

Schweizerische Eidgenossenschaft Confédération suisse Confederazione Svizzera Confederaziun svizra

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

Designs for Efficient Weather & Climate Models

Carlos Osuna¹, Oliver Fuhrer¹, Xavier Lapillonne¹, Mauro Bianco², Paolo Crosetto³, Thomas Schulthess²

¹Federal Institute of Meteorology and Climatology MeteoSwiss ²Swiss National Supercompuing Centre CSCS, Lugano ³Centre for Climate Systems Modeling C2SM, ETH Zurich

Current operational system @ Meteoswiss

ECMWF-Model

U

16 km gridspacing 2 x per day 10 day forecast

COSMO-7

 $\Delta x = 6.6$ km, $\Delta t = 60$ s 393 x 338 x 60 cells 3 x per day 72 h forecast

COSMO-2

 $\Delta x = 2.2$ km, $\Delta t = 20$ s 520 x 350 x 60 cells 7 x per day 33 h forecast 1 x per day 45 h forecast



Next-generation system



Ensemble data assimilation: LETKF

Benefit of high resolution

(18-days for July 9 - 27, 2006)

D



Computational cost $= 40 \times$

(relative to current operational system)

J



Production with COSMO @ CSCS

Cray XE6 (Albis/Lema)

MeteoSwiss operational system Since ~4 years

Next-generation system

Accounting for Moore's law (factor 4)





Co-design: Approach

- Design software, workflow and hardware with the following principles
 - Portability to other users (and hardware)
 - Achieve time-to-solution
 - Optimize energy (and space) requirements
- Collaborative effort between
 - MeteoSwiss, C2SM/ETH, CSCS for software since 2010
 - Cray and NVIDIA for new machine since 2013
 - Domain scientists and computer scientists
- Additional funding from the HPCN Strategy (HP2C, PASC)

New MeteoSwiss HPC system

Piz Kesch (Cray CS Storm)

- Installed at CSCS in July 2015
- Already Operational
- Hybrid system with a mixture of CPUs and GPUs
- "Fat" compute nodes with 2 Intel Xeon E5 2690 (Haswell) and 8 Tesla K80 (each with 2 GK210)
- Only 12 out of 22 possible compute nodes
- Fully redundant (failover for research and development)





•

	Piz Dora	Piz Kesch	Factor
Sockets at required time-to-solution	~16 CPUs	~7 GPUs	2.4 x
Energy per member	6.19 kWh	2.06 kWh	3.0 x

Results Relative to "Old" Code

("Old" = no C++ dycore, double precision)

	Piz Dora	Piz Kesch	Factor
Sockets at required time-	~26 CPUs	~7 GPUs	3.7 x
Energy per member	10.0 kWh	2.06 kWh	4.8 x

Co-design: Software Technologies for Portable Models

- OpenACC
- C++ DSL for PDEs
- GCL for halo exchange communications (MPI based)

Copy to accelerator Boundary conditions || OpenACC port **Physics OpenACC** port C++ / DSL rewrite **Dynamics** Mixed OpenACC / CPU Data assimilation Δt **Communication library (GCL)** Halo-update Diagnostics OpenACC port **Input / Output** Mixed OpenACC / CPU



C++ DSL For Portable and Performance Portable Models?

Separation of Concerns:

Abstract hardware dependent code, underlying programming model from weather model.

Main focus on dynamical cores: composition of complex stencils

Separation of Concerns



Naïve Implementation

- Readable (close to numerical formulation)
- Close to programming agnostic formulation

- Not Portable
- Not Parallel
- Not Optimized

Q

Example kernel for 4th horizontal diffusion in CUDA

```
const int i = threadIdx.x:
const int i = threadIdx.y;
int i h = 0;
int j^h = 0;
if(j < 2)
Ł
  i h = i
 j^{-}h = (j = 0? -1: blockDim.y);
else if(i < 4 \&\&i <= blockDim.y)
 i h = (j = 2? -1: blockDim.x);
 j h = i:
}
for(int k=0; k < kdim; ++k)
  lap(i,j) = -4.0 * phi(i,j,k)
    + phi(i+1,j,k) + phi(i-1,j,k)
```

+ phi(i,i+1,k) + phi(i,i-1,k);

```
if(i h != 0 || i h != 0)
  lap(i h, j h) =
    -\overline{4.0 * phi}(i h, j h, k)
    + phi(i h+1,j h,k) + phi(i h-1,j h,k)
    + phi(i h, j h+1, k) + phi(i h, j h-1, k);
    syncthreads();
flx(i,j,k) = lap(i+1,j,k) - lap(i,j,k);
fly(i,j,k) = lap(i,j+1,k) - lap(i,j,k);
if(i h < 0)
  f\bar{lx}(i h, j h, k) = lap(i h+1, j h, k) -
    lap(i h,j h,k);
if(j h < \overline{0})
  fly(i,j h,k) = lap(i,j h+1,k) -
    lap(i,j h,k);
   syncthreads();
result(i,j) = phi(i,j,k) - alpha(i,j,k)*(
  flx(i,j,k) - flx(i-1,j,k) +
  fly(i,j,k) - fly(i,j-1,k));
```

Example kernel for 4th horizontal diffusion in CUDA



GridTools

 Set of grid tools, including DSL for stencil codes, for solving PDEs on



- Provides separation of concerns: Separates model and algorithm from hardware specific implementation and optimization
- Supports multiple hardware and grid backends.

Encoding Stencil Information in Types

```
struct Laplace
```

```
typedef in_accessor<0, range<-1,1-1,1> > u;
typedef out_accessor<1> lap;
template<typename Evaluation>
```

static void Do(Evaluation const& eval, full_domain)
{

```
eval(lap()) = eval(-4*u() + u(i+1) + u(i-1) + u(j+1) + u(j-1));
```

};

}

ł

Encoding Stencil Information in Types

```
struct Laplace
   typedef in accessor<0, range<-1,1-1,1> > u;
   typedef out accessor<1> lap;
    template<typename Evaluation>
   static void Do(Evaluation const& eval, full domain)
        eval(lap()) = eval(-4*u() + u(i+1) + u(i-1) + u(i-1))
           u(j+1) + u(j-1));
```

Support for Multiple Grids



Co-design: Extending Collaborations



Co-design: Extending Collaborations



Co-design: Extending Collaborations



Ū **Co-design: Extending Collaborations** x-backend x-grid boundary python conditions high level halo language exchanges frontend **Stencil** DSL gridtools INTEL x86 native architecture storage ESCAPE backend storage grids XeonPhi Atlas CUDA Kokkos structured grid octahedral icosahedral C++Std, grid grid **Sandia NVIDIA** unstructured grid

ESCAPE Project



ESCAPE (Energy-efficient Scalable Algorithms for Weather Prediction at Exascale)

Work on weather and climate "dwarfs", explore programming models and adaptations to new computing architectures.

GridTools DSL for solvers on octahedral grids (structured and unstructured)

Support and evaluate multiple architectures (NVIDIA GPU, Xeon Phi)



Summary

- New forecasting system doubling resolution of deterministic forecast and introducing a convection permitting ensemble
- Co-design

(simultaneous code, hardware and workflow re-design in close collaboration with hardware & software vendors, model scientists, ...)

allowed MeteoSwiss to increase computational load by 40x within 4–5 years

- Energy to solution is a factor 3x smaller as compared to a "traditional" CPU-based system
- GridTools:
- Next generation DSL for PDEs: generalize, modular, grid & backend agnostic.



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 http://superfri.org/

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



Energy Measurement

- We use power clamp for comparison
- Measurements from PMDB and RUR were within 1% of clamp

Piz Dora (Cray XC40)

Power clamp

- (external measurement which measures wall consumption including AC/DC conversion, interconnect, but excluding blower)
- 1-2 nodes were down and could not be used (considered in computation)
- **PMDB** (1 Hz, per node)
- **RUR** (total per job)

Piz Kesch (Cray CS Storm) Power clamp

- (external measurement which measures wall consumption including AC/DC conversion, interconnect, but excluding blower)
- Other components (mgmt nodes, extra service nodes, drives) powered down

"Managment summary"

Key ingredients

- Processor performance (Moore's law) ~2.8 x
- Port to accelerators (GPUs)
- Code improvement
- Increase utilization of system
- Increase in number of sockets
- Target system architecture to application

Note Separating hardware investments from software and workflow investments does not make sense!



~2.3 x

~1.7 x





Weak



Gridpoints per node: 128x128x60 or 64x64x60

512x512x60 or 256x256x60 gridpoints total

Performance comparison



- 24h COSMO-1 forecast on 70 nodes of Piz Daint
- Refactoring effort in dynamical core (1.5 x)

Slim vs. fat nodes

Piz Daint (Cay)

- 1 x Sandybridge
- 1 x NVIDIA Tesla K20x GPU
- 220 GB/s bandwidth

OPCODE (Tyan)

- 2 x Sandybridge
- 8 x NVIDIA Tesla K20 GPU
- 185 GB/s bandwidth





1986 s with old version

1402 s

G2G communication



G2G bandwidth

• Bi-directional BW between two GPUs

	Same	CUDA Memcopy	MPI	COSMO	
	PLX	13.1 GB/s	8.2 GB/s	7.1 GB/s	
_	IOH	11.0 GB/s	7.4 GB/s	6.9 GB/s	
	Node	10.3 GB/s	4.2 GB/S	4.1 GB/s	2–8 GPUs)

2	3	4	5	6	7	8
7.0 GB/s	6.3 GB/s	6.5 GB/s	4.0 GB/s	4.5 GB/s	4.6 GB/s	4.6 GB/s



- Pre- / post-processing software!
- Increase of machine utilization [] Error recovery and monitoring



OpenACC vs. STELLA

• Comparison using hor. diffusion and vert. advection

	runtime	occupancy	DRAM throughput read write		shared memory	register usage
non-block	(naive)					
K20X	0.53 ms	0.266	>75.1 GB/s	>35.5 GB/s	0 B	47-53
K20	0.68 ms	0.285	>39.1 GB/s	>26.3 GB/s	0 B	37-44
blocked						
K20X	0.90 ms	0.283	13.9 GB/s	62.9 GB/s	0 B	73
K20	0.69 ms	0.591	12.7 GB/s	63.1 GB/s	4 B	46
shared						
K20	0.54 ms	0.600	15.9 GB/s	16.1 GB/s	4.272 KB	39
shared-3E)					
K20	0.56 ms	0.670	15.4 GB/s	16.1 GB/s	4.272 KB	34
STELLA			Conclusio			
K20X	0.29 ms	0.90		115 Limplomontot	ion in abou	+151
K20	0.35 ms	0.90	 OpenACC code is portable, but not fully performance portable, many manual optimization 			
			-	-	-	-











O Benchmark

COSMO-E

2.2 km gridspacing 582 x 390 x 60 gridpoints 120 h forecast



Details

- operational setup by MeteoSwiss
- Required time-to-solution = 2h
- (333 ms per timestep)
- Fill a full rack with members
- (keeping sockets per member constant)
- COSMO v5.0
- (with additions for GPU porting and C++ dynamical core)
- Single precision
- (both CPU and GPU not fully optimized)

Co-design: Software Design

