



CSCS

Swiss National Supercomputing Centre

Improving the Performance of COSMO-CLM

Matthew Cordery, William Sawyer
Swiss National Supercomputing Centre

Ulrich Schättler
Deutscher Wetterdienst

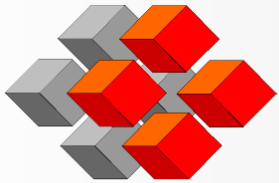


CSCS

What is COSMO-CLM?

Swiss National Supercomputing Centre

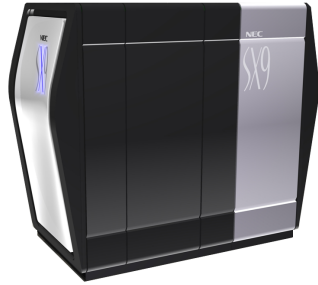
- COSMO is an operational non-hydrostatic meso-to-micro scale NWP system.
 - Used by MeteoSwiss, DWD, etc.
- The COSMO model in CLimate Mode (COSMO-CLM or CCLM) is a non-hydrostatic regional climate model.



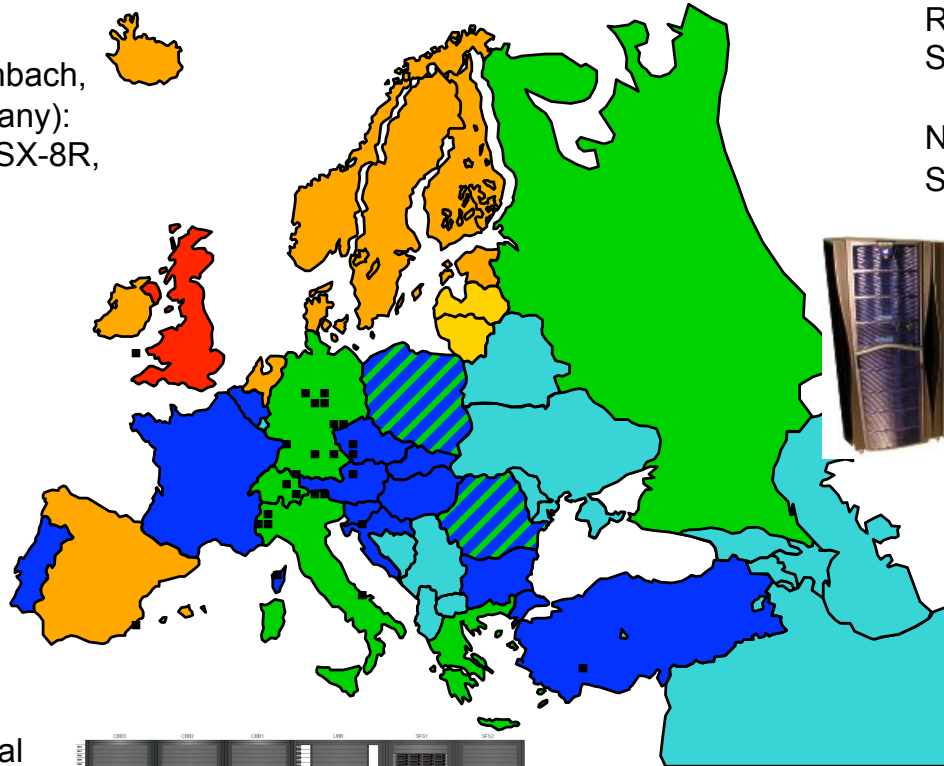
CSCS

COSMO NWP Applications

Swiss National Supercomputing Centre



DWD
(Offenbach,
Germany):
NEC SX-8R,
SX-9



Roshydromet (Moscow, Russia),
SGI

NMA (Bucharest, Romania):
Still in planning / procurement phase



IMGW (Warsawa, Poland):
SGI Origin 3800:
uses 88 of 100 nodes



ARPA-SIM (Bologna, Italy):
IBM pwr5: up to 160 of 512
nodes at CINECA

COSMO-LEPS (at ECMWF):
running on ECMWF pwr6 as
member-state time-critical
application

HNMS (Athens, Greece):
IBM pwr4: 120 of 256 nodes

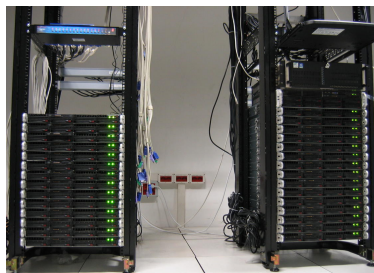
USAM (Rome, Italy):
HP Linux Cluster
XEON biproc quadcore
System in preparation



ARPA-SIM (Bologna, Italy):
Linux-Intel x86-64 Cluster for
testing (uses 56 of 120 cores)



MeteoSwiss:
Cray XT4: COSMO-7 and
COSMO-2 use 800+4 MPI-
Tasks on 402 out of 448 dual
core AMD nodes





CSCS

CLM Community 2010

Swiss National Supercomputing Centre



Afrika:



Tanzania Meteorol. Agency



National Meteorol. Agency of Senegal

America



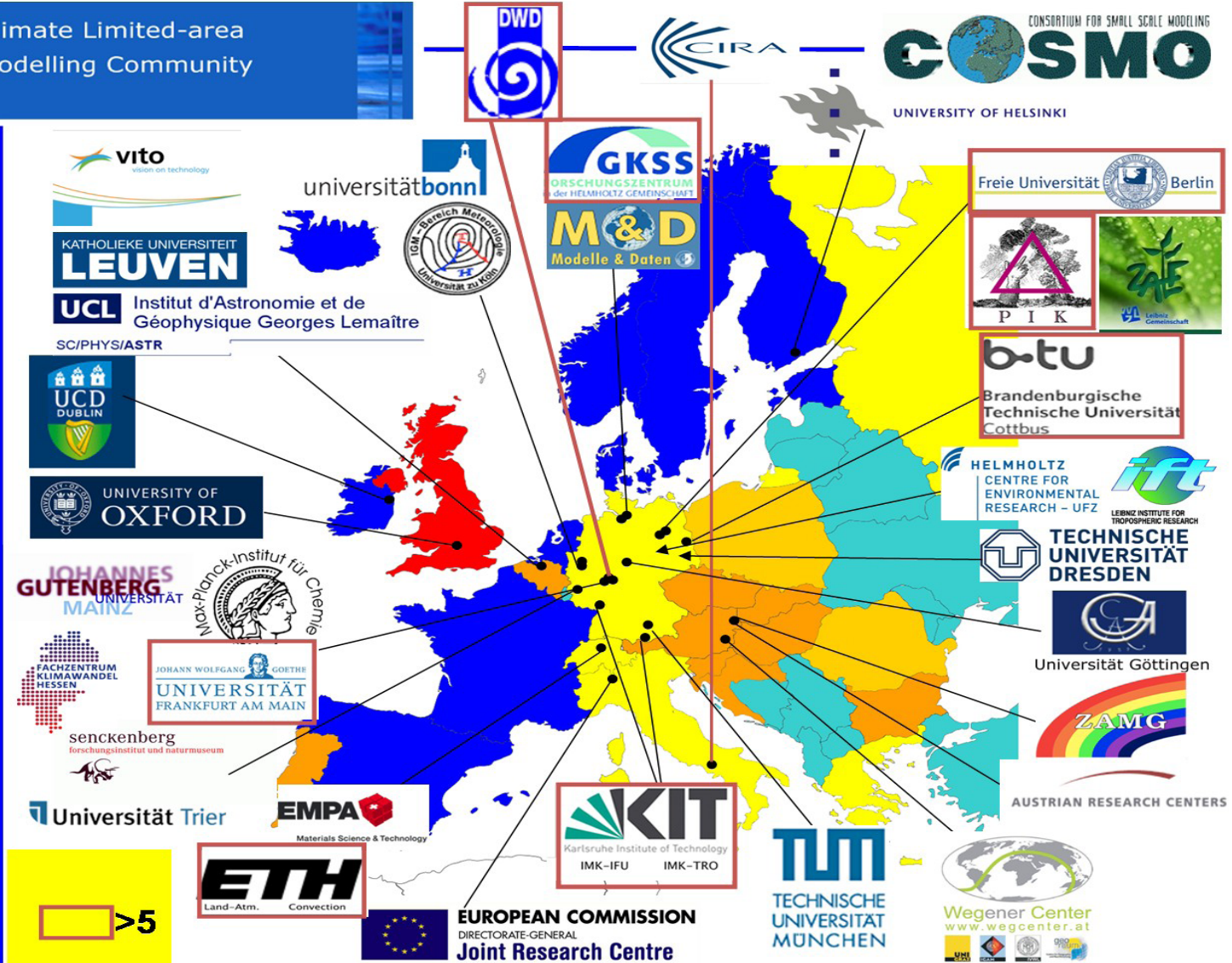
Universidad Nacional Autonoma de Mexico



India



Climate Limited-area Modelling Community



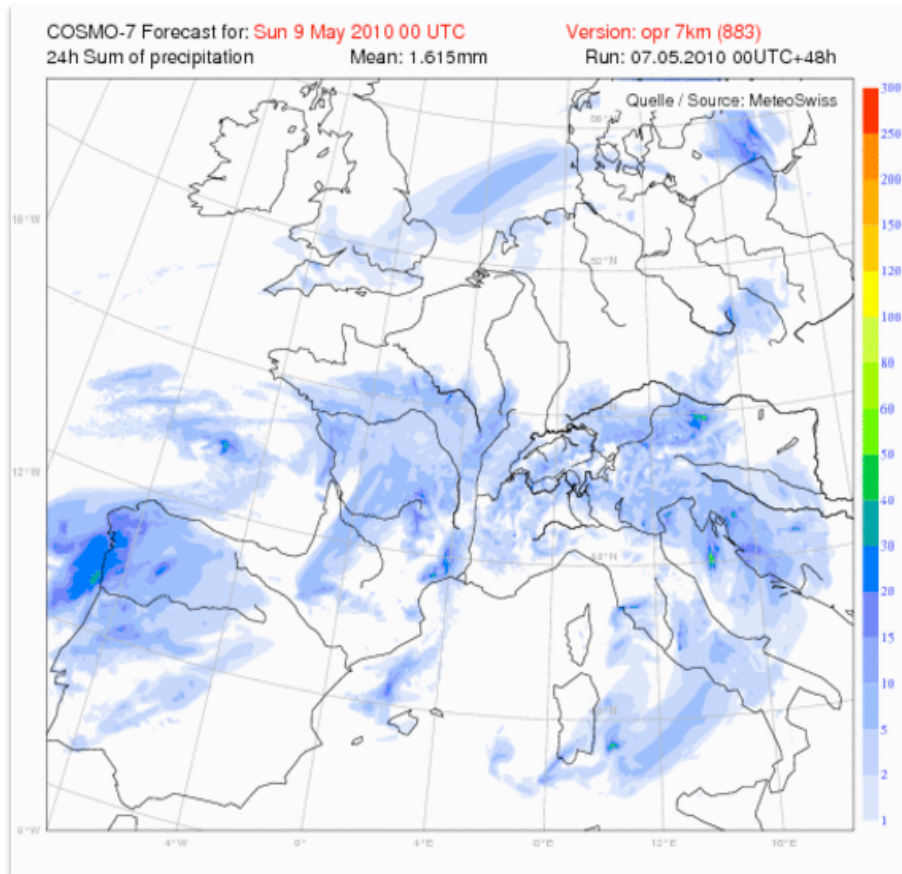
>5



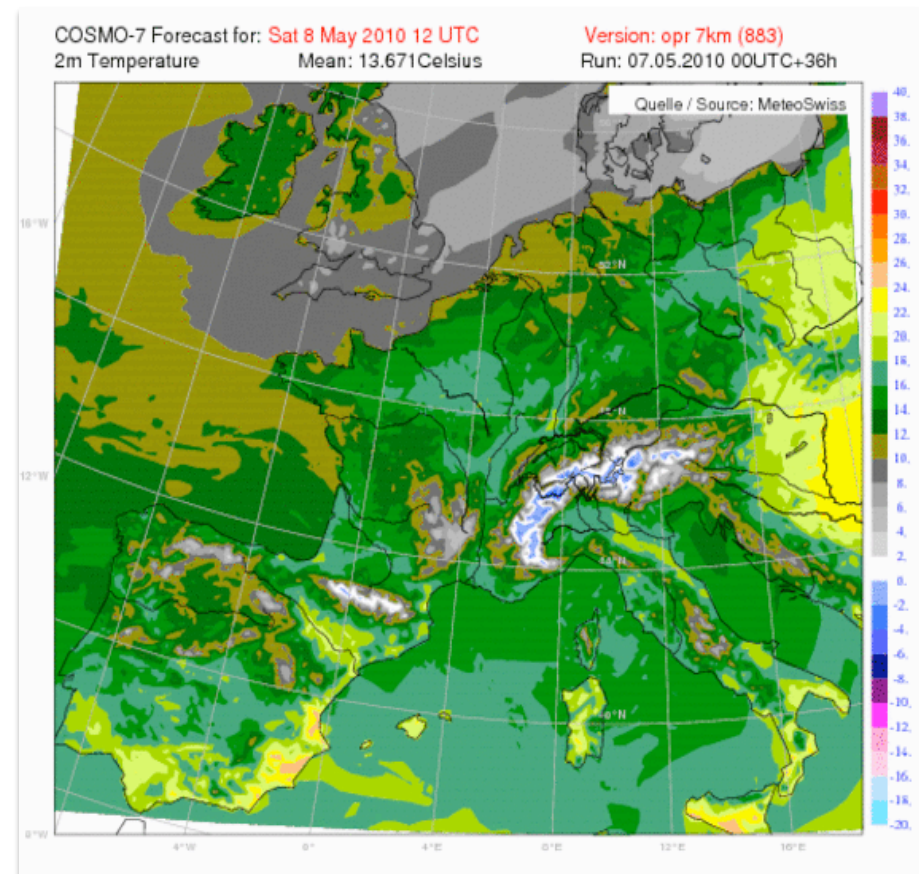
CSCS

Sample NWP output

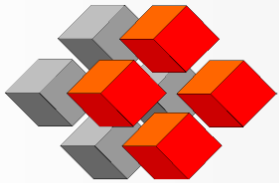
Swiss National Supercomputing Centre



Precipitation



2m Temperature



- An operation NWP prediction forecast for 24 hours make take half an hour of CPU time
- By extension, a 100 year climate simulation on a similar grid would take several months of CPU time.
- Our goal is to examine different hybrid programming schemes that will provide a substantial CPU time savings without requiring extensive use of machine resources.
- First step:
 - Examine mixed programming model using MPI and OpenMP.



CSCS

Swiss National Supercomputing Centre

- The computational grid is a 3D rotated-latitude/longitude structured grid.
- Communications are through 2- or 3-line halo exchanges.
- Many loops are of the form

```
do k = 1, ke
    do j = 1, je
        do i = 1, ie
            ...
        end do
    end do
end do
```

- Some 2D and 4D arrays, but of the same basic structure.



Main computational region is a time stepping loop over

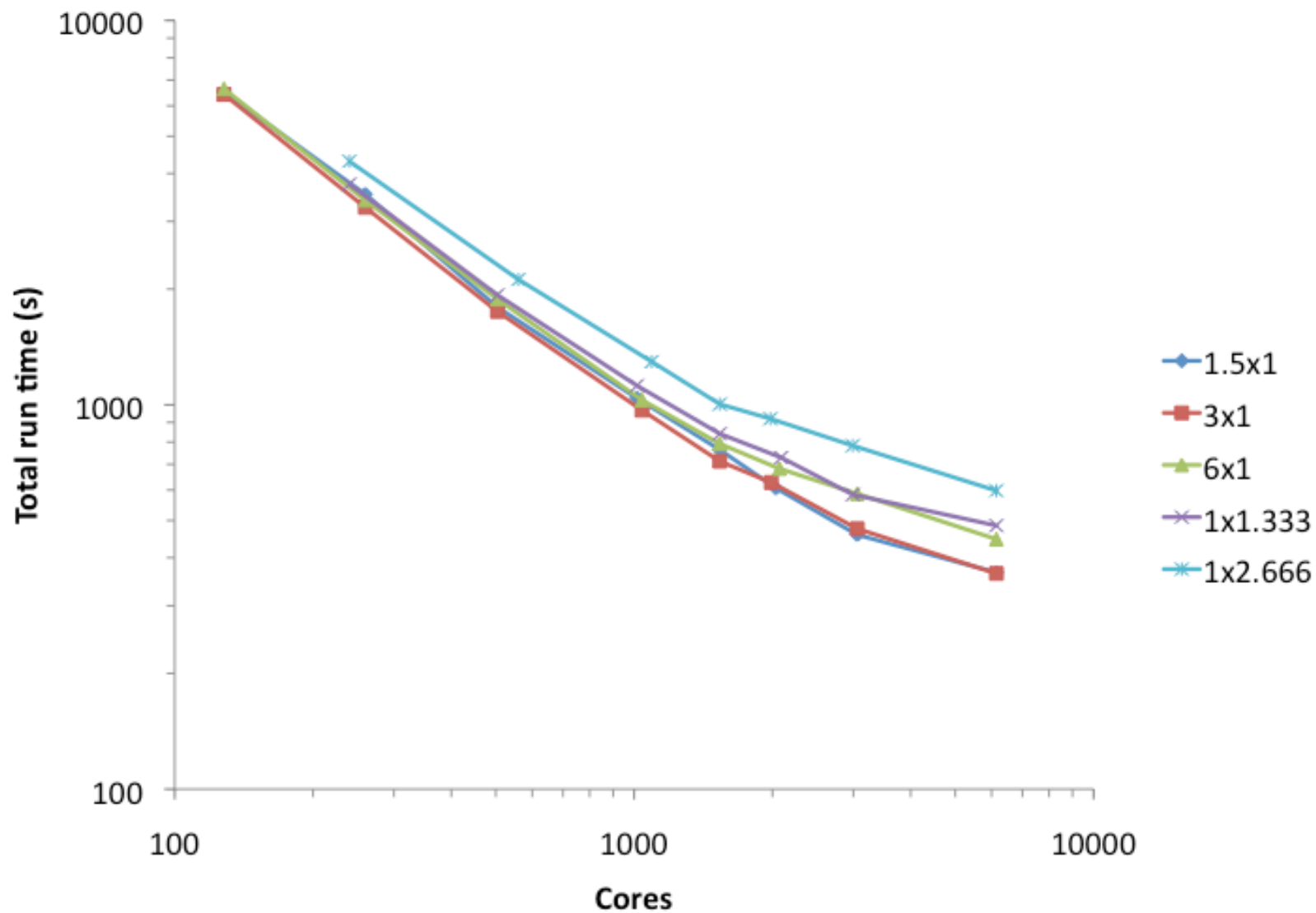
- a dynamical core that solves for fluid motion
 - a physical core that computes radiation transfer, precipitation, etc.
-
- Examine scaling of 1-km to elucidate 'hotspots'.
 - 1142 longitude points
 - 765 latitude points
 - 90 vertical points
 - 3 halo grid points
 - 4 I/O tasks
 - 1 hour of simulated time



CSCS

Swiss National Supercomputing Centre

Effect of Decomposition on Scaling





CSCS

1km observations

Swiss National Supercomputing Centre

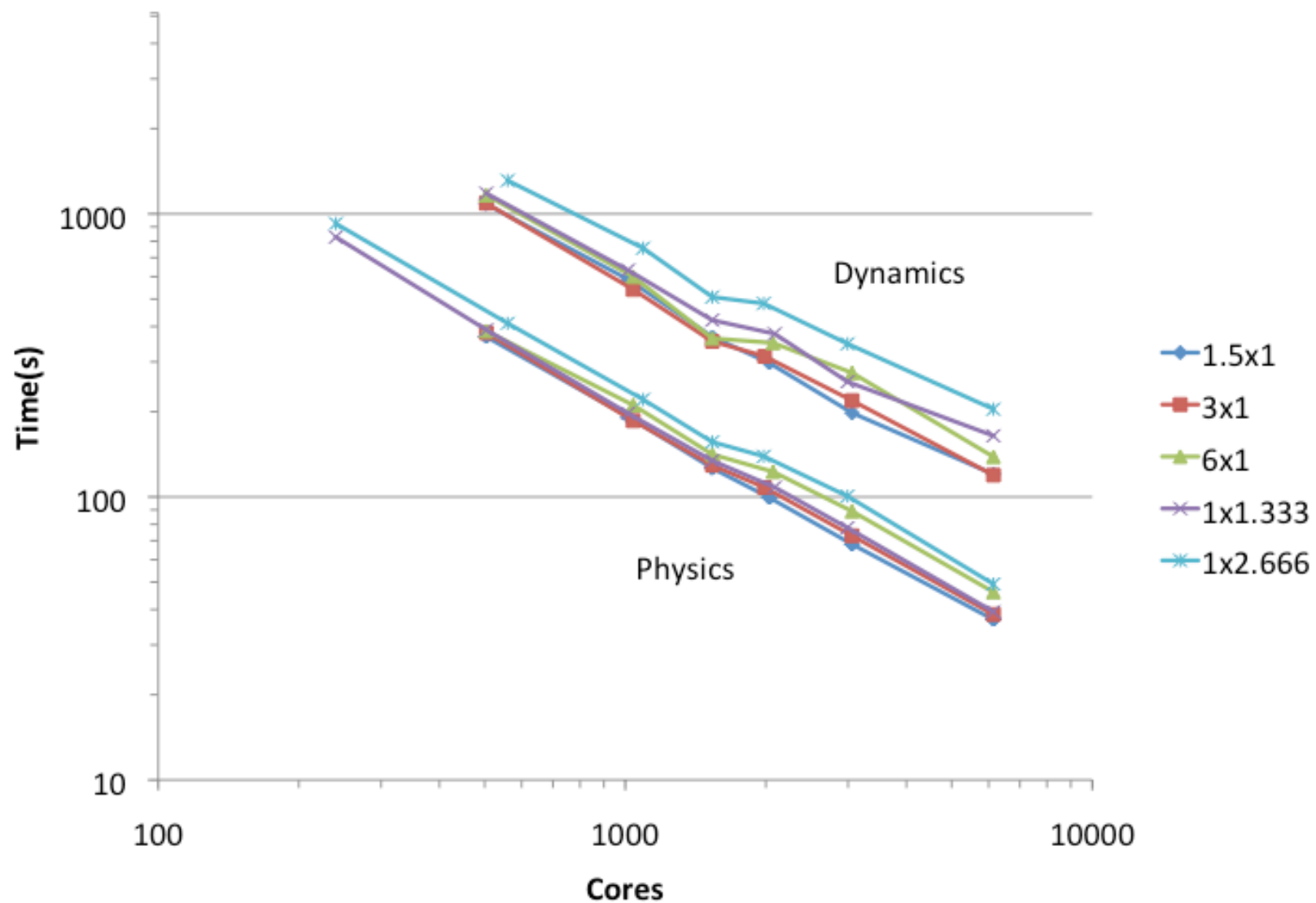
- Computational intensity is low (<1) for most top routines and the cache hit ratio is $> 95\%$.
- Good scaling out to several thousand cores, though performance is rather low.
- Implies algorithms are memory-bandwidth limited.

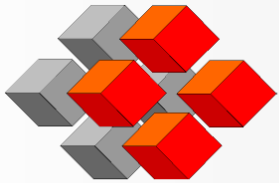


CSCS

Swiss National Supercomputing Centre

Scaling of dynamical and physical cores

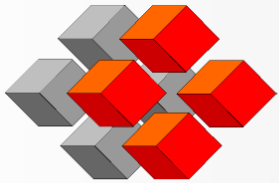




- Physical and dynamical cores scale well out to > 2000 cores
- Dynamical core takes approximately 3x more time than the physical core, regardless of core count.
- I/O and MPI communications are not limiting factors

Profiling indicates that during the main time stepping loop

- I/O and MPI communications are not limiting factors
- The dynamical core requires 3x more time than the physical core, regardless of core count



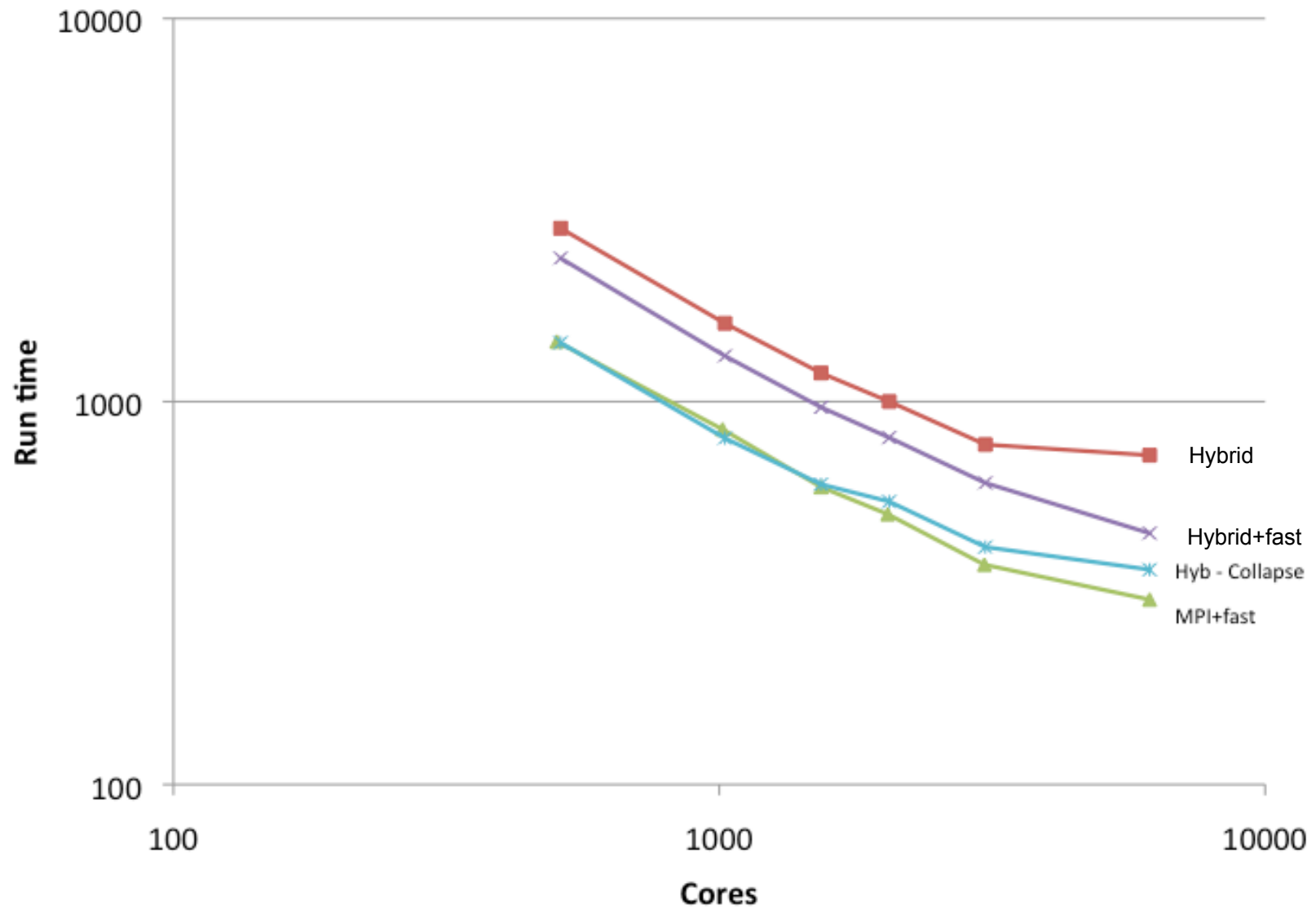
- Most computational work is encapsulated in multiply-nested loops in multiple subroutines that are called from a main driver loop.
- Most outer loops are over the number of levels, inner loops over latitude/longitude.
- Insert OpenMP PARALLEL DO directives on outermost loop
 - Also attempted use of OpenMP 3.0 COLLAPSE directive.
 - Over 600 directives inserted.
 - Also enabled use of SSE instructions on all routines (previously only used on some routines).



CSCS

Swiss National Supercomputing Centre

Hybrid scaling

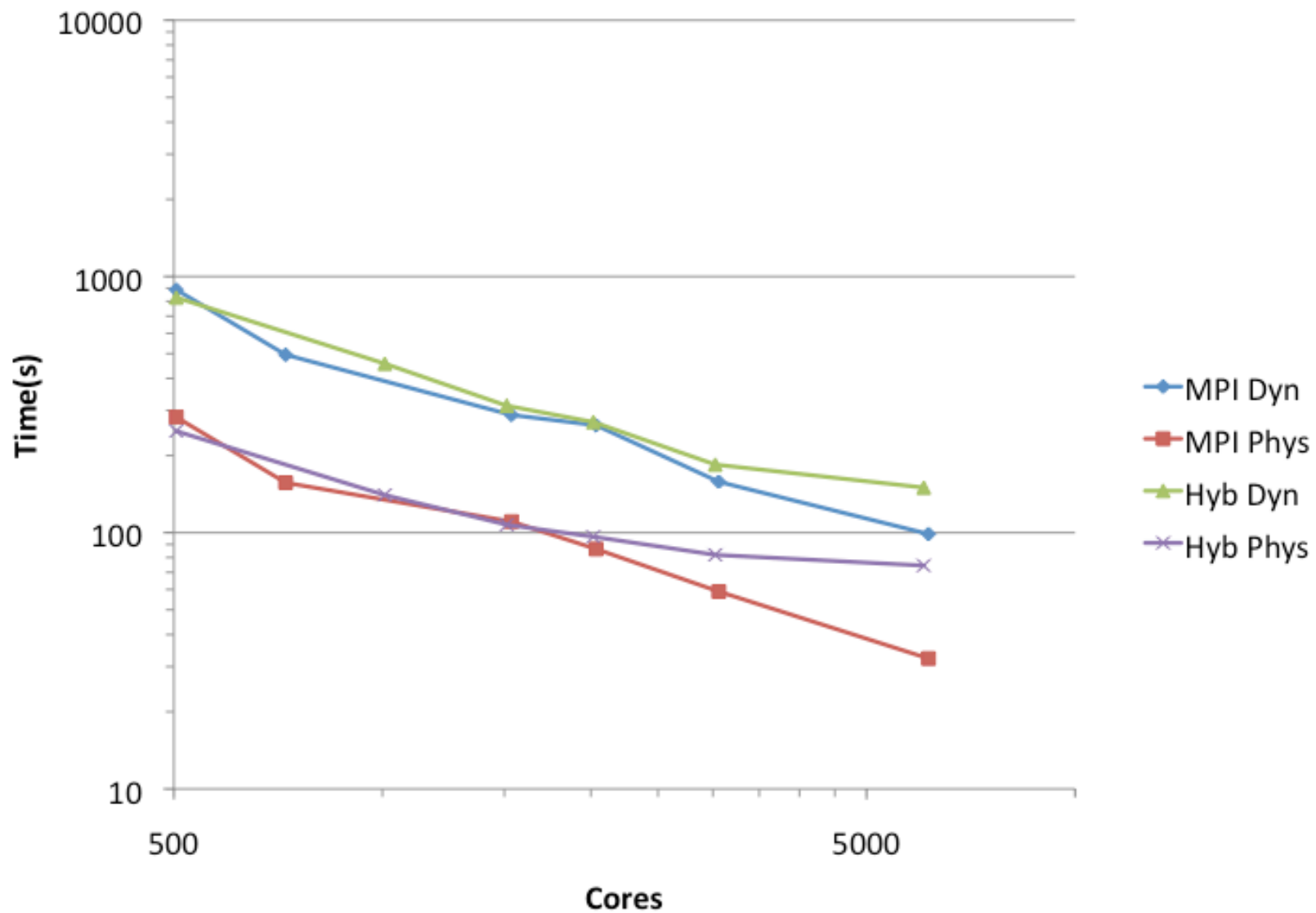




CSCS

Swiss National Supercomputing Centre

Scaling of Dynamics and Physics





- Important to have as much of the code as possible compiled to use SSE instructions.
- Important not to overuse COLLAPSE directive which may interfere with compiler optimization
- Code runs approximately as fast as the MPI only version for most core counts.
- Dynamical and physical cores scale well, though the physical core shows a much more pronounced loss of scaling beyond 4 threads.
- Reducing the number of I/O tasks improves performance (why?) and reduces idle cores.



CSCS

'Climate grid'

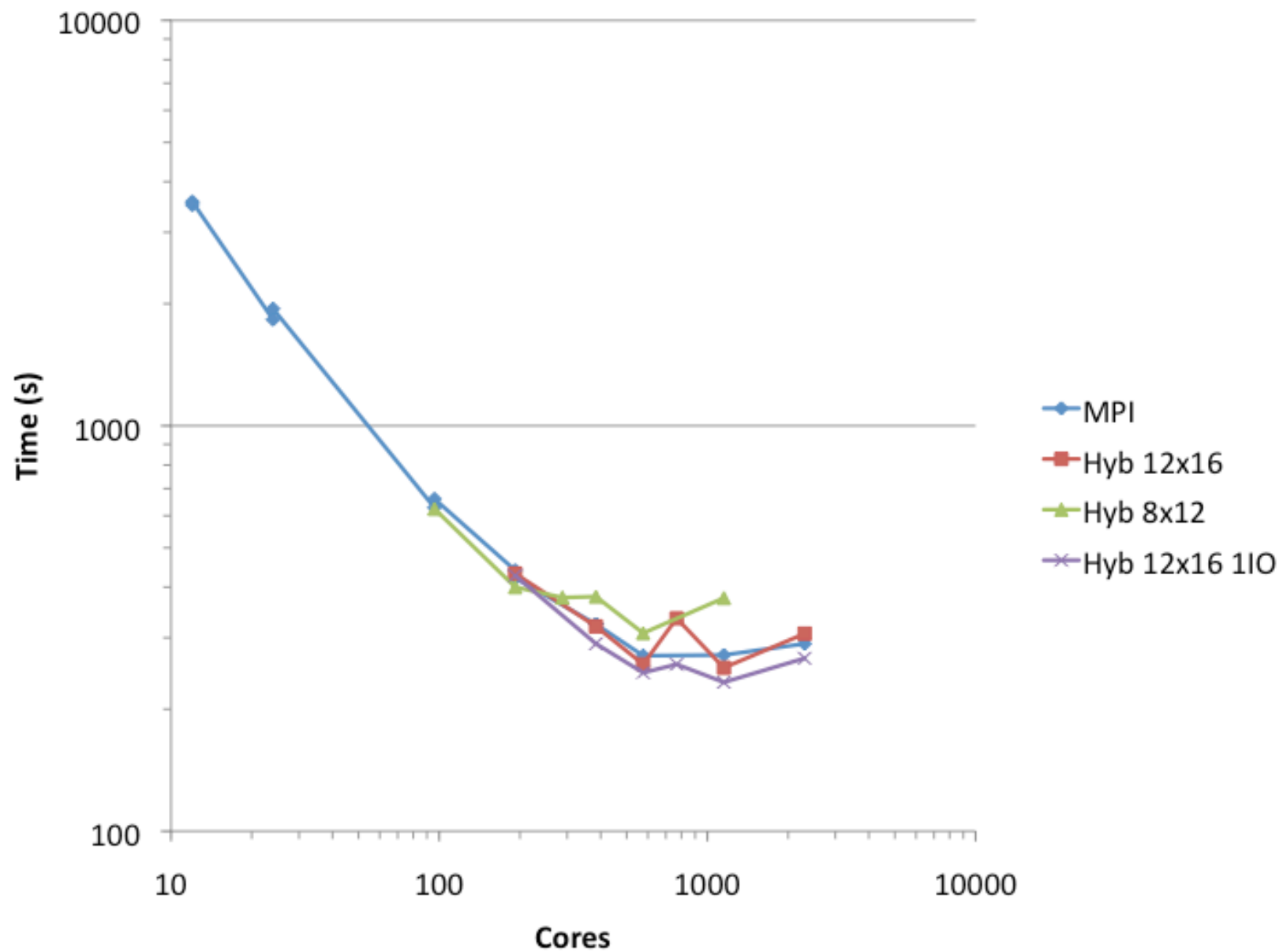
Swiss National Supercomputing Centre

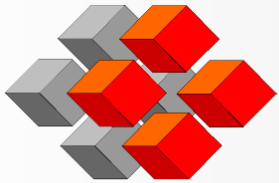
- In order to examine the results more similar to what will be used for a 100 year climate science run:
- 102 latitudes
- 102 longitudes
- 60 height levels
- 1 I/O task
- Run for 24 simulated hours.



CSCS

Swiss National Supercomputing Centre

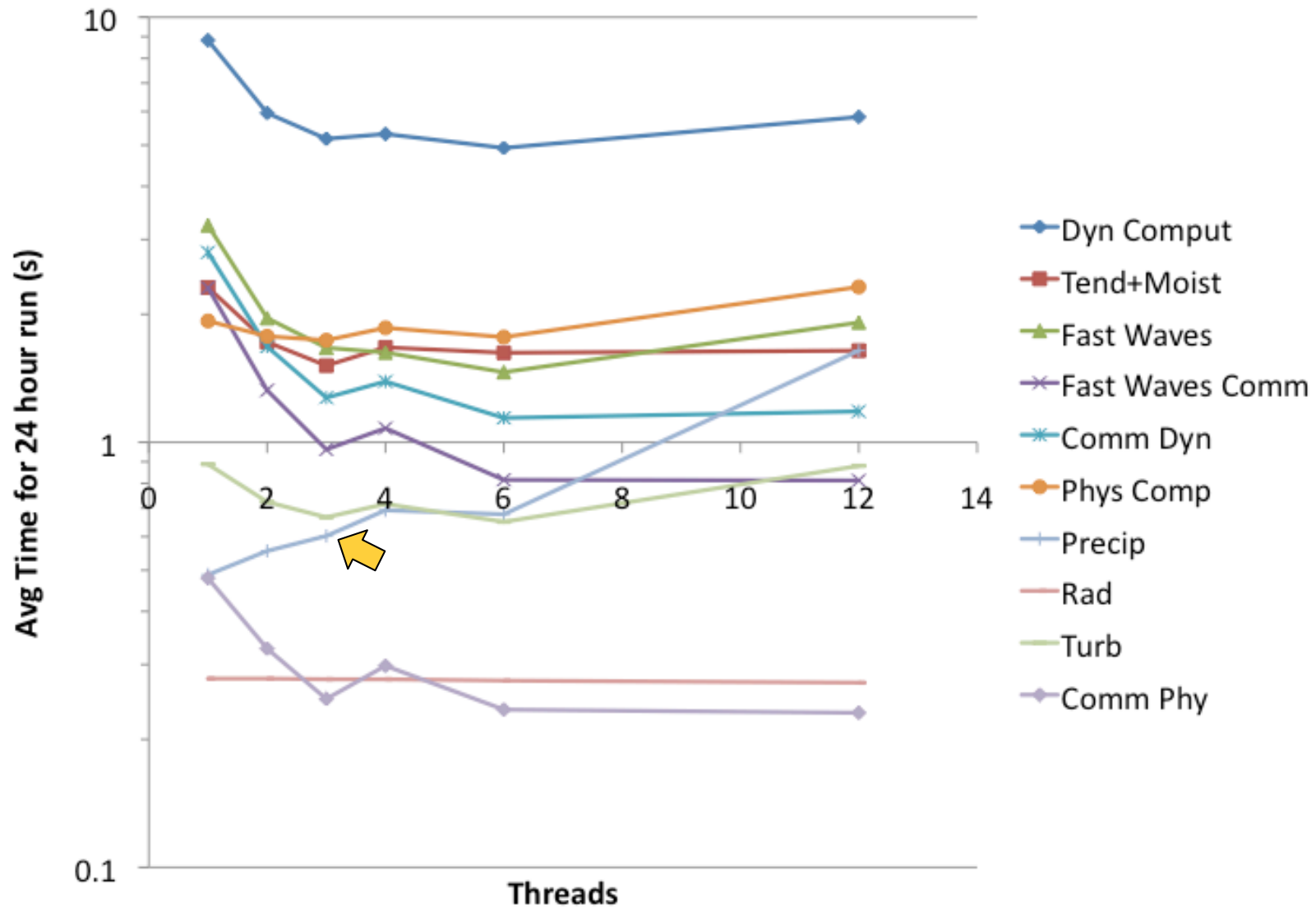


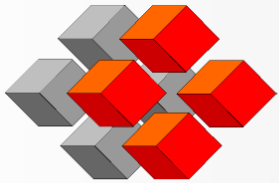


- Relatively weak (<10%) improvement over standard MPI code.
- Best results at 3 or 6 threads.
- Performance decrease for 4 threads where threads from one MPI task will span a node.
- Performance decrease going from six to twelve threads indicates still some performance bottlenecks.



Component Scaling, 12x16 tasks





- Loop level parallelism can achieve some modest performance gains
 - Can require many threaded loops -> OpenMP overhead
 - Can require a lot more software engineering to maintain
 - Introducing new private variables into old loops
 - Introducing new physics that needs to be threaded
 - Can be problematic when dealing with threads that include arrays created using Fortran allocate statement.
- Next task is to examine threading and blocking at a higher level. This will require more extensive work to the code.
 - Improve data re-use
 - Reduce memory footprint
- Investigate new algorithms