Con Con Con

Schweizerische Eidgenossenschaft Confédération suisse Confederazione Svizzera Confederaziun svizra

Swiss Confederation

Federal Department of Home Affairs FDHA Federal Office of Meteorology and Climatology MeteoSwiss

Establishing a baseline for global climate simulations at 1 km on hybrid HPC architectures with COSMO

Carlos Osuna and Oliver Fuhrer, MeteoSwiss

Contributions from T.- Chadha², T. Gysi³, T. Hoefler³, G. Kwasniewski³, X. Lapillonne¹, D. Leutwyler⁴, D. Lüthi⁴, C. Osuna¹, Ch. Schär⁴, T. Schulthess^{5,6}, and H. Vogt⁶

¹MeteoSwiss, ²ITS RI ETH, ³SPCL ETH, ⁴IAC ETH, ⁵ITP ETH, ⁶CSCS

5th HPC Workshop (ESiWACE) - Lecce May 17-18, 2018

Outline

- 1. Why do we want km-scale global weather and climate simulations?
- 2. Why is it hard?
- 3. What can be achieved today using a refactored code on Europe's largest supercomputer?
- 4. Open questions and challenges



The Greyzone



€ to 95% range

IPCC AR5, 2013

RCP8

"Business as Usual" Scenario





Uncertainties are primarily due to uncertainties in the response of clouds. (e.g. Schneider et al. 2017, Nature CC) MeteoSwiss 4 4 4

adapted from Schär, 2017

• What is the difference?

Low climate sensitivity



1. Why do we want km-scale global weather and climate simulations?

- Explicitly resolve important processes (e.g. deep convection, gravity waves, ocean eddies)
- Reduce uncertainty in cloud response in climate model projections



Current state-of-the-art (CMIP5) ∆x ≈ 25 km



 $\Delta x \approx 1 \text{ km}$

How to achieve a factor 10'000?

- Option 1: Money
- Option 2: Wait 15 years

Gordon Moore (1965, Intel co-founder)

"Number of transistors in per inch² in a CPU doubles every 18-24 months at constant cost"



Source: top500.org

	÷	5	5					÷	4			<u>4</u> 2		
÷	55	÷		÷.		4					4		÷	÷
	÷	4	÷			45 C			>	4		<u>ج</u>		
÷	÷			÷		·	수 公		÷		÷	U U	수 순	- ÷
		÷	5	4 4 A	~ ~ ~ ~	÷	4	않 순		C	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	~ ~ ~		
	4						÷ ÷		수 수				÷ ÷	수 수 수
	÷	5 S	100	A. O. L. O. L.	÷	수 · 수 · 수	1	÷ +		÷ ÷	4 4 <i>4</i>		·	÷
53			Mete	OSWISS	4 4	4 4 4	상수		수 수	- ÷	5 th ENES wor	rkshon May 20	184 4	8 🕂
÷	÷	수수 수		0011100		÷ + + + +	ф ф		ኑ ቍ ቍ	· 수수 수		ngiop, may zo	ф с ф	- - - -
÷	수 수수	수수 수수수	፡ ራራ ራ	· ~ ~		· + + +	4 4	4 4		4	-ttt	ኑ _የ የ የ	Ъ Ф	
÷	÷+++	수수수	6666	64 444	ት ት ትትት ት	+ + ++ + +	• + + + +	수수 수석	ንታት ታ	·	ф ф	÷ + + +	÷ ÷ ÷	÷ ÷



Supercomputer: Piz Daint



- 5320 hybrid nodes (Intel Haswell + NVIDIA Tesla P100)
- 1431 "traditional" nodes (Intel Broadwell)

MeteoSwiss

- Currently #3 on Top500 and #1 on Green500
- 90% of compute power from accelerator (GPUs)!!!

5th ENES workshop, May

2. Why is it hard?

- Requires a massive increase in compute power
- Moore's law is dead!
- Requires continually (!) adapting our models to emerging hardware architectures and programming paradigms

Where are we today? COSMO on GPUs

CRAY &



•

Since 1986 - Covering the Fastest in the World and the People Who I



April 1, 2016 Swiss Weather Foreca

Kesch'

John Russell

O



Large investement into software (MeteoSwiss, CSCS, ETH/C2SM)

Altair AMDZ micro

CoolIT

Adapted to run on GPUs

Chelsio

- Operational for weather forecast on a GPU-based system (Piz Kesch)
- Regional climate simulations on Piz Daint

After six months of tweaking – producing a 20 percent reduction in time-to-solution for weather forecasting – MeteoSwiss, the Federal Office of Meteorology and Climatology, today reported its next generation COSMO-1 forecasting system is now operational. COSMO-1 requires 20 times the computing power of COSMO-2 and runs on the hybrid CPU-GPU supercomputer, Piz Kesch, operated by the Swiss National Supercomputing Centre (CSCS) and custom built in collaboration with Cray and NVIDIA.

COSMO-1 was put into service last September (see, Today's

 Bibliotheca Alexandrir Solution to Build Mass
 ASRock Rack to Exhibit



ASRock Rack to Exhib

Visit additional Tabor

5th ENES workshop, May 2018

÷

Baroclinic wave (Jablonowski 2006)





Simulation Setup

- Single precision
- Regular lat/lon grid
- Periodic in i-direction
- 80°S to 80°N
- Covering 98.4% of Earth's surface
- Analytical initial condition
- 10 day simulation

5th ENES workshop May 2018

Minimal I/O

MeteoSwiss





Visualization by Tarun Chadha (C2SM): clouds > 10^{-3} g/kg (white) and precipitation > 4^{-1} 10^{+2} g/kg (blue)



፡፡ ት ት ት ት

Weak scaling

• Increase problem size and computational resources at the



Metric: time-to-solution

- SYPD = Simulated years per wallclock day
- E.g. AMIP (CMIP6)
 - 1979 2014 (36 years)
 - 4 6 months
 - 0.2-0.3 SYPD required



Strong scaling

• Near-global simulations at a fixed horizontal resolution



19

Results (and beyond) Ð

Δx	#nodes	∆t [s]	SYPD	MWh/SY
1.9 km	4,888	12	0.23	97.8
930 m	4,888	6	0.043	596

What would it take to do a 36 year AMIP simulation?

	∆x = 1.9 km	∆x = 930 m
Time to solution	156 days	840 days
Energy to solution	3.5 GWh (150 kCHF)	22 GWh (940 kCHF)
CO2eq* to solution	640 tons	3'800 tons
CHF to solution (2go.cscs.ch)	18 M node hours → 12 MCHF	97 M node hours \rightarrow 68 MCHF

+ 5th ENES workshop, May 2018

3. Where are we today?

- No global weather and climate model ready for production on state-of-the-art hardware
- Global km-scale AMIP-type simulations are feasible (0.23 SYPD at ∆x = 1.9 km)
- Cost of simulation is significant

• How do we get there?



 Many opportunities and challenges on the way, here are some examples...



Software (Domain-specific languages)

- Our models comprise 0.1 1 M lines of legacy code
- Effort to adapt is huge!
- Domain-specific languages (DSLs) can reduce code by 10x
- E.g. "Towards a performance portable, architecture agnostic implementation strategy for weather and climate models " (Fuhrer et al. 2010, Superfri)

```
function avg {
                                                          Example: Coriolis force
     offset off
     storage in
                                                              No loops
                                                              No data structures
     avg = 0.5 * (in(off) + in())
                                                              No halo-updates
    function coriolis force {
                                                              Different hardware backends
     storage fc, in
                                                              (x86 multi-core, NVIDIA Tesla,
     coriolis force = fc() * in()
                                                              Xeon Phi)
    operator coriolis {
     storage u tend, u, v tend, v, fc
     vertical region ( k start , k end ) {
         u tend += avg(j-1, coriolis force(fc, avg(...
         v tend -= avg(i-1, coriolis force(fc, avg(...
MeteoSwiss
                                                                     ENES workshop, May
                                                                                                     22
```

Optimized Algorithms (Artificial Neural Networks)

- Accelerators for ANNs
- Physical parametrizations can be 50% of computational cost
- Inherent uncertainties can be large
- E.g. "Accurate and Fast Neural Network Emulations of Model Radiation for the NCEP Coupled Climate Forecast System: Climate Simulations and Seasonal Predictions" (Krasnopolsky et al. 2010, MWR)



4. Open questions and challenges

- We need to rethink how we formulate and implement our models
- Many research opportunities on the way

C Summary

- 1. Why do we want km-scale global weather and climate simulations?
- 2. Why is it hard?
- 3. What can be achieved today using a refactored code on Europe's largest supercomputer?
- 4. Open questions and challenges







Schweizerische Eidgenossenschaft Confédération suisse Confederazione Svizzera Confederaziun svizra

Swiss Confederation

Federal Department of Home Affairs FDHA Federal Office of Meteorology and Climatology MeteoSwiss

MeteoSwiss Operation Center 1 CH-8058 Zurich-Airport T +41 58 460 91 11 www.meteoswiss.ch

MeteoSvizzera

Via ai Monti 146 CH-6605 Locarno-Monti T +41 58 460 92 22 www.meteosvizzera.ch

MétéoSuisse

7bis, av. de la Paix CH-1211 Genève 2 T +41 58 460 98 88 www.meteosuisse.ch

MétéoSuisse

Chemin de l'Aérologie CH-1530 Payerne T +41 58 460 94 44 www.meteosuisse.ch



Modeling the Earth system



(source https://www.ecmwf.int/sites/default/files/medialibrary/2017-09/atmospheric-physics-754px.jpg)

	÷	5	5					4		÷			25			
÷	5	4		4								÷				4
	÷	4	4			<u>~</u>		÷	÷	4			<u>_</u>			
÷	÷			÷				2	÷		÷			÷		3
		4	5	슈 너	2 4 4	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		수 않 수		4	5	\$ \$ \$	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~			
	4					÷	상 순 순		수 수					÷	~ ÷	÷
	÷	- G - G - G - G - G - G - G - G - G - G	1000	1001000	÷	수 승수	÷	÷	÷	÷			› 슈 슈 슈 ·			÷
53	4	4 4	Mete	oSwiss		유 슈 <i>슈</i>	4 X	ት ት		÷ +	수 <u>5</u> th Fl	VES works	:hon	÷	28 🕂	
÷	÷	수수 수		0011100	· + +	· · · · · · · · ·	4 4 4	4	÷ +	÷ +	·수 · · · ·			÷		
÷	수 수수	수수 수수수		· 수 · 수	· · · ·	6 4 4 4	4 4 4	· 4 (2	4	4	4 A	~ ~ ~	÷	÷ +	
÷	+ + +	÷++ (}+++++++++++++++++++++++++++++++++++++	64 444	÷÷ ÷÷÷	ት <u>ት</u> ት ትት	÷ + ++++	÷ ÷÷	++++++++++++++++++++++++++++++++++++++	÷ ÷	÷++	수 수	ት ት ት ት e	b cb	÷	÷

(source https://www.weather.gov/images/key/Cloud_Chart/Low/Large/L2d.jpg)

Output: Out



(Marvel 2017, doi:10.1038/scientificamerican1217-72)

Cloud response to climate change



likely, observational evidence (e.g. Norris et al. 2016, Nature)

- Small changes in cloud reflectance can have a large impact
- State-of-the-art climate models use $\Delta x = 25$ km and do not explicitly resolve these clouds



Source: ECMWF, 2016

But wait a minute...





• Off-chip memory bandwidth is not increasing at the same rate as FLOP performance





Consequence for atmospheric models

• "Dynamics" code (niter = 48, nwork = 4096000)

U



Machine	Cray XT4	Cray XT5	Cray XE6	Cray XK6
# cores	4	6	12	16
Single core	0.80 s	0.84 s	0.63 s	0.65 s
All cores	0.56 s	0.46 s	0.18 s	0.16 s
Speedup	1.4	1.8	3.4	4.0







HPC Challenges O

Rapid change

MeteoSwiss

- Timescale of HPC system is 3-4 years \rightarrow lagging behind
- New design constraints (not FLOPs!) •
 - Maximize parallelism

- \rightarrow not efficient \rightarrow wrong algorithms
- Minimize data movement and energy consumption
- Minimize synchronizations
- New and disruptive programming models •
 - Emerging "exotic" HW architectures
 - E.g. OpenMP 4.5, Coarray Fortran,
 - CUDA, OpenACC

- - \rightarrow cannot run

Software (Automatic optimization)

- Compilers will not solve the problem!
- DSL-based code can be automatically optimized for a specific hardware target
- E.g. "Design of a Compiler Framework for Domain Specific Languages for Geophysical Fluid Dynamics Models" (Fabian Thüring, MSc thesis)



Software (I/O, data compression)

- We expect data volumes 1-2 PB/year for climate simulations at $\Delta x = 1 \text{ km}$
- Traditional workflow (compute \rightarrow store \rightarrow analyze) will break!
- E.g. "Data compression for climate data" (Kuhn et al. 2016, Superfri) and "Convectionresolving climate modeling on future supercomputing platforms (crCLIM)" (SNF Sinergia project, Lead: Ch. Schär)



Up or down?



Increase level of abstraction

- Hide implementation details
- Can be disruptive

- Decrease level of abstraction
 - Add implementation details
 - Often incremental



DOWN – Decrease level of abstraction



Approaches

- Fortran + MPI + Directives (OpenMP, OpenACC)
- Optimize code for a specific hardware
- Custom implementations (#ifdef) or programming languages





- Example: OpenACC
 - More details
 - data movement
 - data structures
 - Less encapsulation
 - Hardware dependent code
 - #ifdef

• "Easy"

÷.÷.÷.

- Incremental
- Hard to understand / modify
- Hard to maintain,

5th ENES workshop, May 2018

40

DOWN – Challenges



- Is it possible to reach a good compromise?
 - Multiple hardware architectures
 - Good performance
 - Maintainable / modifiable code



UP – Increase level of abstraction



Approaches

- Compilers
- Libraries / Frameworks
- Code generators and source-tosource translators
- Domain-specific languages (DSL)



Source: Thuering, 2017

Domain-specific language (DSL)

 DSL compiler compiles into standard programming language



UP – Challenges



Feasibility

- No turn key solutions.
- Can we achieve good performance on different hardware architectures with a high-level specification of our algorithms?

Acceptance

- Disruptive change
- Can we achieve a community solution?



Example: COSMO



U	Traditional vs. next-generation?							
		CSCS						
		Piz Dora (old code)	Piz Kesch (new code)	Factor				
Socket	S	~26 CPUs	~7 GPUs	3.7 x				
Energy	/	10 kWh	2.1 kWh	4.8 x				

→ Investment into hardware and software!

	Hardware	~ 2-3 MCHF	
Г. су <i>-</i> й	Software (projects, in-kind)	~ 5-7 MCHF	
	ф. ф. с.		ф ф ф
	+ + + + * * + + + + * * * * * * * * * *		
~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~		фф [,] ,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	······································
\$\phi\$         \$\Phi\$<	OSWISS	+ + + + + + + + + + + + + + + + +	p, ⁴ May 2018 ⁺ + 46 + + + + + + + + + + + + + + + + +