



From “Piz Daint” to “Piz Kesch”: the making of a GPU-based weather forecasting system

Oliver Fuhrer and Thomas C. Schulthess

“Piz Daint”



Cray XC30 with 5272 hybrid, GPU accelerated compute nodes

Compute node:

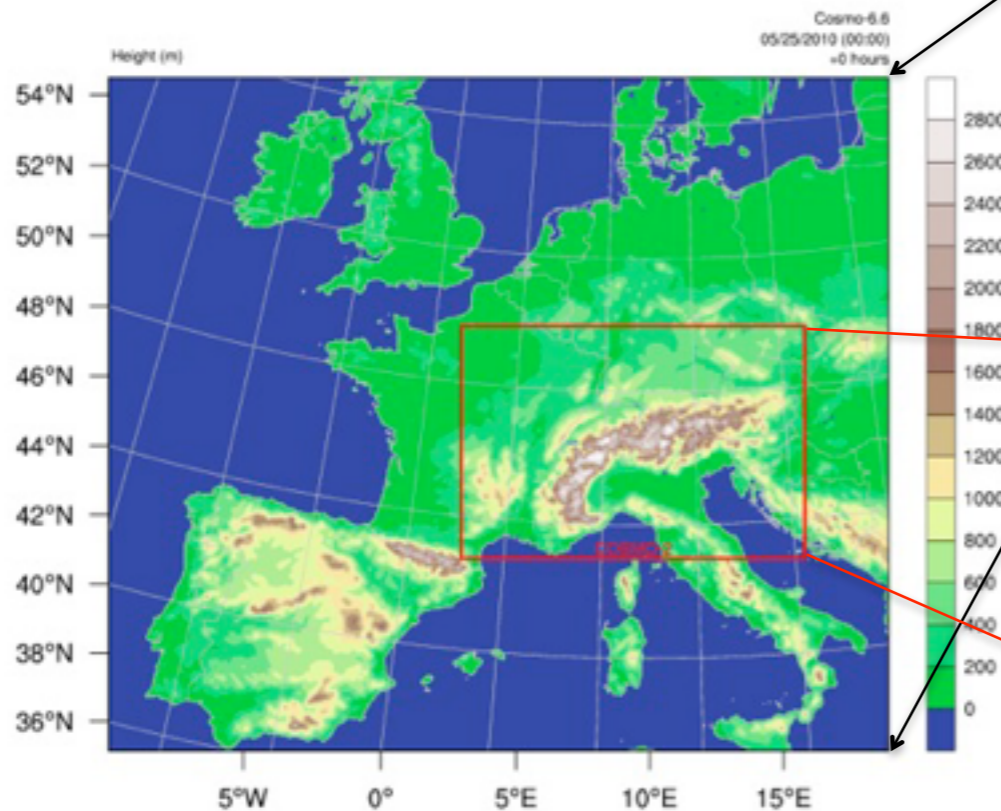
- > Host: Intel Xeon E5 2670 (SandyBridge 8c)
- > Accelerator: One NVIDIA K20X GPU (GK110)

Today's (2015) production suite of Meteo Swiss

COSMO-7

3x per day 72h forecast
6.6 km lateral grid, 60 layers

Orography of COSMO-7

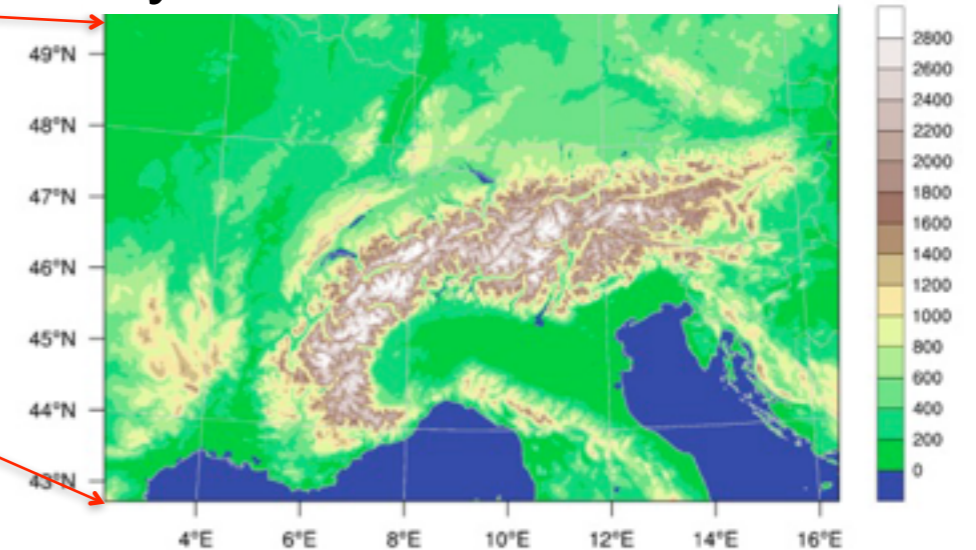


ECMWF

2x per day
16 km lateral grid, 91 layers

COSMO-2

8x per day 24h forecast
2.2 km lateral grid, 60 layers



Some of the products generate from these simulations:

- ▶ Daily weather forecast on TV / radio
- ▶ Forecasting for air traffic control (Sky Guide)
- ▶ Safety management in event of nuclear incidents

“Albis” & “Lema”, CSCS production systems for Meteo Swiss



Cray XE6 procured in spring 2012 based on 12-core AMD Opteron multi-core processors

Cloud resolving simulations

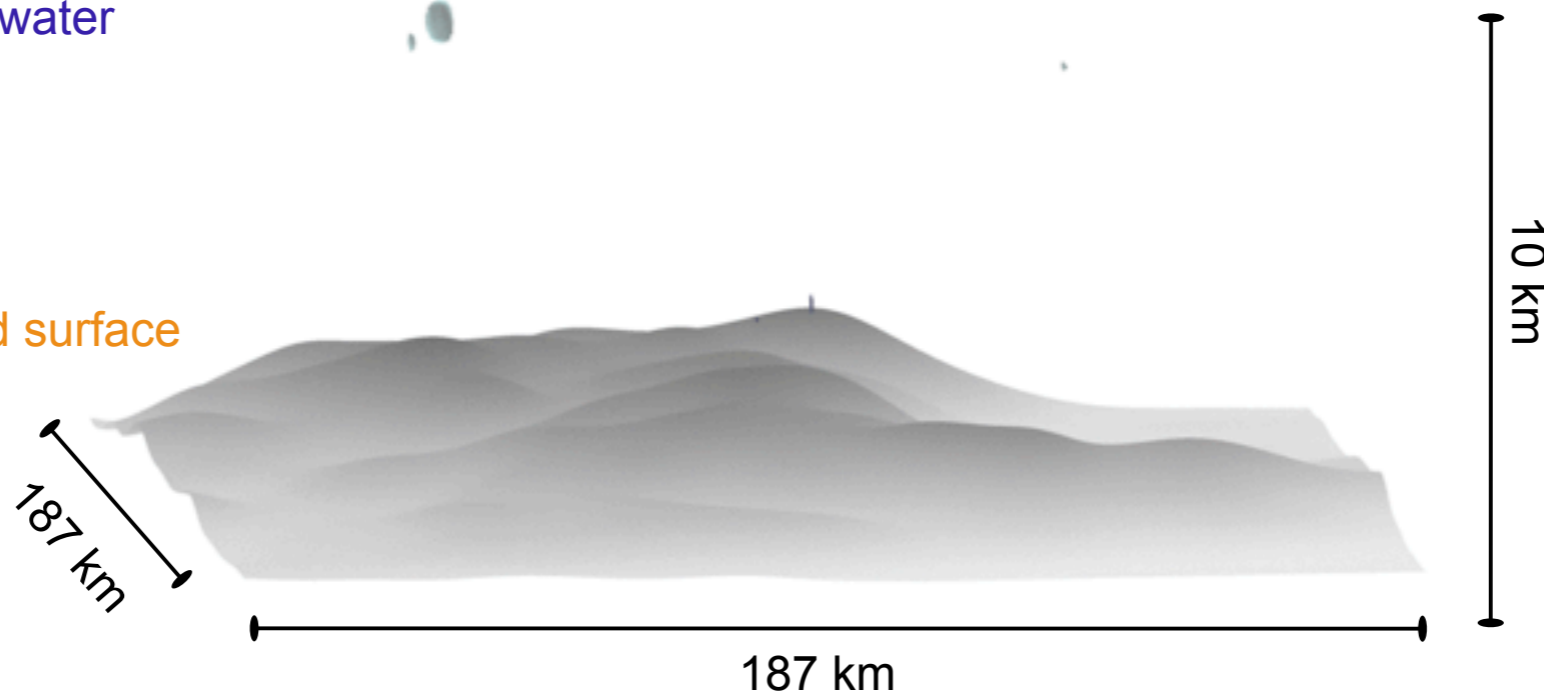
Institute for Atmospheric and Climate Science Study at ETH Zürich (Prof. Schär) demonstrates cloud resolving models converge at 1-2km resolution (at least for convective clouds over the alpine region)

Cloud ice

Cloud liquid water

Rain

Accumulated surface precipitation



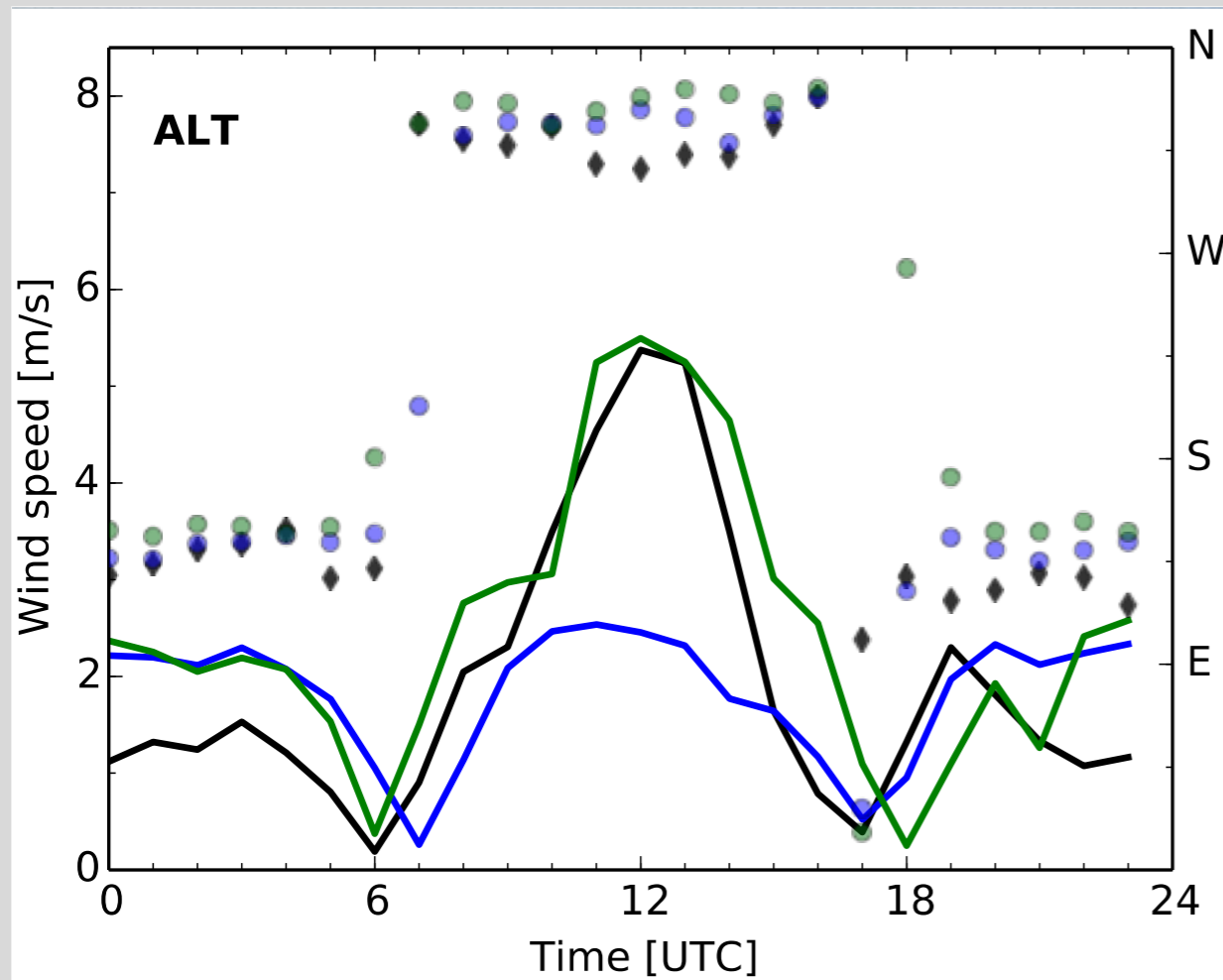
COSMO model setup: $\Delta x=550$ m, $\Delta t=4$ sec Plots generated using INSIGHT

Orographic convection – simulation: 11-18 local time, 11 July 2006 ($\Delta t_{\text{plot}}=4$ min)

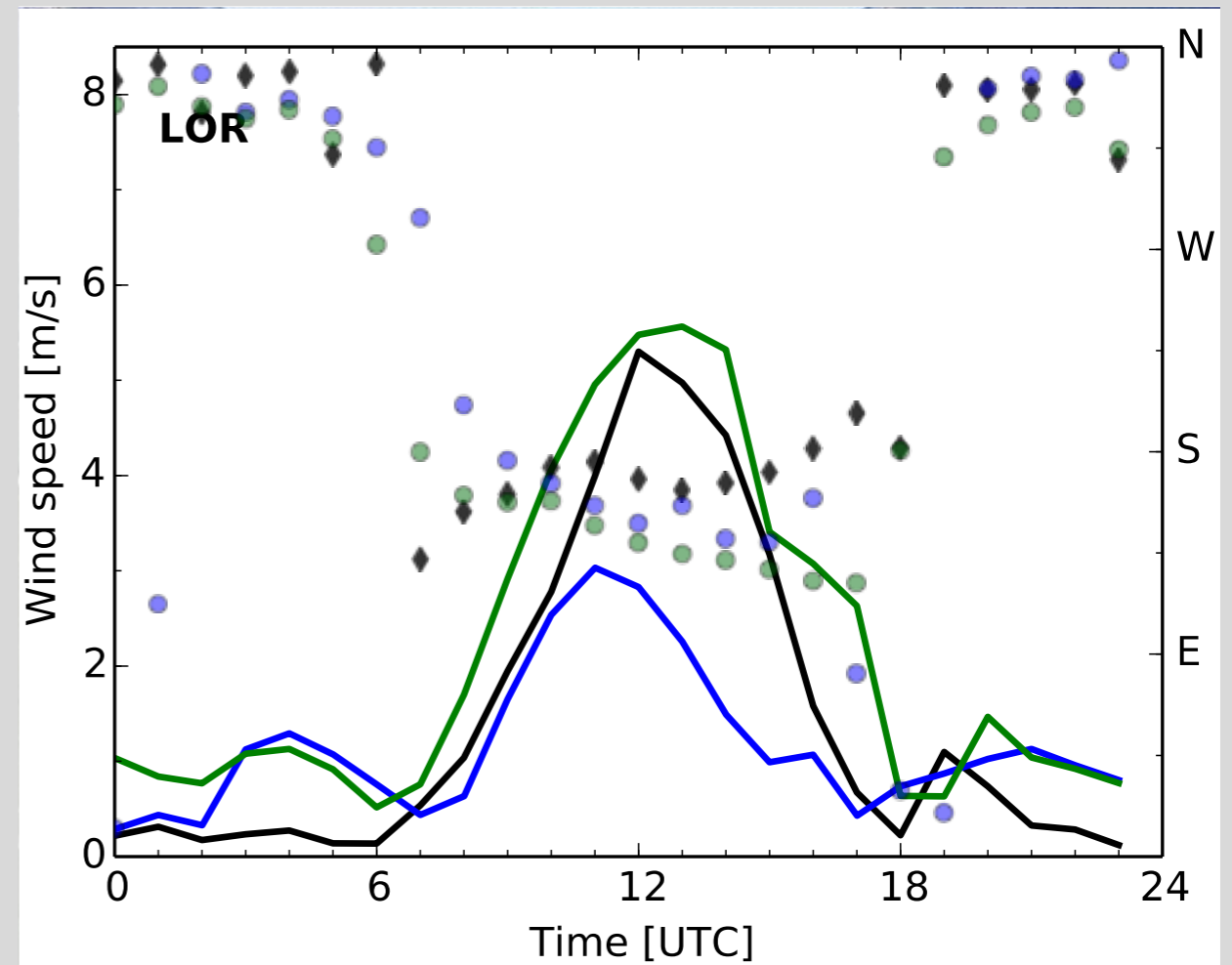
Source: Wolfgang Langhans and Christoph Schär, Institute for Atmospheric and Climate Science, ETH Zurich

Higher resolution is necessary for quantitative agreement with experiment (18 days for July 9-27, 2006)

Aldorf (Reuss valley)



Lodrino (Leventina)



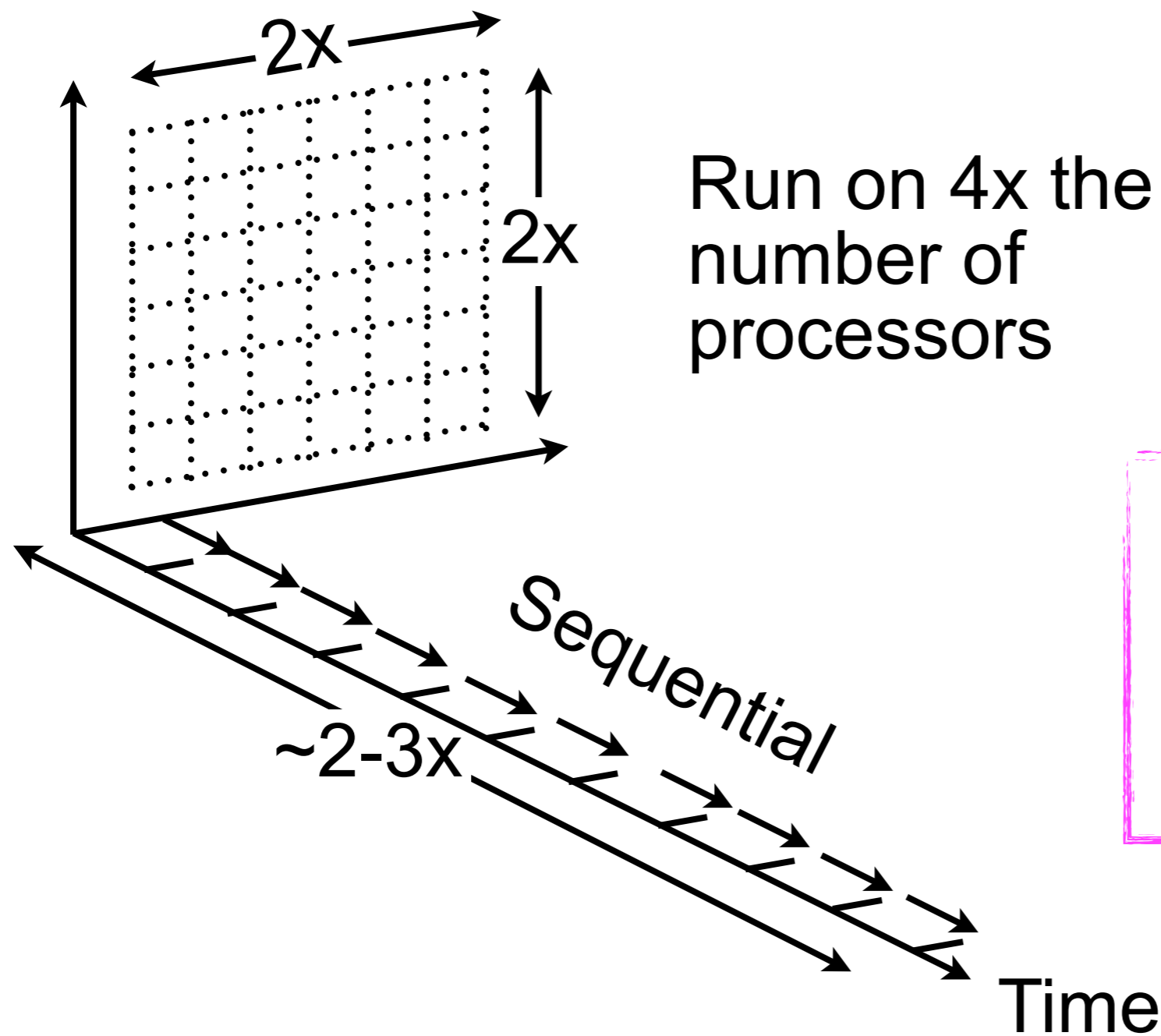
Observation Average wind speed (—) and direction (◇)

COSMO-2

COSMO-1

source: Oliver Fuhrer, MeteoSwiss

Improve resolution of Meteo Swiss model from 2 to 1 km

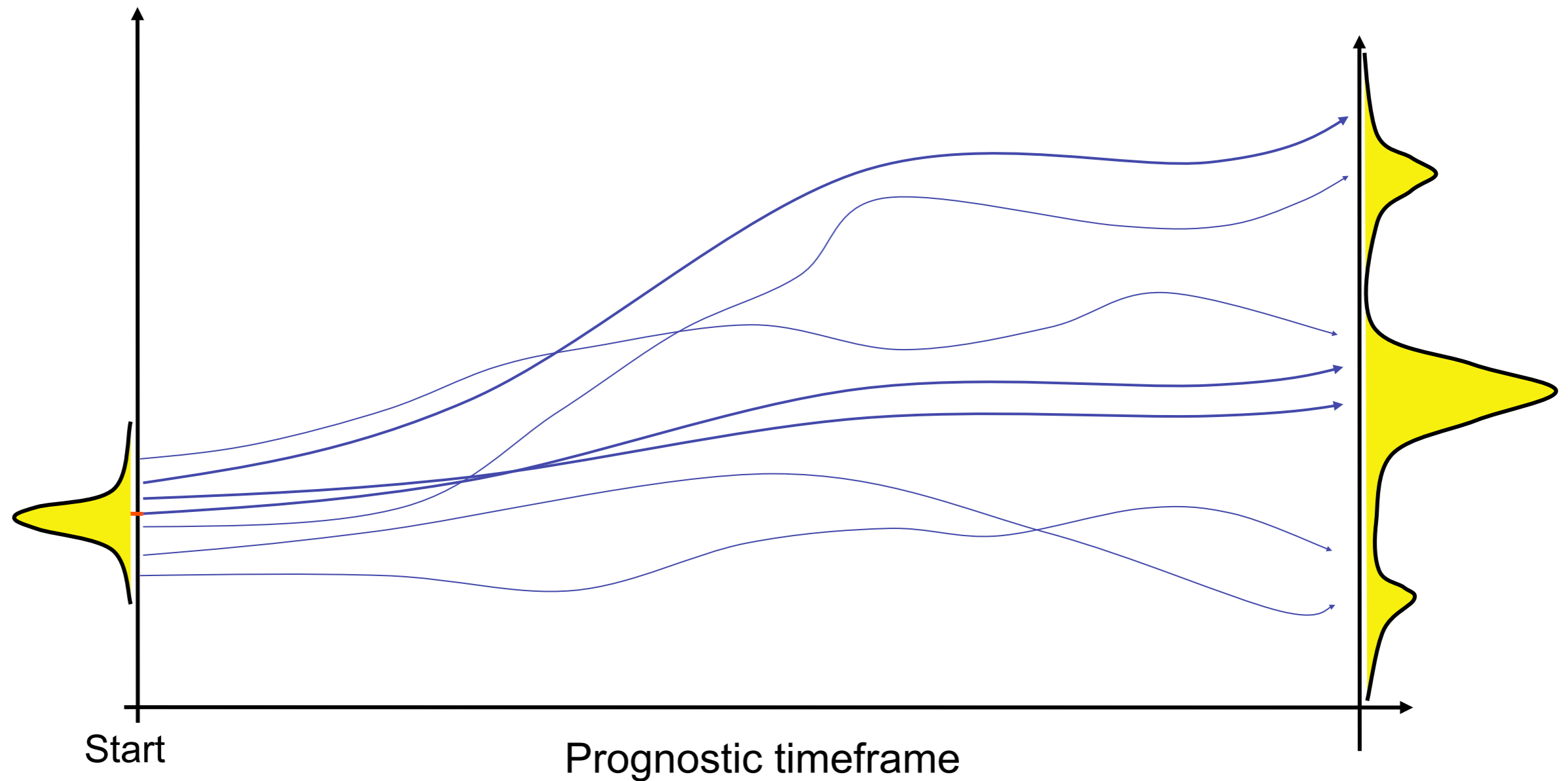


Doubling the resolution requires
 $\sim 10x$ performance increase

Prognostic uncertainty

The weather system is chaotic

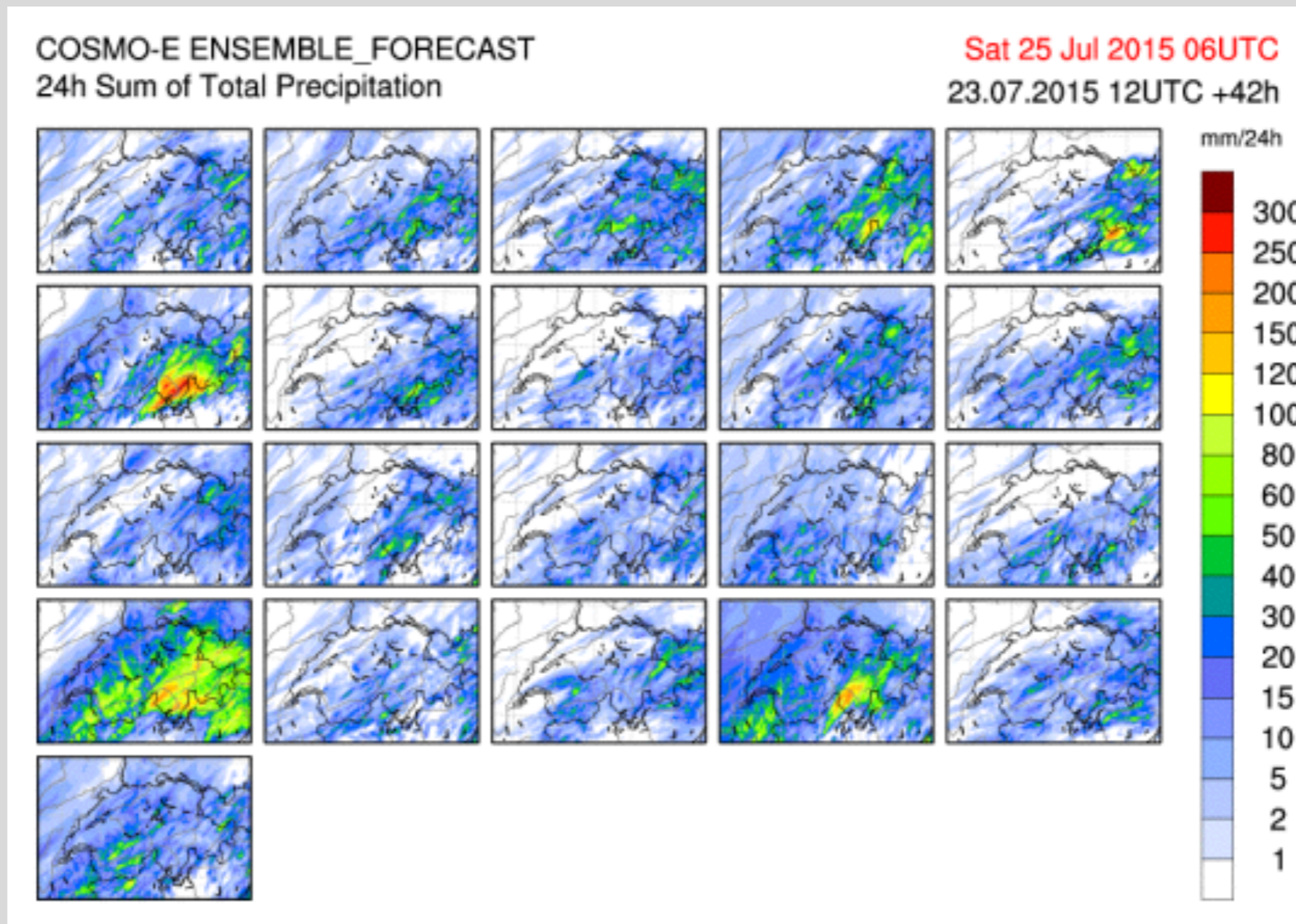
→ rapid growth of small perturbations (butterfly effect)



Ensemble method: compute distribution over many simulations

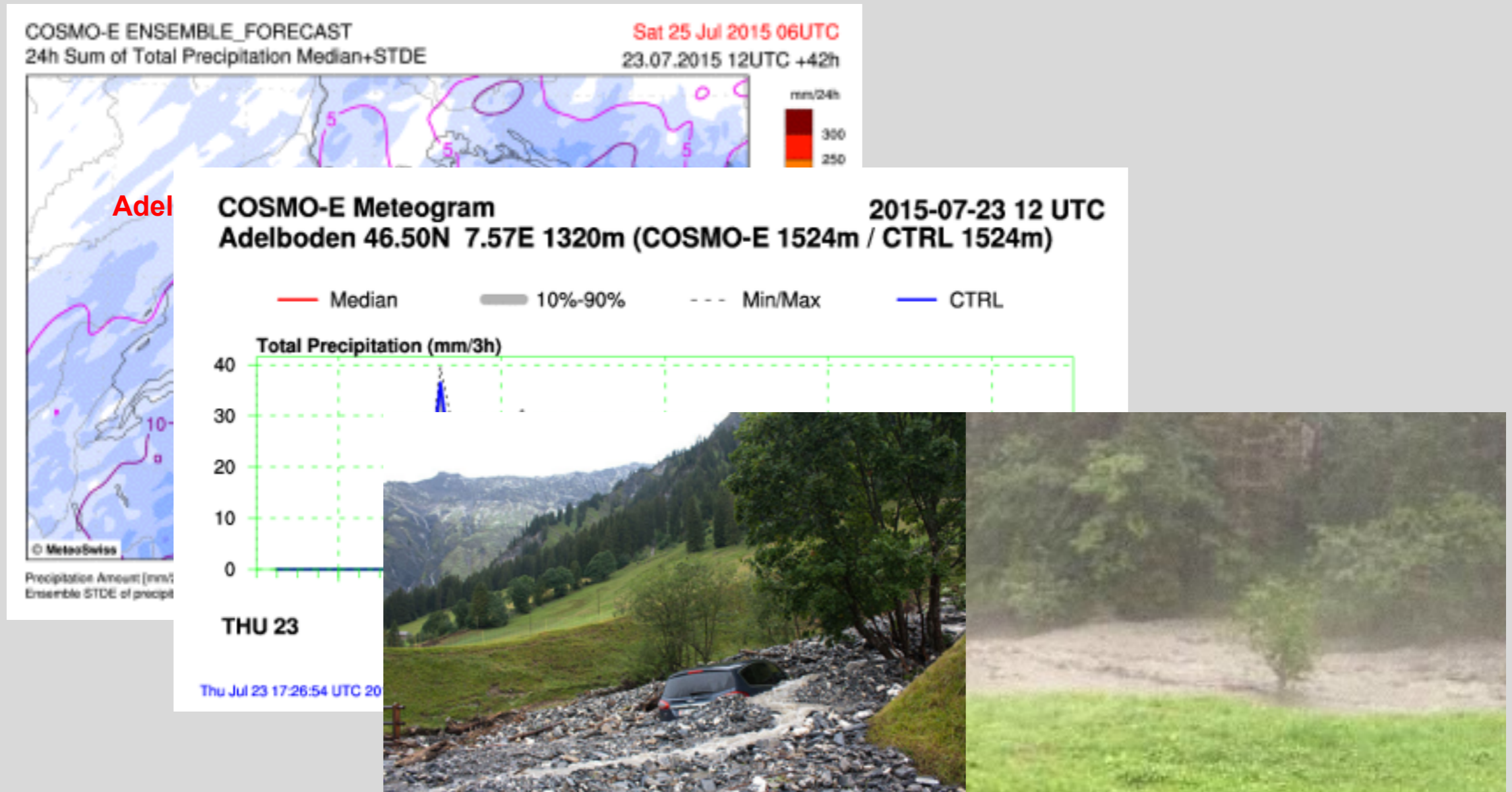
Benefit of ensemble forecast

(heavy thunderstorms on July 24, 2015)



source: Oliver Fuhrer, MeteoSwiss

Benefit of ensemble forecast (heavy thunderstorms on July 24, 2015)

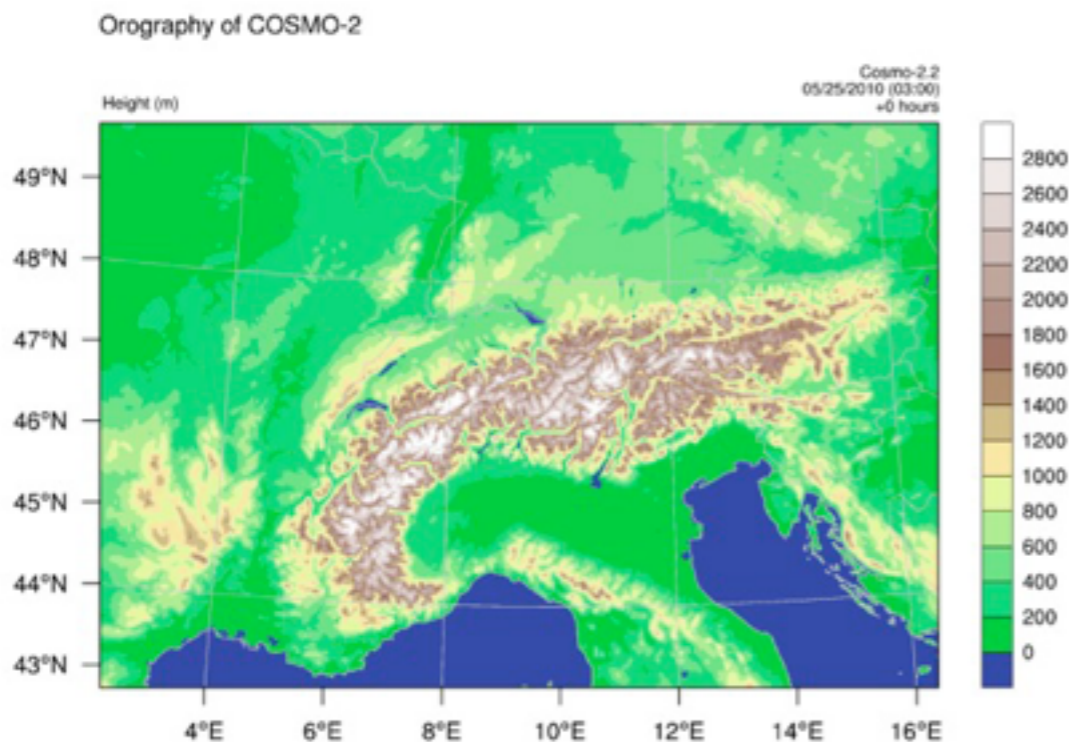


Improving simulation quality requires higher performance – what exactly and by how much?

Resource determining factors for Meteo Swiss' simulations

Current model running through mid 2016

COSMO-2: 24h forecast running in 30 min.
8x per day



New model starting operation on in Jan. 2016

COSMO-1: 24h forecast running in 30 min.
8x per day (~10x COSMO-2)

COSMO-2E: 21-member ensemble, 120h forecast
in 150 min., 2x per day (~26x COSMO-2)

KENDA: 40-member ensemble, 1h forecast
in 15 min., 24x per day (~5x COSMO-2)

New production system must deliver
~40x the simulations performance
of “Albis” and “Lema”

State of the art implementation of new system for Meteo Swiss

Albis & Lema: 3 cabinets Cray XE6 installed Q2/2012

- New system needs to be installed Q2-3/2015
- Assuming 2x improvement in per-socket performance:
~20x more X86 sockets would require 30 Cray XC cabinets

New system for Meteo Swiss if we build it like the German Weather Service (DWD) did theirs, or UK Met Office, or ECMWF ... (30 racks XC)

Current Cray XC30/XC40 platform (space for 40 racks XC)

Thinking inside the box is not a good option!

CSCS machine room

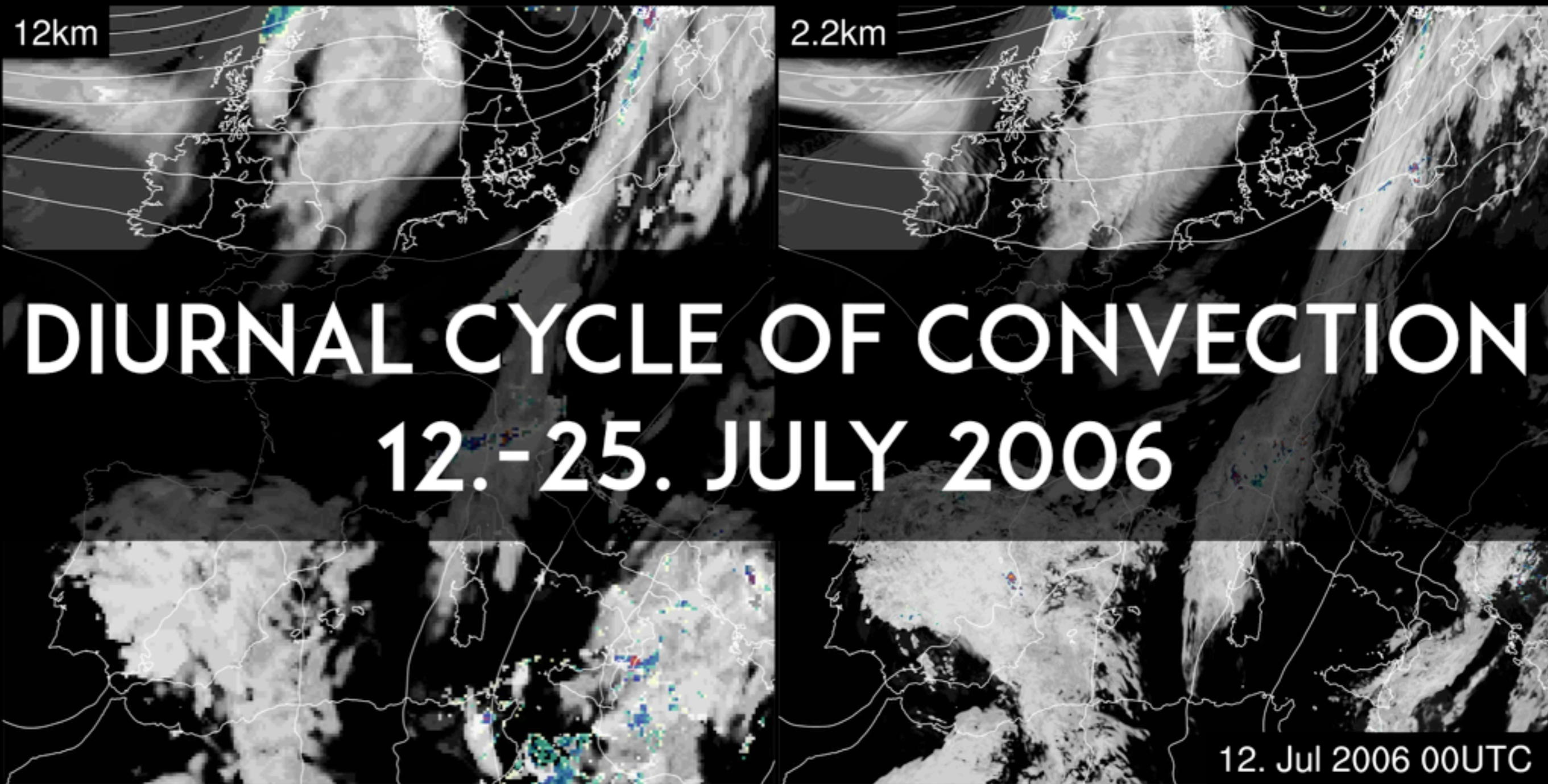
Co-Design our way out?

Potential for co-design

- Time-to-solution driven (specs are clear)
- Exclusive usage
- Only one performance critical application
- Stable configuration (code & system)
- Current code can be improved
- Novel hardware has yet to be exploited

Challenges for making it work

- Community code
 - Large user base
 - Performance portability
 - Knowhow transfer
- Complex workflow
- High reliability required
- Rapidly evolving technology (hardware and software)

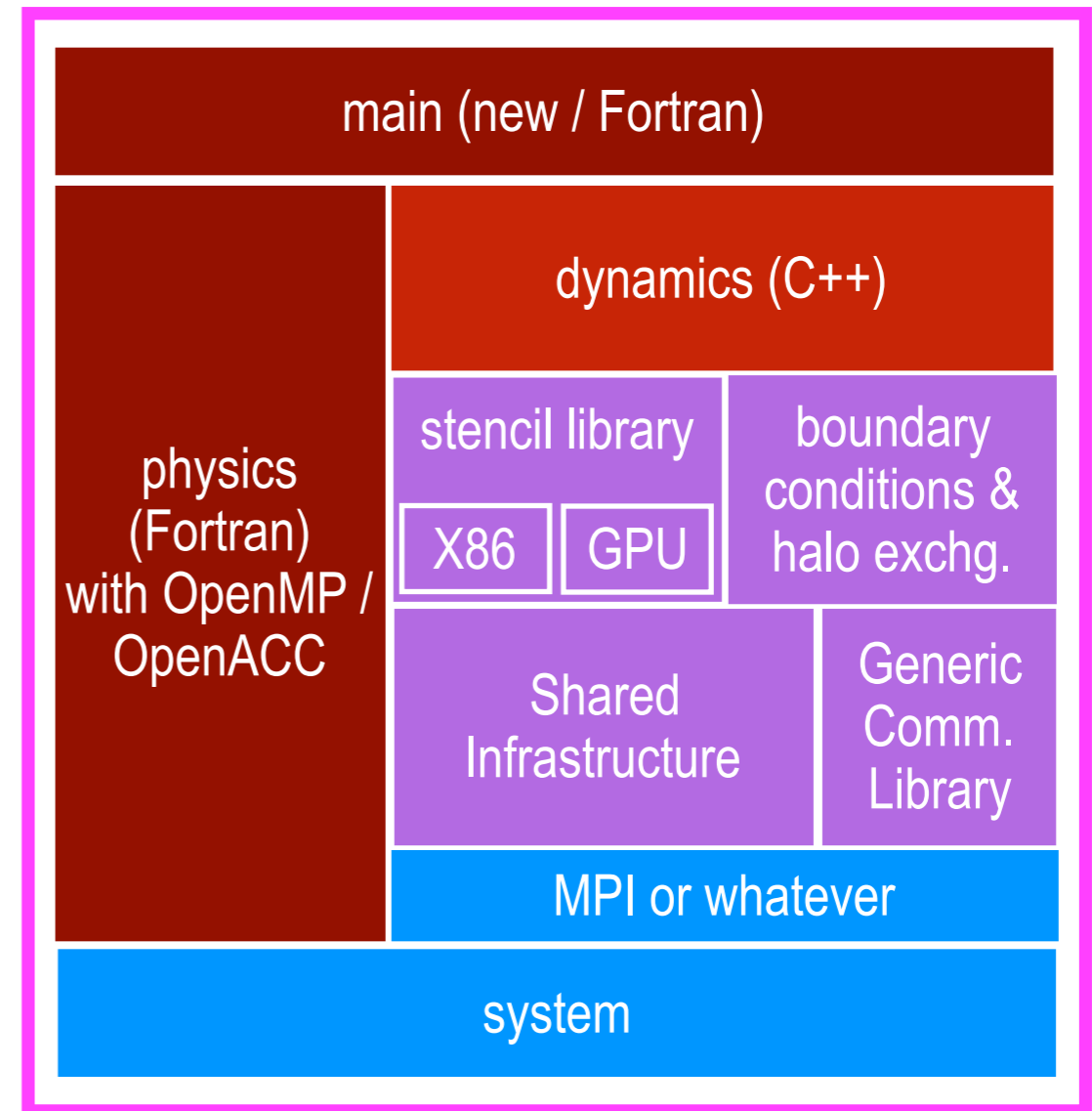
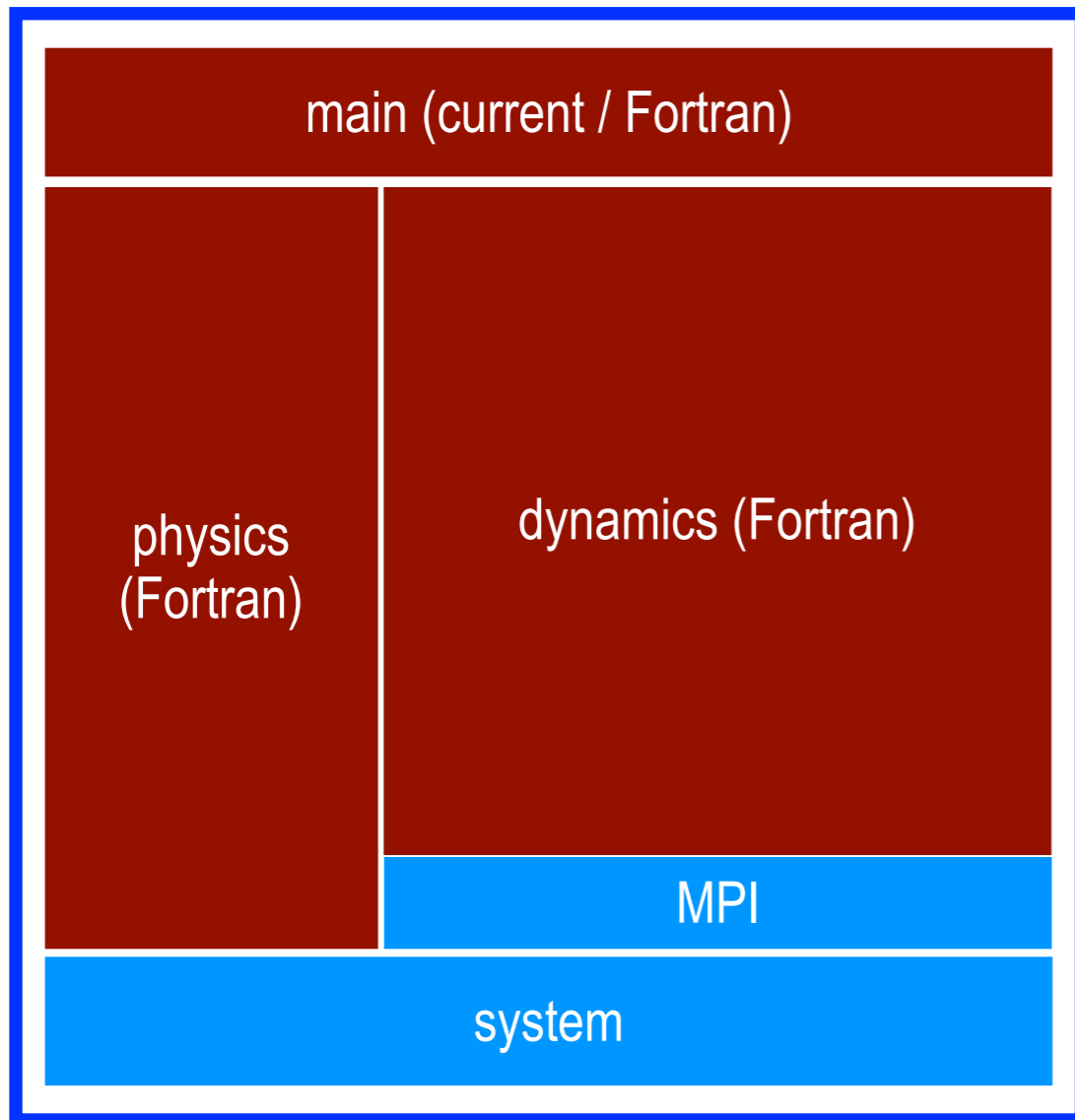


Leutwyler, D., O. Fuhrer, X. Lapillone, D. Lüthi, C. Schär, 2015: Continental-Scale Climate Simulation at Kilometer resolution. ETH Zurich Online Resource, DOI: <http://dx.doi.org/10.3929/ethz-a-010483656>, online video: <http://vimeo.com/136588806>

Co-design approach

- Co-design software / workflow / hardware paying attention to
 - Portability to other users and hardware architectures
 - Achieve specified time-to-solution
 - Optimise hardware footprint and energy
- Several collaboration pre-existed
 - Software development since 2010: MeteoSwiss / C2SM@ETH Zurich / CSCS
 - CSCS with Cray and NVIDIA for development of “Piz Daint” in 2013
 - Domain scientists and computer scientists
- Substantial software investments from HPCN Strategy: HP2C and PASC
- Extreme programming team
 - Oliver Fuhrer (the perfect product owner)
 - Tobias Gysi, Carlos Osuna, Xavier Lapillonne, Mauro Bianco, Andrea Arteaga (not all at the same time)
 - CSCS experts: Ben Cumming, Gilles Fourestey, Guilherme Peretti-Pezzi
 - NVIDIA experts: Peter Messmer, Christoph Angerer

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



A factor 40 improvement with the same footprint

Current production system: Albis & Lema

New system: Kesch & Escha



Piz Kesch / Piz Escha: appliance for meteorology

- Water cooled rack (48U)
- 12 compute nodes with
 - 2 Intel Xeon E5-2690v3 12 cores @ 2.6 GHz 256 GB 2133 MHz DDR4 memory
 - 8 NVIDIA Tesla K80 GPU
- 3 login nodes
- 5 post-processing nodes
- Mellanox FDR InfiniBand
- Cray CLFS Luster Storage
- Cray Programming Environment



Origin of factor 40 performance improvement

Performance of COSMO running on new "Piz Kesch" compared to current production systems



- Current production system installed in 2012
- New Piz Kesch/Escha installed in 2015
 - Processor performance **2.8x** ← Moore's Law
 - Improved system utilisation **2.8x**
 - General software performance **1.7x** ← Software refactoring
 - Port to GPU architecture **2.3x**
 - Increase in number of processors **1.3x**
 - Total performance improvement **~40x**
- Bonus: simulation running on GPU is **3x** more energy efficient compared to conventional state of the art CPU

Outlook

- Continue to invest in software
 - domain specific libraries / embedded languages
 - improve scientist's productivity through Python bindings
 - refactor entire software toolchain
- Continued performance improvements for climate / meteorology simulations
 - hardware-software co-design
 - improved memory performance
 - processor performance / explore new architectures
- Longterm investment in new model with even higher resolution

References and Collaborators

- Peter Messmer and his team at the NVIDIA co-design lab at ETH Zurich
- Teams at CSCS and Meteo Suisse, group of Christoph Schaer @ ETH Zurich
- O. Fuhrer, C. Osuna, X. Lapillonne, T. Gysi, B. Cumming, M. Bianco, A. Arteaga, T. C. Schulthess, “Towards a performance portable, architecture agnostic implementation strategy for weather and climate models”, Supercomputing Frontiers and Innovations, vol. 1, no. 1 (2014), see superfri.org
- G. Fourestey, B. Cumming, L. Gilly, and T. C. Schulthess, “First experience with validating and using the Cray power management database tool”, Proceedings of the Cray Users Group 2014 (CUG14) (see arxiv.org for reprint)
- B. Cumming, G. Fourestey, T. Gysi, O. Fuhrer, M. Fatica, and T. C. Schulthess, “Application centric energy-efficiency study of distributed multi-core and hybrid CPU-GPU systems”, Proceedings of the International Conference on High-Performance Computing, Networking, Storage and Analysis, SC’14, New York, NY, USA (2014). ACM
- T. Gysi, C. Osuna, O. Fuhrer, M. Bianco and T. C. Schulthess, “**STELLA: A domain-specific tool for structure grid methods in weather and climate models**”, to be published in Proceedings of the International Conference on High-Performance Computing, Networking, Storage and Analysis, SC’15, New York, NY, USA (2015). ACM

paper at SC15: 11/18 @ 1:30-2PM room 18AB

Join us @ the PASC16 Conference

PASC16 provides an opportunity for scientists and practitioners to discuss key issues in the use of High Performance Computing (HPC) in branches of science that require computer modelling and simulations. The scientific program will offer invited lectures, minisymposia, contributed talks and poster presentations. The active participation of graduate students and postdocs is strongly encouraged.

PASC16

Platform for Advanced Scientific Computing Conference

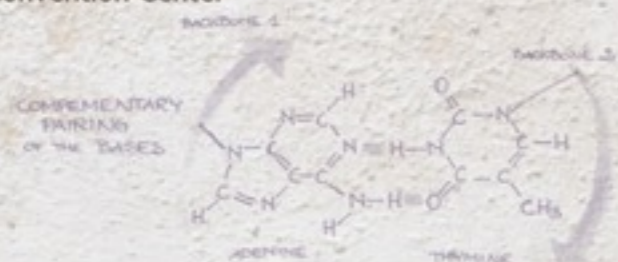
Lausanne Switzerland | 08-10 June 2016



Conference information, registration and submission www.pasc16.org

Queries may be addressed to pasc16@pasc-ch.org

Venue
EPF Lausanne
Swiss Tech Convention Center
Lausanne
Switzerland



NAVIER-STOKES EQUATION

$$\rho \left(\frac{\partial \mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{v} \right) = -\nabla p + \mu \nabla^2 \mathbf{v} + \mathbf{f}$$

\mathbf{f} = body forces (gravity or centrifugal)



quicksort(A, i, k)
if i < k
p = partition(A, i, k)
quicksort(A, i, p-1)
quicksort(A, p+1, k)

Contributions

Researchers from the academic and from the corporate world are invited to participate and present their research area in the form of minisymposia, contributed talks and/or poster presentations. PASC16 welcomes submissions in the following scientific fields:

- CLIMATE & WEATHER
- SOLID EARTH
- LIFE SCIENCE
- CHEMISTRY & MATERIALS
- PHYSICS
- COMPUTER SCIENCE & MATHEMATICS
- ENGINEERING
- EMERGING DOMAINS

METROPOLIS ALGORITHM

```

initialize  $x_0, n$  and  $s$ 
for  $i = 1: (n-1)$  do
  while  $x_{i+1}$  not assigned do
    draw  $z \in [0,1]$  and  $u_i \in [-1,1]$ 
     $x_{new} = x_i + u_i s$ 
    if  $f(x_{new})/f(x_i) \geq z$  then  $x_{i+1} = x_{new}$ 
  end while
end for
    
```

Abstracts should describe original, interesting, and solid scientific content that is relevant to computational sciences and HPC. Cross-disciplinary approaches are highly encouraged.

Organization Committee

Maria Grazia Giuffreda, ETH Zurich / CSCS

Jan Hesthaven, EPFL

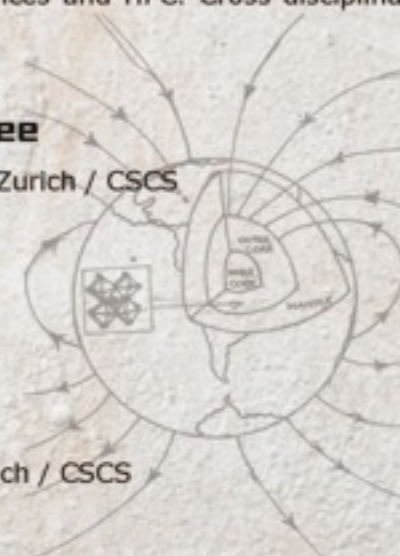
Torsten Hoefler, ETH Zurich

David Keyes, KAUST

Nicola Marzari, EPFL

Olaf Schenk, USI

Thomas Schulthess, ETH Zurich / CSCS



POISSON'S EQUATION

$$\Delta g = f$$

Δ = LAPLACE OPERATOR
 $f = g$ IDEAL OF COMPLEX-VALUED FUNCTIONS

$$\nabla^2 g = f$$

IN THREE-DIMENSIONAL CARTESIAN COORDINATES:

$$\left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} \right) g(x,y,z) = f(x,y,z)$$

when $f=0$ we obtain LAPLACE'S EQUATION

EULER EQUATION

$$\frac{\partial \rho}{\partial t} + \sum_{i=1}^3 \frac{\partial (\rho u_i)}{\partial x_i} = 0$$

$$\frac{\partial (\rho u_j)}{\partial t} + \sum_{i=1}^3 \frac{\partial (\rho u_i u_j)}{\partial x_i} + \frac{\partial p}{\partial x_j} = 0$$

$$\frac{\partial E}{\partial t} + \sum_{i=1}^3 \frac{\partial ((E+p)u_i)}{\partial x_i} = 0$$

i, j label the three Cartesian components
 $(x_1, x_2, x_3) = (x, y, z)$ and
 $(u_1, u_2, u_3) = (u, v, w)$

