

# Some computational aspects of HARMONIE

Ulf Andræ, SMHI

Trygve Aspelien, Ole Vignes, Dag Bjørge met.no

Torgny Faxén, NSC

Jacob Weissman Poulsen DMI

ALADIN 23<sup>nd</sup> WS - HIRLAM ASM

Reykjavik

15-19<sup>th</sup> April 2013

## Content

- Benchmark activities
- Diagnosing the IO performance of cy38
- Portability
- Surface assimilation lesson learned

## Benchmarking

- Several HIRLAM countries are in the process of upgrading their HPC.
- A benchmarking package have been created
  - harmonie-38h1.alpha.2, no netcdf dep, wrgp2fa.F90 update, OpenMP fixes,
  - Simple sample scripts
  - Required background data
  - 1h boundaries up to 6h for several domains with 2-5km and 65 vertical levels

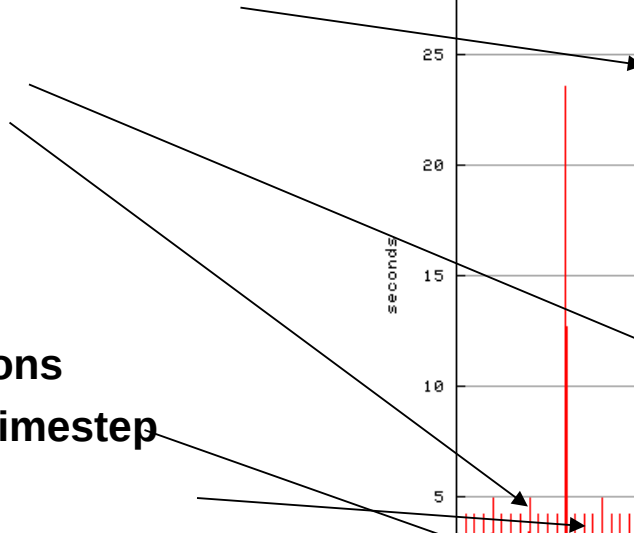
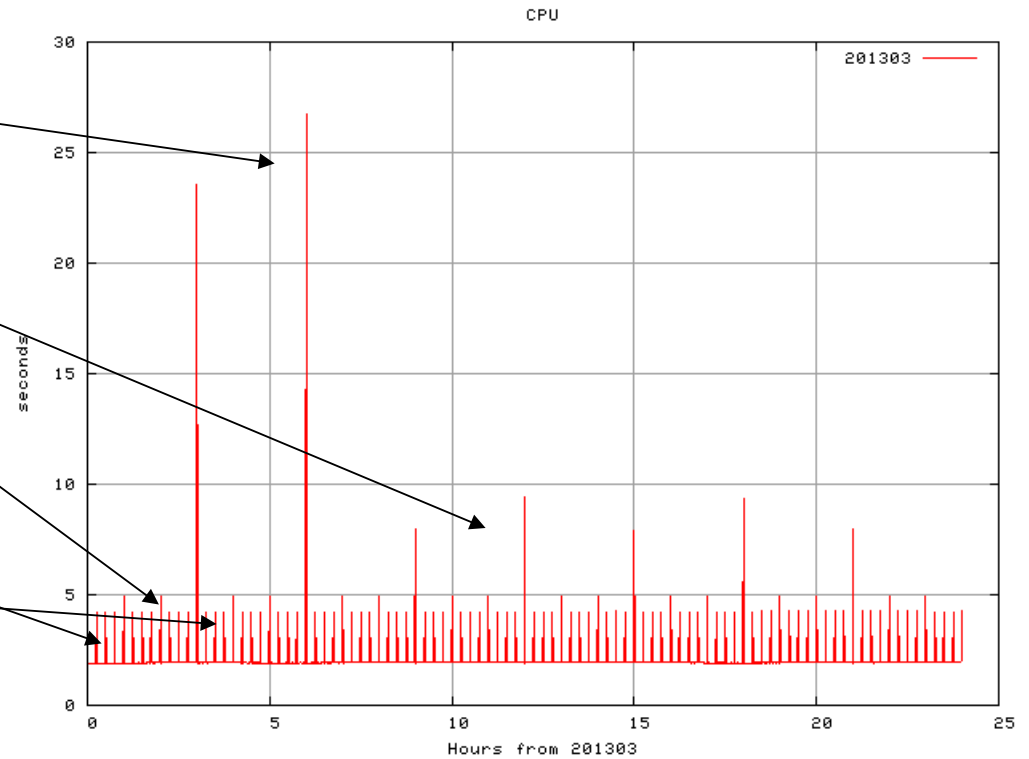
XS: 50x50, M:384x400, L:750x960,  
XL:1200x1200, XXL:1600x1600
- <https://hirlam.org/trac/wiki/HarmonieSystemDocumentation/HarmonieBenchMark>

## Benchmarking, some properties

- **Runs with EDKF ( EDMFM had problems)**
- **Tested for IBM, gfortran, intel**
- **MPI reproducible, different decompositions**
- **OpenMP reproducible for different number of threads**
- **Reproducibility issues with MKL libraries ( intel ) but good performance and reproducibility with other blas/lapack libraries.**
- **Forecast model only. Assimilation still considered second of importance ( or too complicated to deal with )**

# The typical cost of a forecast

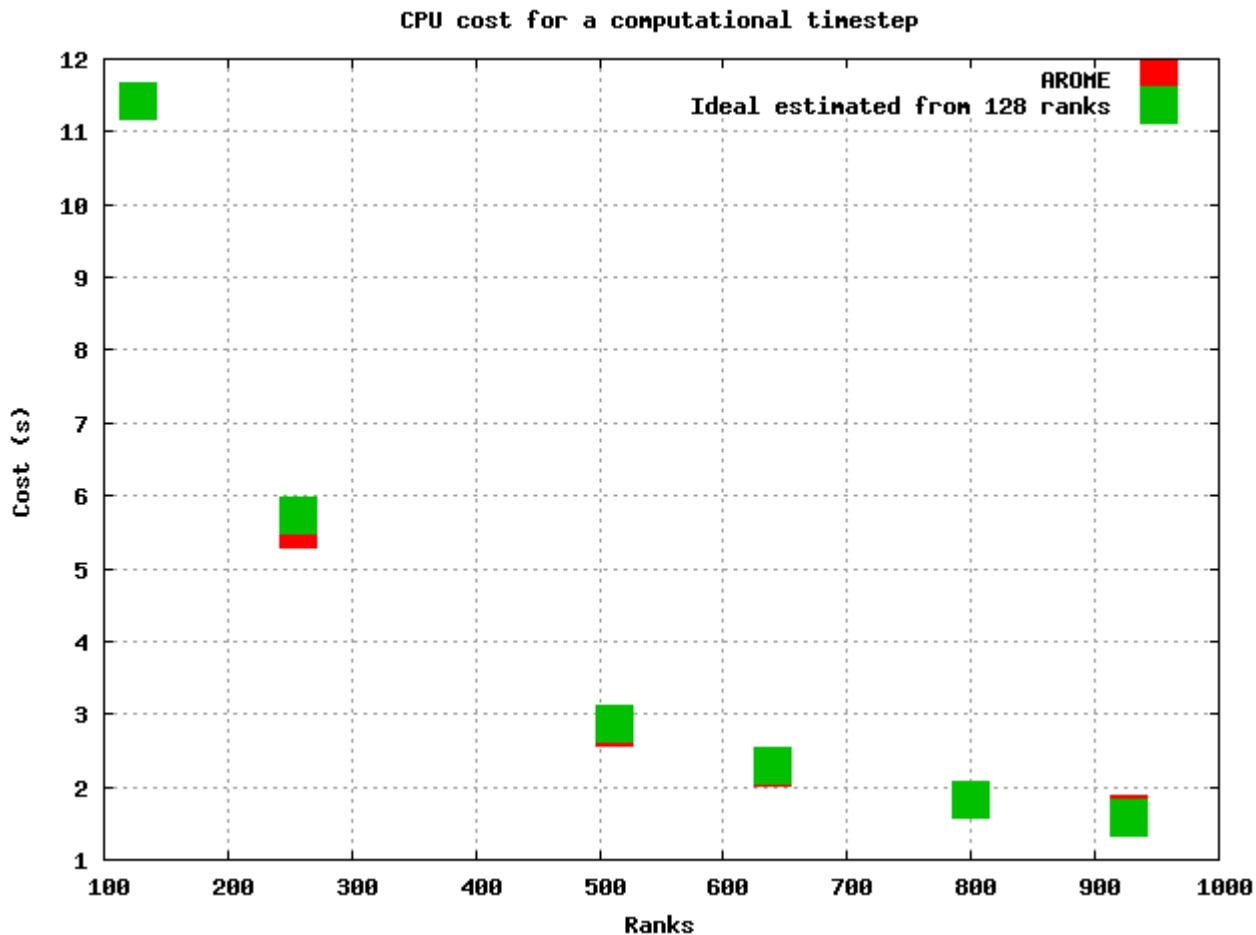
- SURFEX lfi output
- FA output
- BD input
  
- Computations
- Radiation timestep



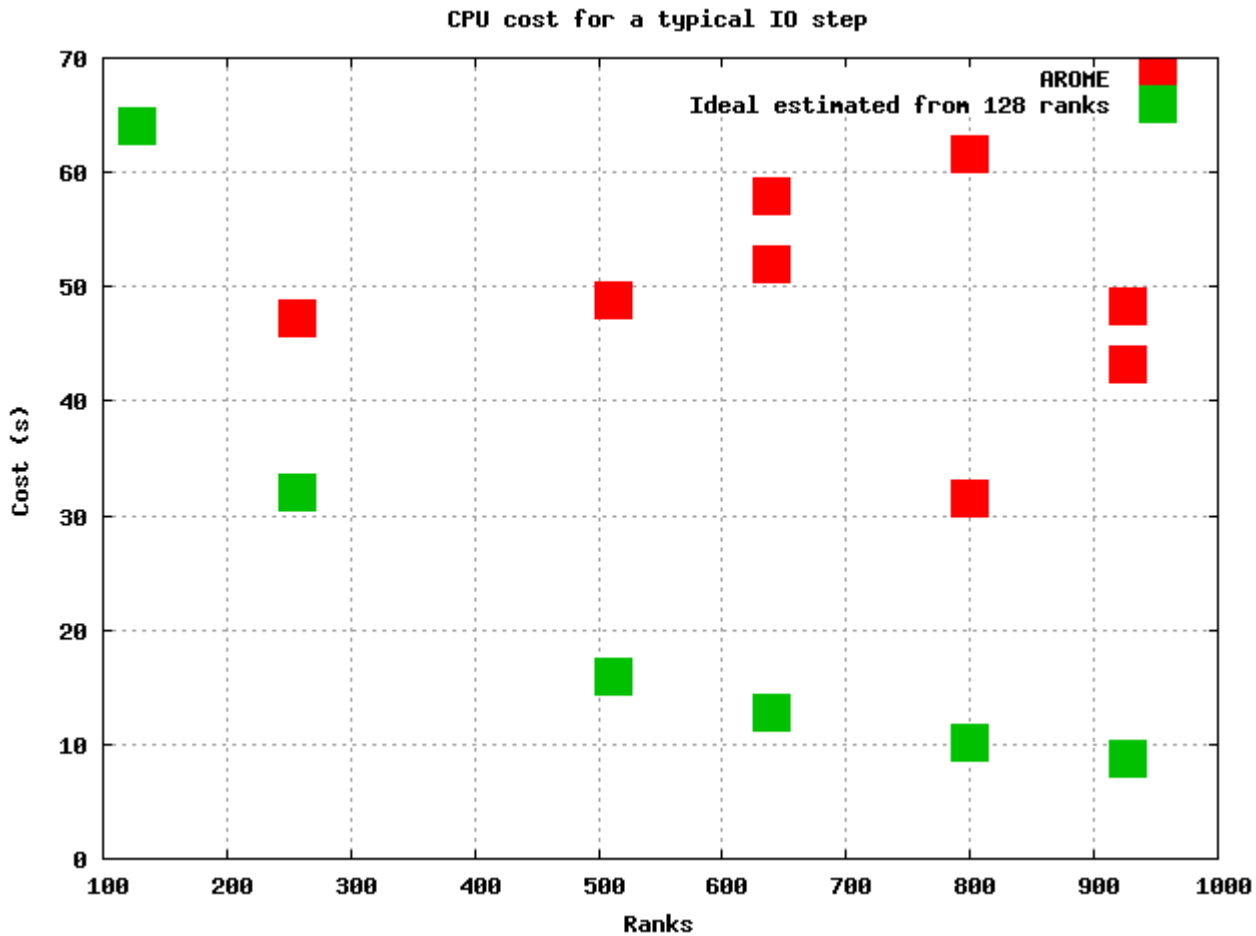
# Scalability on a 1200x1200x65 domain

## Sandybridge 2.2GHz 16 core nodes

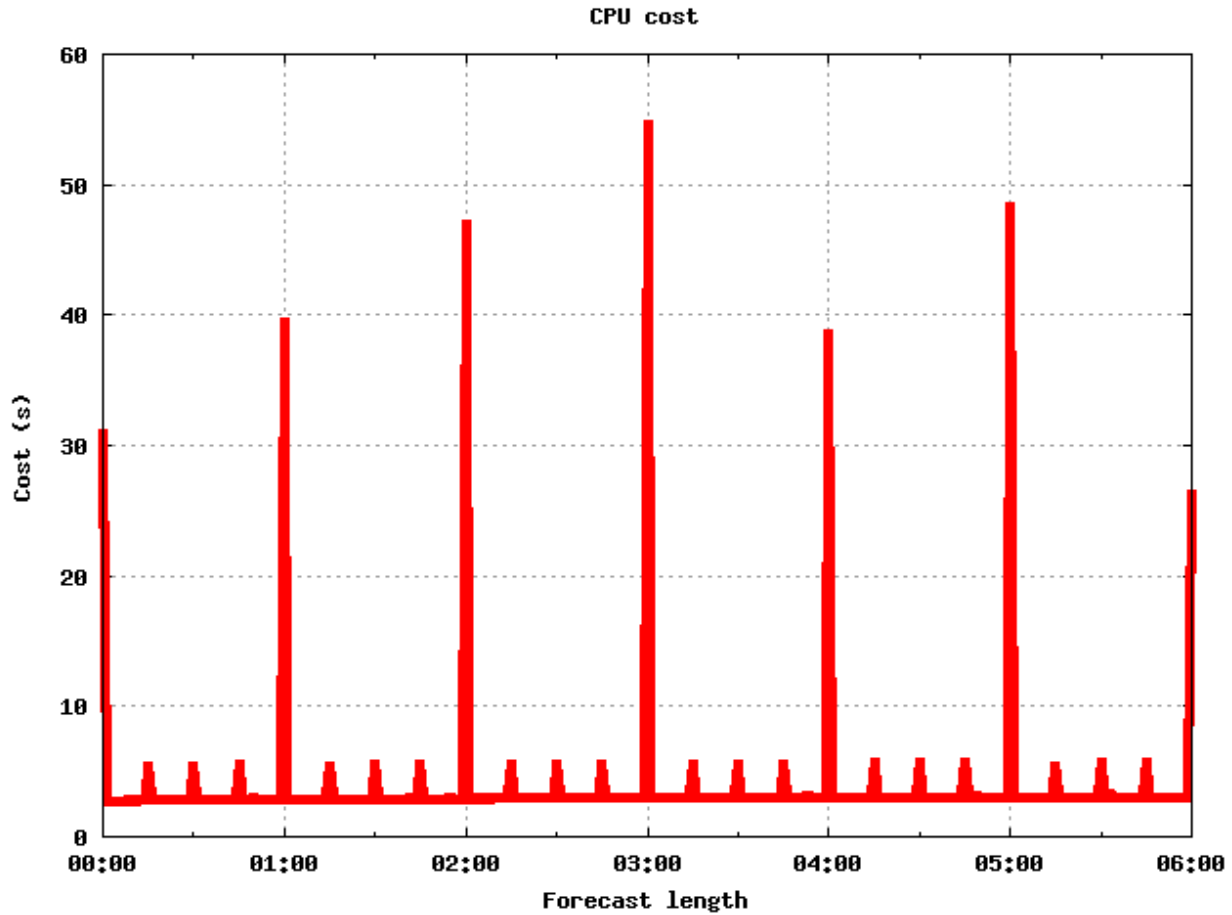
## Mellanox infiniband



# IO step only

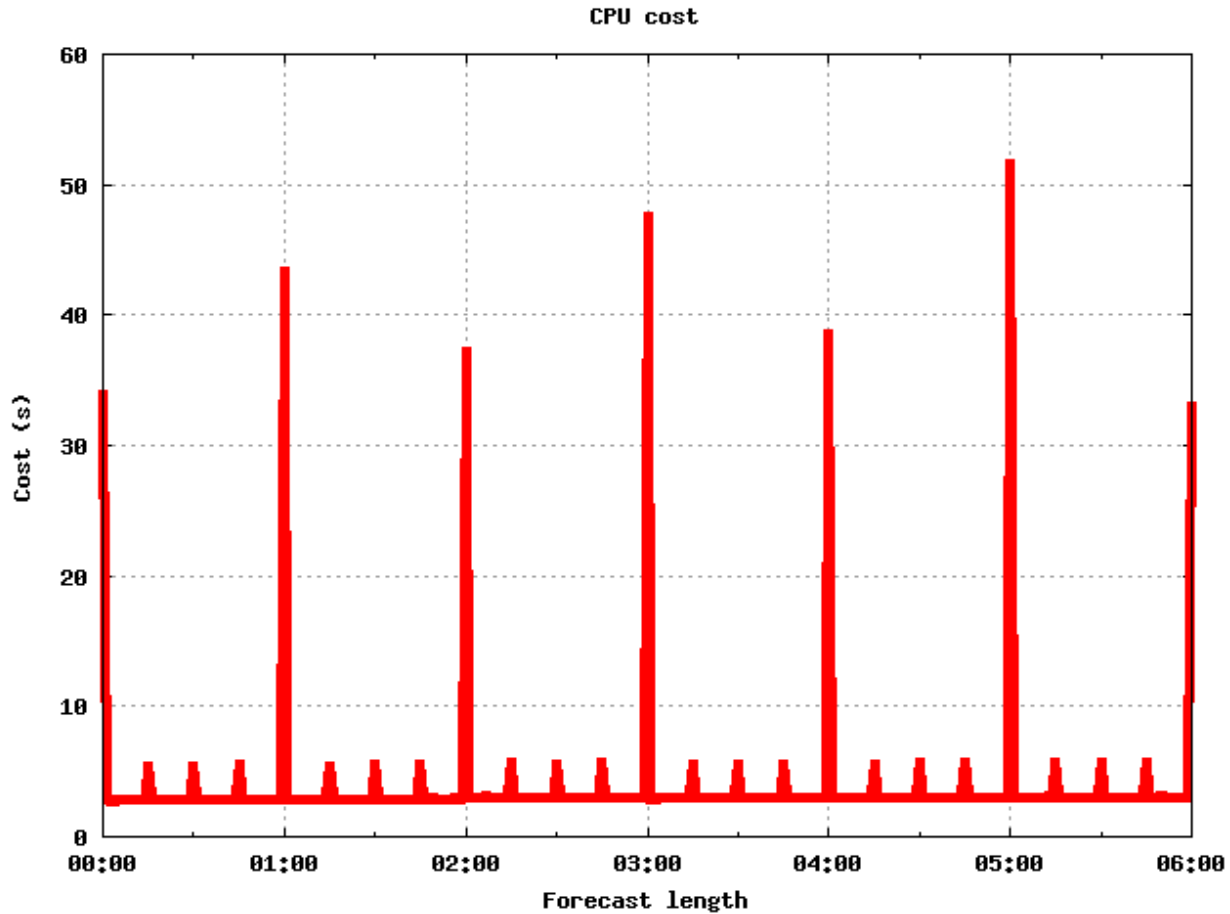


# IO FA + LFI

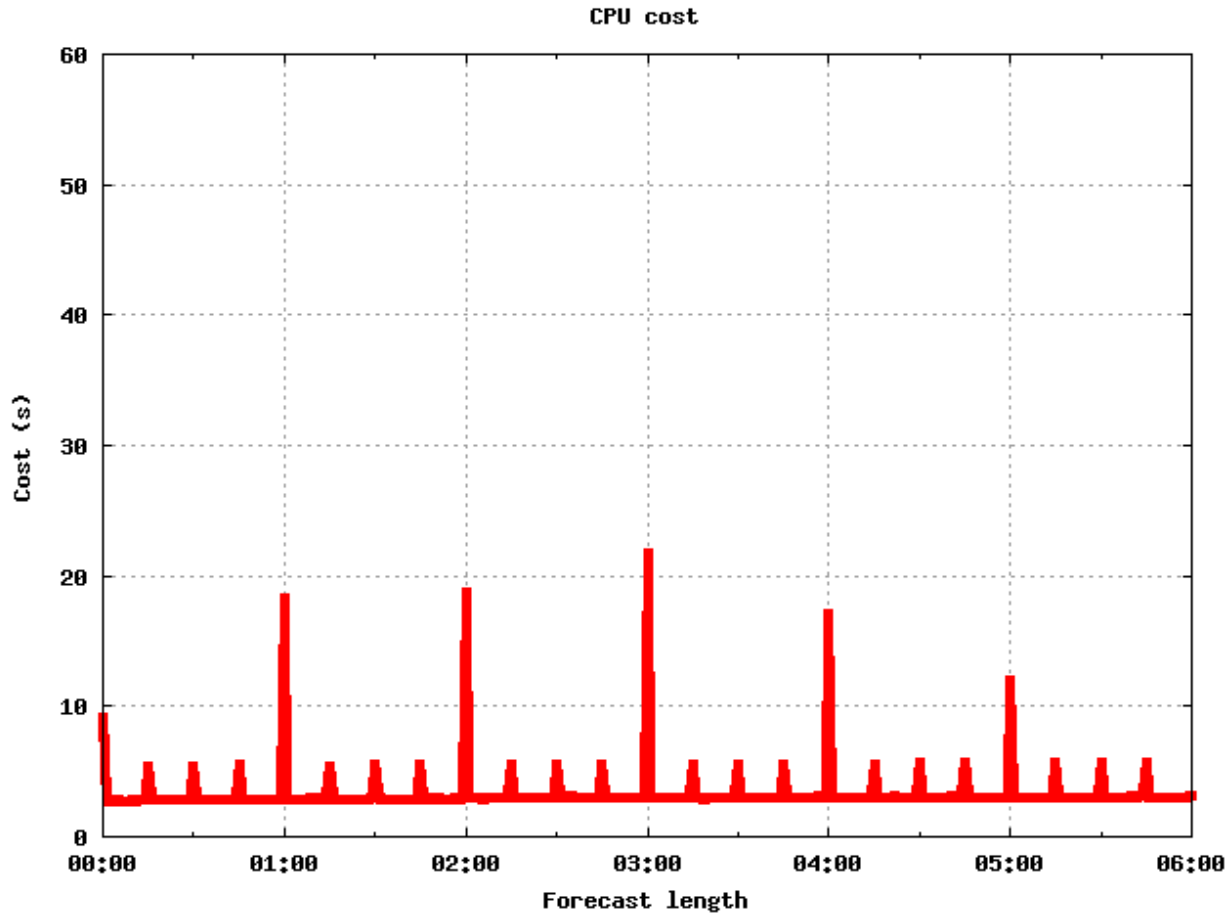




# IO FA+SURFEX as FA



# IO FA + SURFEX as FA + IO SERVER



# A more careful look on the IO steps (without IO server)

**Runs done on Lustre file system,  
 maximum BW for single file  
 120 MB/sec.**

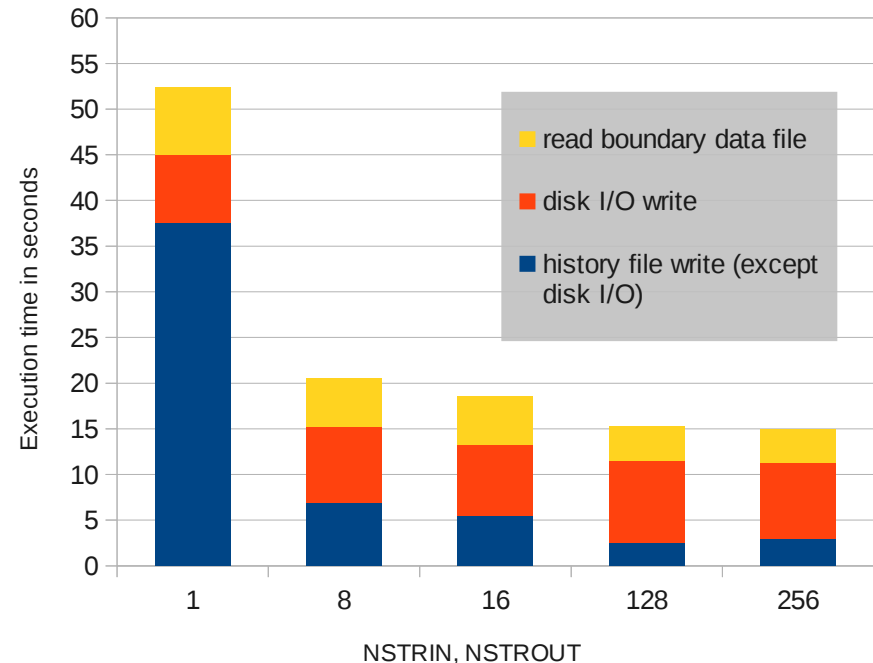
**“Excluding compute” means  
 subtracting the time for an  
 ordinary time step or radiation  
 time step.**

**NSTROUT important to maximize.**

**NSTRIN has rather small effect on  
 execution time and can vary a  
 lot depending on what else is  
 going on and where file  
 resides. Factor of three slower  
 when NSTRIN=1 → NSTRIN=8  
 observed in one case!**

Arome cy38 execution time for one "I/O step", 16 nodes, 256 ranks  
 excluding compute time.

Area: 750x960x65. 2.5km. Without IO-server.



*Courtesy Torgny Faxen NSC*

# IO performance with IO server.

**IO-server works fine. Time to concatenate output files not included though. Maybe reading of output data can be done in parallel instead?**

**With IO-server it seems like reading of boundary data is now the largest time consuming routine.**

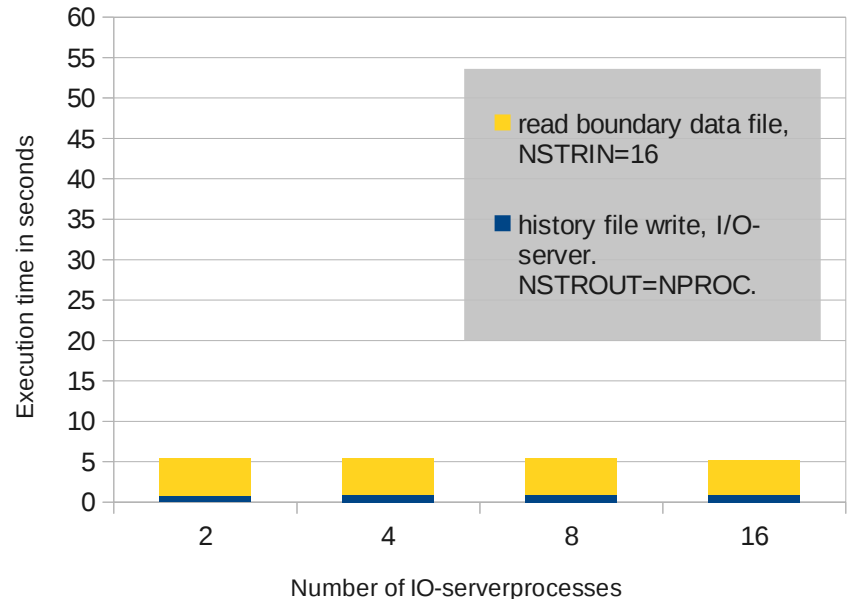
How to improve reading of boundary files?

Asynchronously through the IO-server should be possible?

Modify the actual READ?

Arome cy38 execution time for one "I/O step", 16 nodes, 256 ranks, excluding compute time. NSTROUT=NPROC.

Area: 750x960x65. 2.5km. With IO-server.

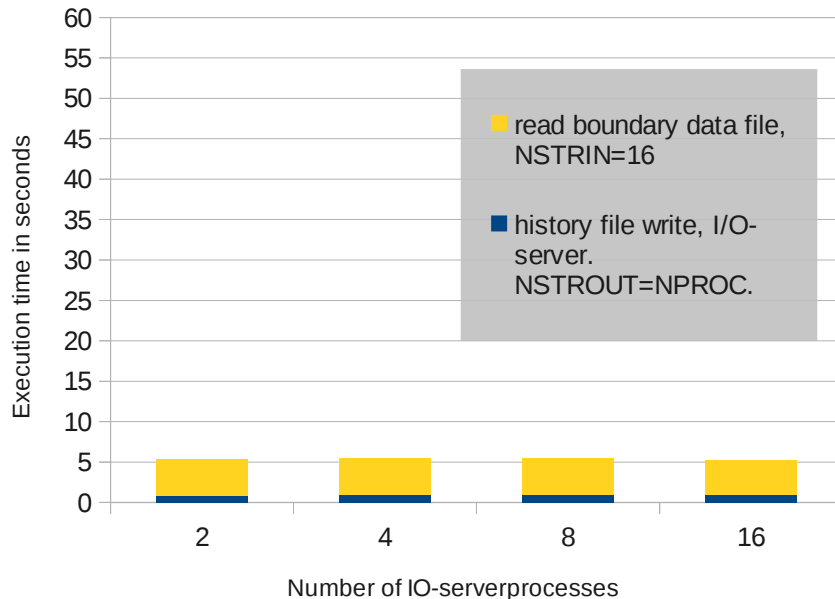


# Better performing by prestaging the input file ( just an example )

## Default

Arome cy38 execution time for one "I/O step", 16 nodes, 256 ranks, excluding compute time. NSTROUT=NPROC.

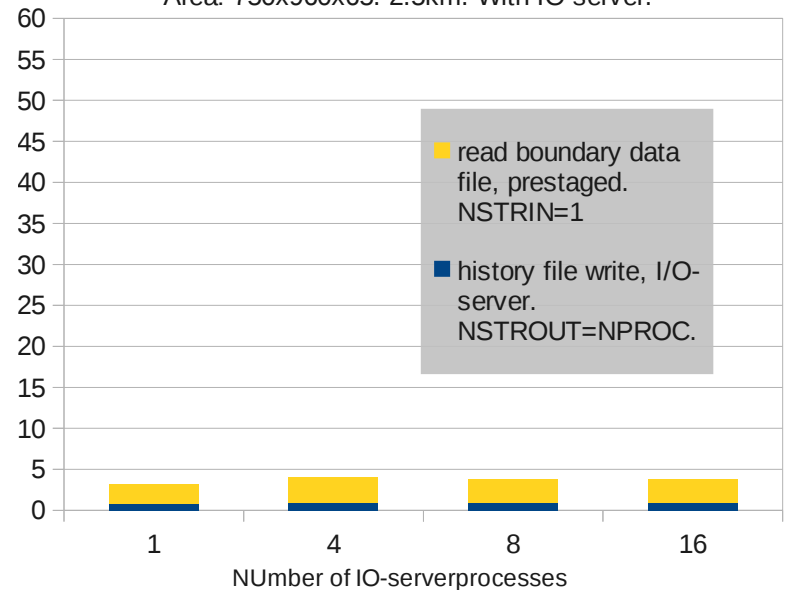
Area: 750x960x65. 2.5km. With IO-server.



## Prestage

Arome cy38 execution time for one "I/O step", 16 nodes, 256 rank, Prestaged boundarydata input file. NSTROUT=NPROC.

Area: 750x960x65. 2.5km. With IO-server.



So we believe we have a reasonably well working benchmark package!

Well....

**Mixed MPI  
OpenMP**

mpi	1 thread	2 threads	12 threads
pathscale	<b>X</b>	<b>X</b>	<b>X</b>
intel	<b>X</b>	<b>X</b>	<b>X</b>
cray	<b>X</b>	<b>X</b>	<b>F</b>
gfortran	<b>X</b>	<b>X</b>	<b>F</b>
pgi	<b>F</b>		

**Pure MPI with  
different  
compiler  
options**

MPI	Default	ieee	stack	bound
pathscale	<b>X</b>	<b>X</b>	<b>F</b>	<b>F</b>
gfortran	<b>X</b>	<b>X</b>	<b>X</b>	<b>F</b>
cray	<b>X</b>	<b>X</b>	<b>X</b>	<b>F</b>
intel	<b>X</b>	<b>F</b>	<b>F</b>	<b>F</b>
pgi	<b>F</b>			

## Compilation warnings and interface problems some examples

- Fortran pointer variable "FOO" is being used before being pointer assigned or allocated (2)
- Variable "FOO" is used before it is defined (95)
- Dummy argument "FOO" has the INTENT(OUT) attribute, but is never assigned a value or used as an actual argument (73)
- "FOO" is used but never set (31)
- Argument type differ from declaration (49)

(some of these were sent as corrections to cy39t1.)

Some warnings are more harmful than others

**It's natural that "real" errors are dealt with first, but how can we do better here? (back to the cycling strategy)**



## How to speedup your code

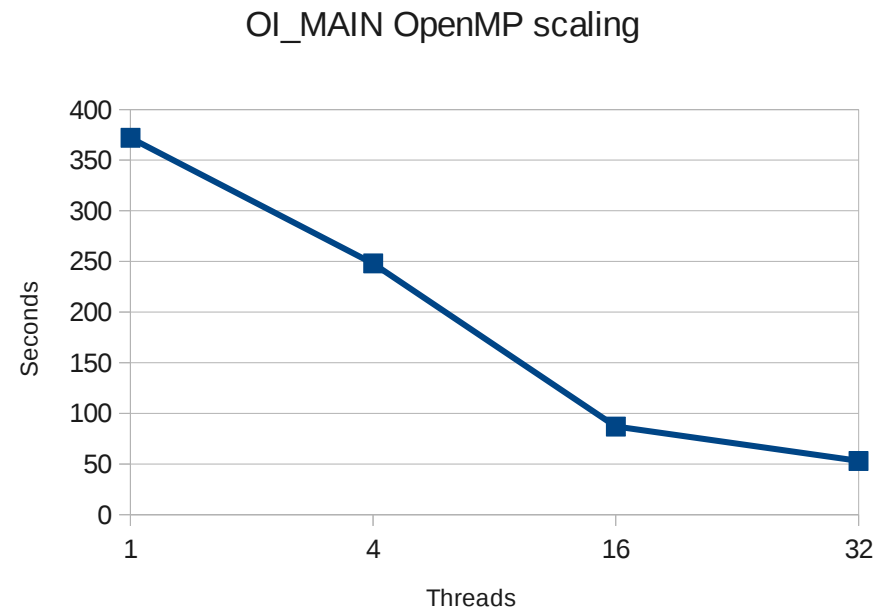
- **Nothing beats doing less**
- **If you have to do it, do it better**
- **Share your work**
  - **OpenMP, loop, single node parallelisation**
  - **MPI, distributed computations**

## Fighting with the surface assimilation

- On the way to 37h1.2 we modified OI\_MAIN and introduced SODA in parallel.
  - Increased the cost ~10 times compared to 37h1.1
  - OI main 1
  - The bad setup ) a
  - Could de extrapolate
- In cy38h1 we have OI main comes for
  - HARMON inside C/
  - Extrapolations not reproducible in the full environment

## Back to OI\_main

- **Share the work, add OpenMP directives to the painful part in the extrapolatio**
- **Scaling example for MetCoOp domain : reasonable number**
- **All tricks tried**
  - Doing less**
  - Share with OpenMP**
  - Doing it right: Reasonable number of threads**
  - reason for extrapolations!**



*Courtesy Ole Vignes met.no*

## Conclusions

- We have a benchmark package for the forecast model in cy38h1
- The IO server works well for output. Will be optional in HARMONIE together with SURFEX FA/LFI output. More work needed for input
- Several versions of surface assimilation exists with different computational and meteorological properties. Convergence discussion started!
- Our system is sem-implicit, semi-lagrangian, semi-portable and semi-fortran standard compliant

Thanks again for your attention  
Questions?