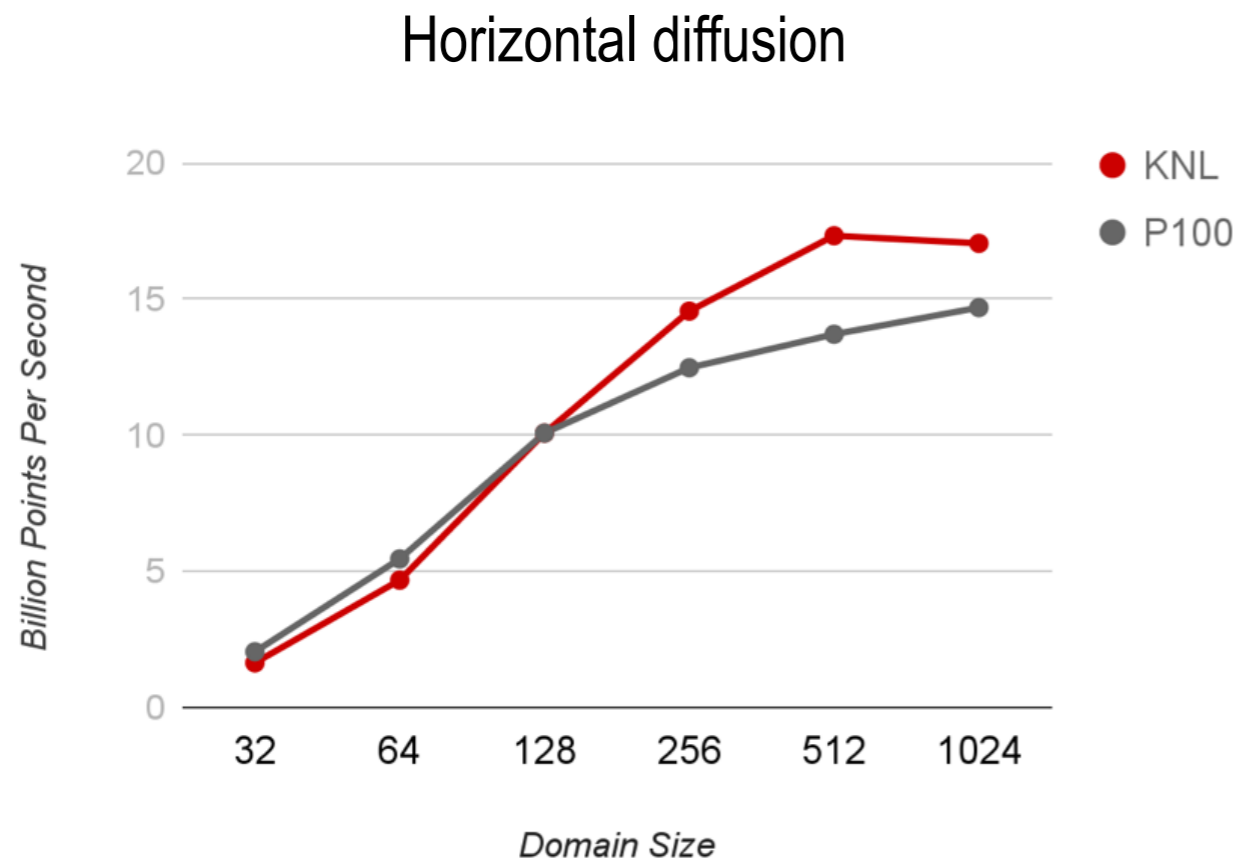
This depends on the architecture

COSMO @ MeteoSwiss

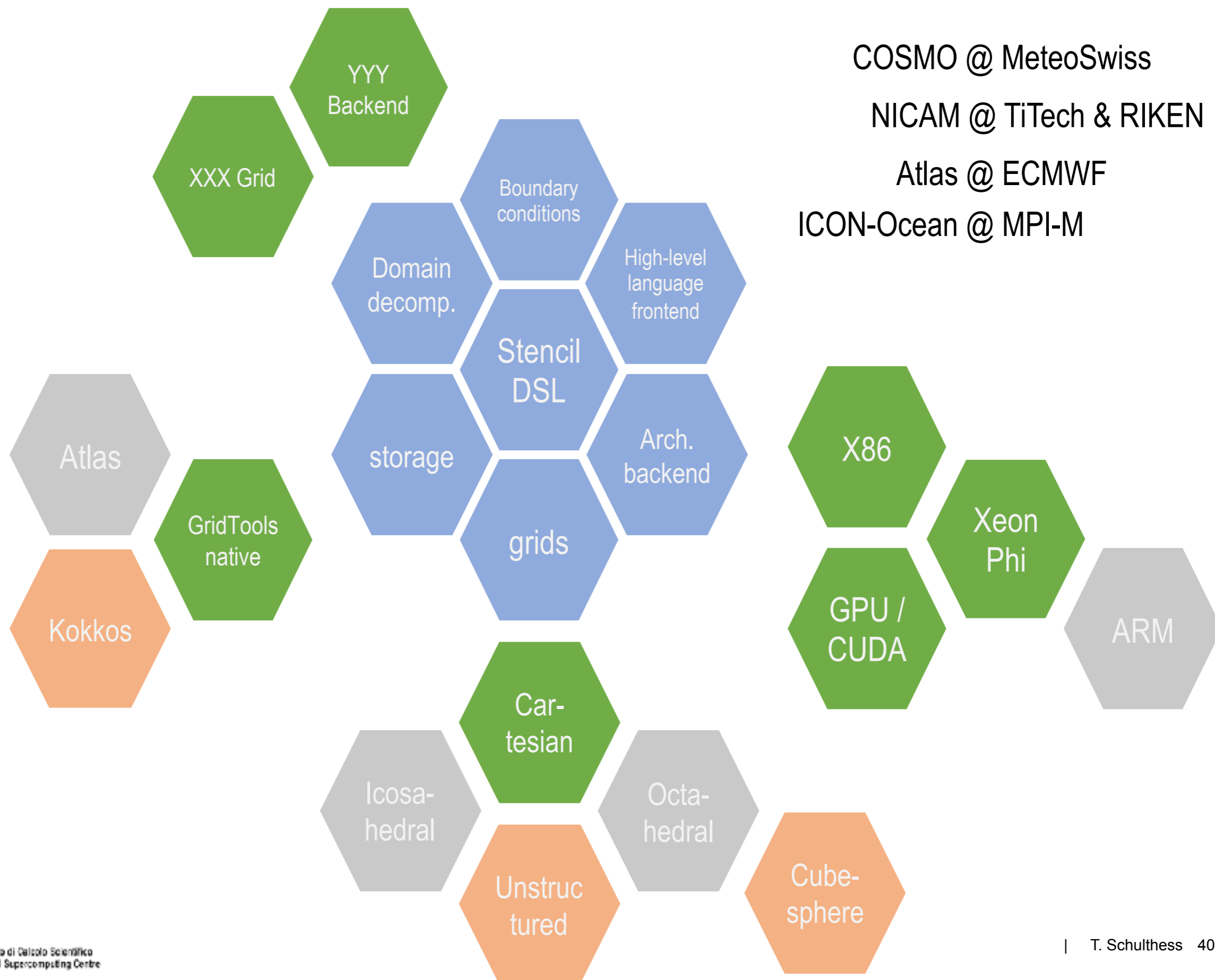


Besides storage order, there are several other optimisation techniques (e.g. loop/stencil fusion) where details depend on architecture specifics

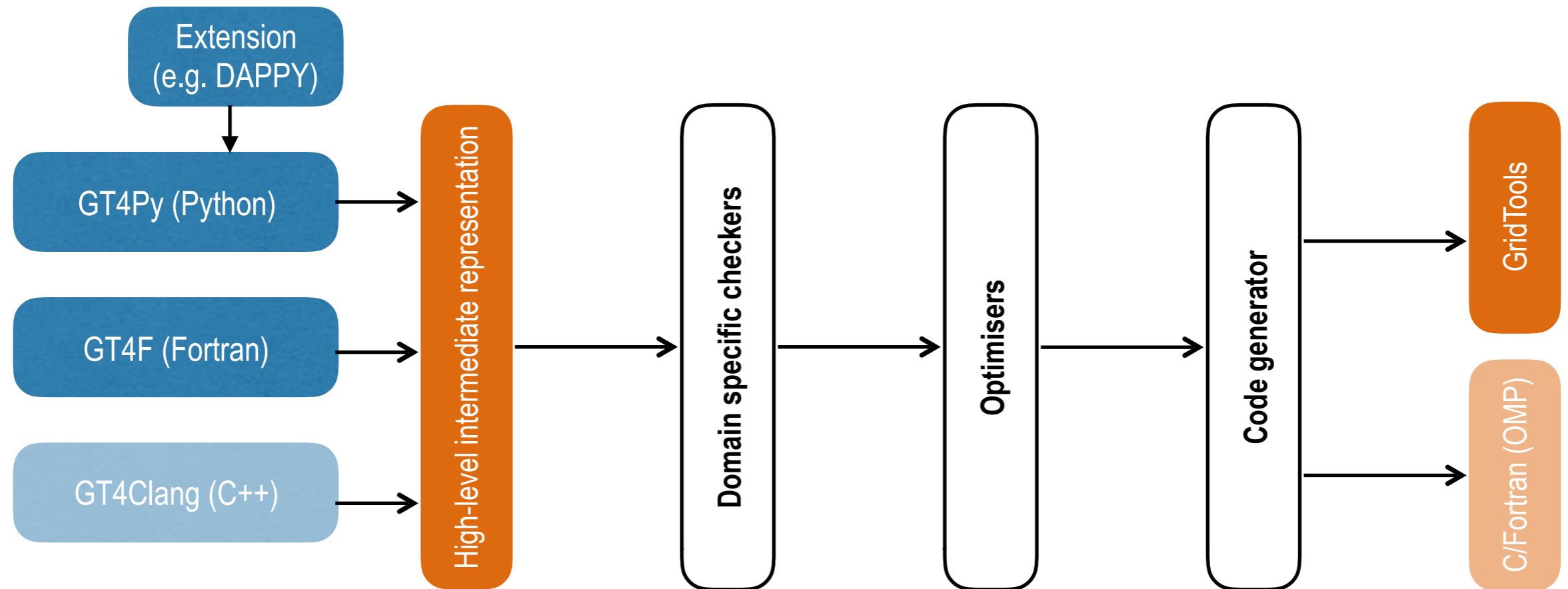
Exploring Intel Xeon Phi (KNL) and NVIDIA's P100



Source: Felix Thaler, CSCS



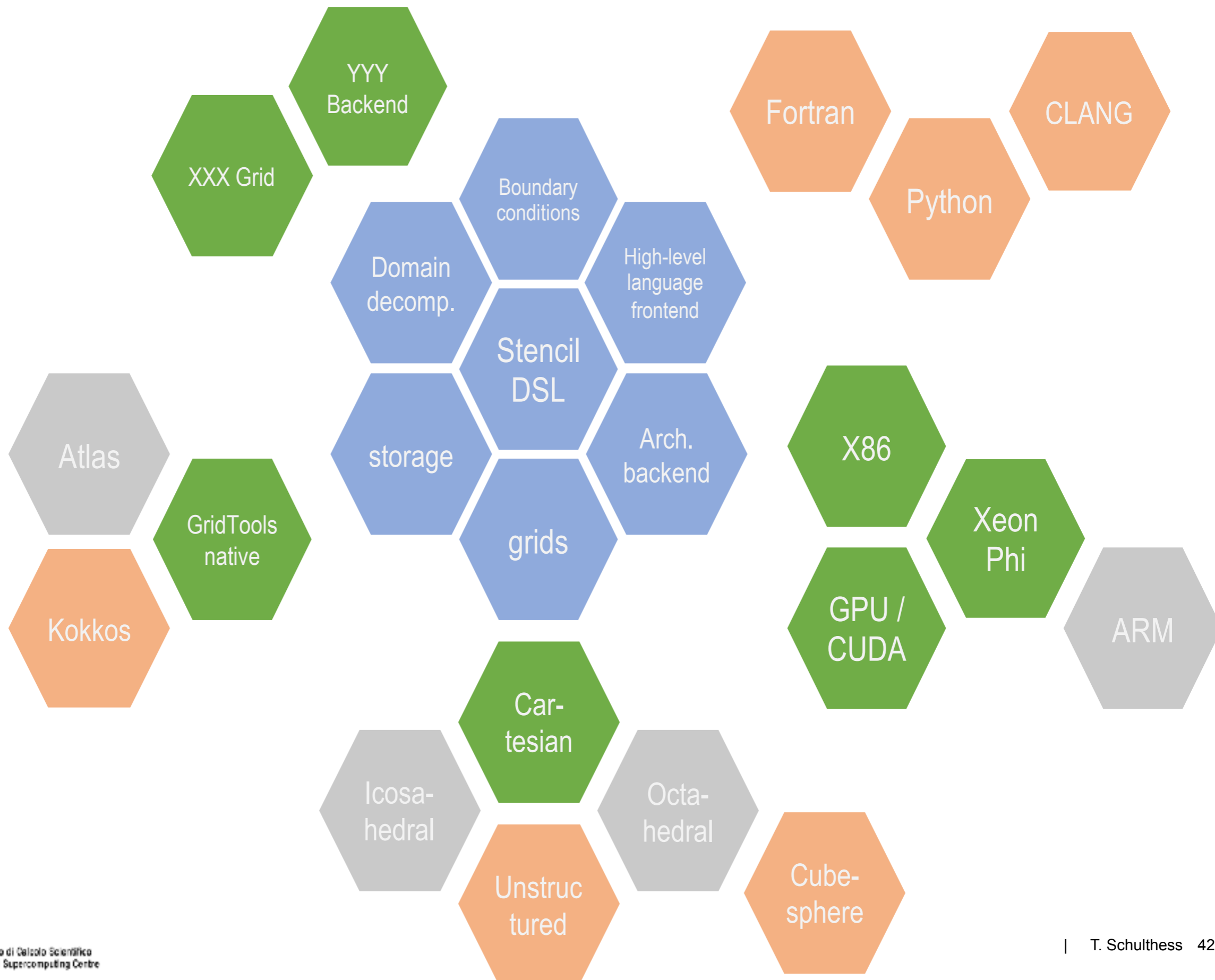
Toolchain



- read before write
- missing boundary update
- data dependency race conditions
- out of bounds stencil access

- Software managed caches
- Full vertical parallelisation
- Stage fusion
- Data locality exploit
- Strong/weak scaling optimiser

- Native C/C++/Fortran
- Optimised GridTools Generator (C++)



The GridTools Team



Anton Afanasyev
Software Architect



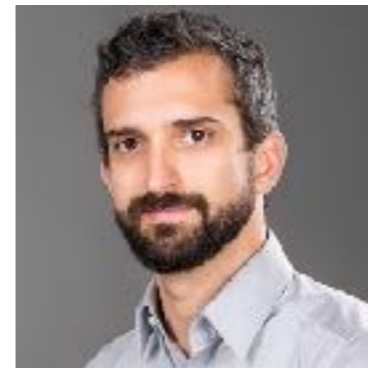
Carlos Osuna
Other Grids, CUDA
Product Owner



Christopher Bignamini
Dycore (no full time)



Felix Thaler
KNL Backend



Hannes Vogt
Dycore Lead, CUDA
Scrum Master



Mauro Bianco
Stencil Composition
Technical Lead



Nora Abi Akar
ARM Backend



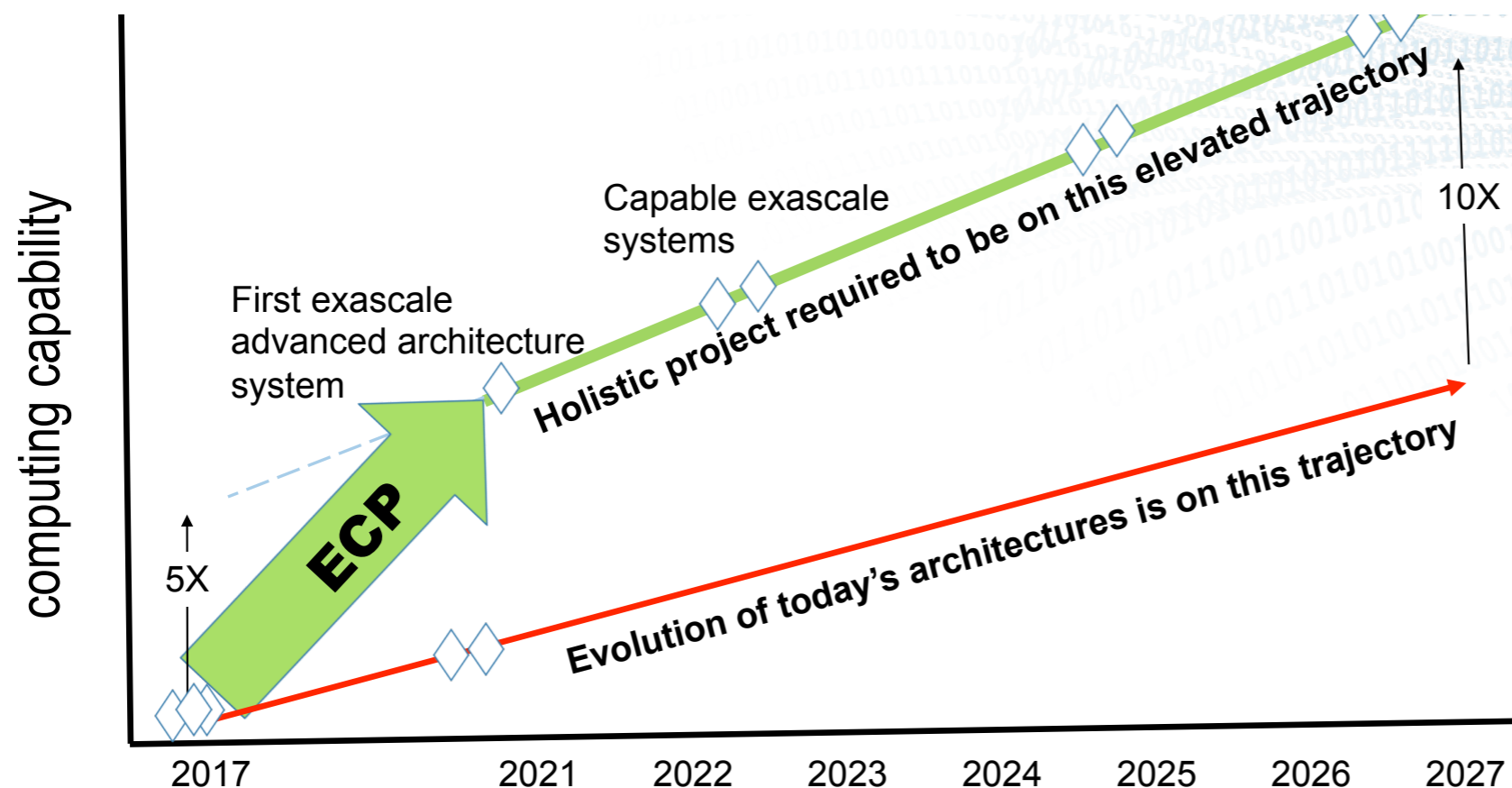
Stefan Moosbrugger
KNL, Storage

What exactly is “exascale” computing?

Exascale is not just $\sim 10^3$ PF

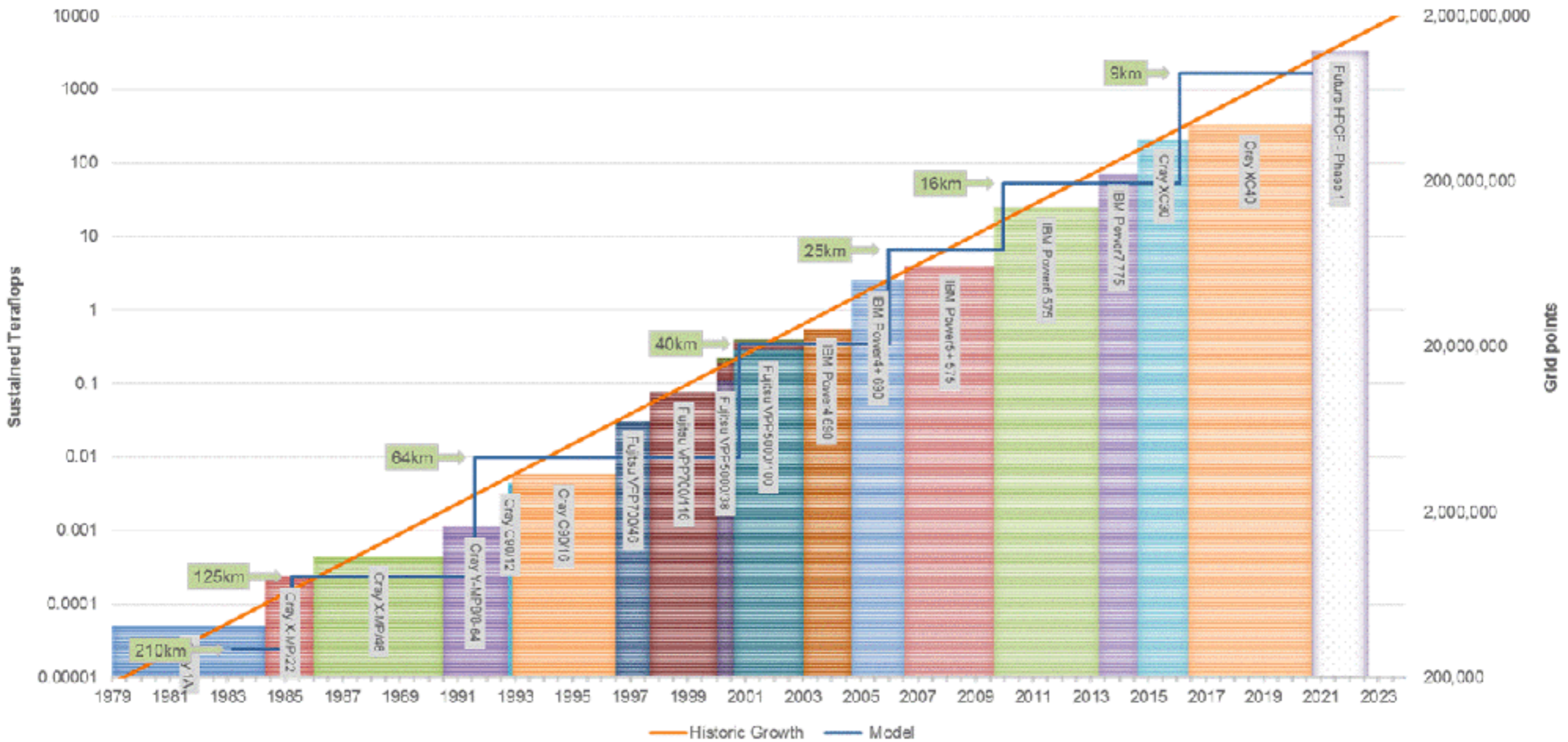
A capable exascale computing system requires an entire computational ecosystem that:

- Delivers 50x the performance of today's ~ 20 PF systems, supporting application that deliver high-fidelity solutions in less time and address problems of greater complexity
- Operates in a power envelope of 20-30 MW
- Is sufficiently resilient (perceived full rate: ≤ 1 /week)
- Includes a software stack that supports a broad spectrum of applications and workloads

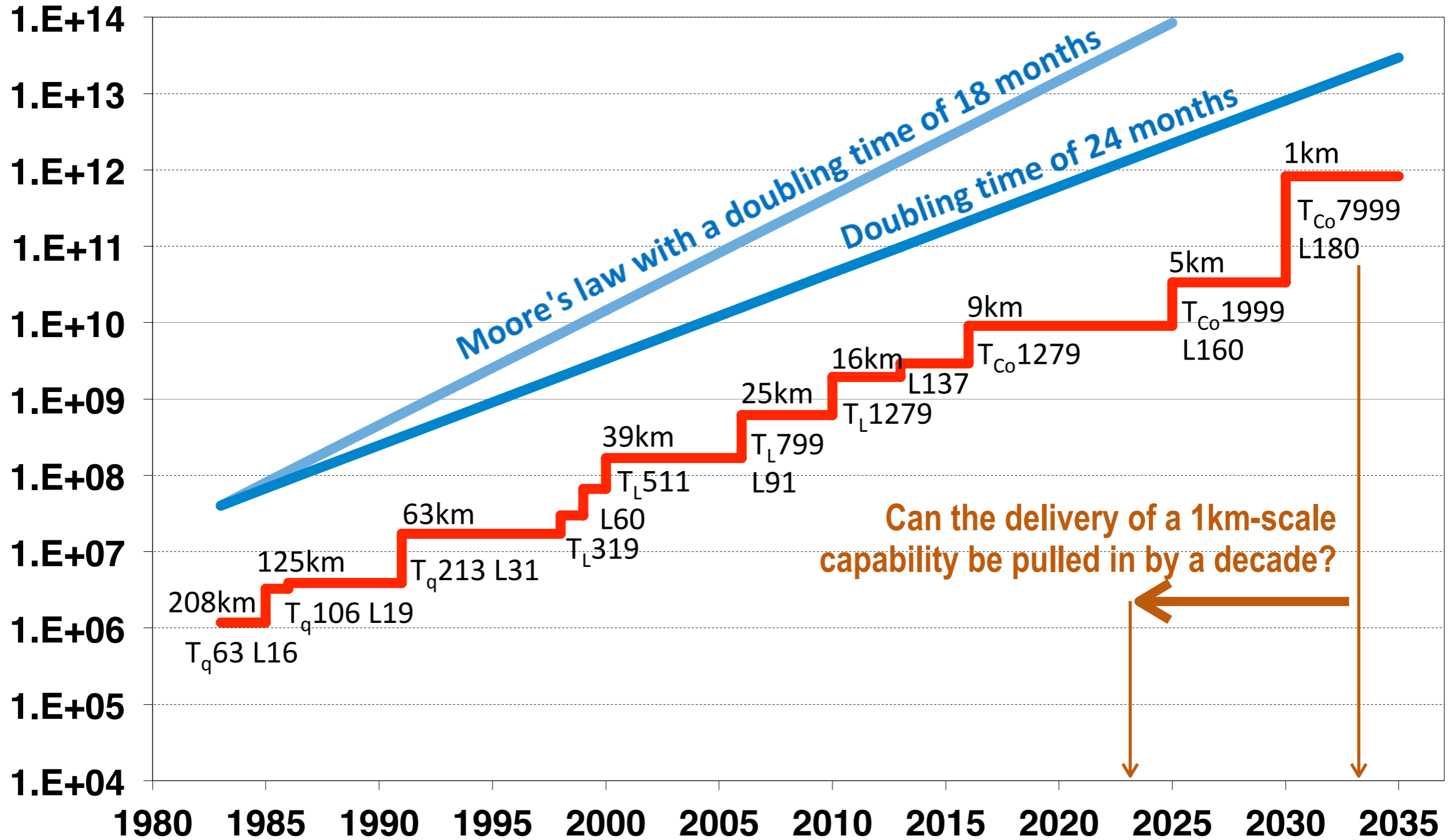


Transition to higher trajectory with advanced architecture

But what is the goal for exascale computing, and the baseline?



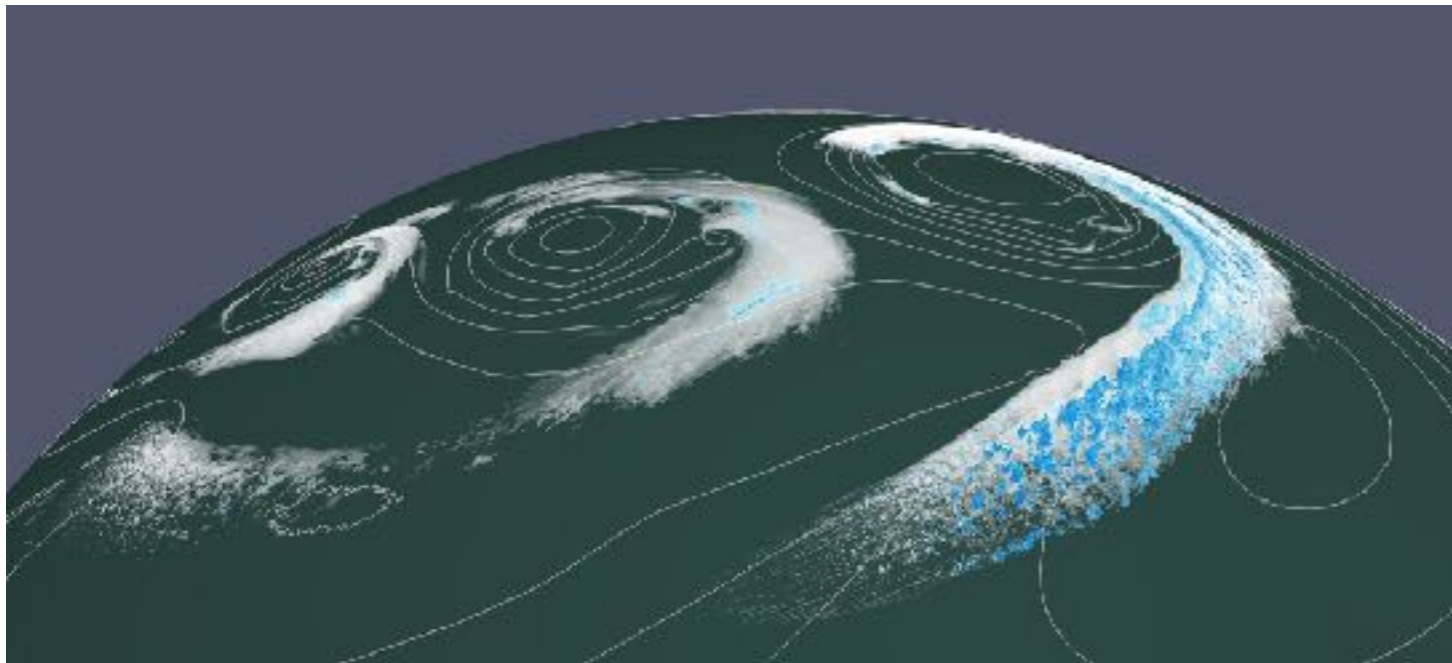
Computational power drives spatial resolution



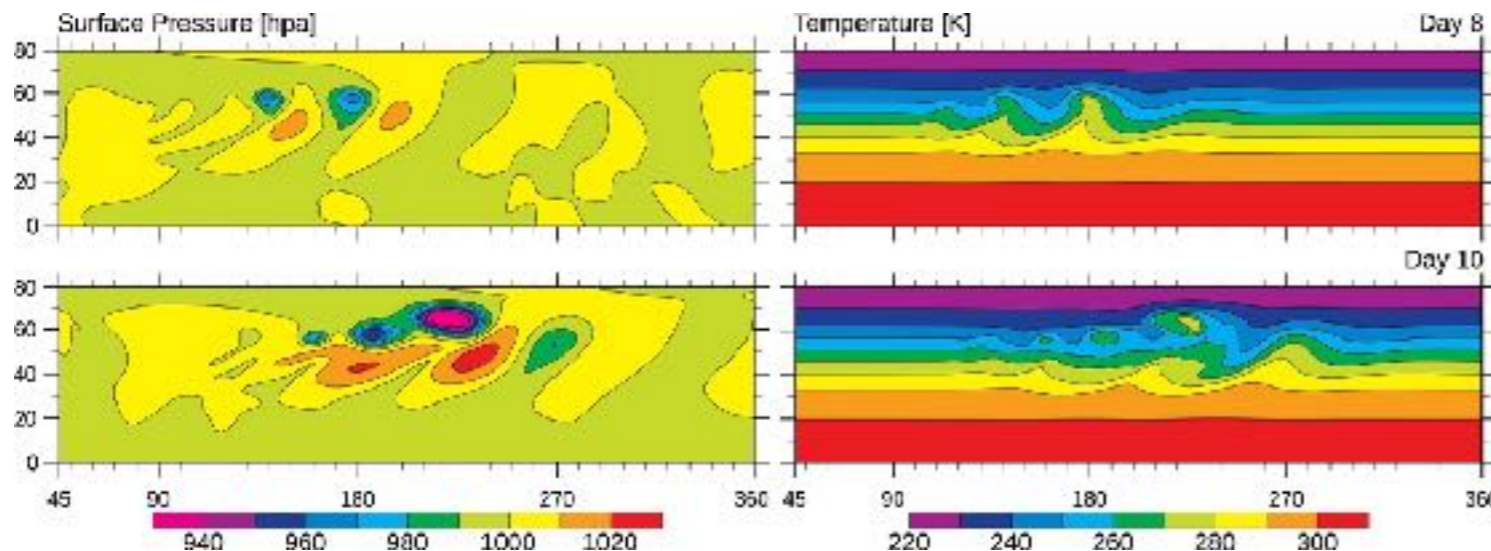
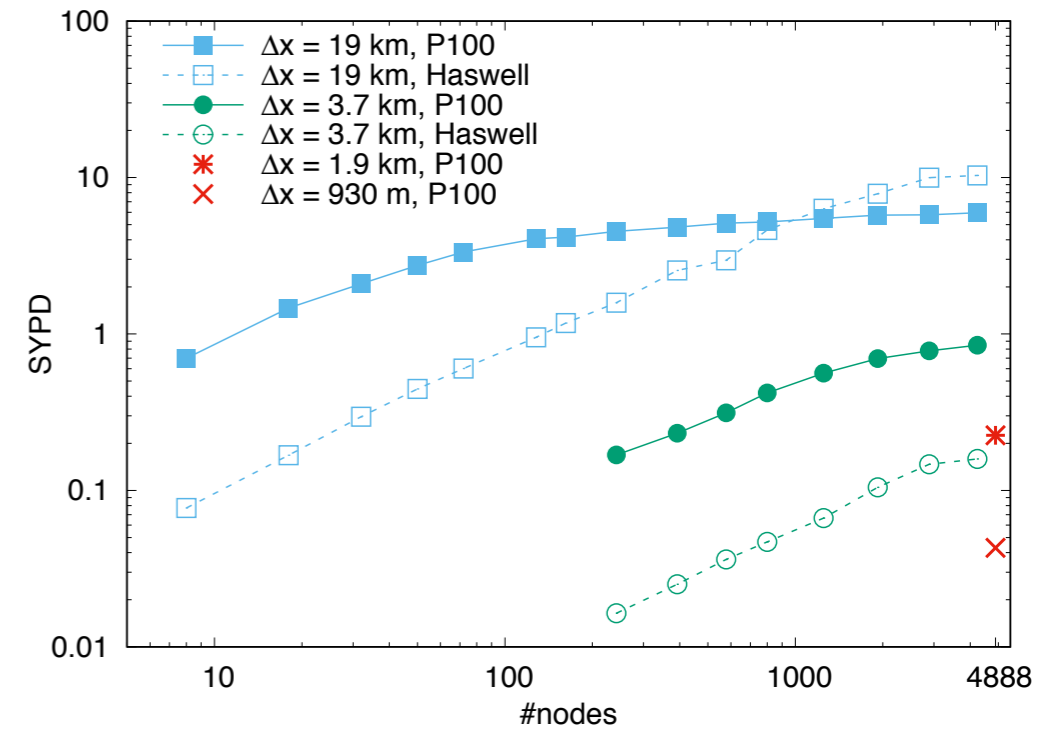
**Let's assume for a moment we can
build on the CSCS-MCH experience**

Near-global climate simulation at 1km resolution: establishing a performance baseline on 4888 GPUs with COSMO 5.0

Fuhrer et al., Geosci. Model Dev. Discuss., <https://doi.org/10.5194/gmd-2017-230>, published 2018



Metric: simulated years per wall-clock day



$\langle \Delta x \rangle$	#nodes	Δt [s]	SYPD	MWh/SY	gridpoints
930 m	4,888	6	0.043	596	3.46×10^{10}
1.9 km	4,888	12	0.23	97.8	8.64×10^9
47 km	18	300	9.6	0.099	1.39×10^7

(c) Time compression (SYPD) and energy cost (MWh/SY) for three moist simulations. At 930 m grid spacing obtained with a full 10d simulation, at 1.9 km from 1,000 steps, and at 47 km from 100 steps

2.5x faster than Yang et al.'s 2016 Gordon Bell winner run on TaihuLight!

“Exascale” goal for global weather and climate runs

Horizontal resolution	1 km (globally quasi-uniform)
Vertical resolution	180 levels (surface to ~100 km)
Time resolution	0.5 minutes
Coupled	Land-surface/ocean/ocean-waves/sea-ice
Atmosphere	Non-hydrostatic
Precision	Single or mixed precision
Compute rate	1 SYPD (simulated years per wall-clock day)

The baseline for COSMO-global and IFS

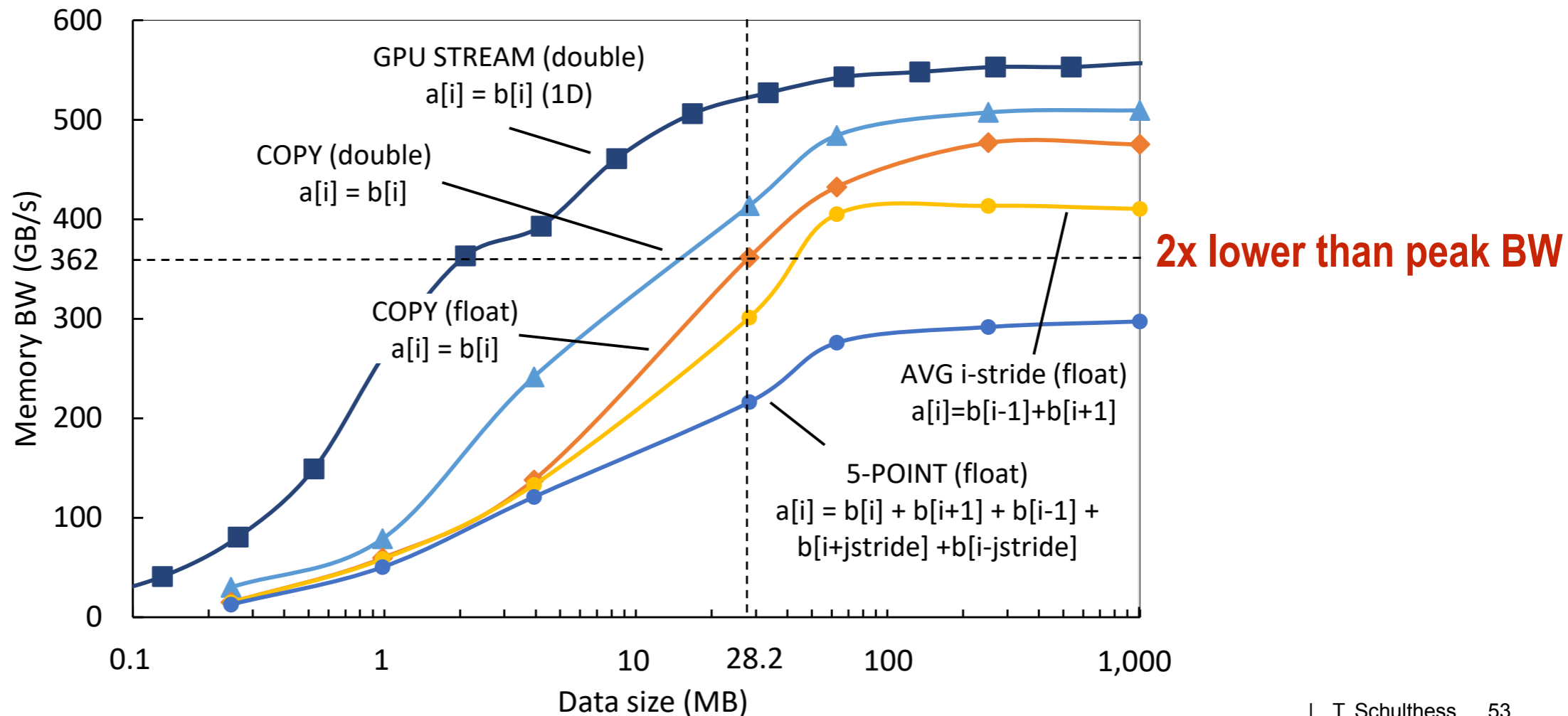
	Near-global COSMO [Fuh2018]		Global IFS [Wed2009]	
	Value	Shortfall	Value	Shortfall
Horizontal resolution	0.93 km (non-uniform)	0.81x	1.25 km	1.56x
Vertical resolution	60 levels (surface to 25 km)	3x	62 levels (surface to 40 km)	3x
Time resolution	6 s (split-explicit with sub-stepping)*	-	120 s (semi-implicit)	4x
Coupled	No	1.2x	No	1.2x
Atmosphere	Non-hydrostatic	-	Non-hydrostatic	-
Precision	Double	0.6x	Single	-
Compute rate	0.043 SYPD	23x	0.088 SYPD	11x
Other (e.g. physics, ...)	microphysics	1.5x	Full physics	-
Total shortfall		60x		247x

Memory use efficiency

Fuhrer et al., Geosci. Model Dev. Discuss., <https://doi.org/10.5194/gmd-2017-230>, published 2018

$$MUE = \text{I/O efficiency} \cdot \text{BW efficiency} = \frac{Q}{D} \frac{B}{\hat{B}} = 0.67$$

Necessary data transfers → Q (0.88) ← Achieved BW
Actual data transfers → D (0.76) ← Max achievable BW



How realistic is it to overcome 65-fold shortfall of a grid-based implementation like COSMO-global?

1. Icosahedral grid (ICON) vs. Lat-long/Cartesian grid (COSMO)

2x fewer grid-columns

Time step of 10 ms instead of 5 ms

4x

2. Improving BW efficiency

Improve BW efficiency and peak BW
(results on Volta show this is realistic)

1.5x

3. Weak scaling

4x possible in COSMO, but we reduced
available parallelism by factor 2

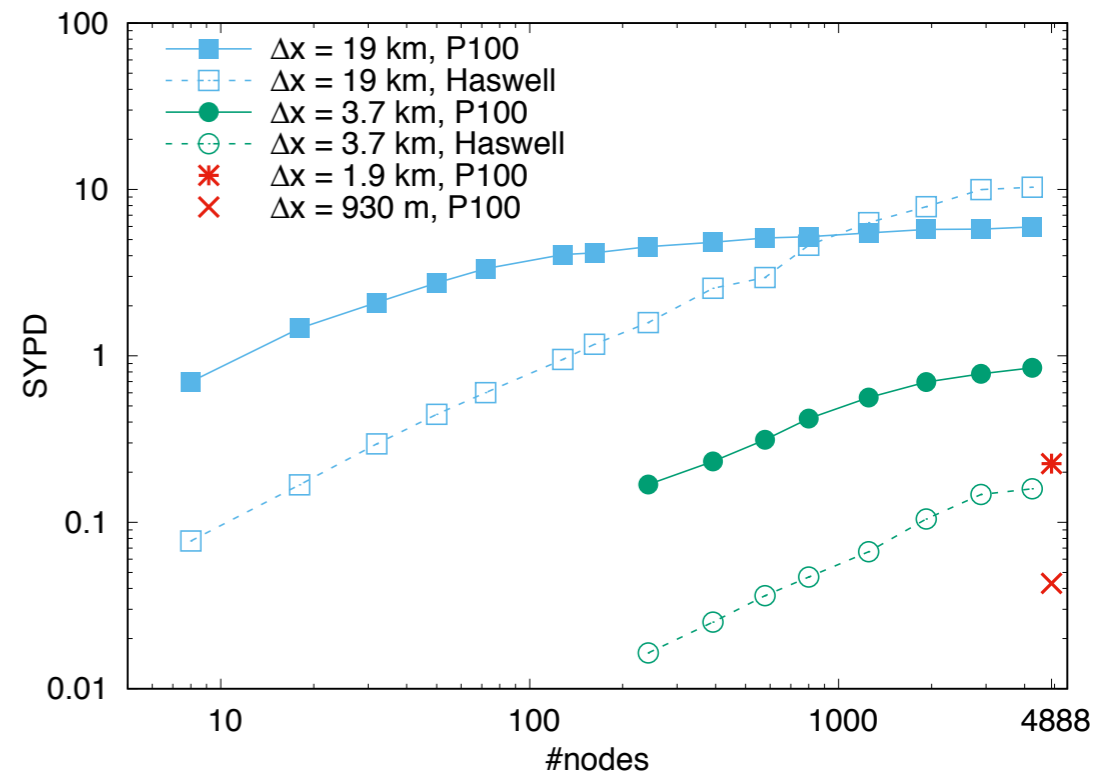
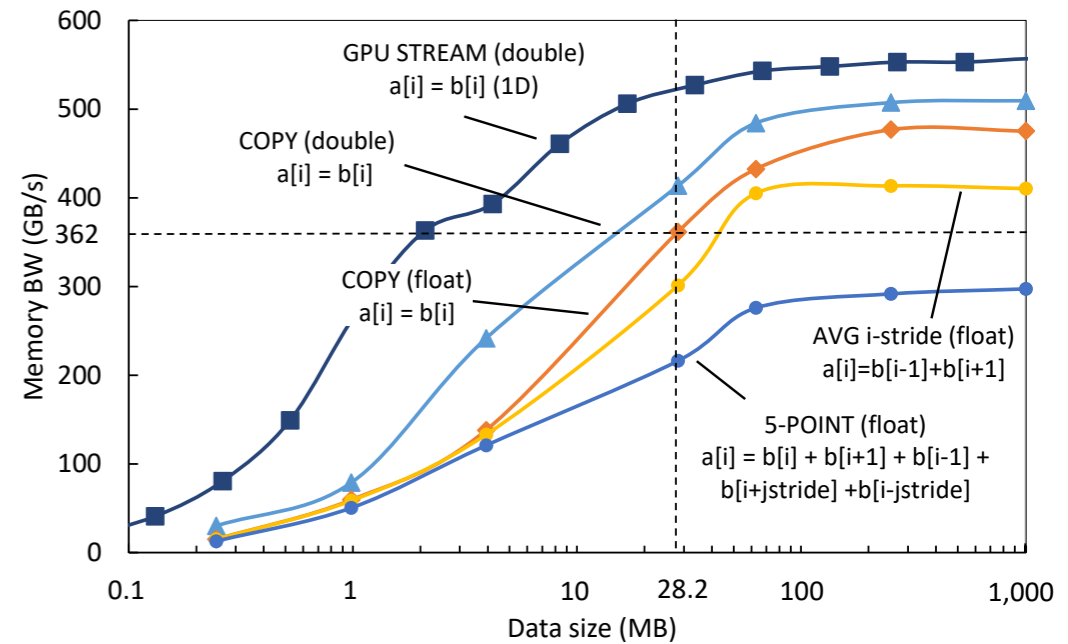
2x

4. Remaining reduction in shortfall

Numerical algorithms (larger time steps)

Further improved processors / memory

5x



But we don't want to increase the footprint of the 2021 system beyond "Piz Daint"

The main conclusions

- Change is nothing new to HPC, nor is the reluctance to adapt to change
 - “Killer micros”, memory wall, end of Dennard Scaling and multi-core, GPU
- CMOS scaling tapering due to constraints in device physics and fabrication
- Architectural improvements & diversity seem a good option to improve performance
- New opportunities for materials science and device physics?
- Fundamental challenge to software / application development
- Domain specific libraries and frameworks are a way out
 - GridTools framework with successful demonstration to COSMO @ MeteoSwiss
- “Exascale” computing, if properly defined and pursued, could give us ~1km scale horizontal resolution in simulation with good throughput

Great motivations to clean up our software stack!

Collaborators



Tim Palmer (U. of Oxford)



Bjorn Stevens (MPI-M)



Peter Bauer (ECMWF)



Oliver Fuhrer (MeteoSwiss)



Nils Wedi (ECMWF)



Torsten Hoefler (ETH Zurich)



Christoph Schar (ETH Zurich)