# Performance profiling of COSMO on three generations of Cray supercomputers

**Jean-Guillaume Piccinali, CSCS**

**Anne Roches, C2SM**

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Goal

Analyse performances:

- of main configurations of COSMO (running currently on CSCS machines)

- with different domain decompositions

- on several platforms

In order to:

- understand the code and its bottlenecks

- identify potential quick wins

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

2

# Method (I)

Systematic procedure:

- Performance overview

- Computation scaling

- Communication scaling

- Memory usage

- I/O

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Method (II)

| TESTS 1 (no tools) | # of repet | CONFIG | DECOMP | # of runs |
|---|---|---|---|---|
| DOLE | 2 | 3 | 1 | 6 |
| ROSA | 2 | 3 | 3 | 18 |
| PALU | 2 | 3 | 3 | 18 |

| TESTS 2 (craypat) | Sampling + Tracing | CONFIG | DECOMP | # of runs |
|---|---|---|---|---|
| DOLE | 2 | 3 | 1 | 6 |
| ROSA | 2 | 3 | 3 | 18 |
| PALU | 2 | 3 | 3 | 18 |

| TESTS 3 ( Tau+ Scalasca) | Sampling + Tracing | CONFIG | DECOMP | # of runs |
|---|---|---|---|---|
| PALU | 2*2 | 3 | 1 | 12 |

| TESTS 4 | Nodes Occupation (6,12,24) | CONFIGS | DECOMP | # of runs |
|---|---|---|---|---|
| PALU | 3 | 3 | 1 | 9 |

# COSMO configurations

- ## Based on RAPS5.0

  - ➜ All runs : 24 h simulation

- ## 2 operational NWP versions:

  - ➜ COSMO2
  - ➜ COSMO7

- ## 1 climate (RCM) version:

  - ➜ "IPCC"

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Configurations: grid, time

## Space and time characteristics:

| | COSMO2 | COSMO7 | "IPCC" |
|---|---|---|---|
| nx | 520 | 393 | 441 |
| ny | 350 | 338 | 429 |
| nz | 60 | 60 | 40 |
| $\Delta x$, $\Delta y$ (°) | 0.02 | 0.06 | 0.11 |
| dt (s) | 20 | 60 | 100 |

# Configurations: obs

## Use of observations:

| | COSMO2 | COSMO7 | "IPCC" |
|---|---|---|---|
| luseobs | T | T | F |
| Use of satellite data | T | T | F |
| Nudging | T | T | F |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Domain decompositions

- Usual decomposition for each configuration
  + 1 "small" decomposition
  + 1 "big"    decomposition

|  | COSMO2 | COSMO7 | "IPCC" |
|---|---|---|---|
| usual | 28 x 35 + 4 (984) | 28 x 35 + 4 (984) | 30 x 30 (900) |
| small | 20 x 24 + 4 (484) | 20 x 24 + 4 (484) | 24 x 24 (576) |
| big | 48 x 60 + 4 (2884) | 48 x 60 + 4 (2884) | 53 x 53 (2809) |

(Number of cores)

=> Keep aspect ratio ~ cst

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Decompositions

## Aspect ratios:

|        | COSMO2 | COSMO7 | "IPCC" |
|--------|--------|--------|--------|
| usual  | 0.8    | 0.8    | 1      |
| small  | 0.83   | 0.83   | 1      |
| big    | 0.8    | 0.8    | 1      |

## Grid points per core:

|        | COSMO2 | COSMO7 | "IPCC" |
|--------|--------|--------|--------|
| usual  | ~185   | ~135   | ~210   |
| small  | ~376   | ~274   | ~328   |
| big    | ~63    | ~46    | ~67    |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Platforms

- ## 3 generations of Cray machines

  - ### Dole: Cray XT4

  - ### Rosa: Cray XT5

  - ### Palu: Cray XE6

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# Platforms

| | Dole (XT4) | Rosa (XT5) | Palu (XE6) |
|---|---|---|---|
| Pes (nodes * cores/node) | 172 x 4 = 688 | 1844 x 12 = 22128 | 176 x 24 = 4224 |
| CPU Frequency | 2.3 GHz | 2.4 GHz | 2.1 GHz |
| CPU type (AMD Opteron) | Barcelona | Istanbul | Magny-Cours |
| Interconnect injection bandwidth | 2 GB/s | 2 GB/s | 5 GB/s |
| Memory per node | 8 GB | 16 GB | 32 GB |
| L1 cache | 4 × 128 KB | 6 × 128 KB | 12 × 128 KB |
| L2 cache | 4 × 512 KB | 6 × 512 KB | 12 × 512 KB |
| L3 cache | 2 MB | 6 MB | 2 × 6 MB |
| Total memory | 5.5 TB | 354 TB | 135 TB |
| Peak performance | 6 Tflops | 212 TFlops | 35 TFlops |

http://fr.wikipedia.org/wiki/AMD_Opteron

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Timings (1 day)



Chart 1 legend: PALU PGI c7, ROSA PGI c7, DOLE PGI c7
Values: 601, 513, 421, 377, 421, 1410
X-axis (cores): 484, 984, 2884
Y-axis: Time (sec)

Chart 2 legend: PALU PGI C2 notool, ROSA PGI C2 notool, DOLE PGI C2 notool
Values: 2367, 2229, 1998, 1598, 1314, 1938, 1295
X-axis (cores): 484, 984, 2884

Chart 3 legend: ROSA PGI IPCC notool, PALU PGI IPCC notool, DOLE PGI IPCC notool
Values: 331, 309, 261, 260, 217, 230, 163
X-axis (cores): 576, 900, 2809
Y-axis: wallclock time (sec)

# Computational intensity (xpat)



rosa : cosmo2  palu : cosmo2  rosa : ipcc  palu : ipcc

- Computational intensity (ops/ref) should be >= 1
- D1 cache miss ratio should be ~= 0%

Legend:
- Peak (%)
- ops/cycle
- ops/ref
- d1 cache miss (%)

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Profile by regions (%)

# Profile by regions (sec)



palu : cosmo2

palu : ipcc

ETH
Eidgenössische Technische Hochschule Zürich
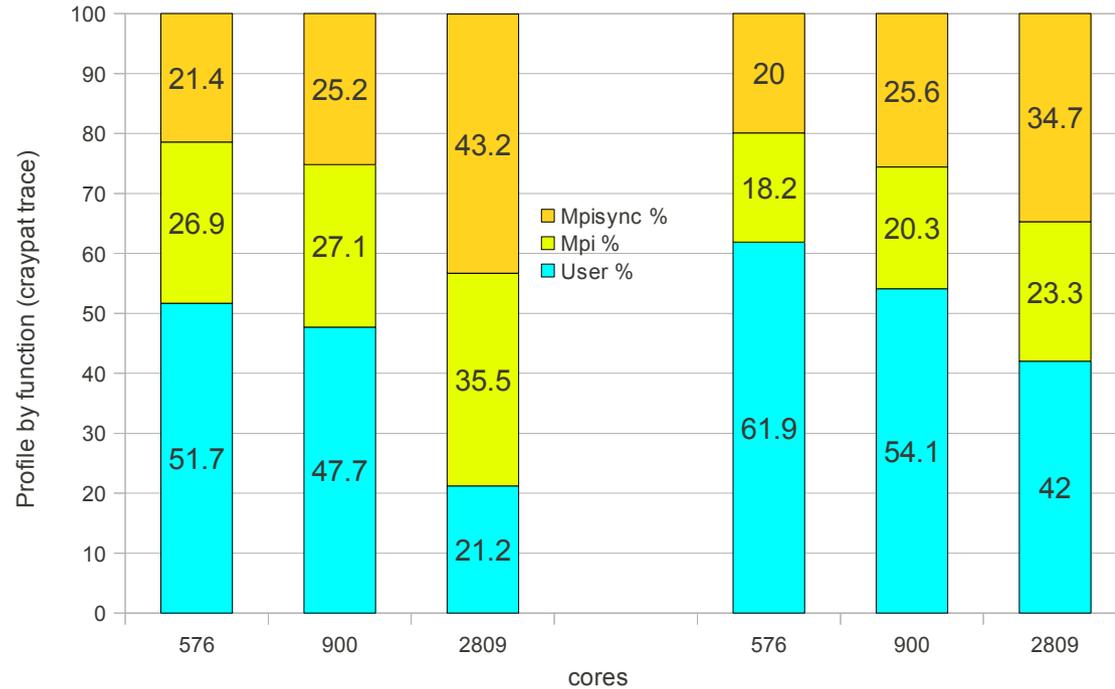Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Profile by MPI calls (ipcc, xpat)



cosmo2 : rosa, palu

ipcc : rosa, palu

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Swiss National Supercomputing Centre

# Profile by user functions (c2)

cosmo2 : rosa, palu

■ fast_waves_runge_kutta  ■ org_runge_kutta  ■ other user calls



cosmo2

rosa / 484c

■ FW_RK (21%)
♦ ORG_RK (12%)
▼ GATHER2D

cosmo2

palu / 484c

■ FW_RK (22%)
♦ ORG_RK (9%)
▼ GATHER2D

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Profile by user functions (ipcc)

ipcc : rosa, palu

## ipcc

### rosa / 576c

- ■ FW_RK (18%)
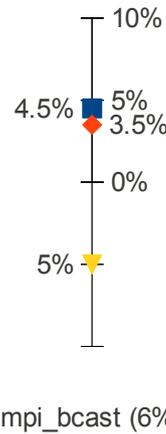- ◆ ORG_RK (8%)
- ▼ mult_fill_DDI (%)

mpi_sendrecv (20%)



4.5% ◆ — 5%

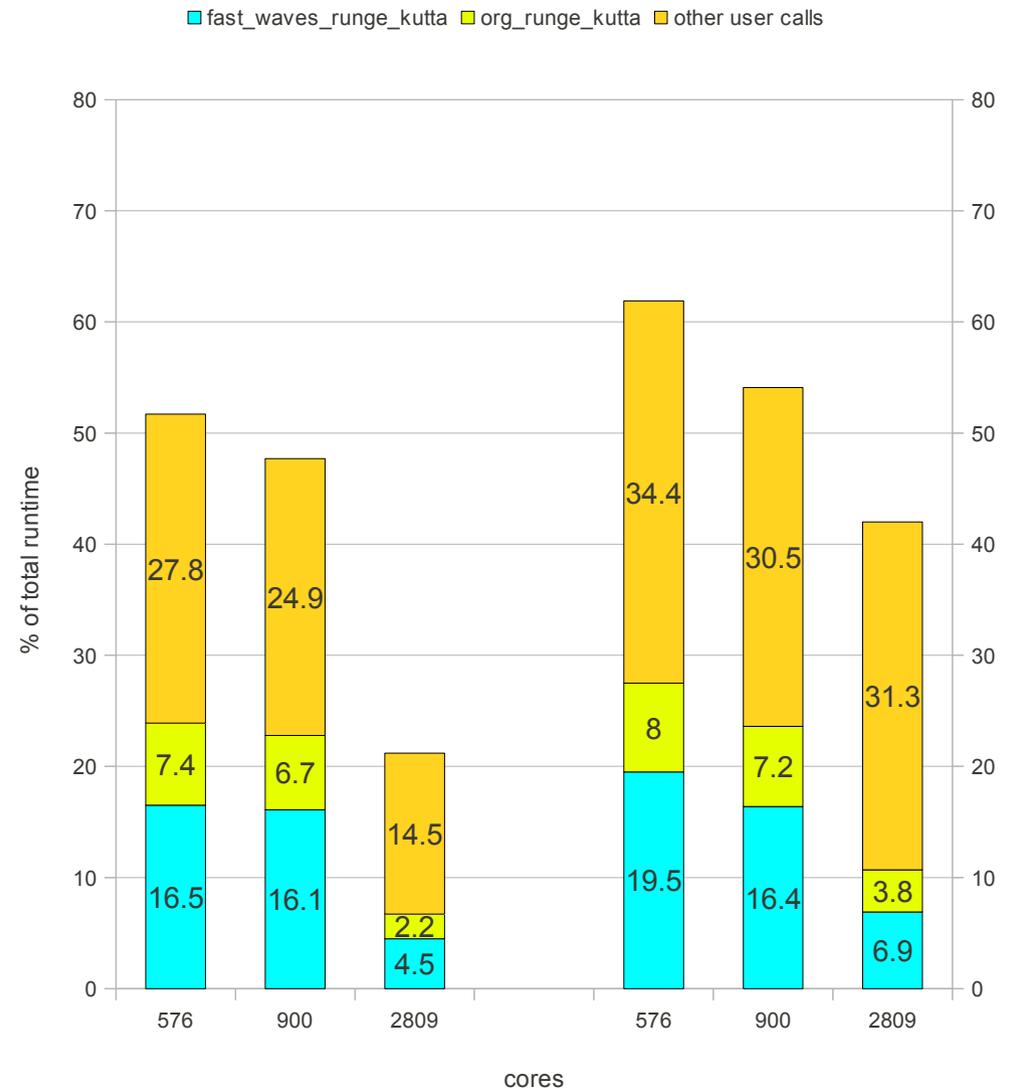1.1% ■ — 0%

3.5% ▼

mpi_allreduce (7.5%)

## ipcc

### palu / 576c

- ■ FW_RK (22%)
- ◆ ORG_RK (9%)
- ▼ org_output (%)

mpi_sendrecv (13%)



4.5% ■ — 5%
◆ 3.5%

0%

5% ▼

mpi_bcast (6%)

**Legend:** ■ fast_waves_runge_kutta  ■ org_runge_kutta  ■ other user calls



% of total runtime

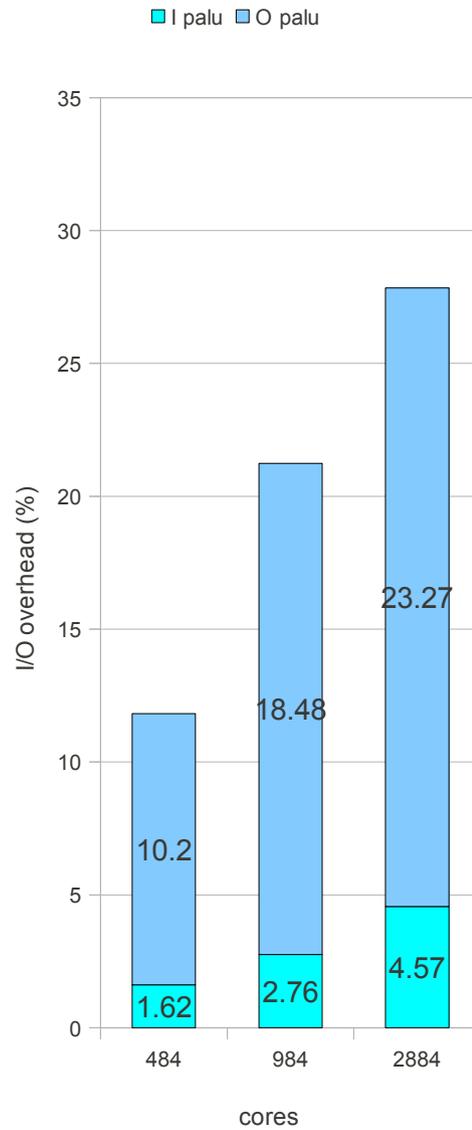| cores | 576 | 900 | 2809 | 576 | 900 | 2809 |
|-------|-----|-----|------|-----|-----|------|
| other user calls | 27.8 | 24.9 | 14.5 | 34.4 | 30.5 | 31.3 |
| org_runge_kutta | 7.4 | 6.7 | 2.2 | 8 | 7.2 | 3.8 |
| fast_waves_runge_kutta | 16.5 | 16.1 | 4.5 | 19.5 | 16.4 | 6.9 |

Swiss Federal Institute of Technology Zurich

# I/O overhead (yutiming)



cosmo2

ipcc

# MPI message stats (cosmo2, xpat)



rosa=palu : cosmo2

~960 000 msgs

palu : cosmo2 (no O)

85 000 msgs

ETH Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Configurations: I/O

I/O:

| | | COSMO2 | COSMO7 | "IPCC" |
|---|---|---|---|---|
| I/O | asynchronous | T | T | F |
| I/O | format | Grib1 | Grib1 | NetCDF |
| I/O | nincwait,nmaxwait | >0 | >0 | 0 |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Configurations: I/O

I:

| | | COSMO2 | COSMO7 | "IPCC" |
|---|---|---|---|---|
| I | hincbound | 1 | 3 | 6 |
| I | hnewbcdt | 3 | 3 | 0 |
| I | lana_qr_qs | T | F | F |
| I | llb_qr_qs | T | F | F |
| I | lana_rho_snow | T | F | F |
| I | lan_rho_snow | T | T | F |

# 4 more 2D fields read for COSMO2 as for IPCC

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Configurations: I/O

**O:**

| | | COSMO2 | COSMO7 | "IPCC" |
|---|---|---|---|---|
| O | # output blocks | 2 | 1 | 7 |
| O | on model levels | 109 x 2D and 13 x 3D @ 1 h freq<br>-----------------<br>19 x 2D and 7 x 3D @ 30min freq on a reduced domain | 99 x 2D and 12 x 3D @ 1 h freq | 5 x 2D @ 1 h freq<br>--------------------<br>30 x 2D and 7 x 3D @ 3 h freq<br>--------------------<br>17 x 2D and 6 x 3D @ 6 h freq<br>--------------------<br>1 x 3D @ 24 h freq |
| O | on pressure levels | 7 x 3D @ 1 h freq | 8 x 3D @ 1 h freq | - |
| O | on z-levels | 1 x 2D and 8 x 3D @ 1 h freq | 7 x 3D @ 1 h freq | - |
| O | On sat. levels | 1 x 3D @ 1 h freq | - | - |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Swiss National Supercomputing Centre

# Conclusion and outlook (I)

- Systematic profiling runs performed

- Lot of profiling data

- First overview of the results

- Need more analysis in order to reach the goals

  - Understand the code

  - Identify quick wins

# Conclusion and outlook (II)

- Other potential next steps:

  - Test external Lustre

  - Perform weak scaling

  - Try different compilers

  - Analyse additional configurations (ART, M7, …)

  - (Other platforms → done by the colleagues)

# THANK YOU

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre