

efficient computing

at the

netherlands eScience center

Dr. Jason Maassen
October 2nd 2017

My background

Computer Science Dept. of VU University Amsterdam
.... in Henri's group!:

1993-1998 MSc: Fast Parallel Java

1998-2003 PhD: Method Invocation Based Prog. Models for Parallel Prog. in Java

2002-2012 Postdoc: Ibis + (GridLab, StarPlane, VL-e, PROMM-Grid, ...)

Netherlands eScience Center:

2012-2013: eScience Research Engineer

2013-2016: eScience Coordinator

2016-now: Technology Lead – Efficient Computing



Jason Maassen

netherlands eScience center ?



Research institute founded by NWO and SURF in 2011

“Enabling digitally enhanced research through efficient utilization
of data, software and e-infrastructure”

“Big data” in the real world



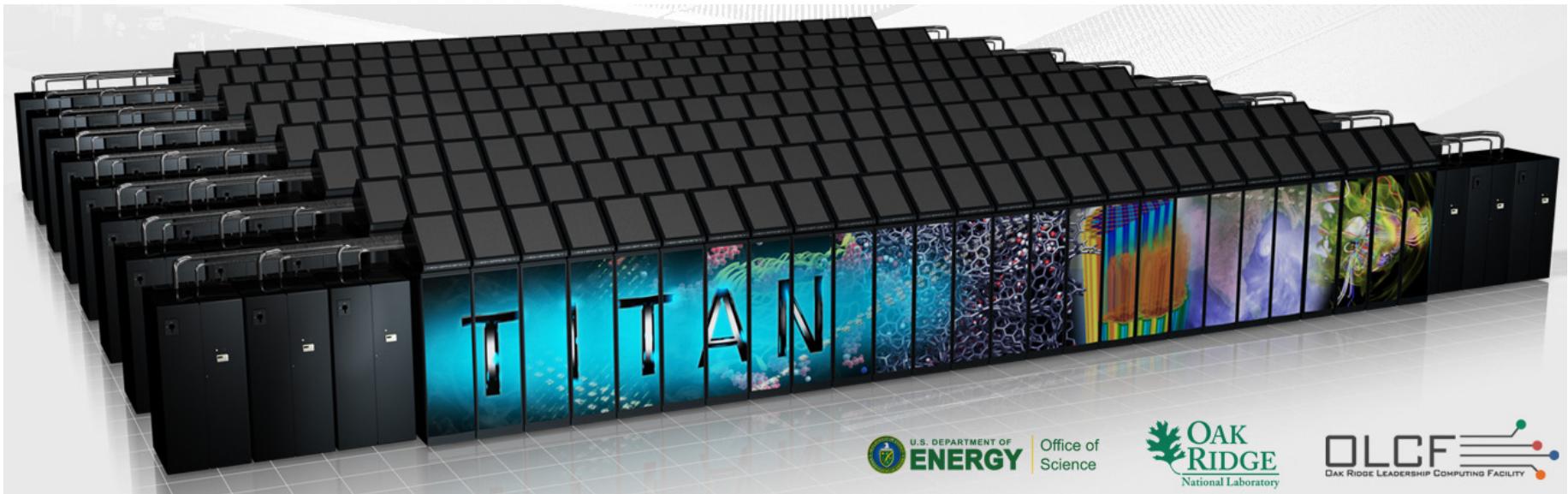
For **a few** scientists “big data” is a 20 PB storage cluster
(like this JUST cluster at Julich Supercomputing Center)

“Big data” in the real world



For **most** scientists “big data” is simply more than fits on their USB stick

“Big compute” in the real world



A **few** scientists run their computations on big machines.
(like this 560640 core, 17 PFlop/s Titan at Oak Ridge National Laboratory)

“Big compute” in the real world



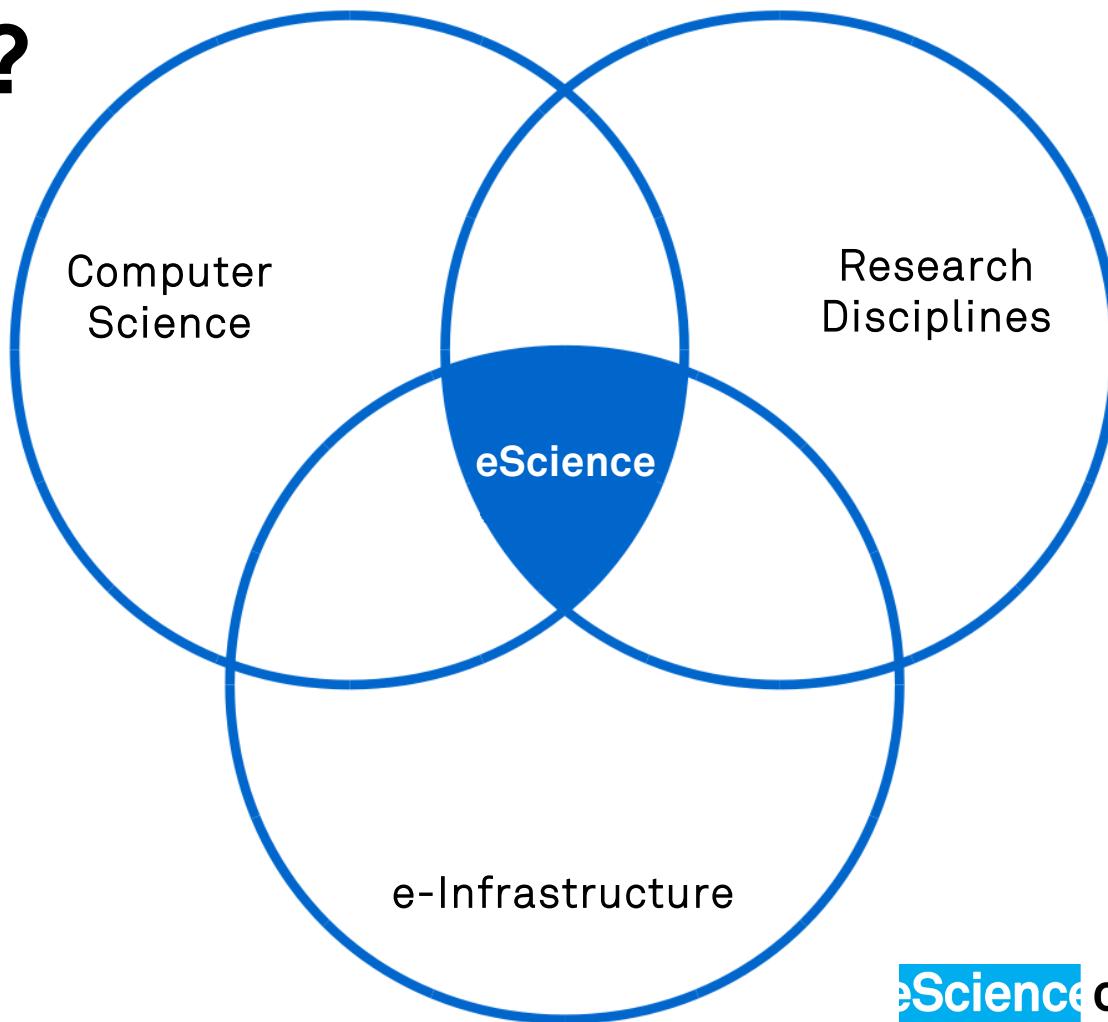
Most scientists simply stick to whatever their laptop can compute over the weekend.

Why eScience



Bridging the gap between:
science and infrastructure
science and computer science
scientific domains

eScience?



~~35~~

30

~~41~~

**Broadly oriented
scientists**

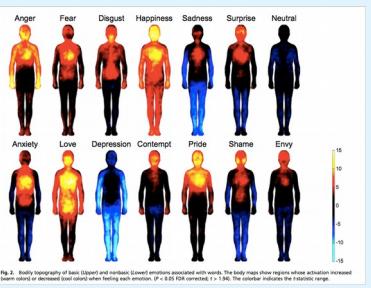
at the interface of
research and ICT

**Close collaboration
with researchers**

to implement
eScience projects
and technology

**Developing usable &
sustainable tools**

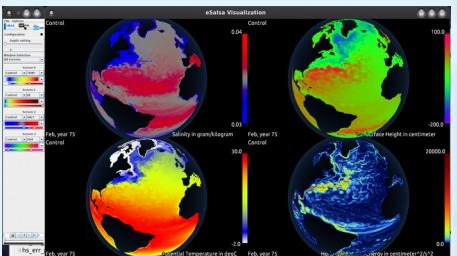
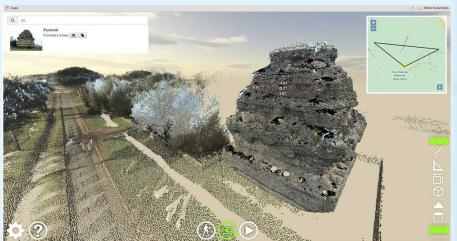
suitable for a broad
range of users



so far: 94 projects (on many different topics)

Humanities & Social Sciences

incl. SMART cities,
text analysis, creative technologies



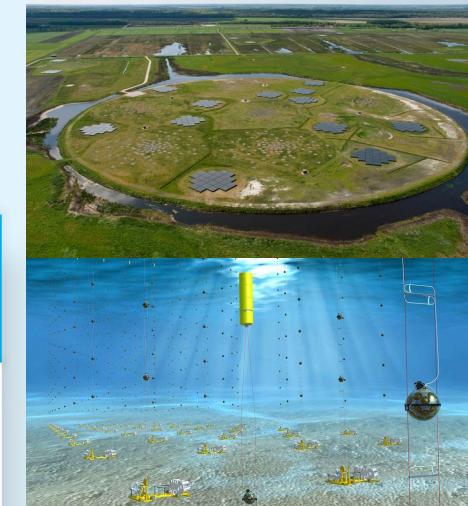
Sustainability & Environment

incl. climate, ecology, energy, logistics, water management



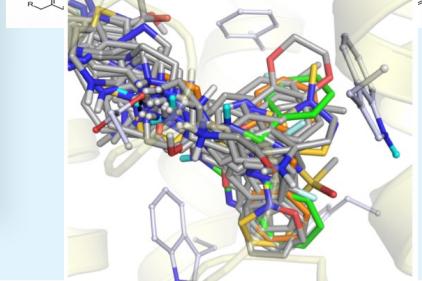
Physics & Beyond

incl. astronomy,
high-energy physics,
advanced materials



Life Sciences & eHealth

incl. bio-imaging,
next generation sequencing, molecules



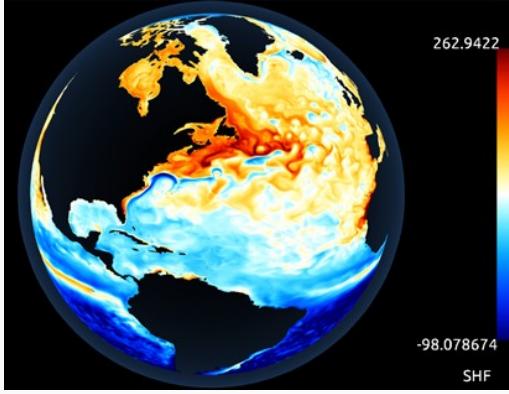
efficient computing ?

sheer number of
computations

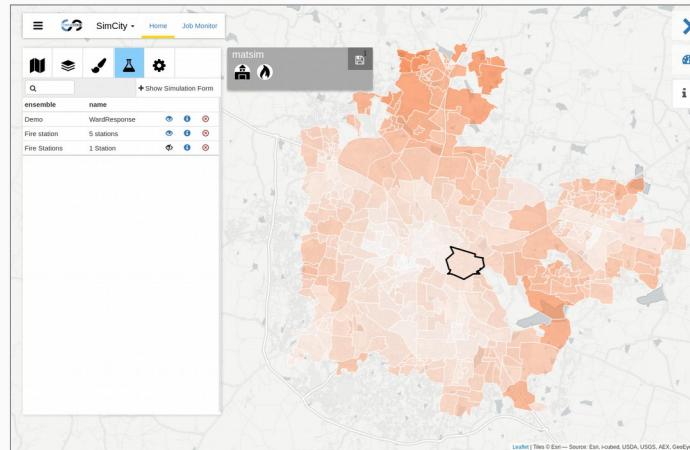
speed & volume of
incoming data

efficient computing ?

ease of use



efficient computing ?



high performance
computing

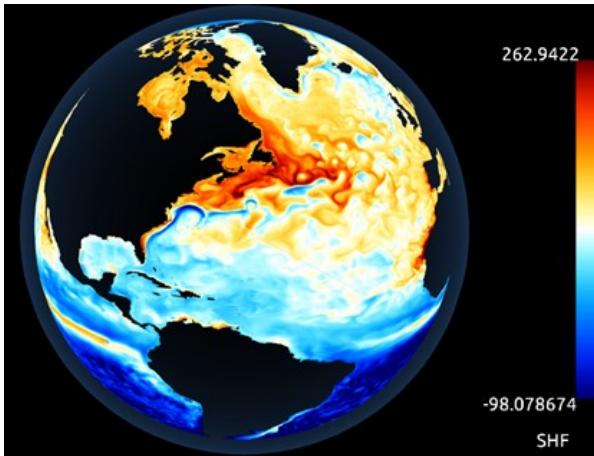
real time streaming
sensor data

efficient computing ?

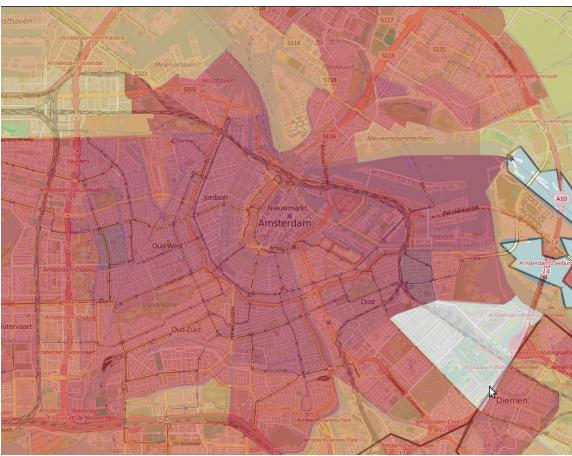
deployment and
coupling

A selection of
projects
using
**high performance
(distributed) computing**

The quest for resolution



eSalsa



Summer in the city



eWaterCycle

Many of our “traditional HPC” projects have a climate focus. They need to increase the resolution of their simulations, couple models, integrate observation data, after which they have trouble with load balancing or the large amounts of data they need to store

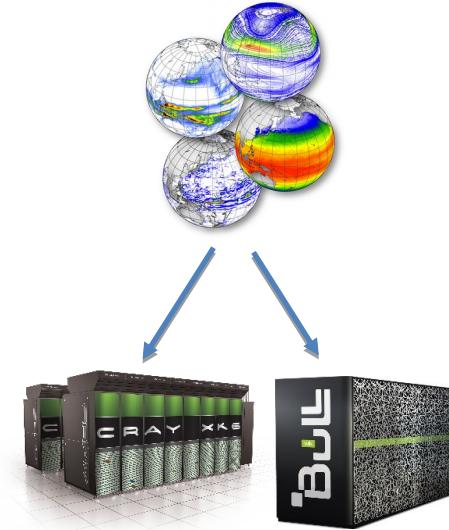
Potential solutions

```
211  ENDDO WALL_LOOP  
218  
219 !$OMP PARALLEL DO SCHEDULE(STATIC) ← Parallel Begins  
DO K=1,KBAR  
  DO J=1,JBAR  
    DO I=1,IBAR  
      IF (SOLID(CELL_INDEX(I,J,K))) CYCLE  
      FRHO(I,J,K) = -FRHO(I,J,K)  
      + (FX(I,J,K,θ)*UU(I,J,K)*R(I)-FX(I-1,J,K,θ)*UU(I-1,J,K)*  
      + (FY(I,J,K,θ)*VV(I,J,K)-FY(I,J-1,K,θ)*VV(I,J-1,K)  
      + (FZ(I,J,K,θ)*WW(I,J,K)-FZ(I,J,K-1,θ)*WW(I,J,K-1)  
    ENDDO  
  ENDDO  
231 !$OMP END PARALLEL DO ← Parallel Ends  
232  
233 SPECIES LOOP: DO N=1,N TRACKED SPECIES
```

parallelization



GPUs



distributed computing

GPUs deliver more performance per node than parallelization, but require a significant effort to rewrite existing code. Distributed work over multiple machines is easier, but results in less performance improvement and requires additional tools.

The eSalsa Project

Gain insight into **regional** sea-level changes (caused by climate change) by simulating the oceans with an unprecedented level of detail.

26% to 55%
below
sea level



Universiteit Utrecht



COMMIT/

Source: Actueel Hoogtebestand Nederland



Sea levels are changing...

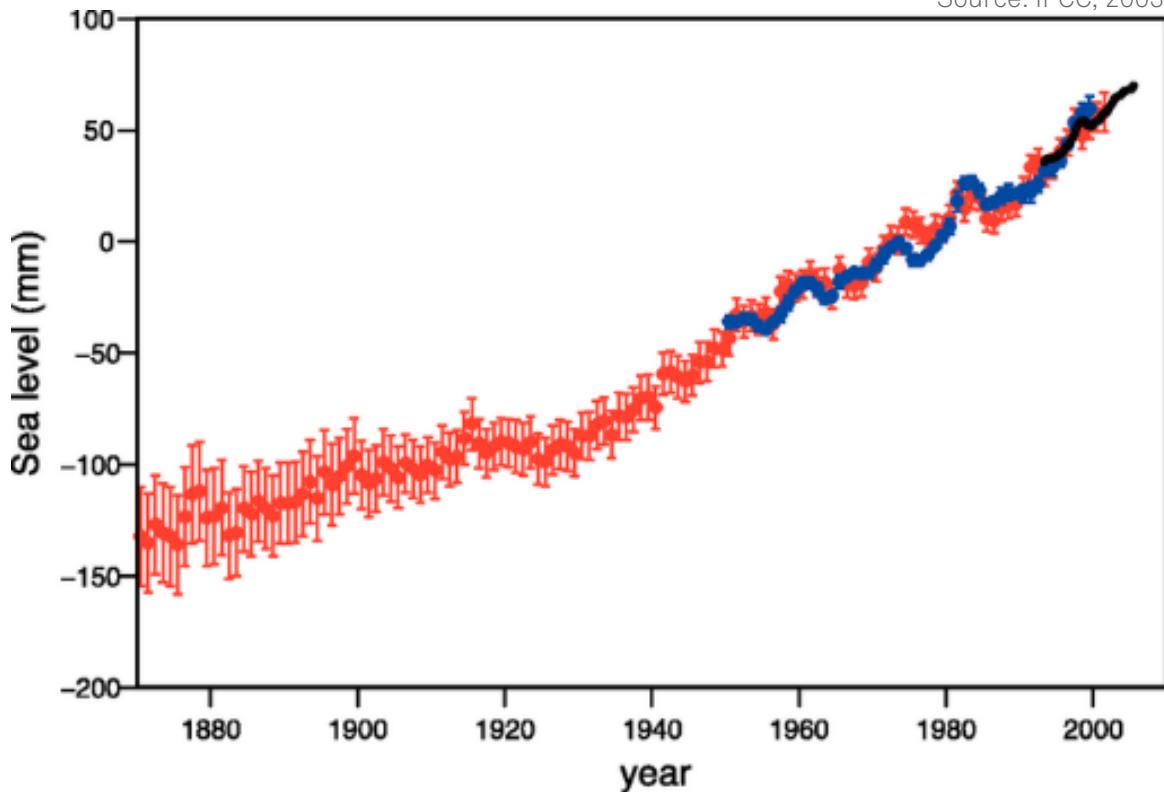
Source: IPCC, 2003

red -- historical records

blue -- tidal gauges

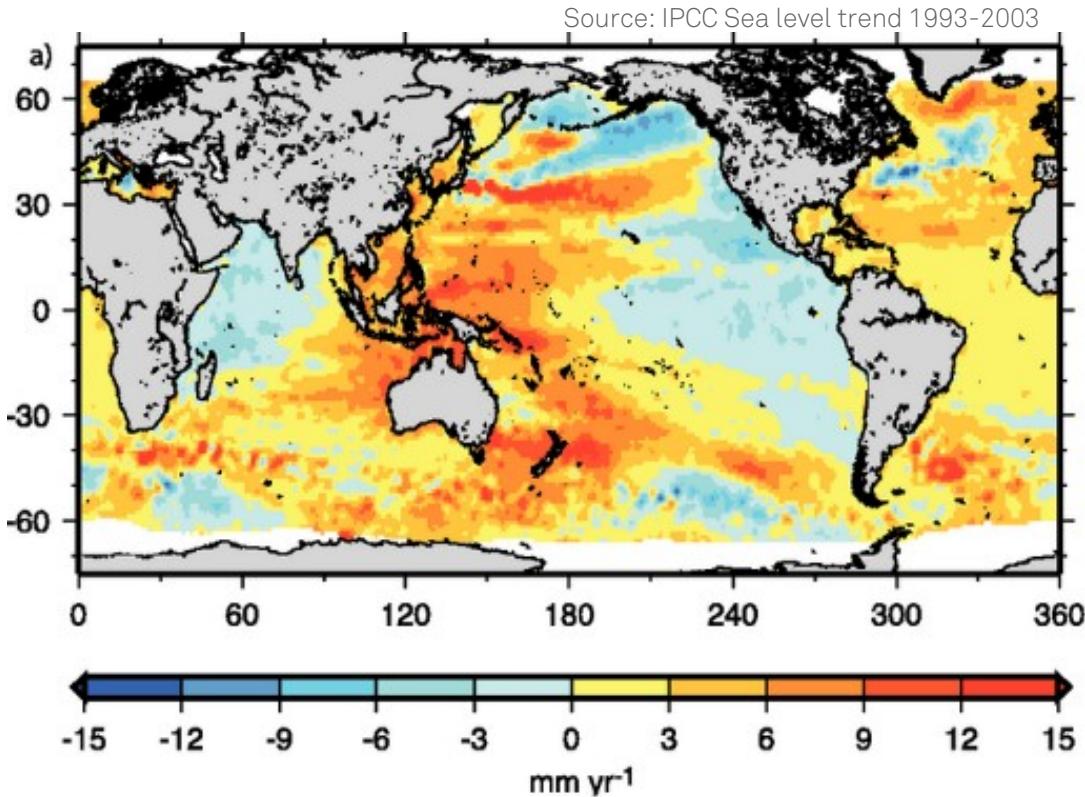
black -- satellite observations

All data show an upward trend!



...but the change is not uniform!

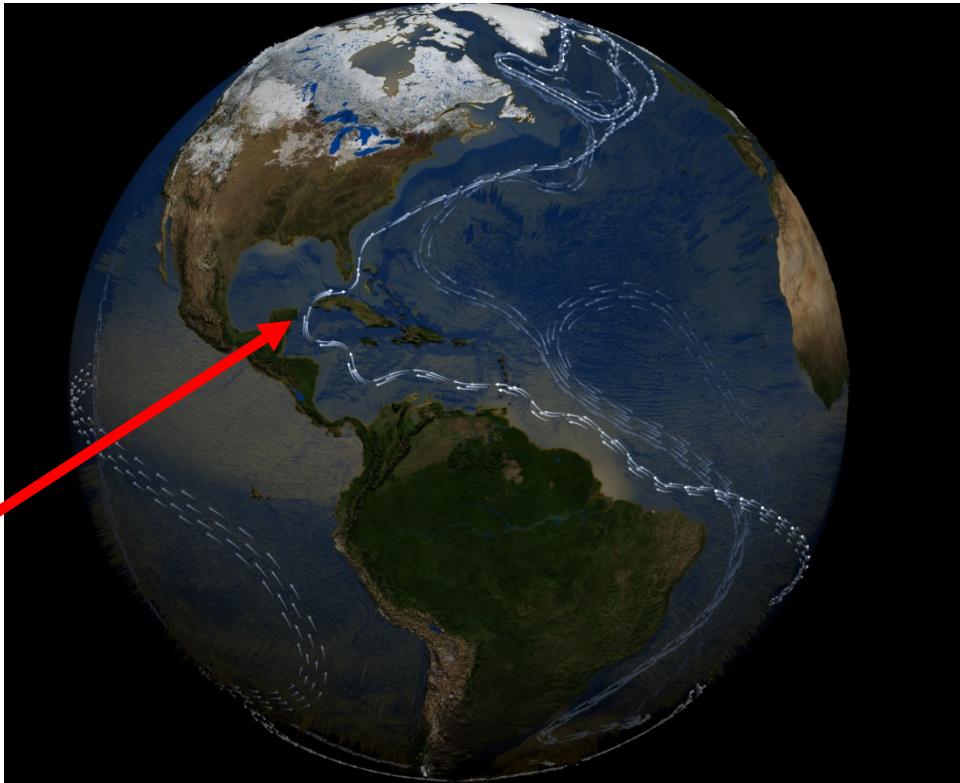
Satellite observations show large regional variations in sea-level change.



Sea level varies regionally

Caused by **large ocean currents** which are driven by temperature, and salinity differences and wind.

Meridional
Overturning
Circulation
(MOC)



Source: NASA/Goddard Space Flight Center Scientific Visualization Studio

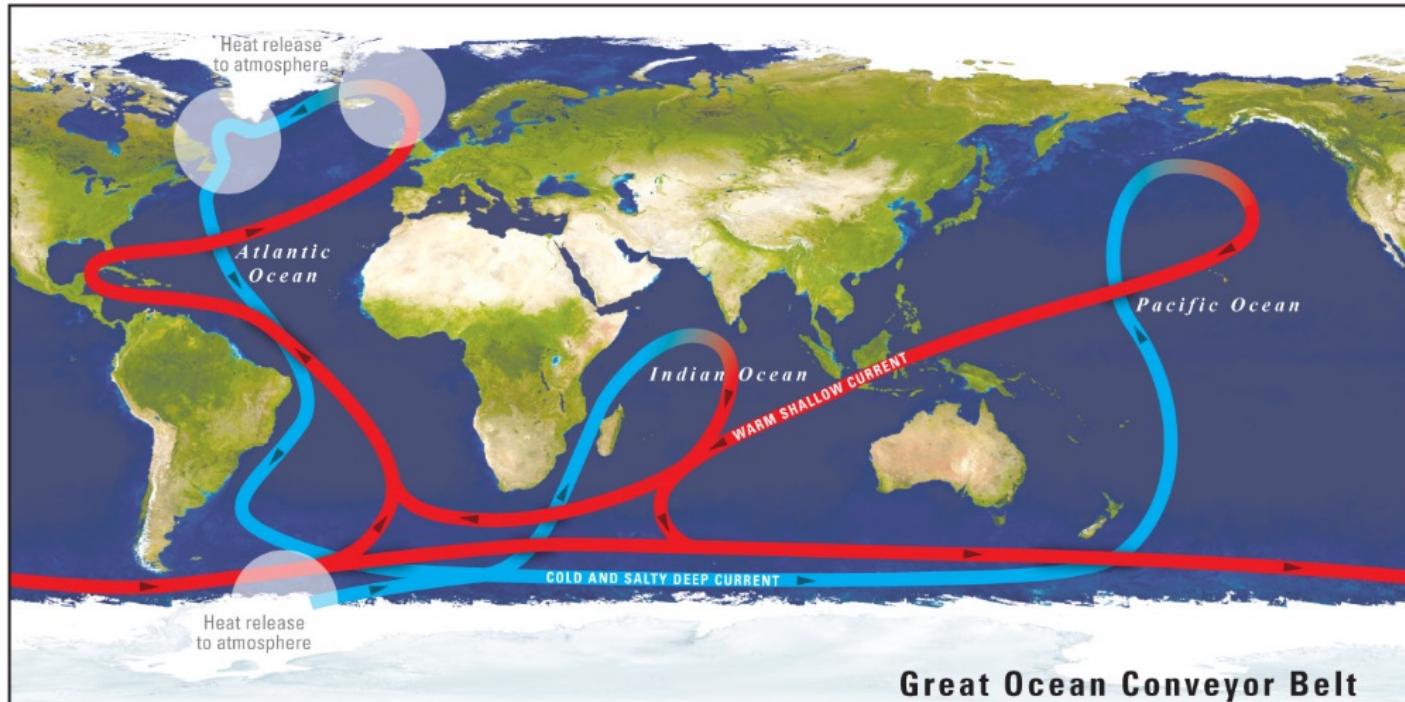
Meridional Overturning Circulation

Water transport:
20 billion liters/sec.

Heat release:
500 GigaWatt

What is the effect
of climate change
on the MOC?

Use simulations to
gain more insight



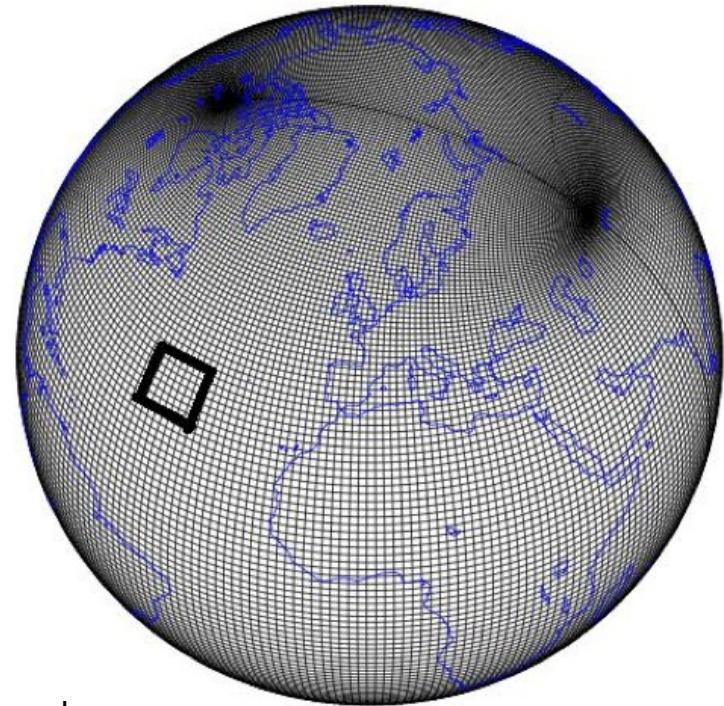
Source: IPCC, 1996

Parallel Ocean Program (POP)

“The POP ocean model is a level-coordinate ocean general circulation model that solves the three-dimensional primitive equations for ocean dynamics”

Resolution is important for the results:
1° resolution (100x100 km) was the norm.
0.1° resolution (10x10 km) is **eddie permitting**

Direct relation between resolution and compute time!

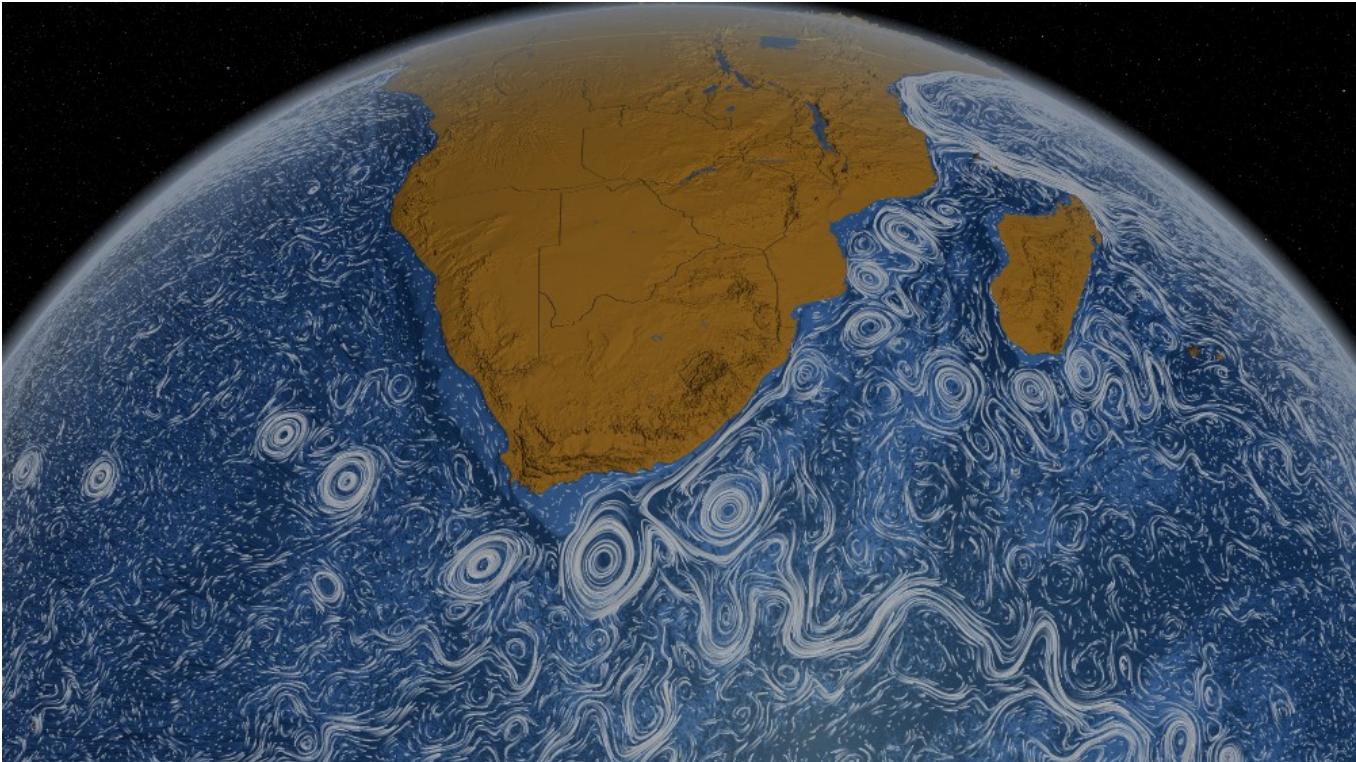


Source: Los Alamos National Laboratory

What are eddies ?

'Whirlpools', up to 300 km in diameter and 4 km deep.

They have a large effect on ocean behavior.



Source: NASA/Goddard Space Flight Center Scientific Visualization Studio

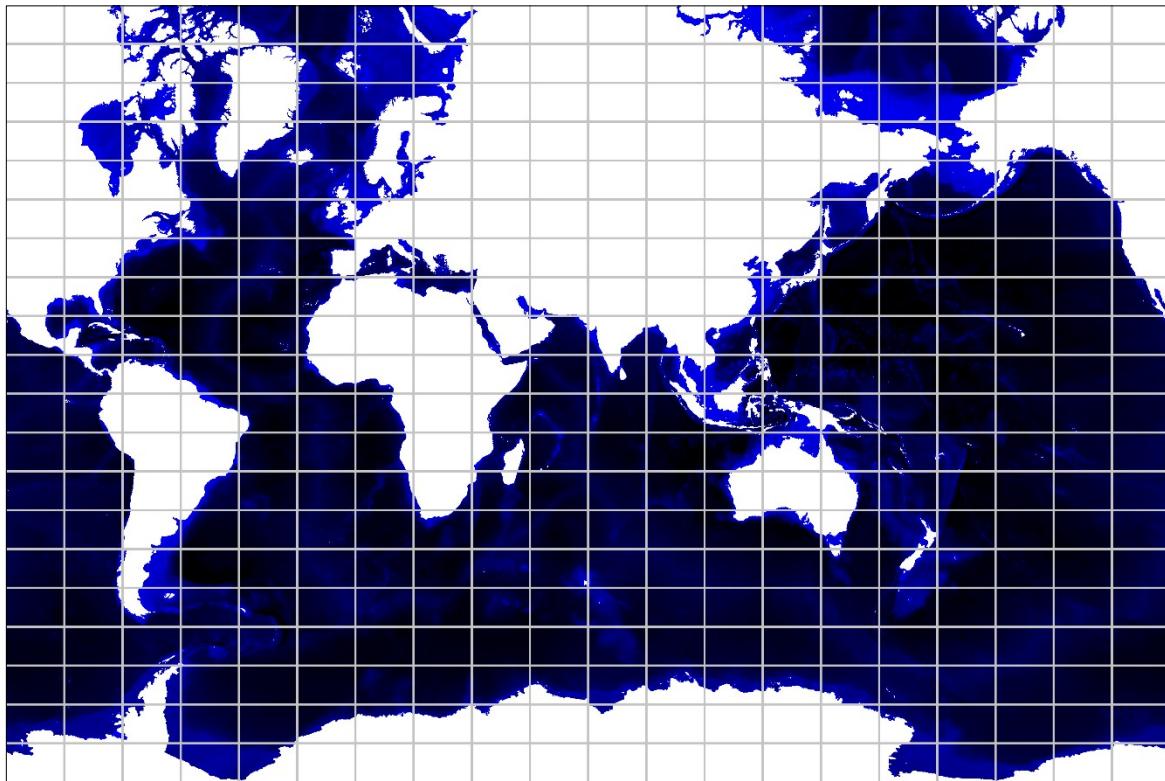
How does POP work ?

Fortran/MPI application (1992)
25 years old!!!

POP divides the world into a grid,
which is divided into blocks.

These blocks are **distributed**
over many **processes** (= cores)
using **MPI**.

Traditionally a **cartesian**
distribution is used that
assigns one block to each
MPI process.



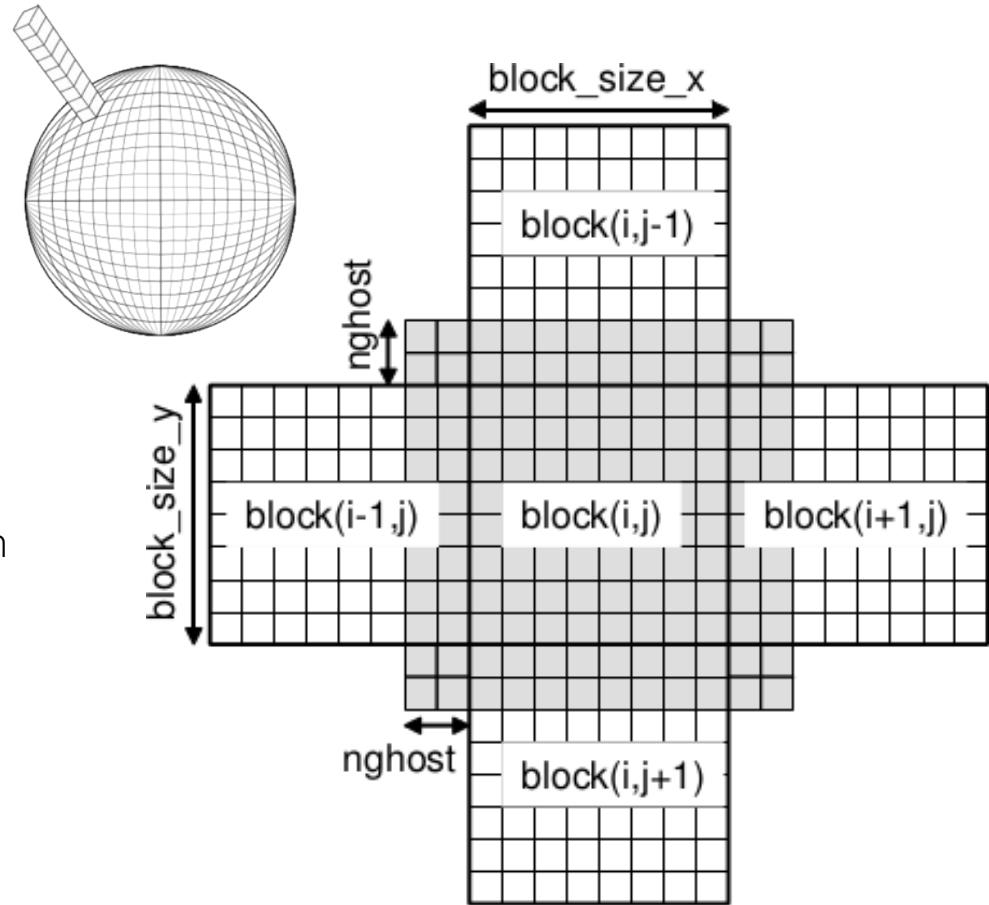
Halo exchange

Each MPI process computes the state of many variables (water temperature, salinity, velocities, etc) in its local 3D column for each time step.

To do this every grid cell needs information from its neighbors.

Therefore, each block needs to exchange data with its neighbors many times during each time step (3D halo exchange).

This is the main source of MPI communication in POP



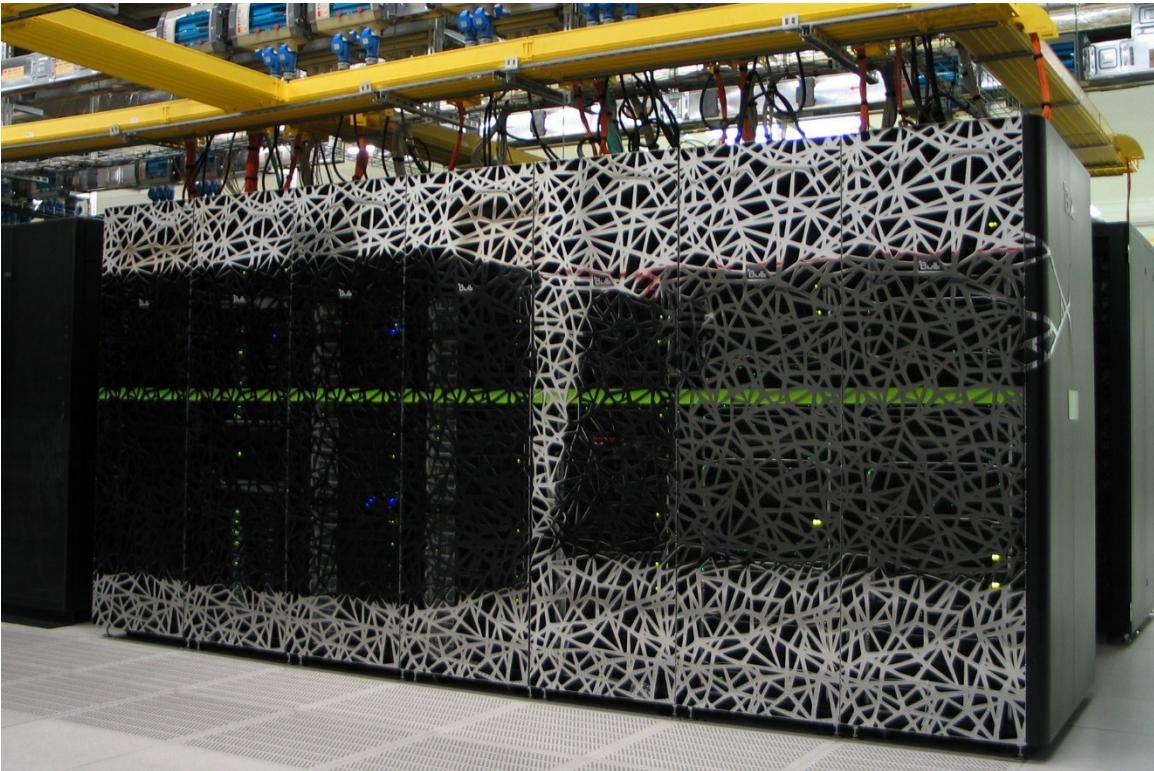
Source: A Distributed Approach to Improve the Performance of the Parallel Ocean Program
Ben van Werkhoven et al., Geoscientific Model Development, 7, 257-281, 2014

How we run our ocean simulations?

SURFsara Cartesius
40960 cores
117 TB memory
1.0 PFlop/s

1 simulation of 100 years
at 0.1° resolution (10x10 km)
takes **20 days** on O(1000)
cores and produces
10+ TB output.

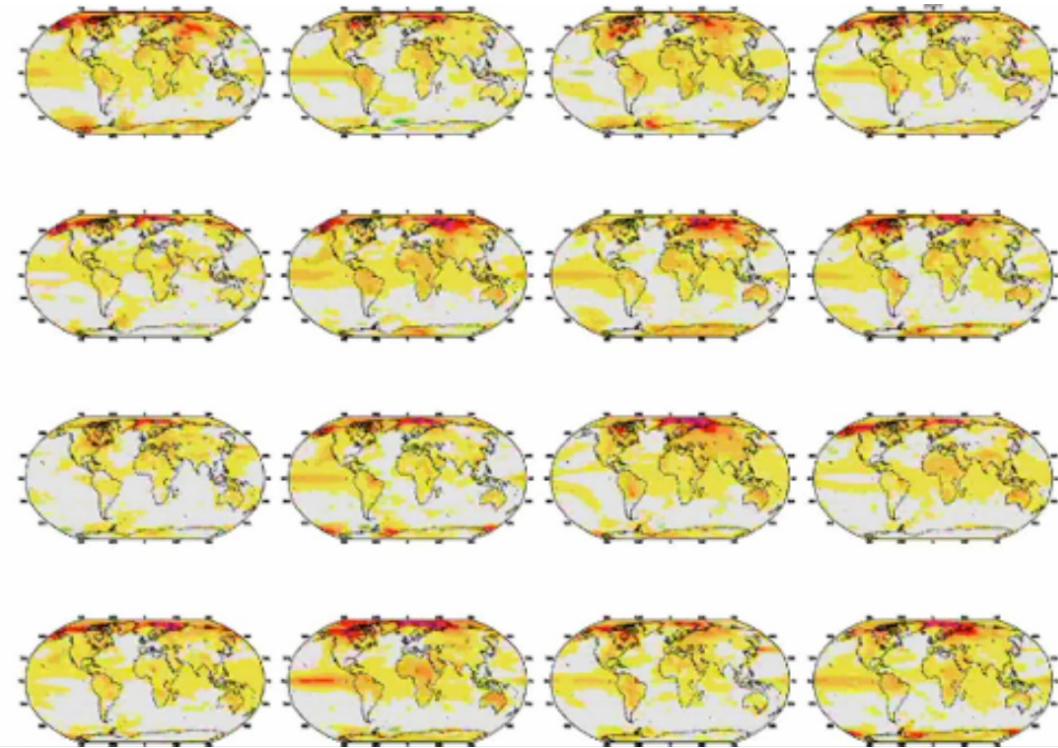
(but there is more!)



Source:SURFSara

Ensembles

We don't run 1 simulation but
an **ensemble** of 16, each using
a slightly different forcing.



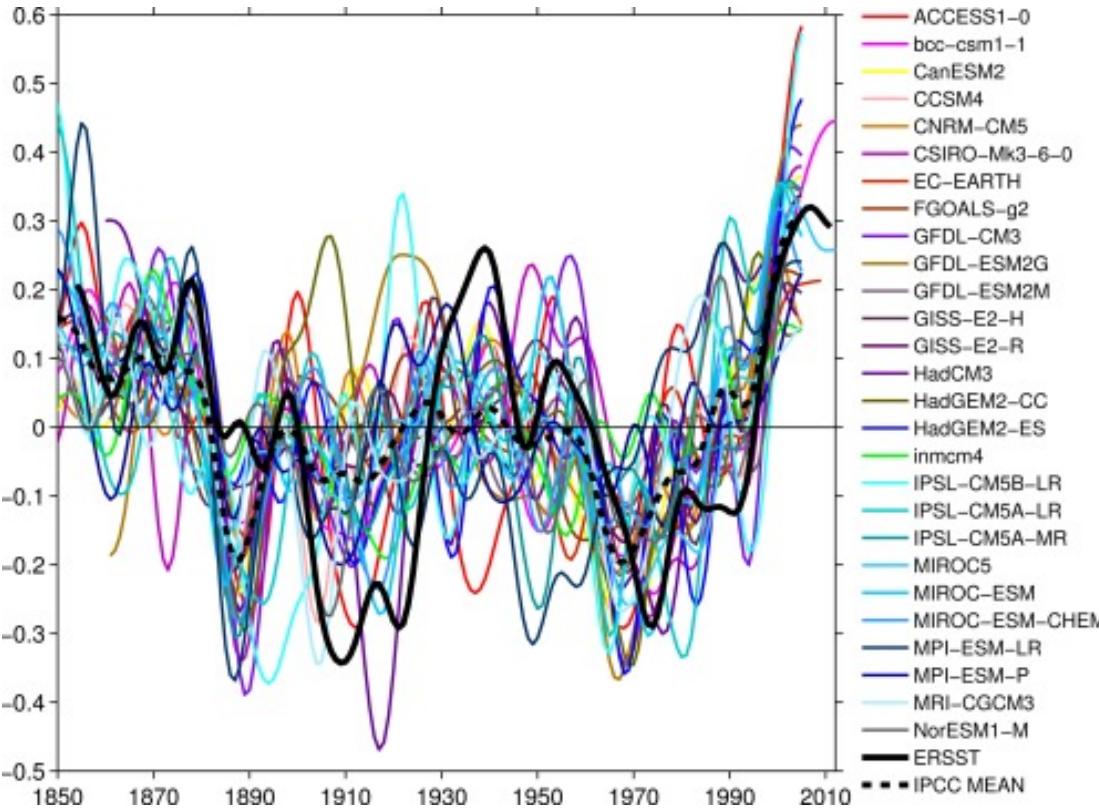
Source: IPCC 1996

16x increase in compute time

Why ensembles?

Climate is a chaotic system:
a small change in forcing, model
or starting conditions may change
the outcome significantly.

By running many simulations
and/or different models, we get
many results and do statistics on
them to determine the certainty
of the results.



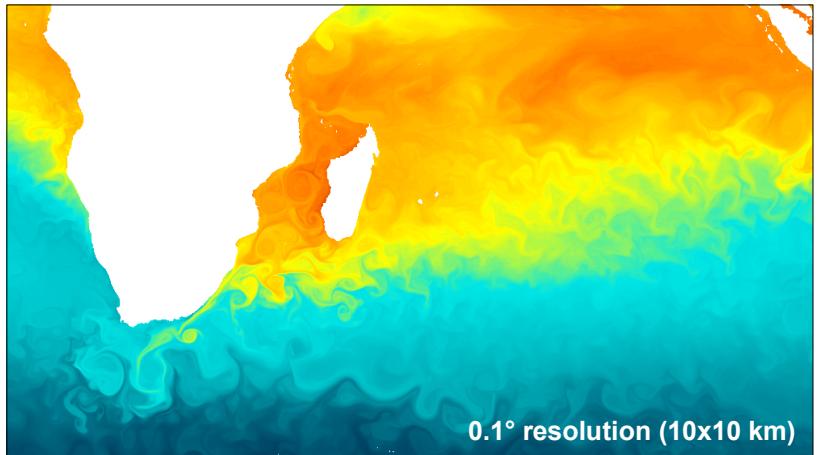
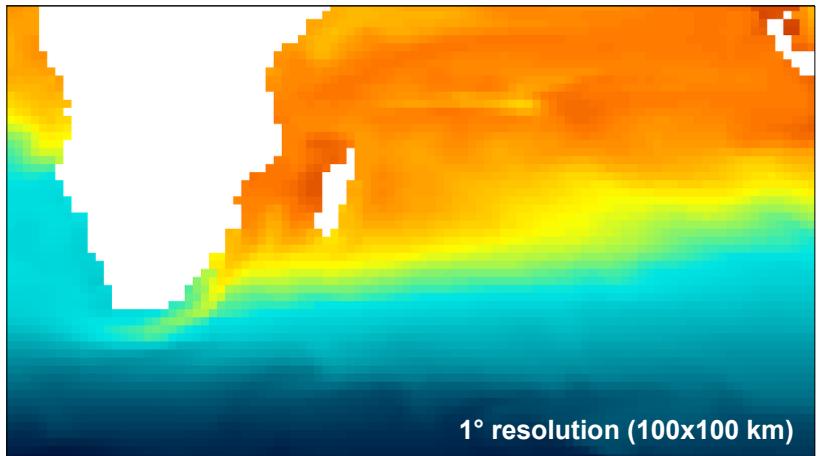
Source: L. Zhang, C. Wang, DOI: 10.1002/jgrc.20390

Higher Resolution

0.1° resolution (10x10 km) is only the start! We want to increase the model resolution even further to get more detailed results.

Ultimate goal (last time I asked):
0.01° resolution (1x1 km)
(fully eddie resolving)

100x increase in compute time!



Source: eSalsa results

We need more compute power

Needed for (a combination of):

higher resolution + large ensembles

We looked at the following solutions:

Distributed task farming - for running ensembles

Accelerators - for more compute power per node

Distributed computing - for more overall compute power

Visualization - for fast evaluation of the results

distributed computing

(the simple way)

Distributed Task Farming

Running ensembles takes a long time and quite a bit of effort!

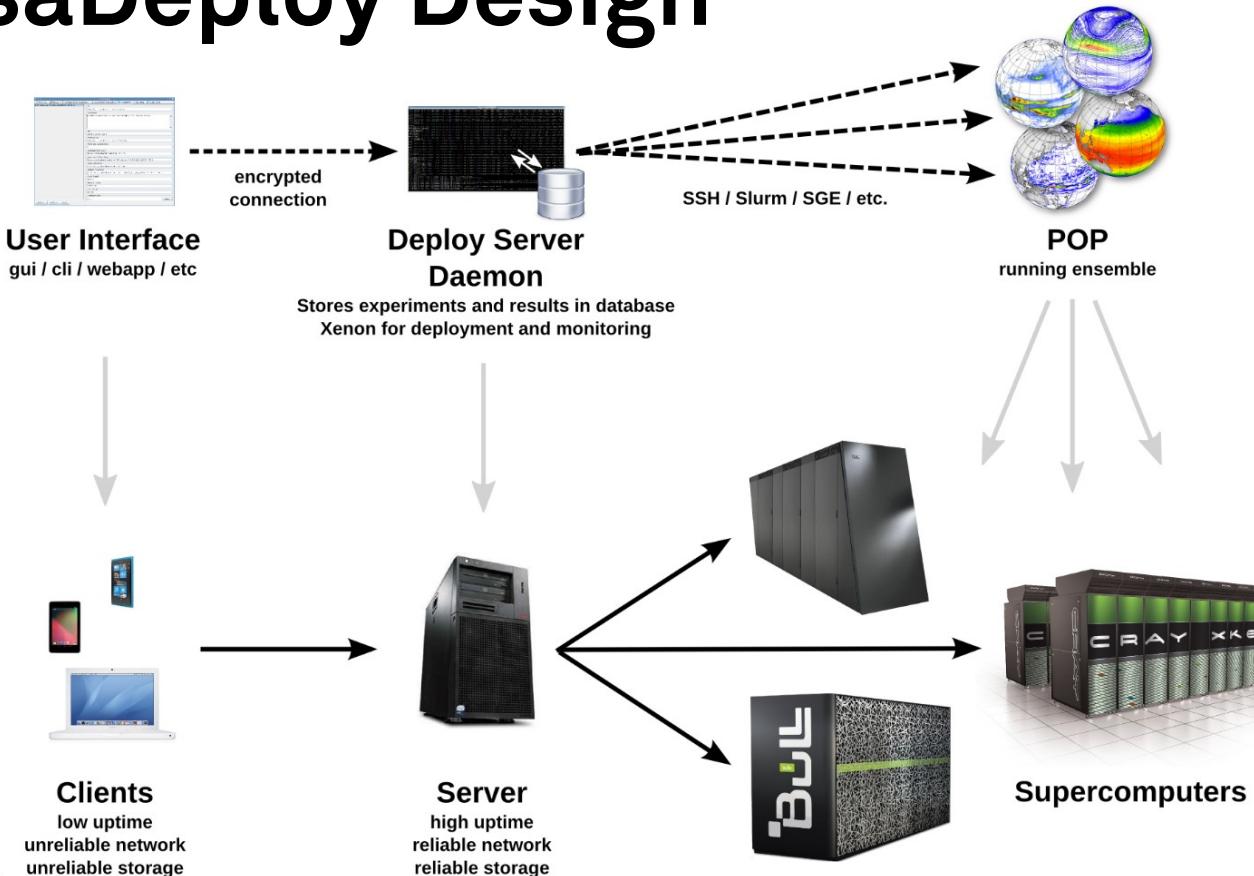
We need to:

- create submission scripts for each target supercomputer
- (possibly) adapt submission scripts for each ensemble member
- create POP configurations for each ensemble member
- need to ‘babysit’ each ensemble run (which may take up to a month).

eSalsaDeploy help the scientist by:

- automatically running **ensemble members** on **multiple supercomputers**
- generating the necessary config files
- managing the input files, logs files and results.
- babysit the application

eSalsaDeploy Design



eSalsaDeploy GUI

The screenshot displays the eSalsaDeploy GUI interface. On the left, a sidebar lists workers: DAS4.Astron, 1-degree.cartesian, 32-cores; DAS4.LU1, 1-degree.cartesian, 32-cores; DAS4.MN, 1-degree.cartesian, 32-cores; DAS4.TUD, 1-degree.cartesian, 32-cores; DAS4.UVA, 1-degree.cartesian, 32-cores; and DAS4.VU, 1-degree.cartesian, 32-cores. The main window shows an experiment template configuration for 'DAS4.VU, 1-degree.cartesian, 32-cores'. It includes fields for ID, Comment, URI (ssh://fs0.das4.cs.vu.nl), Template Directory (/home/jason/esalsa/templates/pop_1d_c_32c), Experiment Directory (/home/jason/esalsa/experiments), Input Directory (/var/scratch/jason/esalsa/experiments/input), and Output Directory (/var/scratch/jason/esalsa/experiments/output). Additional properties include clinic_distribution_type (cartesian), nprocs_tropic (32), tropic_distribution_type (cartesian), distribution_file (unknown), and nprocs_clinic (32). Buttons for Refresh, Delete, Save, and Clear are at the bottom. To the right, two windows show configuration files and logs. The top window shows a configuration file with various parameters like domain_nml, nprocs_clinic, and distribution_file. The bottom window shows a log titled 'Deploy Log' with entries from 2014-10-14 11:13:02 to 2014-10-14 11:13:18, detailing the preparation of remote directories and files, starting the job, and changing states to STAGE_IN and STAGE_IN_COMPLETE.

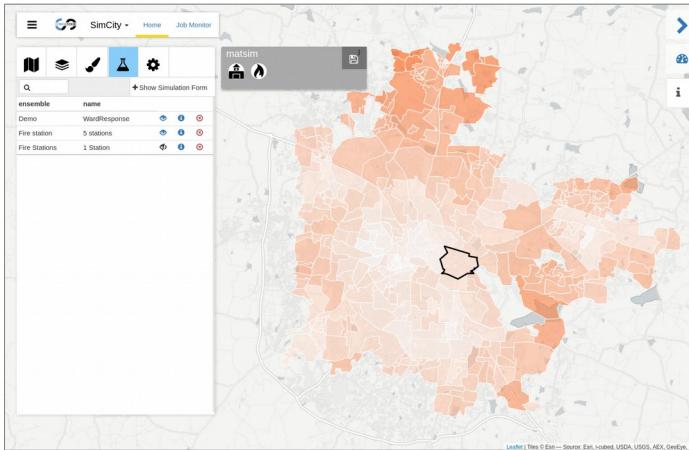
```
Configuration
1 &domain_nml
2 nprocs_clinic = ${worker.nprocs_clinic}
3 tropic_distribution_type = ${worker.tropic_distribution_type}
4 clinic_distribution_type = ${worker.clinic_distribution_type}
5 tropic_distribution_type = ${worker.tropic_distribution_type}
6 nprocs_clinic = ${worker.nprocs_clinic}
7 ns_boundary_type = "closed"
8 distribution_file = ${worker.distribution_file}
9
10
11 &run
12 num_tasks=1
13 redirect_stdout = true
14 log_filename = ${generated.log}
15 use_pointer_files = true
16 pointer_filename = ${generated.experimentDir}/pointer
17
18 &time_manager.nml
19 runid = ${generated.runID}
20 stop_offset = "mmmonth"
21 stop_count = 1
22 time_min_opt = "avgfit"
23 fit_freq = 1
24 fit_req = 17
25 dt_option = "steps_per_day"
26 dt_count = 23
27 dt_min = 23
28 laccel = false
29 accel_file = "unknown_accel_file"
30
31 year=0
32 allow_leapyear = false
33 year0=0
34 hour0=1
35 iday0=1
36 hour0=0
37 min0=0
38 second0=0
39 date_separator = ""

OK Check Cancel
```

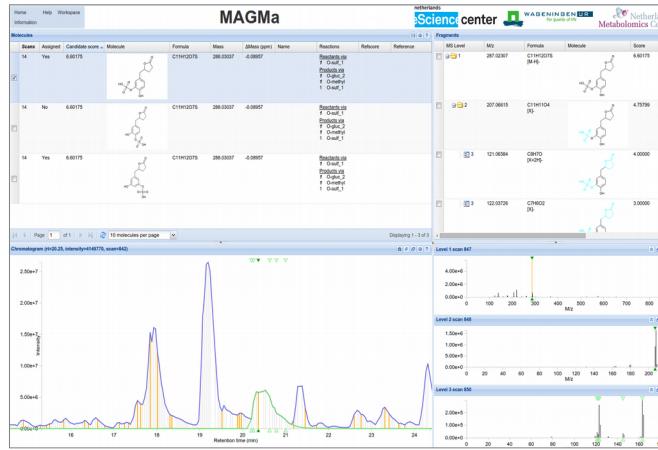
```
Deploy Log
1 2014-10-14 11:13:02 Preparing remote directories and files.
2 2014-10-14 11:13:07 Starting file transfer to STAGE_IN-
3 2014-10-14 11:13:11 File transfer completed.
4 2014-10-14 11:13:11 Remote job submitted successfully.
5 2014-10-14 11:13:11 [FileTransferService] Not copying ssh://fs0.das4.cs.vu.nl/var/scratch/jason/pop/input/small/gravity.grd.xl -> ssh://fs0.das4.cs.vu.nl/nih
6 2014-10-14 11:13:11 [FileTransferService] Not copying ssh://fs0.das4.cs.vu.nl/var/scratch/jason/pop/input/small/gravity/graphics.xl -> ssh://fs0.das4.cs.vu.nl/nih
7 2014-10-14 11:13:11 [FileTransferService] Not copying ssh://fs0.das4.cs.vu.nl/var/scratch/jason/pop/input/small/config/transport_contents.xls -> ssh://fs0.das4.cs.vu.nl/nih
8 2014-10-14 11:13:11 [FileTransferService] Not copying ssh://fs0.das4.cs.vu.nl/var/scratch/jason/pop/input/small/config/twist_angles.xls -> ssh://fs0.das4.cs.vu.nl/nih
9 2014-10-14 11:13:11 [FileTransferService] Not copying ssh://fs0.das4.cs.vu.nl/var/scratch/jason/pop/input/small/config/twist_angles.xls -> ssh://fs0.das4.cs.vu.nl/nih
10 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/monitor.sh -> ssh://fs0.das4.cs.vu.nl/nih
11 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/monitor.sh -> ssh://fs0.das4.cs.vu.nl/nih
12 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/pop -> ssh://fs0.das4.cs.vu.nl/nih
13 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/pop -> ssh://fs0.das4.cs.vu.nl/nih
14 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/pop -> ssh://fs0.das4.cs.vu.nl/nih
15 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/stop.sh -> ssh://fs0.das4.cs.vu.nl/nih
16 2014-10-14 11:13:11 [FileTransferService] Copied ssh://fs0.das4.cs.vu.nl/home/jason/esalsa/templates/pop_1d_c_32c/stop.sh -> ssh://fs0.das4.cs.vu.nl/nih
17 2014-10-14 11:13:11 [FileTransferService] Remote directories and files are ready.
18 2014-10-14 11:13:11 Changed state to STAGE_IN_COMPLETE.
19 2014-10-14 11:13:11 Remote directories and files are ready.
20 2014-10-14 11:13:11 Submitting remote job.
21 2014-10-14 11:13:18 Reading output of start script.
22 2014-10-14 11:13:18 stderr:
23 2014-10-14 11:13:18 stdout: OK 4823602 Your job 4823602 ("pop.benchx1") has been submitted
24 2014-10-14 11:13:18 Job 4823602 has been submitted. JobID: 4823602
25 2014-10-14 11:13:18 Remote job submitted successfully.
26 2014-10-14 11:13:24 Reading output of monitor script.
27 2014-10-14 11:13:24 stderr:
28 2014-10-14 11:13:24 Job 4823602 is OK RUNNING 4823602.055617 pop.benchx jason r 10/14/2014 11:13:13 all.q@node007.cm.cluster 32
29 2014-10-14 11:13:24 Changed state to RUNNING .
30 2014-10-14 11:14:24 Reading output of monitor script.
31 2014-10-14 11:14:24 stderr:
32 2014-10-14 11:14:24 Job 4823602 is OK RUNNING 4823602.055617 pop.benchx jason r 10/14/2014 11:13:13 all.q@node007.cm.cluster 32
33 2014-10-14 11:15:25 Reading output of monitor script.
34 2014-10-14 11:15:25 stderr:
35 2014-10-14 11:15:25 Job 4823602 is OK RUNNING 4823602.055617 pop.benchx jason r 10/14/2014 11:13:13 all.q@node007.cm.cluster 32
36 2014-10-14 11:16:25 Reading output of monitor script.
37 2014-10-14 11:16:25 stderr:
38 2014-10-14 11:16:25 Job 4823602 is OK RUNNING 4823602.055617 pop.benchx jason r 10/14/2014 11:13:13 all.q@node007.cm.cluster 32
39 2014-10-14 11:17:25 Reading output of monitor script.
40 2014-10-14 11:17:25 stderr:
```

Edit machine descriptions, configuration and experiment templates, remotely start experiments, inspect log files, etc...

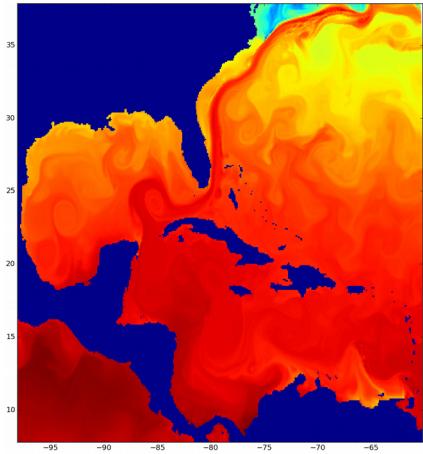
Compute at the touch of a button



SIMCITY



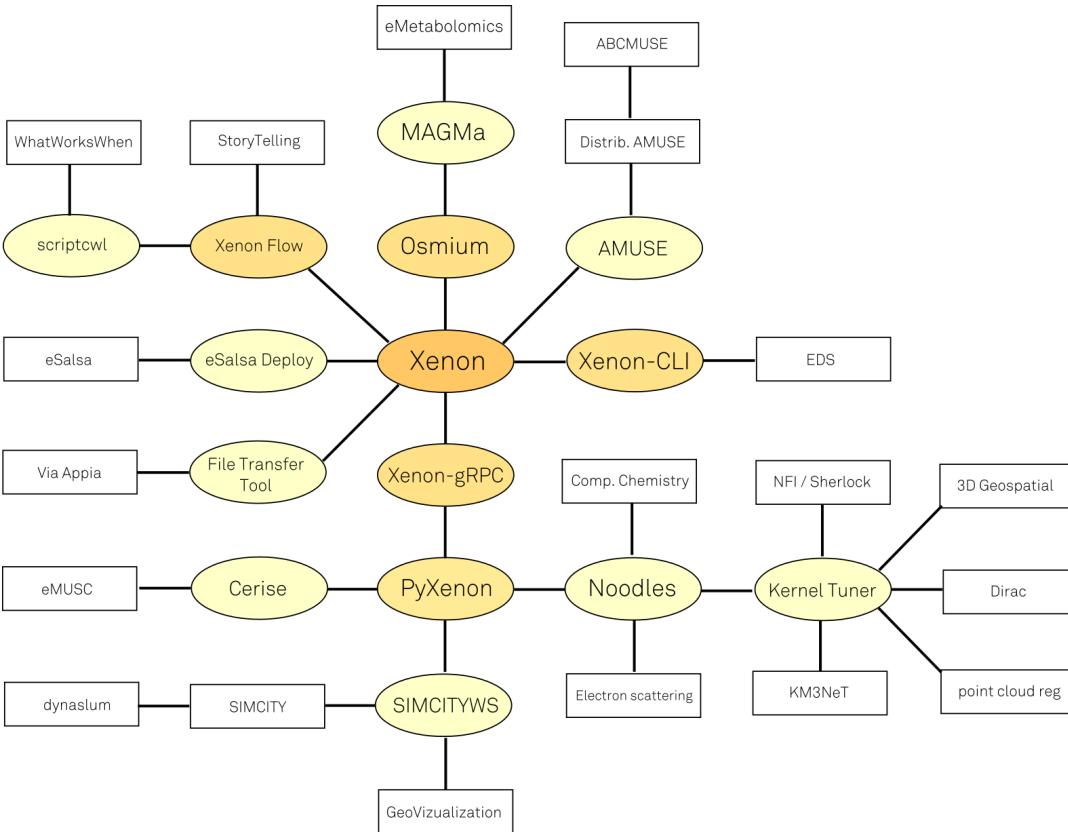
MAGMa



OMUSE

This “deployment and coupling” turned out to be a recurring theme in many of our projects:
easy access to compute and storage and babysitting applications

Xenon and friends



Xenon is a software library that provides easy access to compute and storage.

It is currently used in 14 other tools, which are used in 19 different projects.

This makes generic tools like Xenon more sustainable, as we spread their development efforts over many different projects.

Descendant from JavaGAT developed here at the VU!

accelerators

Accelerators

Accelerators (such as GPUs) offer a **huge** increase in compute power.



NVidia P100 PCI version

3584 cores, 4.7 Tflop/s

300 Watt

Fastest Nvidia GPU available in 2017

\$6000



ASCI Red

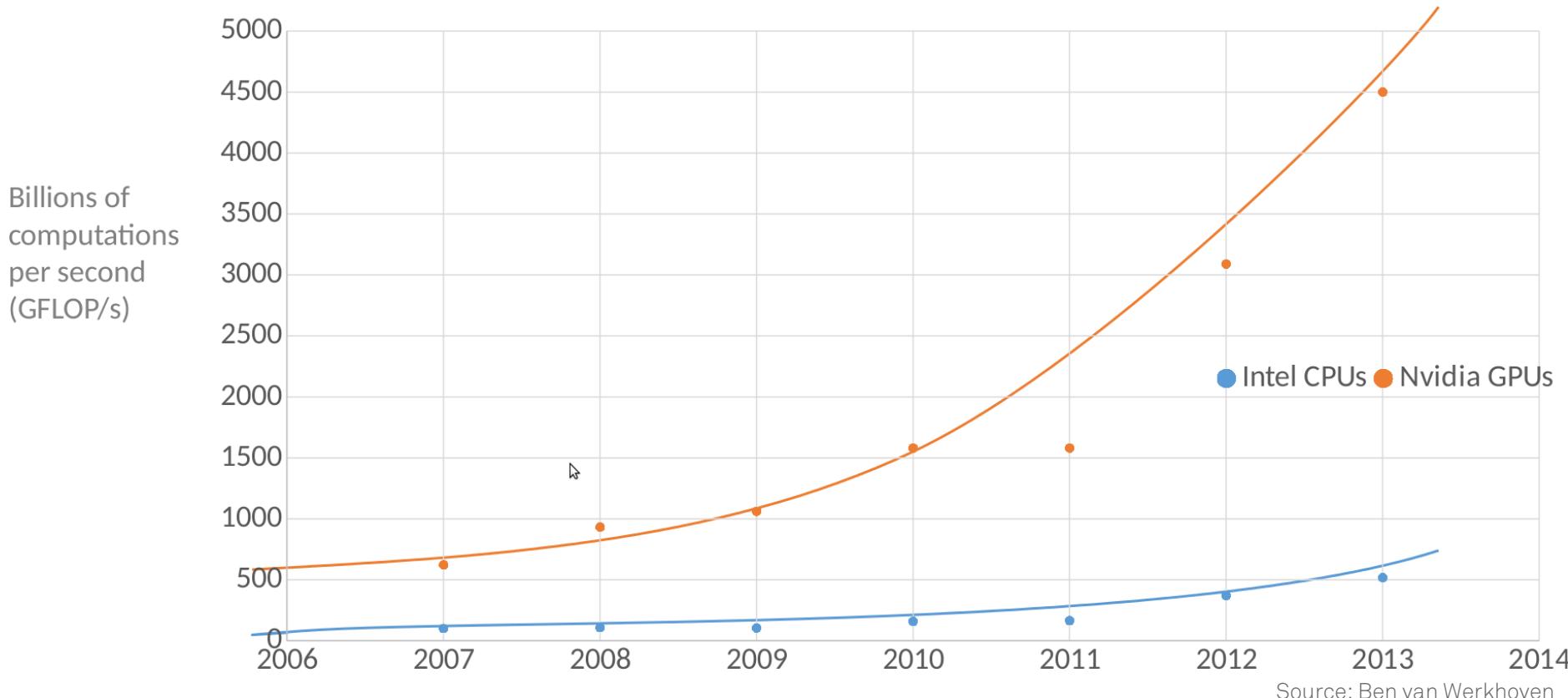
9632 cores, 2.4 Tflop/s

850.000 Watt

Fastest super computer in the world 1997-2000

\$46.000.000

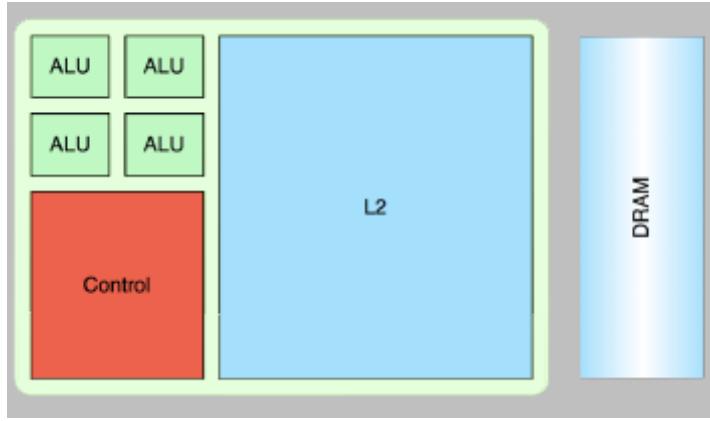
Performance of GPUs and CPUs



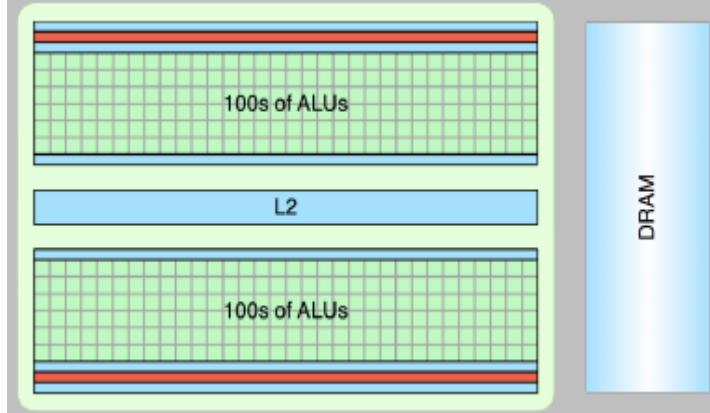
GPUs

Massively parallel: work is divided across thousands of extremely light-weight threads (usually one thread per element)

Functions that run on the GPU are called kernels



Typical CPU vs GPU design

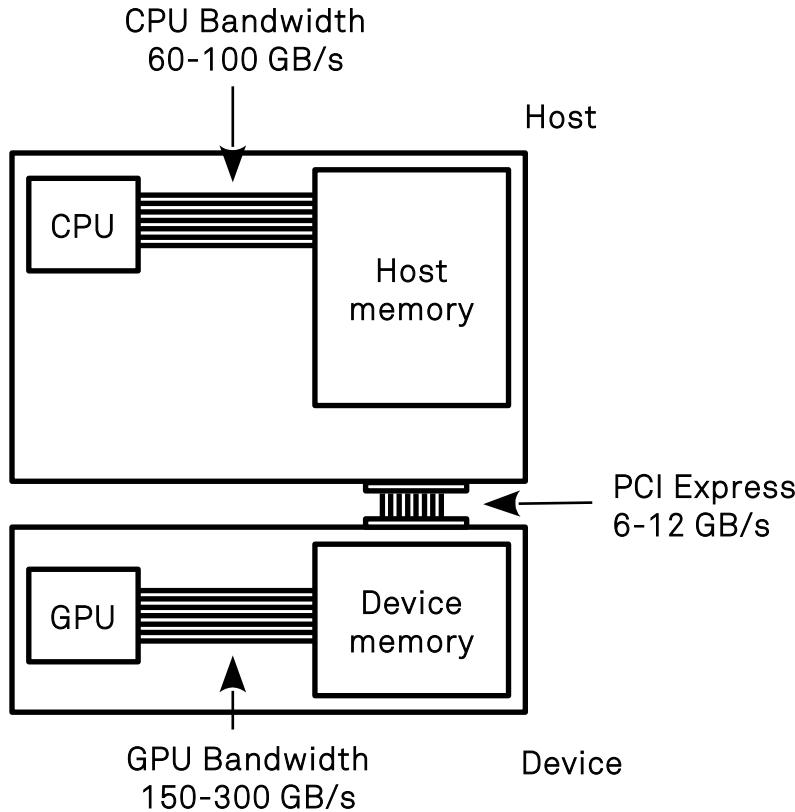


GPUs

Requires special programming languages
(CUDA, OpenCL)

Separate device in computer

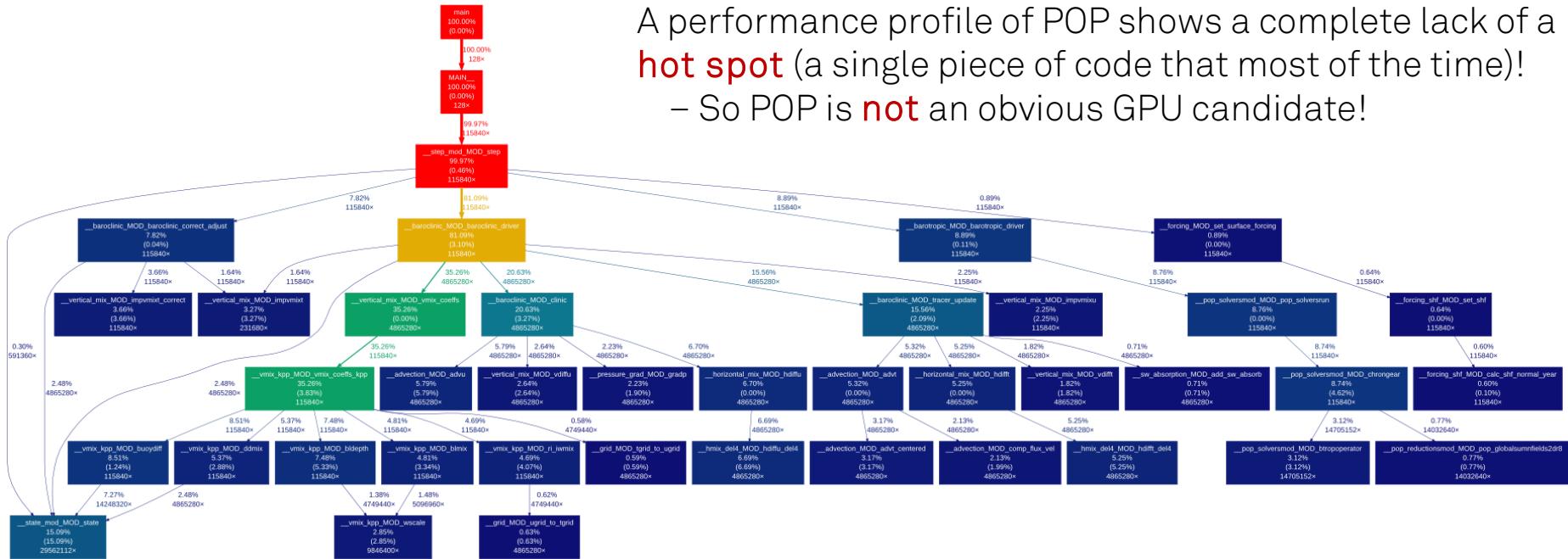
Transferring data from the CPU memory to GPU memory is part of the application and can be a major bottleneck!



Source: Ben van Werkhoven

Is POP suitable for GPUs ?

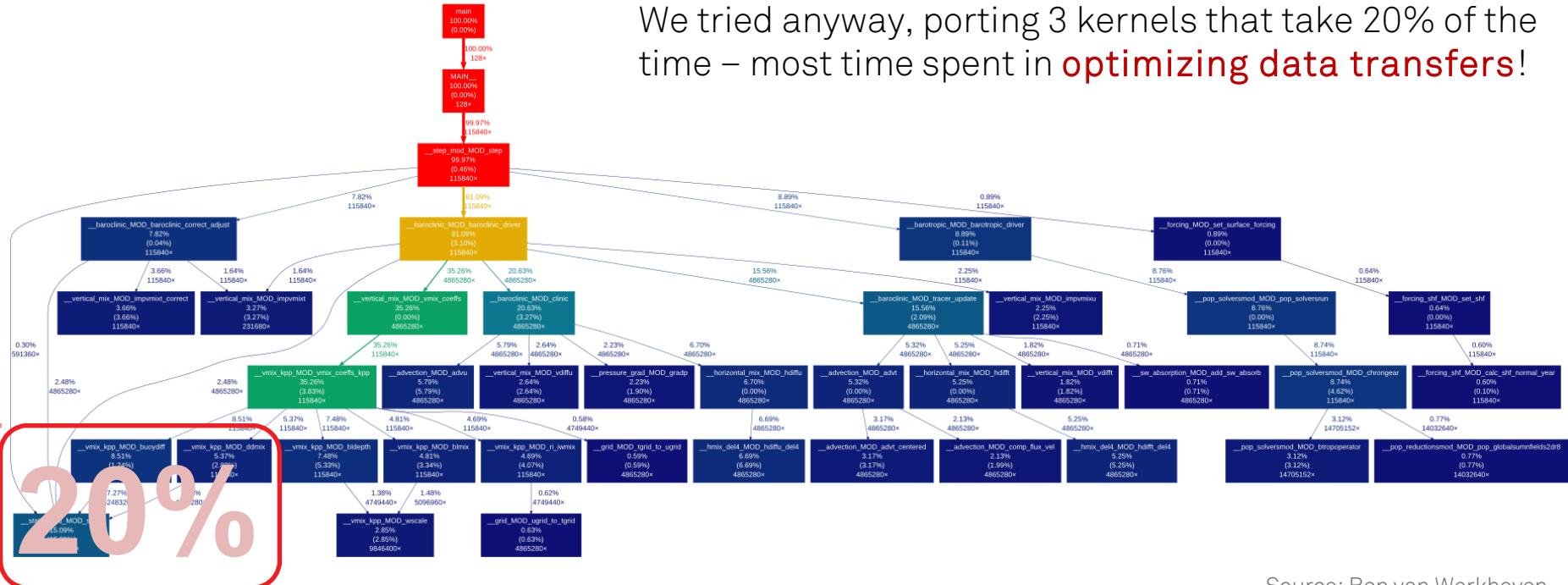
A performance profile of POP shows a complete lack of a **hot spot** (a single piece of code that most of the time)!
– So POP is **not** an obvious GPU candidate!



Source: Ben van Werkhoven

Is POP suitable for GPUs ?

We tried anyway, porting 3 kernels that take 20% of the time – most time spent in **optimizing data transfers!**



Ported 20% of code to CUDA

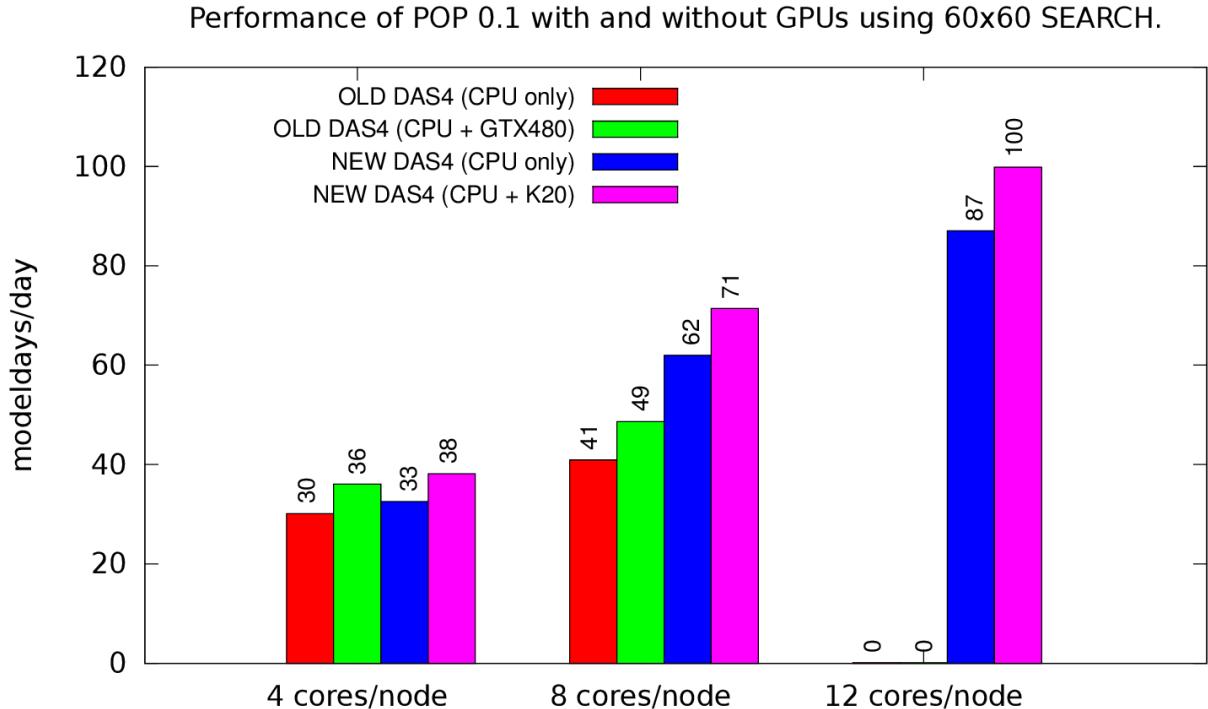
Source: Ben van Werkhoven

Initial results

Porting **20%** of POP to GPUs, improves performance with **13%**

Sharing 1 GPU with 12 CPU cores is not a problem: the GPU has flops to spare!

CPU-GPU data transfer is the bottleneck, not compute!



Source: A Distributed Approach to Improve the Performance of the Parallel Ocean Program
Ben van Werkhoven et al., Geoscientific Model Development, 7, 257-281, 2014

Performance Models

This research also lead to the paper:

Performance Models for CPU-GPU data transfers,
Ben van Werkhoven et al, CCGrid 2014
(best paper nominee)

Source: Performance Models for CPU – GPU Data Transfers, Ben van Werkhoven et al.,
14th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID),
best paper nominee, Chicago, IL, May 2014.

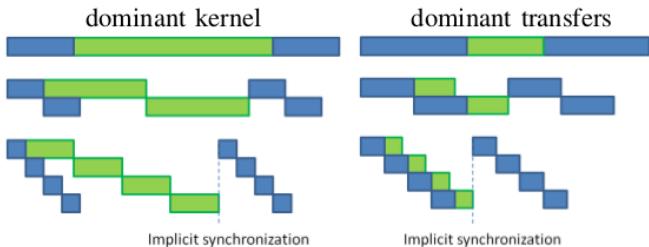


Fig. 1. Diagrams showing possible overlap for devices with implicit synchronization and 1 copy engine when using 1, 2, or 4 streams.

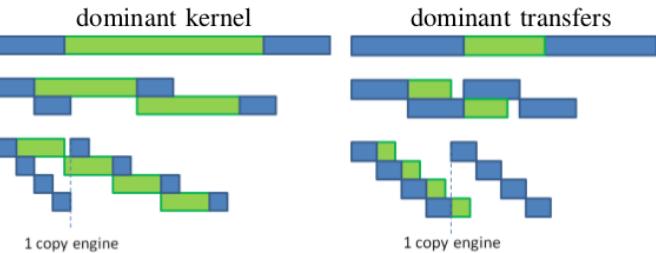


Fig. 2. Diagrams showing possible overlap for devices with no implicit synchronization and 1 copy engine when using 1, 2, or 4 streams.

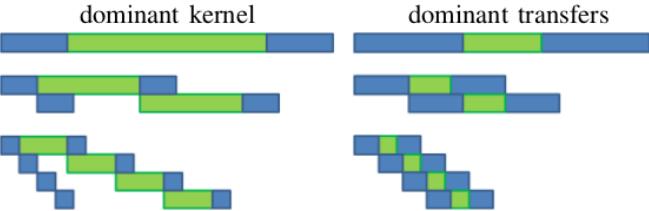


Fig. 3. Diagrams showing possible overlap for devices with no implicit synchronization and 2 copy engines when using 1, 2, or 4 streams.

Performance Models

Ben's performance models can be used to estimate the performance of each of the data transfer options.

Implicit synchronization and 1 copy engine.

Dominant kernel

$$time = L + o + \frac{B_{hd}}{nStreams} \times G_{hd} + t_E + L + o + B_{dh} \times G_{dh} + g \times (nStreams - 1) \quad (1)$$

Dominant transfers

$$time = L + o + B_{hd} \times G_{hd} + g \times (nStreams - 1) + \frac{t_E}{nStreams} + L + o + B_{dh} \times G_{dh} + g \times (nStreams - 1) \quad (2)$$

No implicit synchronization and 1 copy engine.

Dominant kernel

$$time = L + o + \frac{B_{hd}}{nStreams} \times G_{hd} + t_E + L + o + \frac{B_{dh}}{nStreams} \times G_{dh} \quad (3)$$

Dominant transfers

$$\begin{aligned} time &= \max(L + o + B_{hd} \times G_{hd} + g \times (nStreams - 1) + L + o + B_{dh} \times G_{dh} + g \times (nStreams - 1), \\ &\quad L + o + B_{hd} \times G_{hd} + g \times (nStreams - 1) + \frac{t_E}{nStreams} + L + o + \frac{B_{dh}}{nStreams} \times G_{dh}, \\ &\quad L + o + \frac{B_{hd}}{nStreams} \times G_{hd} + \frac{t_E}{nStreams} + L + o + B_{dh} \times G_{dh} + g \times (nStreams - 1)) \end{aligned} \quad (4)$$

No implicit synchronization and 2 copy engines.

$$\begin{aligned} time &= \max(L + o + B_{hd} \times G_{hd} + g \times (nStreams - 1) + \frac{t_E}{nStreams} + L + o + \frac{B_{dh}}{nStreams} \times G_{dh}, \\ &\quad L + o + \frac{B_{hd}}{nStreams} \times G_{hd} + t_E + L + o + \frac{B_{dh}}{nStreams} \times G_{dh}, \\ &\quad L + o + \frac{B_{hd}}{nStreams} \times G_{hd} + \frac{t_E}{nStreams} + L + o + B_{dh} \times G_{dh} + g \times (nStreams - 1)) \end{aligned} \quad (5)$$

Device-mapped host memory.

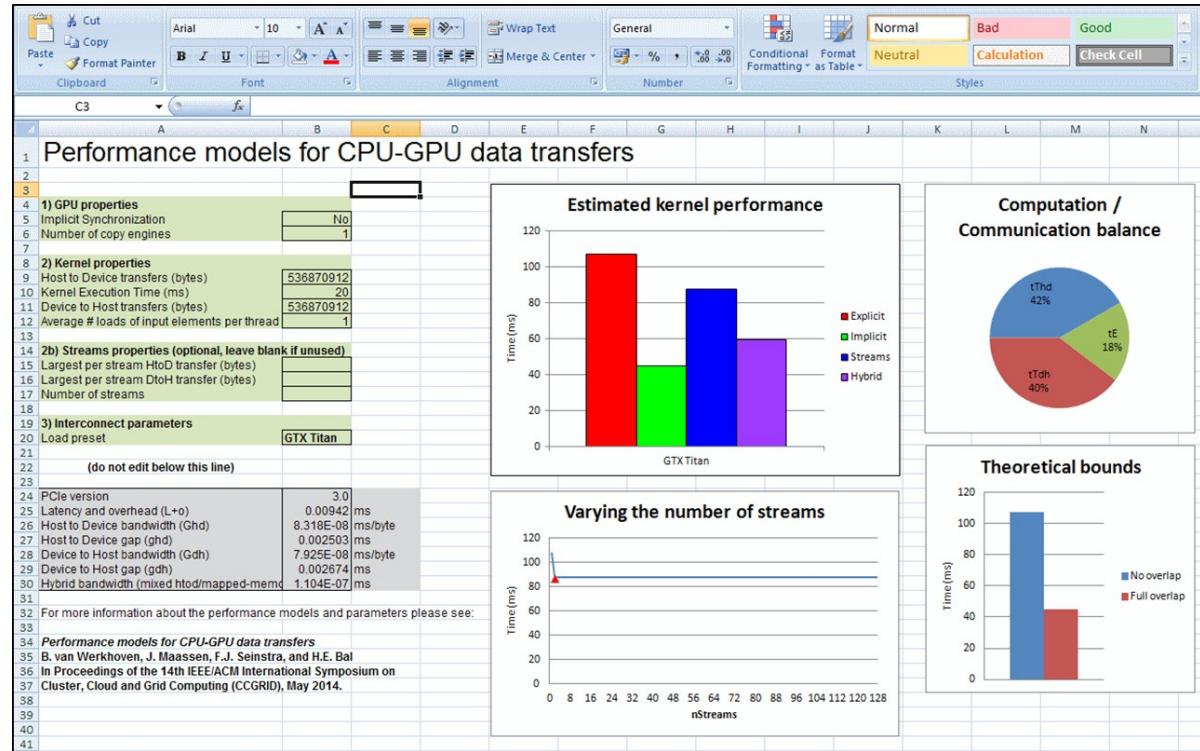
$$time = \max(L + o + B_{hd} \times G_{hd} + L + o, L + o + t_E + L + o, L + o + L + o + B_{dh} \times G_{dh}) \quad (6)$$

Fig. 7. Overview of performance models. B_{hd} is the number of bytes transferred from host to device. B_{dh} is the number of bytes transferred from device to host.

Source: Ben van Werkhoven

Performance Models

Using a simple excel sheet
You can test what the
optimal transfer method
is for your code!



Source: Ben van Werkhoven

Performance Models

Results show that the performance estimates are quite accurate.

Estimates for **implicit** implementation are most accurate (max error 3.85%).

Estimates for **hybrid** implementation have largest error (max error 10.75%).

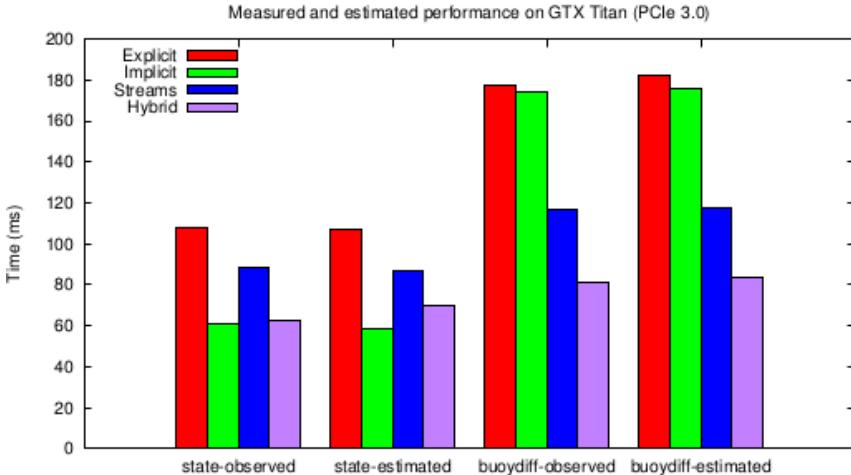


Fig. 9. Performance of state and buoydiff kernels on GTX Titan.

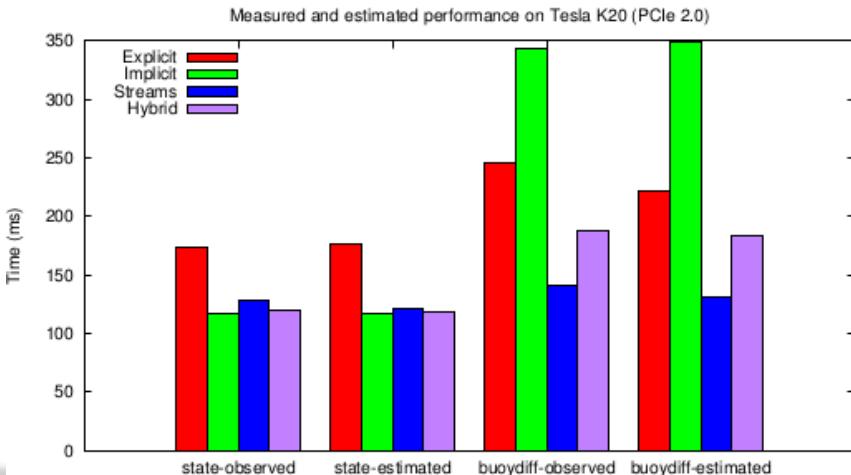


Fig. 10. Performance of state and buoydiff kernels on Tesla K20.

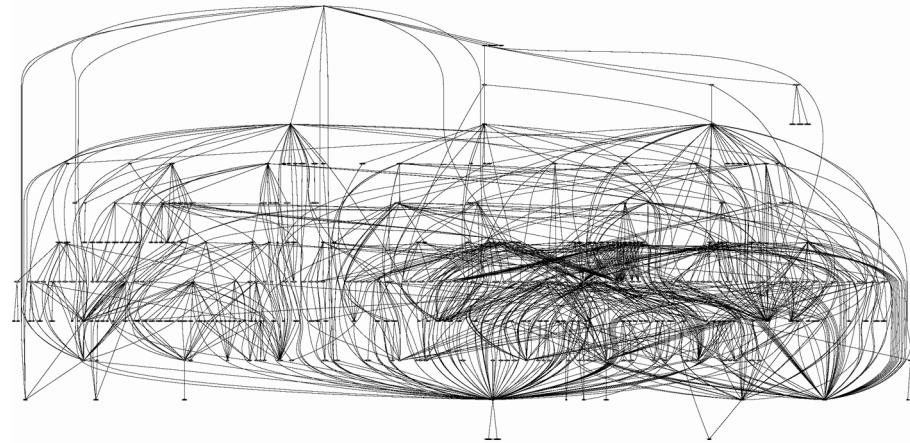
The joy of porting...

Porting 20% of POP to CUDA is hard!

It is hard get insight into the inner workings of a Fortran application that has been worked on for 25 years!

In addition, a lot of 'boilerplate code' is required to connect the Fortran and CUDA parts.

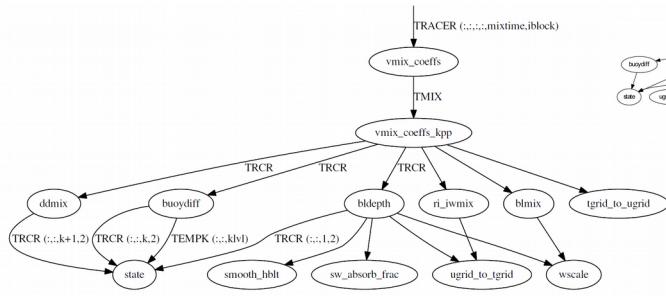
We need more powerful tools!



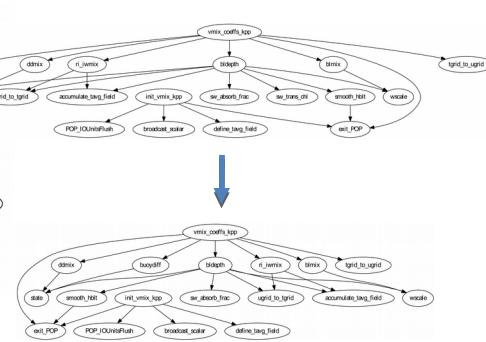
Callgraph of the POP

Marver: semi-automatic translation

source code analysis



transformation



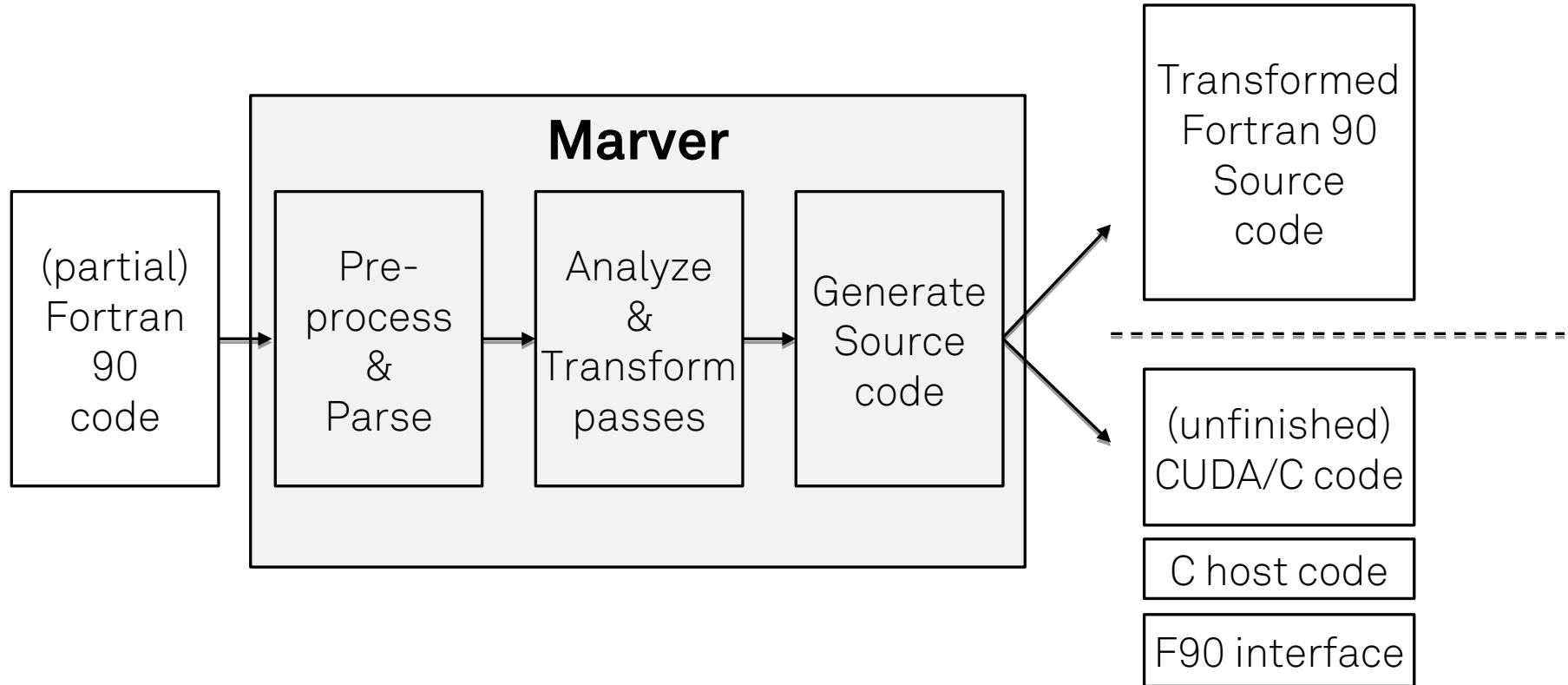
source-to-source translation

$VVC(:,:,k) = \text{merge}(\text{WORK2}, c0, (k < \text{KMU}(:,:,bid)))$

$VVC(i,j,k) = ((k < \text{KMU}(i,j,bid)) ? \text{WORK2} : c0);$

Source: Ben van Werkhoven

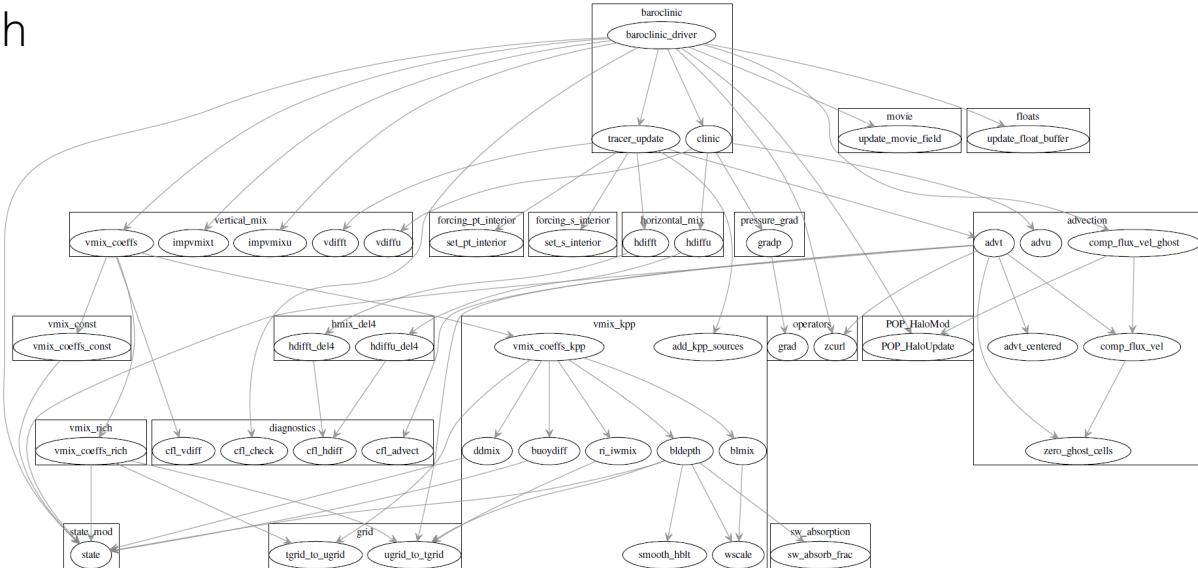
Marver workflow



Source: Ben van Werkhoven

Marver is semi-automatic

Marver generates a call graph

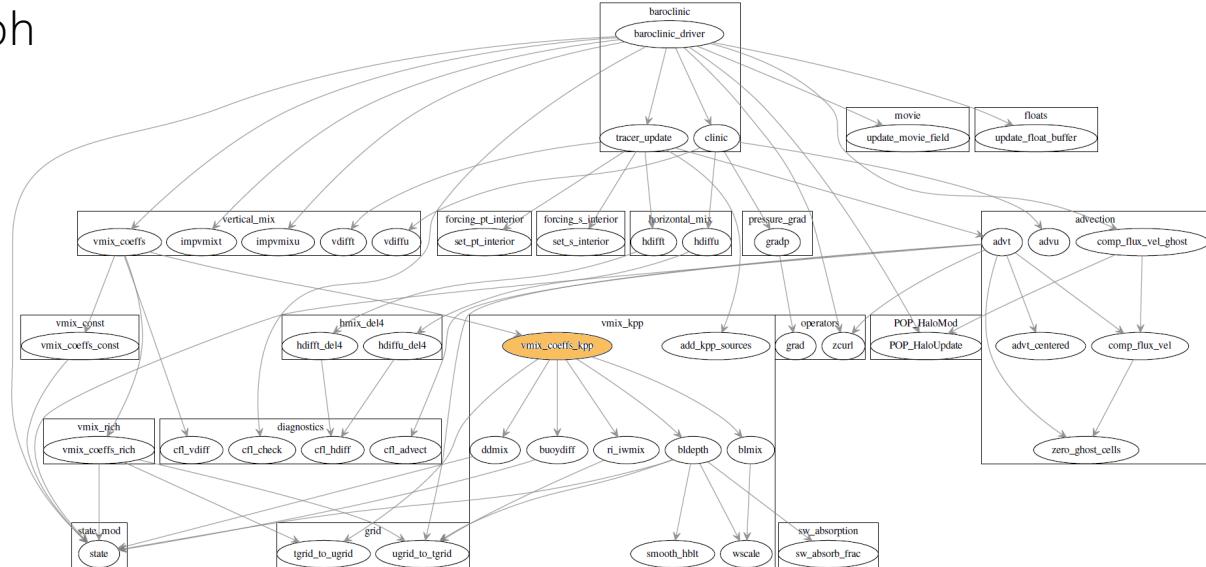


Source: Ben van Werkhoven

Marver is semi-automatic

Marver generates a call graph

User needs to identify
entry point in call-graph
and parallelism



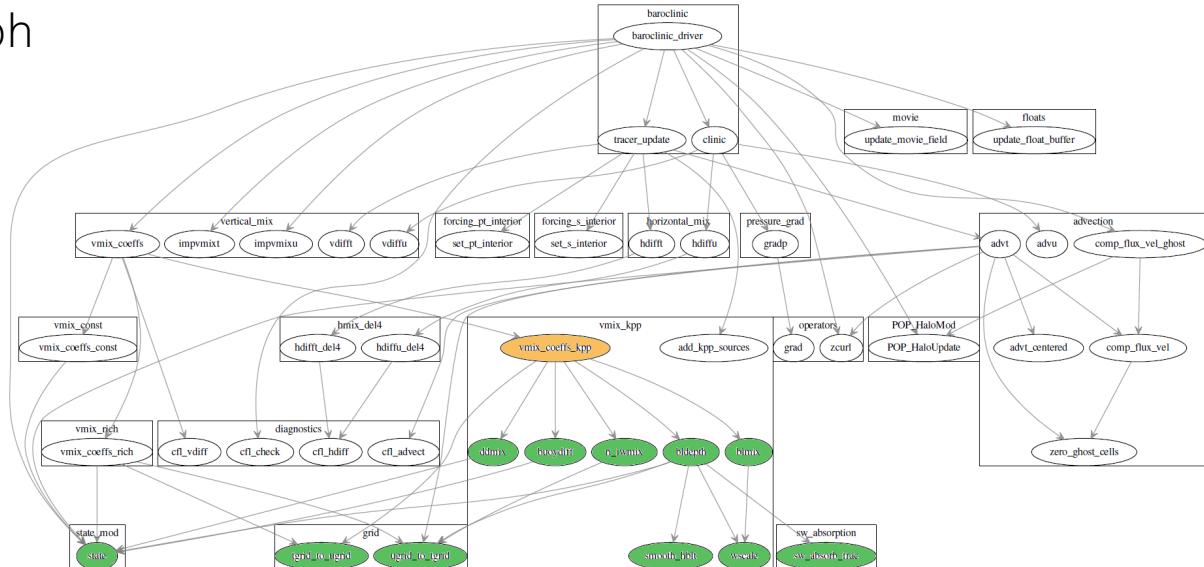
Source: Ben van Werkhoven

Marver is semi-automatic

Marver generates a call graph

User needs to identify entry point in call-graph and parallelism

Marver traverses the call graph and generates CUDA and boilerplate code



Source: Ben van Werkhoven

Marver produces human readable code

User needs to finish the CUDA manually

It is important that Marver generates human readable CUDA code than has a clear link to the original Fortran

Fortran 90 of code:

```
do j=1,ny_block
do i=1,nx_block
  if ( k >= KBL(i,j) ) then
    VISC(i,j,k) = VISC(i,j,k) + convect_visc * FCON(i,j)
    VDC(i,j,k,1) = VDC(i,j,k,1) + convect_diff * FCON(i,j)
    VDC(i,j,k,2) = VDC(i,j,k,2) + convect_diff * FCON(i,j)
  endif
end do
end do
```

CUDA/C kernel code:

```
if ( k >= KBL(i,j)-1 ) {
  VISC(i,j,k) = VISC(i,j,k) + convect_visc * FCON;
  VDC(i,j,k,0) = VDC(i,j,k,0) + convect_diff * FCON;
  VDC(i,j,k,1) = VDC(i,j,k,1) + convect_diff * FCON;
}
```

Source: Ben van Werkhoven



Marver results

Finishing code took about 3 days, plus 5 day to fix the bugs that the user introduced.

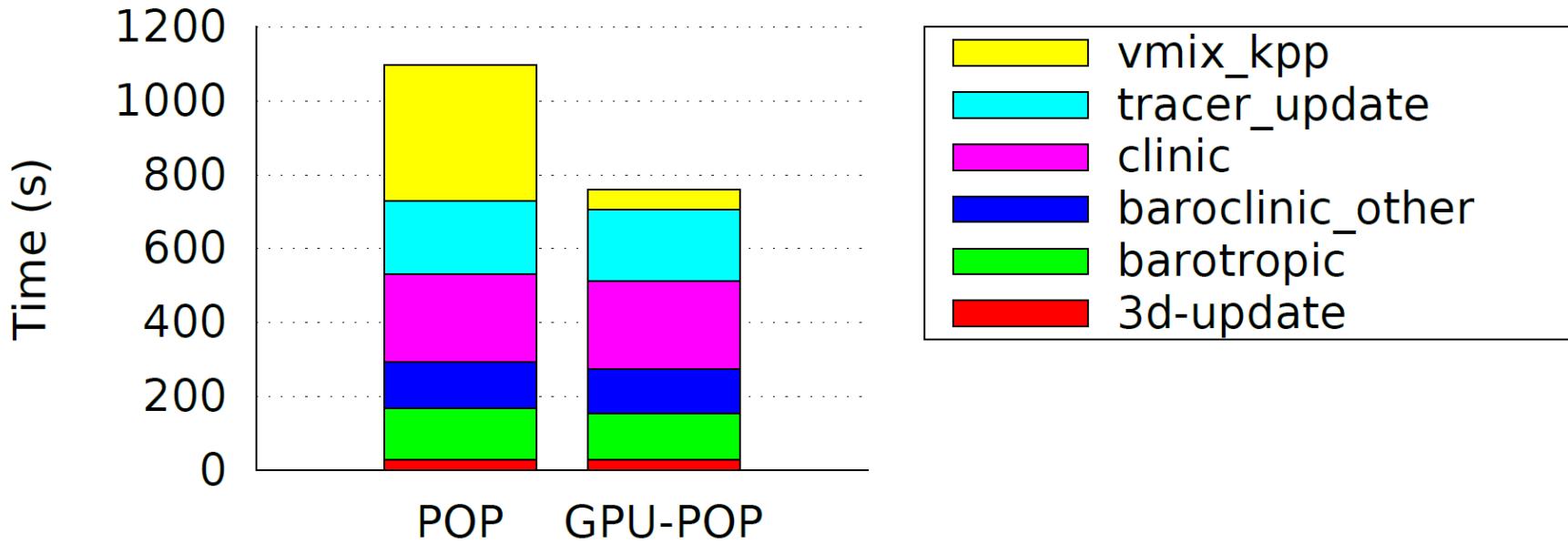
Porting effort is reduced significantly. It still requires good GPU skills, but saves a lot of the boring work!

original name	input	output	split	final
vmix_coeffs_kpp	152	232	vmix_coeffs_kpp_gpu_entry	299
buoydiff	102	84	buoydiff_kernel	71
			compute_vshearu	29
ri_iwmix	305	191	ri_iwmix_kernel	177
ddmix	129	100	ddmix_kernel	131
			compute_ustaru	15
			ugrid_to_tgrid	47
			vshearu_bldepth	21
bldepth	596	325	bldepth_kernel	297
			stability_and_buoyancy	34
blmix	434	279	blmix_kernel	280
interior_convection	231	143	interior_convection	80
			interior_part2	85
wscale	108	94	wsscale	29
			wmscale	29
smooth_hblt	182	100	smooth_hblt	100
ugrid_to_tgrid	58	41	ugrid_to_tgrid	40
tgrid_to_ugrid	58	41	tgrid_to_ugrid	29
sw_absorb_frac	75	61	sw_absorb_frac	53
state	410	131	mwjf_numerator	3
			mwjf_denominator	6
			compute_drhodt	12
			compute_drhods	9
Total	2840	1822		1876

Lines of code in the input (Fortran 90), output (CUDA/C), and final code

Source: Ben van Werkhoven

Marver results



Running POP on 8 GPU nodes of Cartesius (16 CPU cores, 2 K40 GPUs) shows a significant speedup in the yellow ported part

Source: Ben van Werkhoven

distributed computing

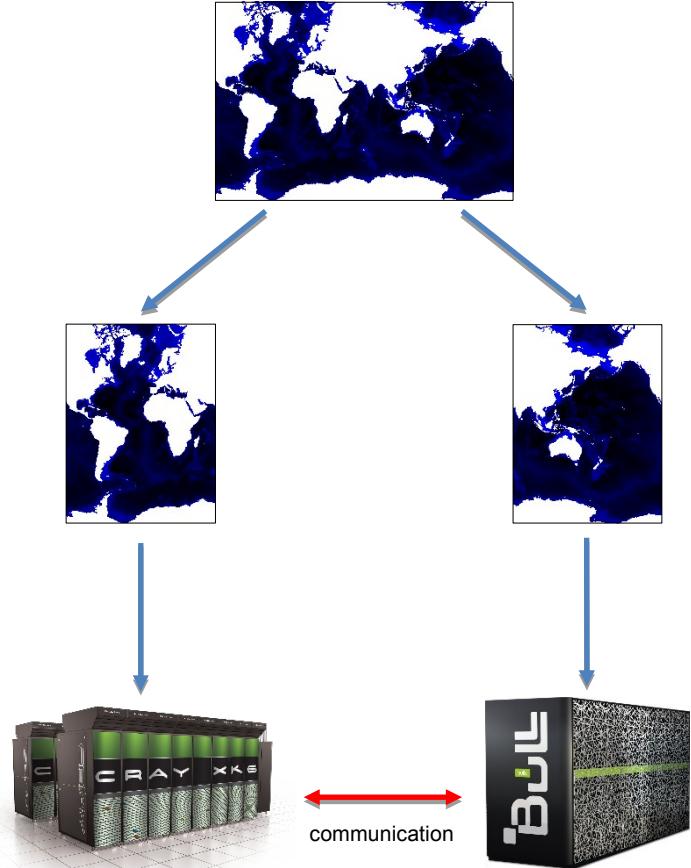
(the hard way)

Distributed Computing

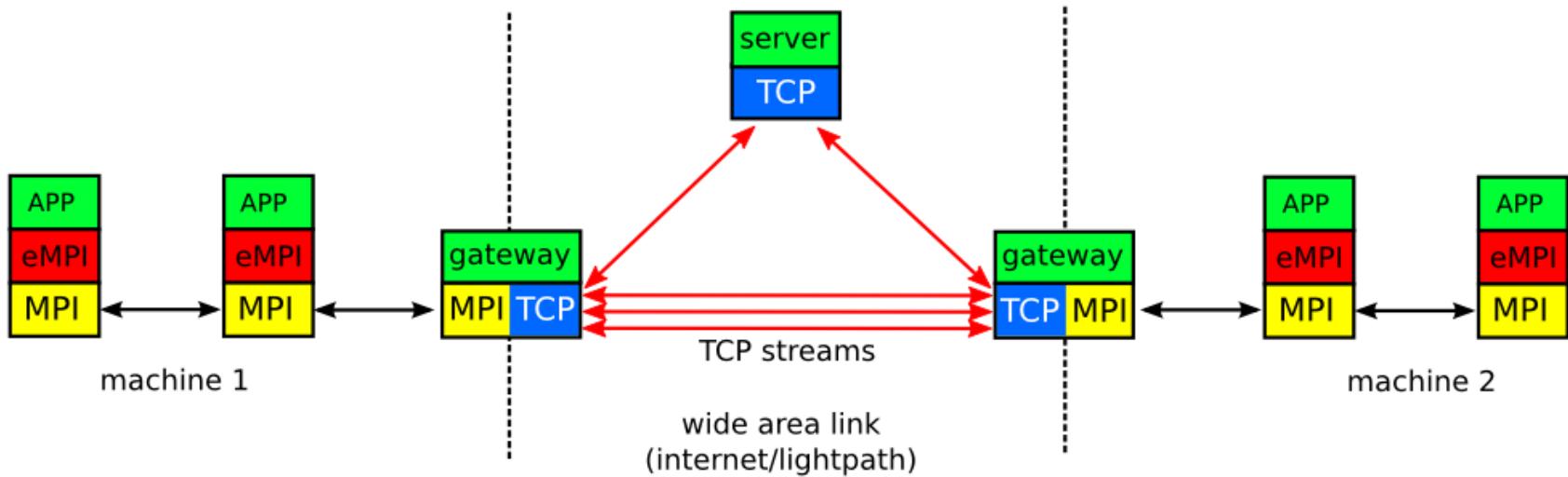
By combining supercomputers we can double / triple / ... the amount of available compute power.

Example: split POP into 2 pieces and run each piece on a different supercomputer.

However: these supercomputers need to communicate (a lot) during the model run!



eSalsa-MPI connects supercomputers



eSalsa-MPI is an **MPI wrapper library** that can combine two MPI runs on different supercomputers, and presents it as a **single set of processes** to the application.
No change to the application is needed!

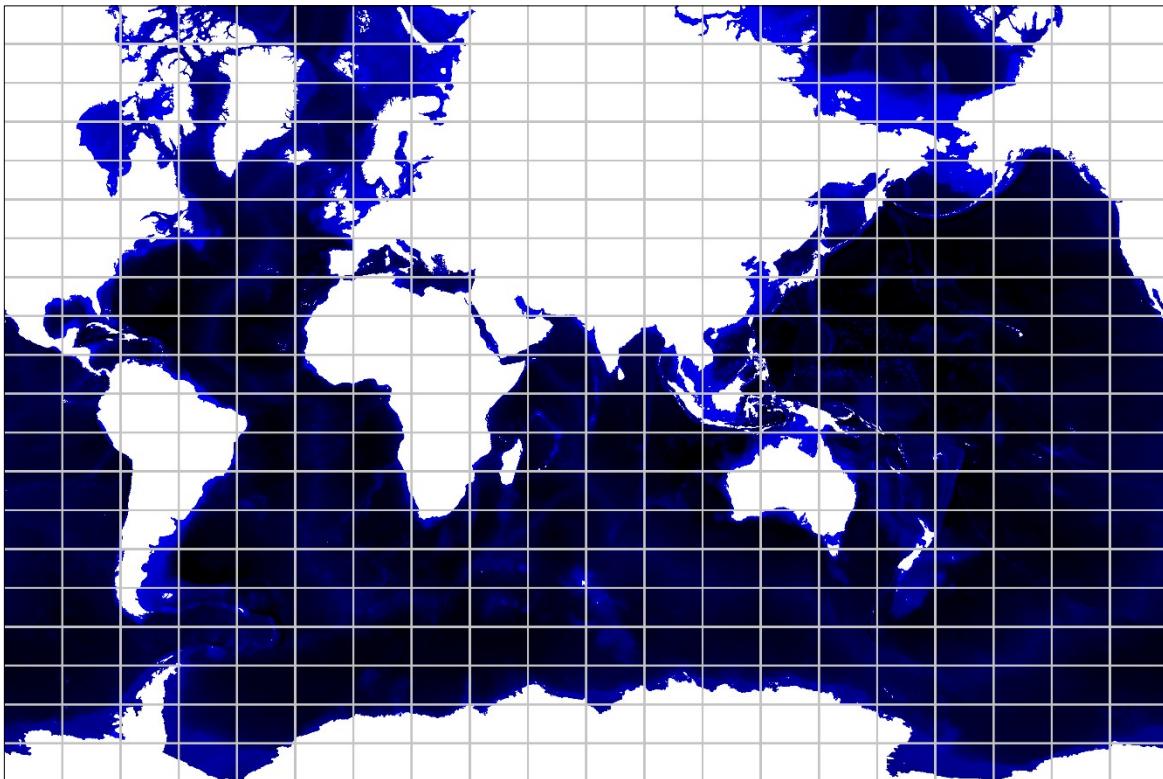
How to split POP ?

Remember:

Within one supercomputer
POP also divides the world
into blocks.

But there is a lot of
communication between
neighboring blocks!

How do you split the world in
such a way that communication
is minimized ?

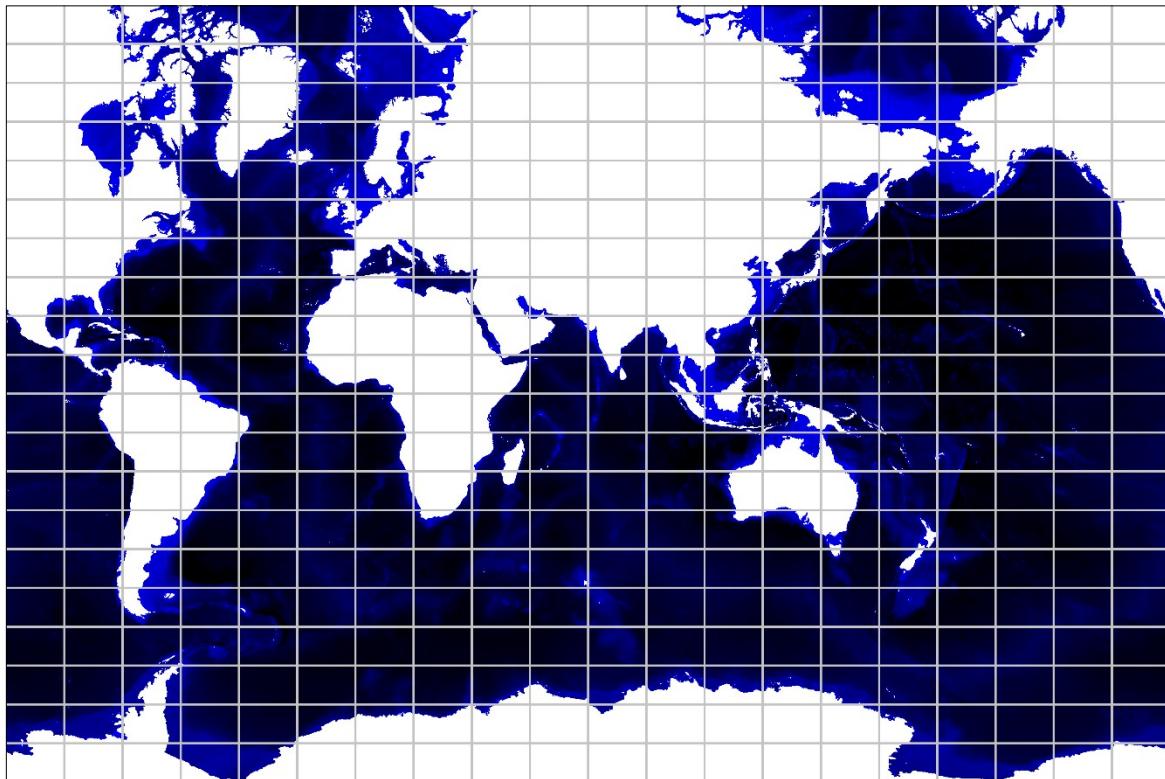


Discarding land-only blocks

Land-only blocks may be discarded.

Mixed ocean + land blocks are fully computed.

Example: 20x20 blocks (400) of 180x120x42 grid points each.



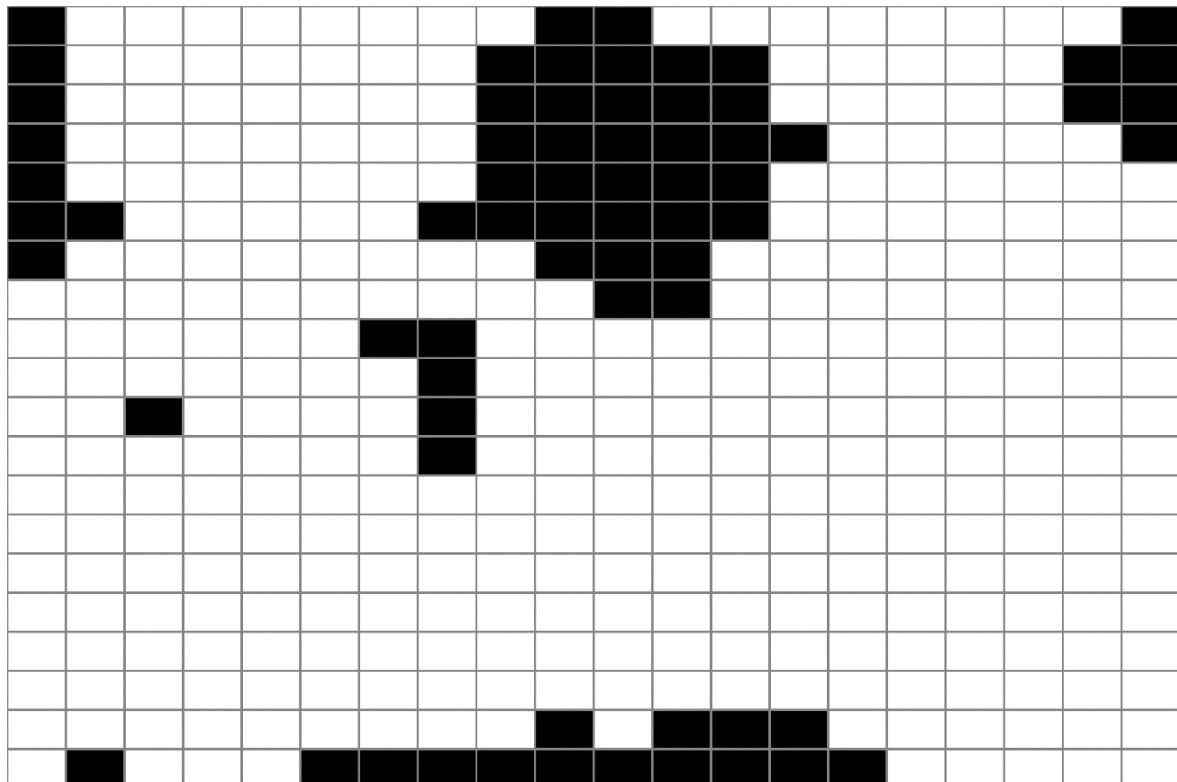
Discarding land-only blocks

Land-only blocks may be **discarded**.

Mixed ocean + land blocks are **fully computed**.

Example: 20x20 blocks (400) of 180x120x42 grid points each

Only **331** blocks are computed!



Using smaller blocks

To reduce the amount of useless computation, we can split the world into **smaller blocks**, and assign **multiple blocks** to each processor.

This is hard as it leads to **load imbalance** and extra **communication overhead!**



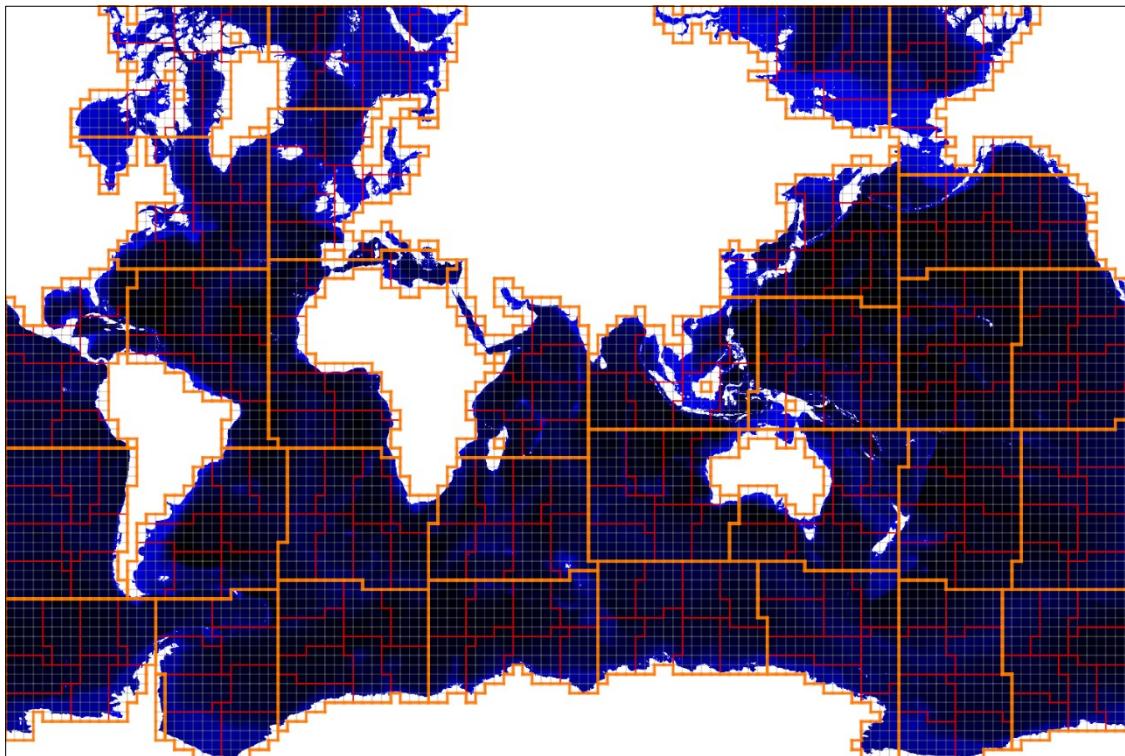
Hierarchical work distribution

To solve this problem we designed a hierarchical work distribution algorithm that recursively splits the world into blocks (**top-down**).

Example:

Split each result into 32 (nodes)
Split each result into 16 (cores)

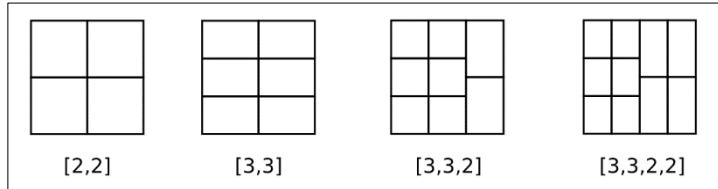
Maximize load balance and
minimize communication at each level!



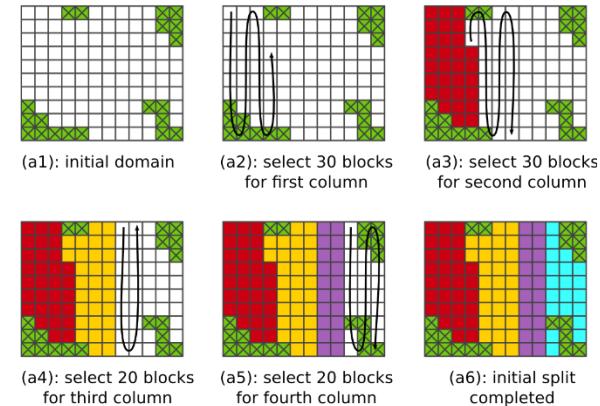
Source: A Distributed Approach to Improve the Performance of the Parallel Ocean Program
Ben van Werkhoven et al., Geoscientific Model Development, 7, 257-281, 2014

Algorithm

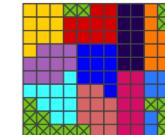
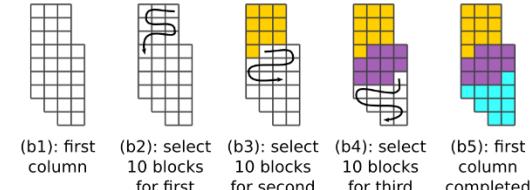
```
f := floor(sqrt(N));  
c := ceiling(sqrt(N))  
  
if (f = c) We have found a square grid of [f x f]  
if (f*c = N) We have found a rectangular grid of [f x c]  
if (N < f*c) We have found a rectangular grid of [f x c] - (f*c-N)  
if (N > f*c) We have found a square grid of [c x c] - (c*c-N)
```



Step 1: Split domain column wise



Step 2: Split column selections row-wise
(only first column selection is shown)

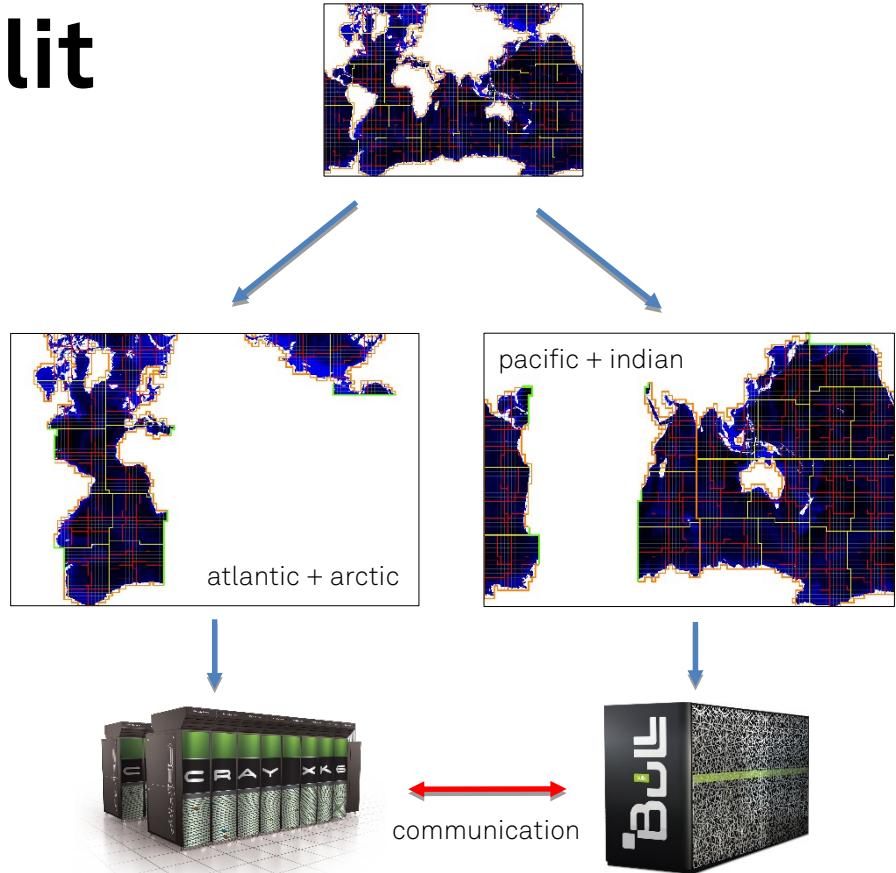


(c): final result

Geography aware split

We can easily extend the hierarchical work distribution scheme to include multiple supercomputers.

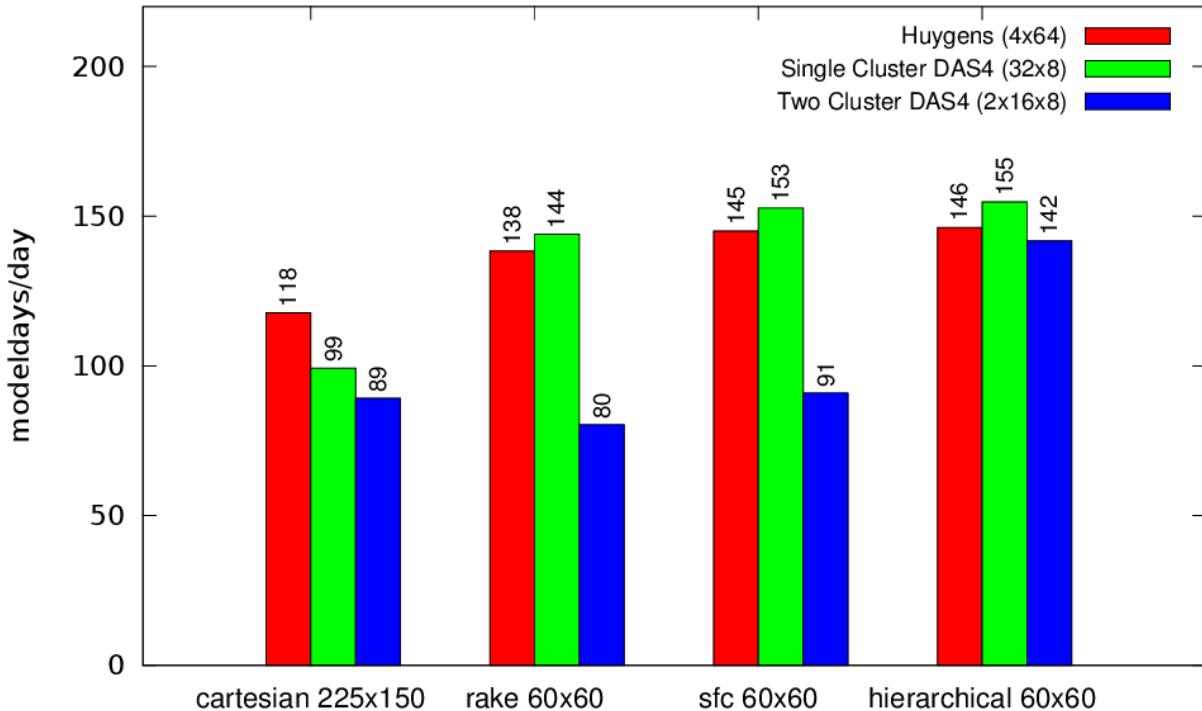
This results in a **geography aware** distribution of work, which **significantly reduces communication** between machines!



Performance of Distributed POP

Hierarchical work distribution performs well on two DAS4 cluster (256 cores) with **only** a shared 1 Gbit/s link between them.

Produces 17% to 76% less communication between sites than other distributions.



Source: A Distributed Approach to Improve the Performance of the Parallel Ocean Program
Ben van Werkhoven et al., Geoscientific Model Development, 7, 257-281, 2014

Communication volume

Table 2. Configuration of the Cartesian, rake, and space-filling curve, and hierarchical distributions.

algorithm	block size	blocks per core (min/max)	blocks discarded	communication per task (min/avg/max)	communication between clusters (messages/volume)
Cartesian	225×150	1/1	0 (of 256)	0/1267.4/2408	22.3 M/99.0 GB
rake	60×60	5/8	628 (of 2400)	748/1940.5/3936	77.9 M/337.4 GB
space-filling curve	60×60	6/7	628 (of 2400)	1007/1707.7/2960	41.0 M/212.7 GB
hierarchical	60×60	6/7	628 (of 2400)	504/1394.9/2584	20.0 M/82.5 GB

Hierarchical work distribution reduces both message count and message volume between supercomputers.



Visualization

We visualize the output of the simulations to get a quick insight into the results.



Presenting eSalsa to King Willem-Alexander

distributed computing

(the next step)

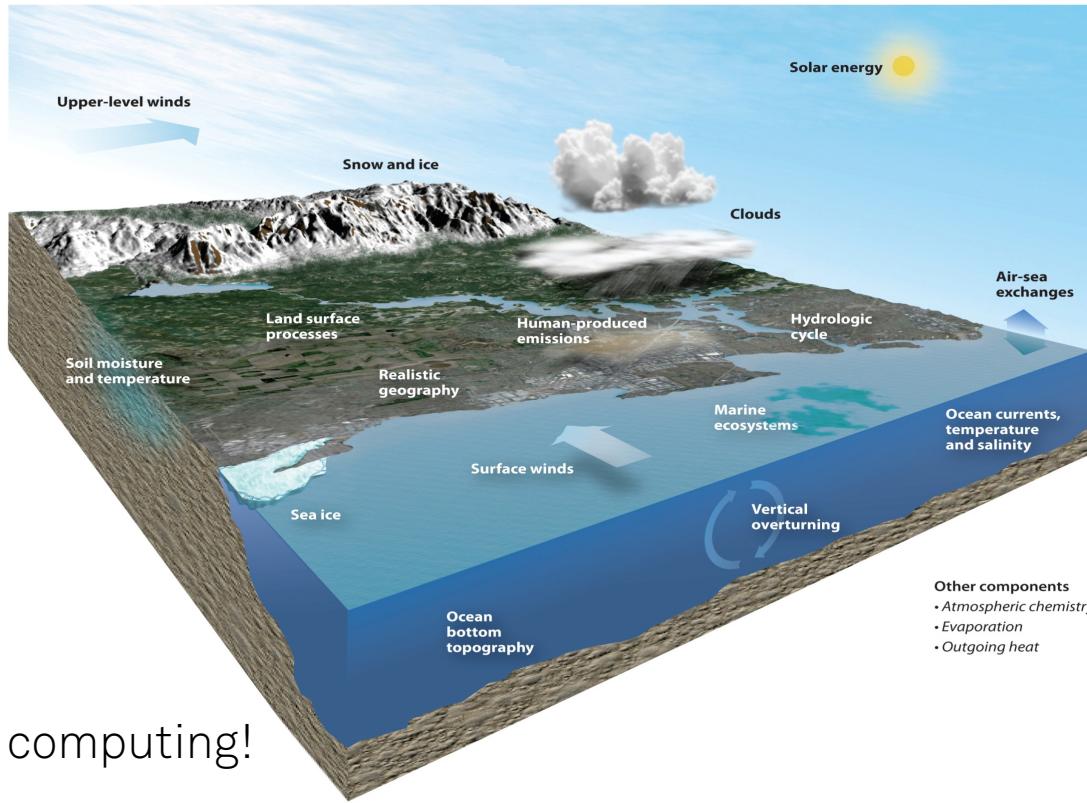
Earth System Modeling

Climate is more than just ocean!

We also want to **combine** atmosphere, ocean, land and ice models to simulate the interactions.

2~3x increase in compute time.

More opportunities for distributed computing!

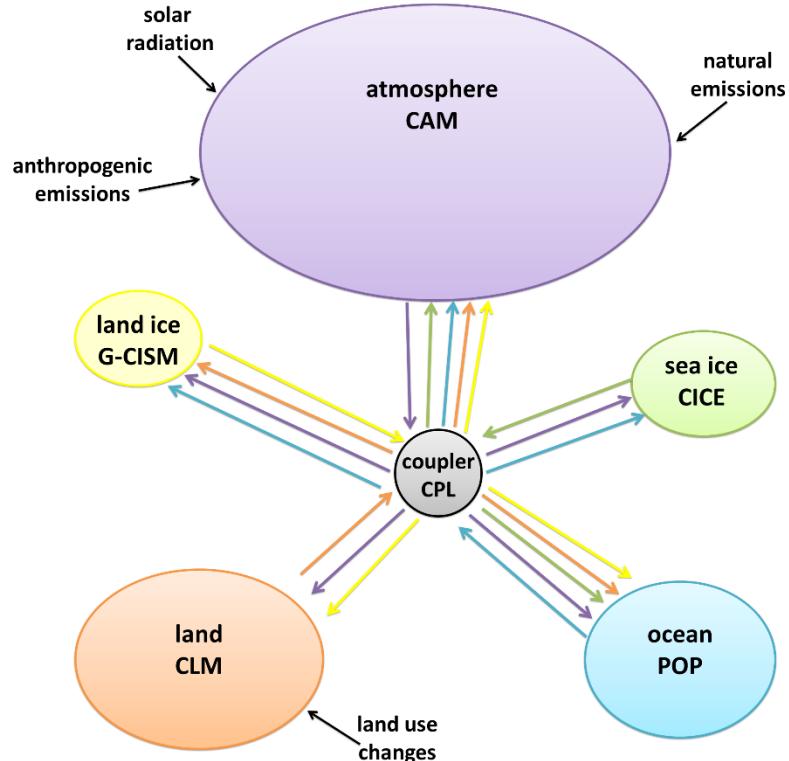


Community Earth System Model

CESM combines 5 simulations each simulating a different aspect of the climate (ocean, atmosphere, land, etc).

~1 million lines of Fortran / MPI

Models regularly **exchange data** to model the influence on each other.

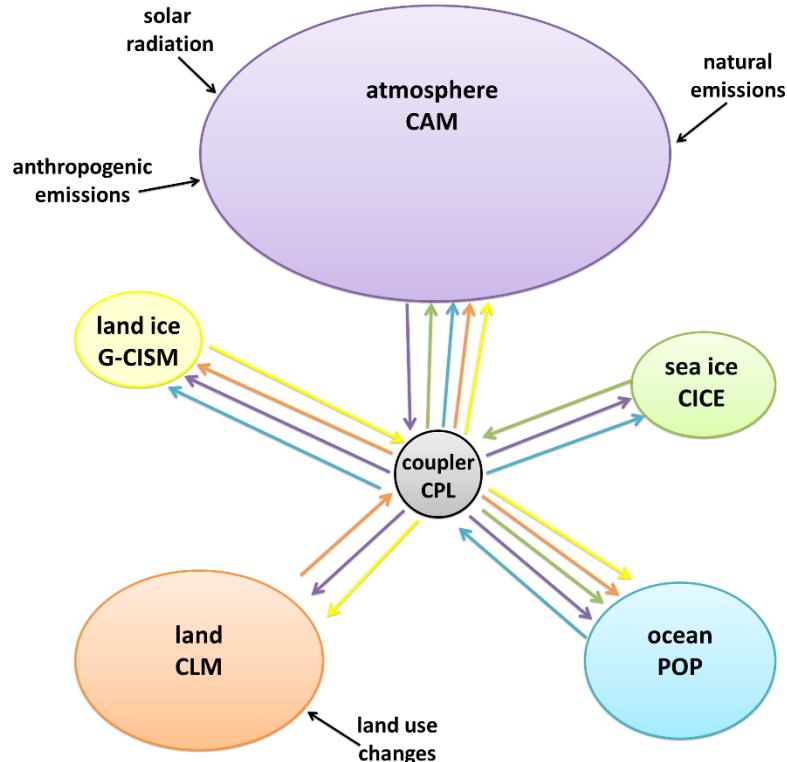


Source: climatesight.org

Community Earth System Model

Using **distributed computing** we run different submodels on different supercomputers.

This allows us to **scale** each submodel to **higher resolutions** and combine different **architectures** (for example a traditional machine + a GPU machine).

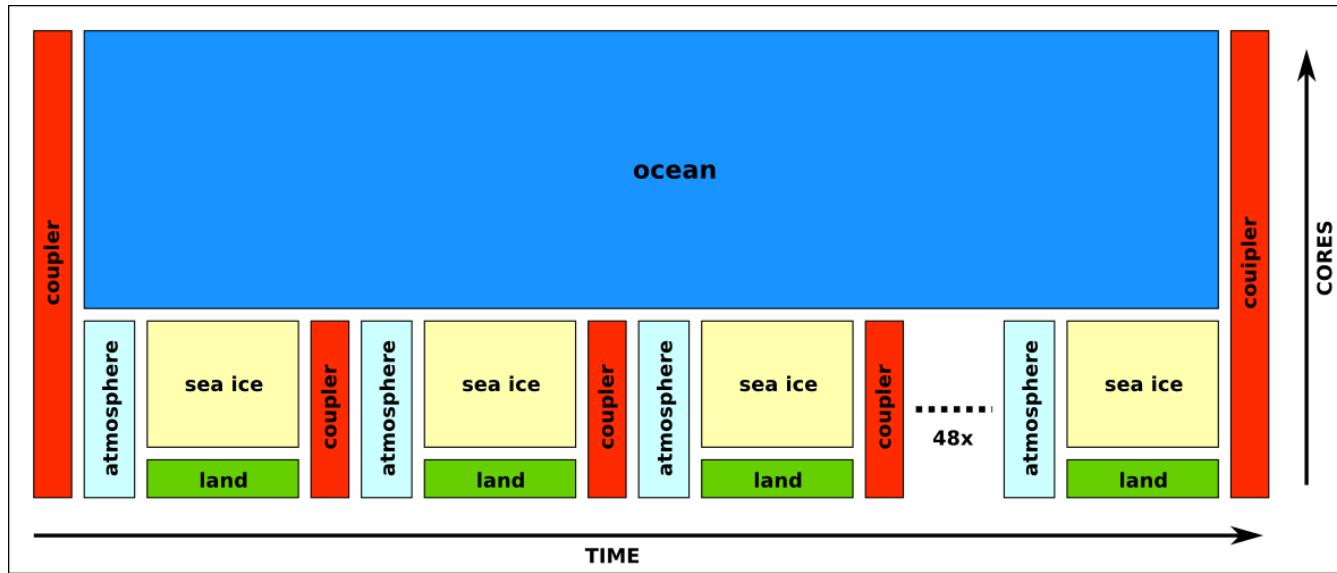


Source: climatesight.org

CESM model coupling

Data dependencies
allow some models
to run concurrently,
while others must
run in sequence.

(balancing their run times
is black magic)

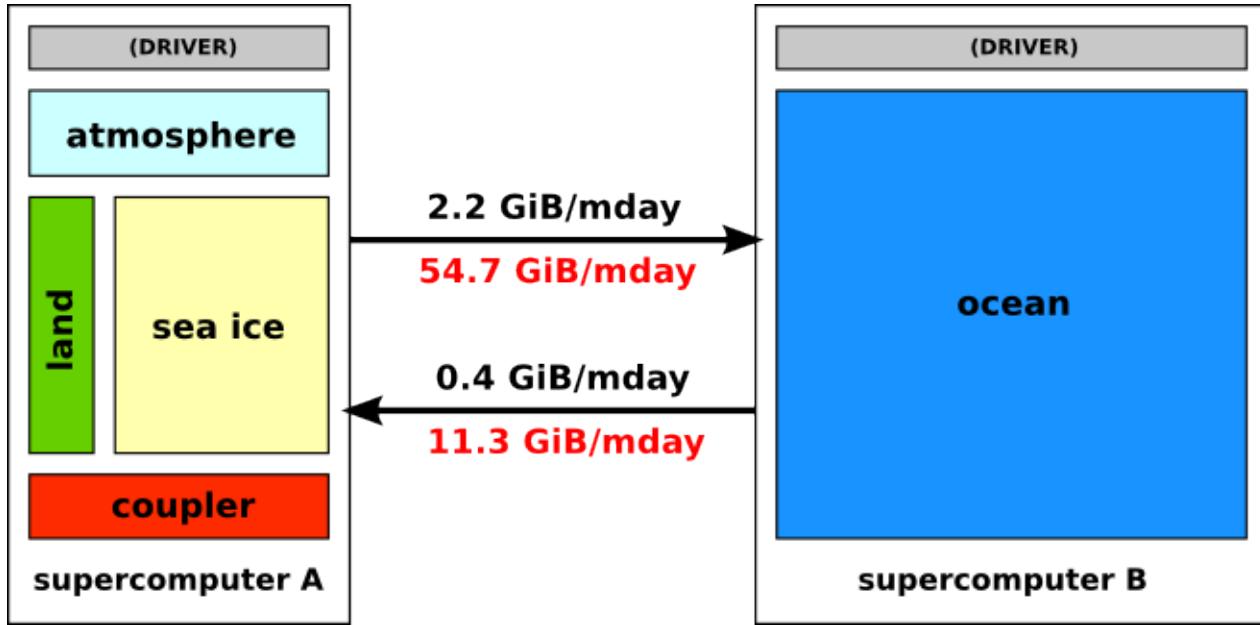


Atmosphere, land and sea ice exchange data every 30 model minutes, but ocean only 1x to 4x each model day!

Can we distribute CESM ?

We can run the ocean model **separately**

Coupling occurs 1~4x per simulated day
(once every **30~120 sec**) and **data volume** is limited.

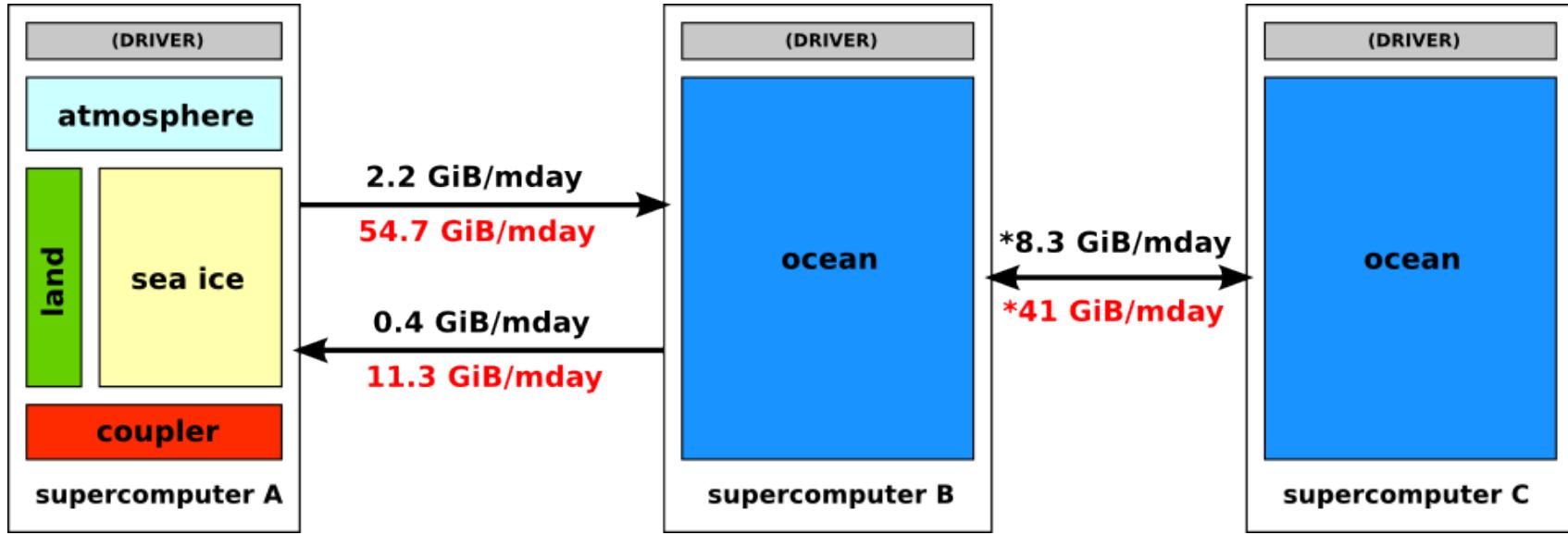


2.6 Gigabyte per exchange for 10 km resolution

66 Gigabyte per exchange for 2 km resolution

bursty traffic, not latency sensitive!

Can distribute CESM even more ?



Split ocean model using **geography aware load balancing** (GMD paper)

8.3 Gigabyte per exchange for 10 km resolution
41 Gigabyte per exchange for 2 km resolution
streaming traffic, latency sensitive!

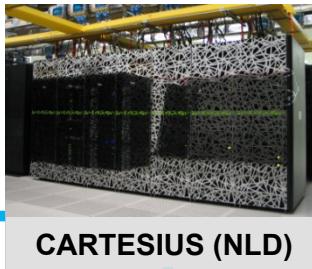


Top 500 #23
2.9 PFlop/s



SUPERMUC (GER)

Top 500 #69
1.0 PFlop/s



CARTESIUS (NLD)

10G



EMERALD (UK)

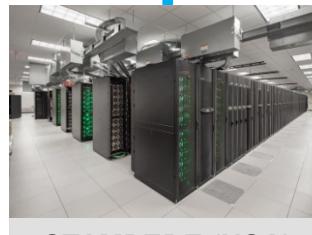
10G

GPUs

GPUs

Idea: combine 2 or 3 supercomputers from EU and the US to run an **Earth System Model** at 0.02° (2x2 km) resolution.

(this required quite a bit lot of cooperation...)



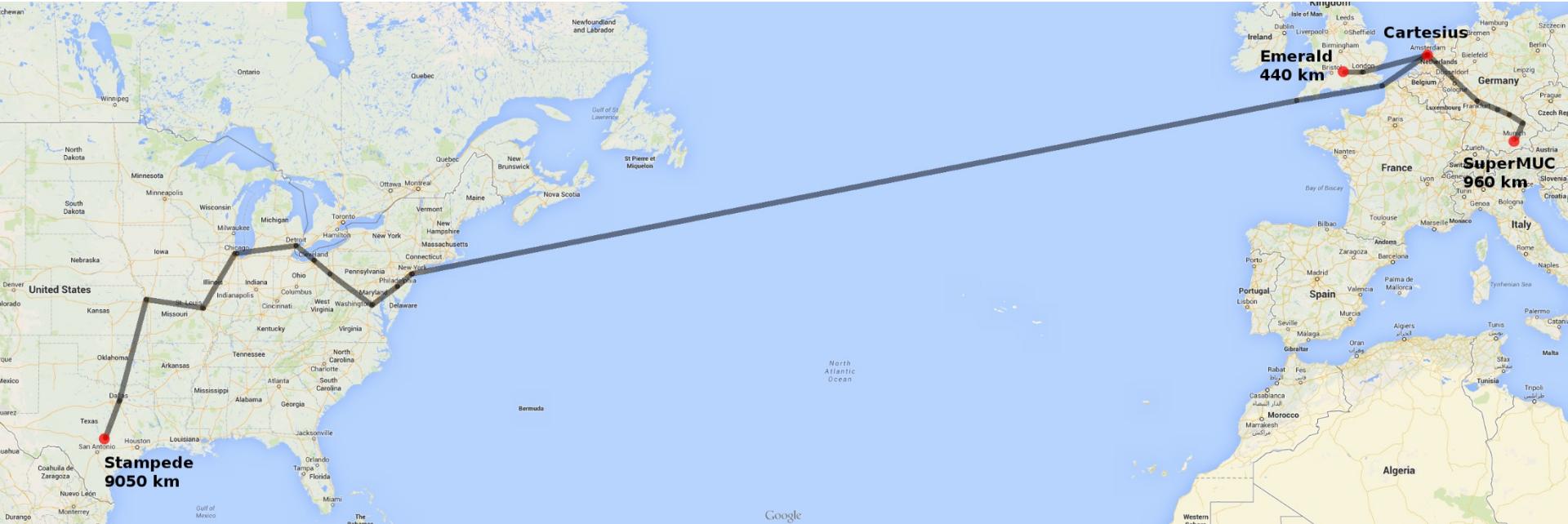
STAMPEDE (USA)

Xeon Phi

Top 500 #10
5.1 PFlop/s



Our prize: World wide 10G network links!



FUNET

janet

INTERNET²

SURF NET

NRENs provide the **essential** 10G network links between the sites.

Supers are difficult to connect!

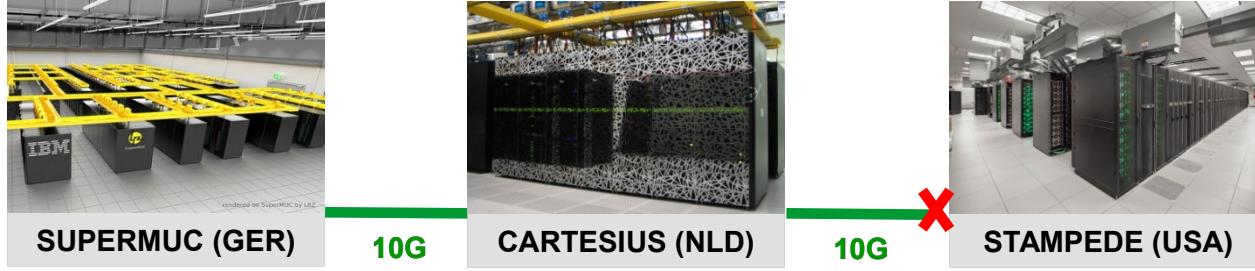
Arranging the infrastructure was quite hard!

Interestingly, the software and the long distance links are not the hard part, but the **endpoints** are!

For eSalsa-MPI we need an extra machine connected to the wide area link and the internal network used by MPI – this was very hard to organize!



Too much red tape....

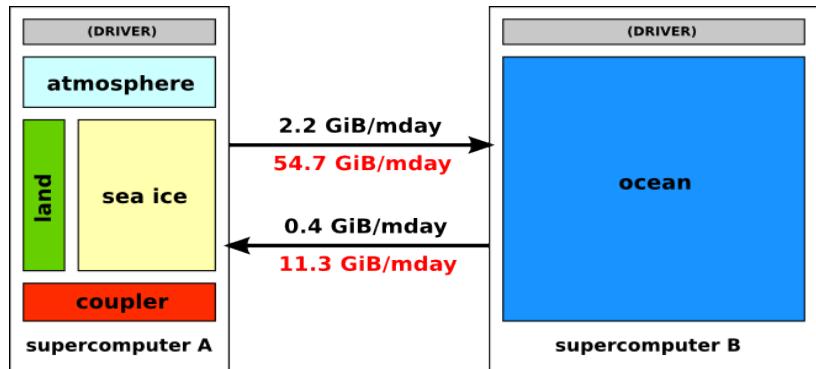


Supercomputers are usually not connected in this way, so we needed to cut through quite a bit of red tape. In the end we only managed to get the link to Germany up.

Initial results

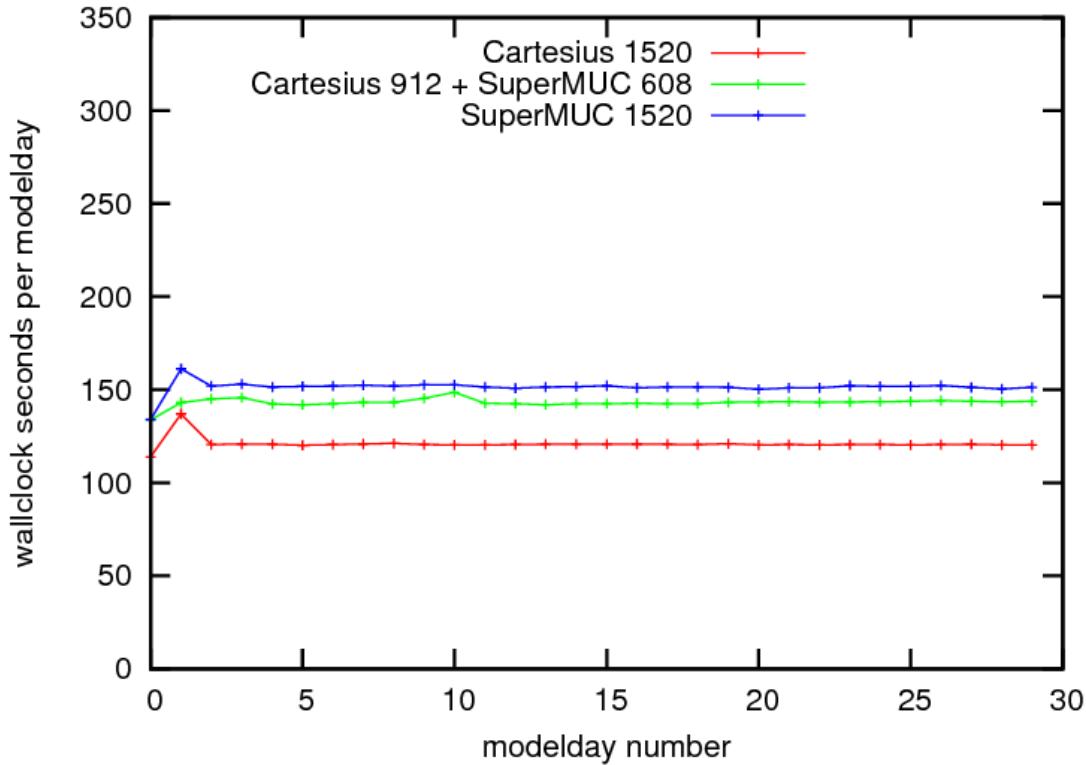
Initial experiments are running on
912 cores at SuperMUC plus
608 cores at Cartesius (1520 total)

These weird core counts are determined
by the limitations of the models, and
their relative performance on both
machines.



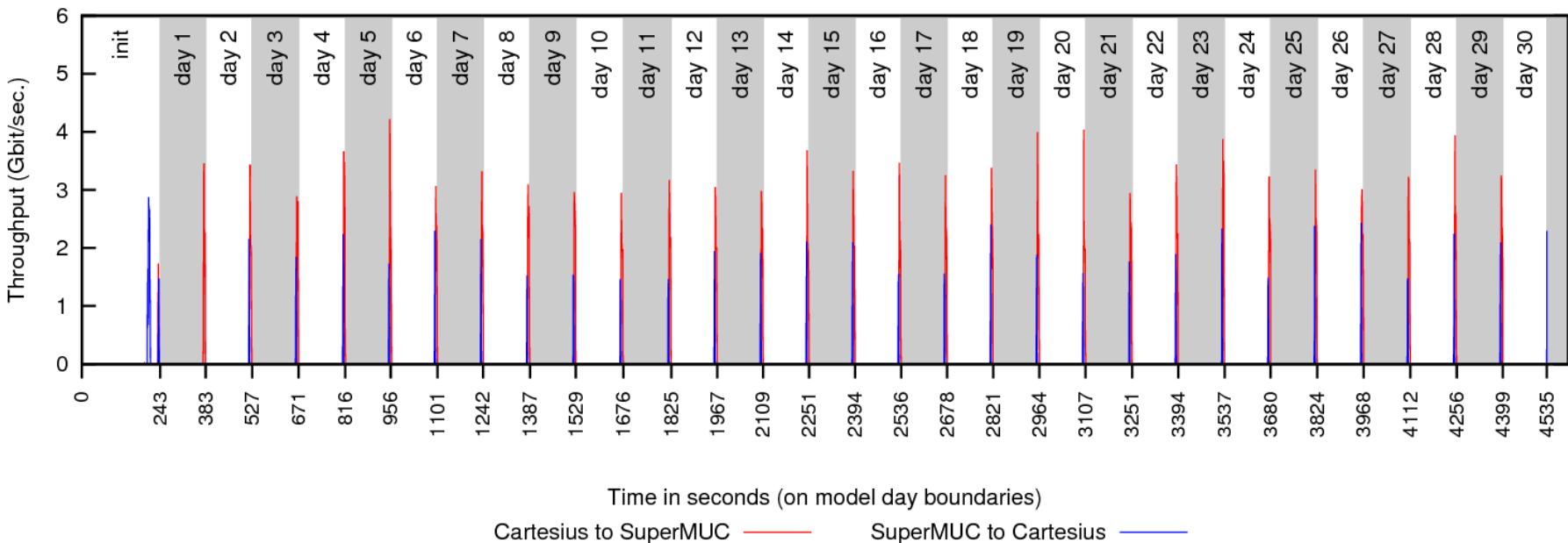
Results

Performance of combined system is “in between” the individual systems at the same core count!



Wide area bandwidth

Cartesius - SuperMUC



Wide area communication on model day boundaries reaches 3~4 Gbit/s & 2 Gbit/s

eSalsa overview

In this project we used:

MPI, GPUs, 3 flavors of distributed computing, load balancing, optical wide area networks (lightpaths), performance models, semi-automatic translation from Fortran to CUDA, visualization of large data sets ...

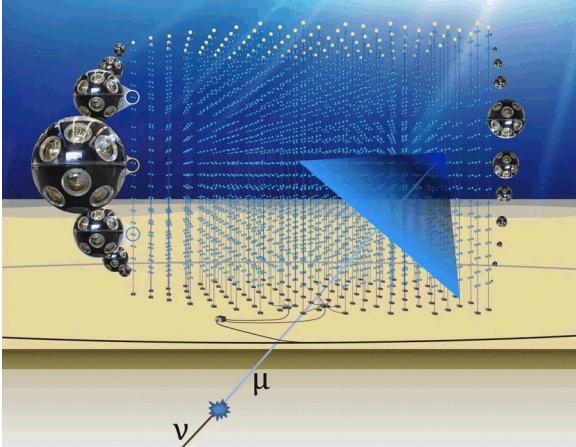
... we may have overdone it a little!

A few words
on
**streaming sensor
data**

Keeping up with sensor data



LOFAR
25 GB/s



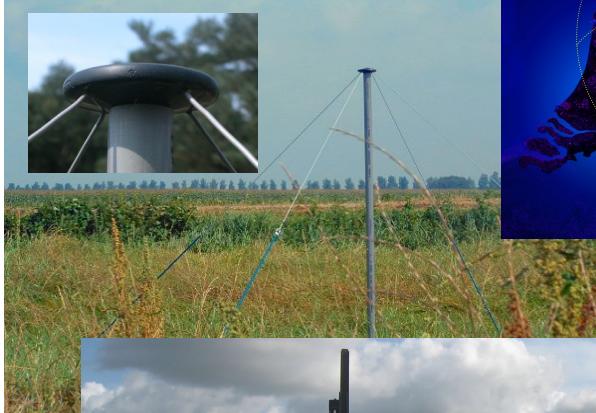
KM3NET
55 GB/s



AMBER
7 GB/s

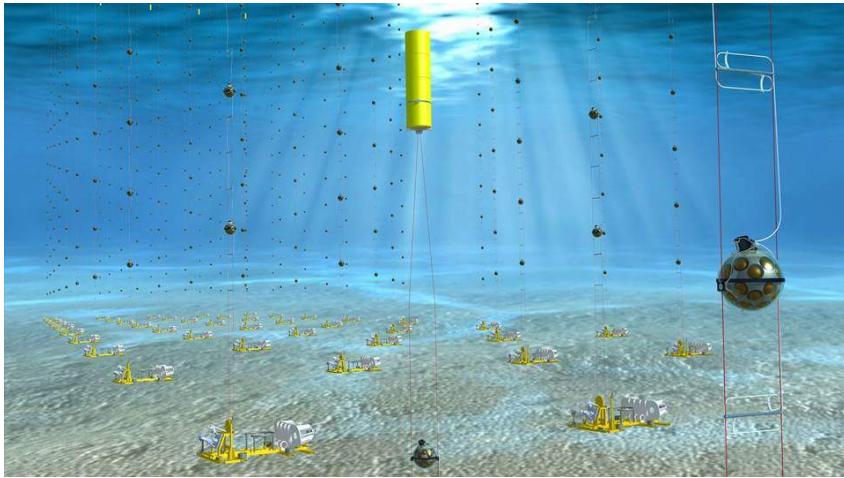
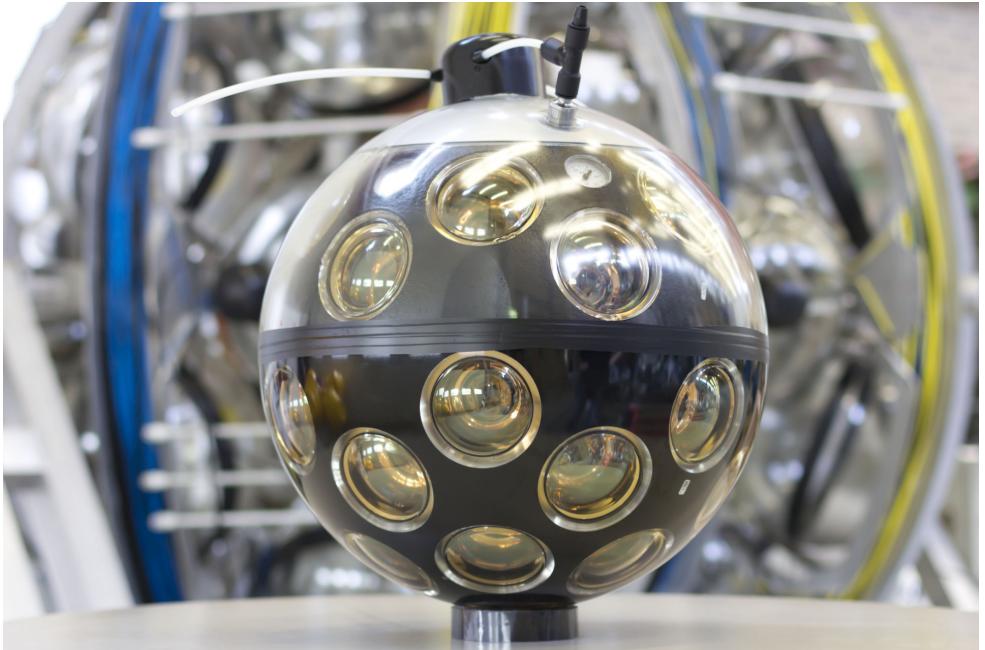
Many scientific instruments produce large amounts of data
that needs to be processed in real-time.

eAstronomy



LOFAR radio telescope – 88.000 antennas in 64 stations producing 25 GB/s, processed real time to remove radio frequency interference (cell phones, radio, lightning, etc), and to combine the signals into a single view of the sky.

Real-time detection of neutrinos from the distant Universe



KM3NeT neutrino detector:
12.000 spheres with 31 photo detectors
on 600 strings in the Mediterranean.
Process 55 GB/s in real time to remove
interference from potassium decay and
bio-luminescence. Produce output when a
hit is likely (22 MB/s).

Real-time SAR processing



AMBER Synthetic Aperture Radar:
8 sending and 24 receiving antennas, whose
signals are combined in software. Requires a
quite a bit of processing, but on a limited energy
budget of about 50 watts.

GPUs are essential

(of all shapes and sizes)



Only GPUs deliver the raw performance **and** performance per watt needed for these projects.
The same expertise is required for the big telescopes and the small drones!

Straightfoward algorithms

(LOFAR correlator)

```
for (time = 0; time < integrationTime; time++) {  
    sum += samples[ch][station1][time][pol1] *  
          samples[ch][station2][time][pol2];  
}
```

Straightfoward algorithms

(LOFAR RFI mitigation)

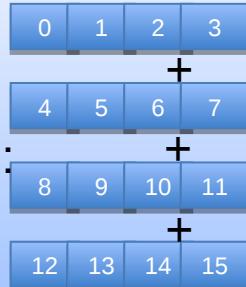
Stream of samples



Period 8:



Period 4:



Straightfoward algorithms

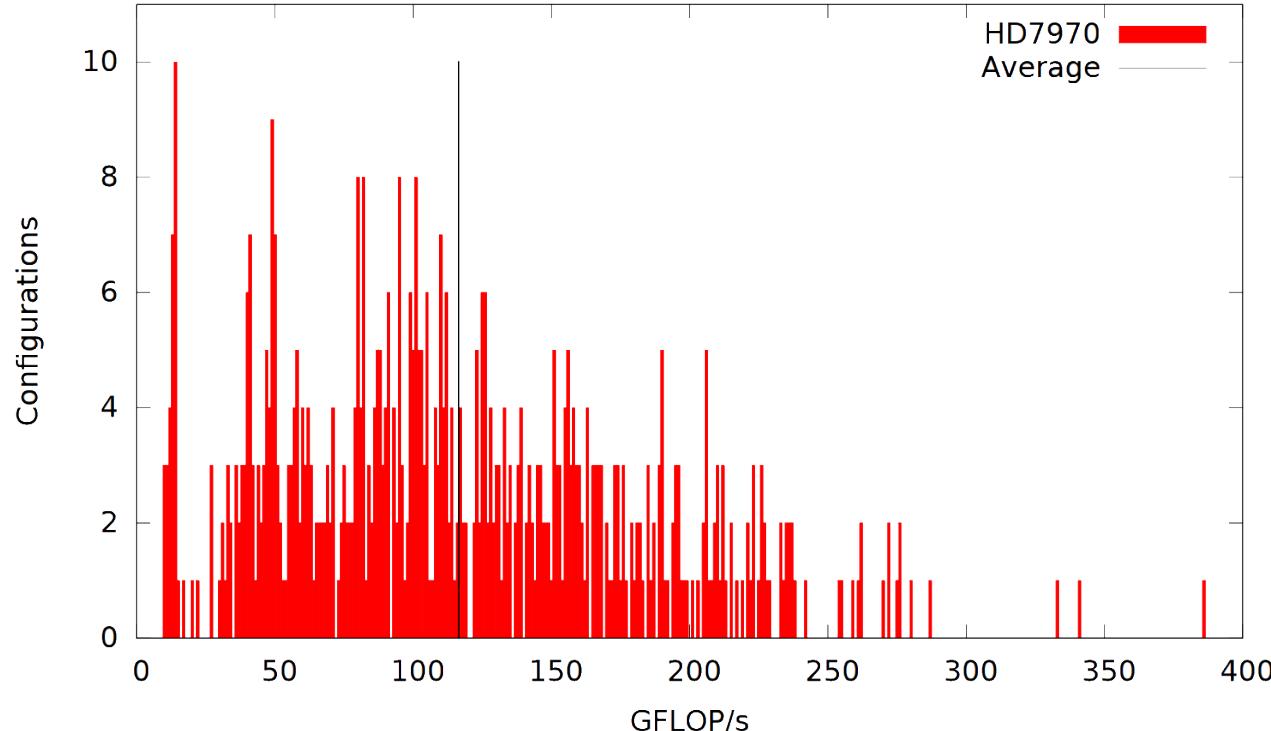
(KM3NeT hit correlation)

For every pair of detected hits (x_1, y_1, z_1, t_1) and (x_2, y_2, z_2, t_2) compute:

$$(t_1 - t_2)^2 \times c < (x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2$$

Tuning GPUs codes is a lot of work

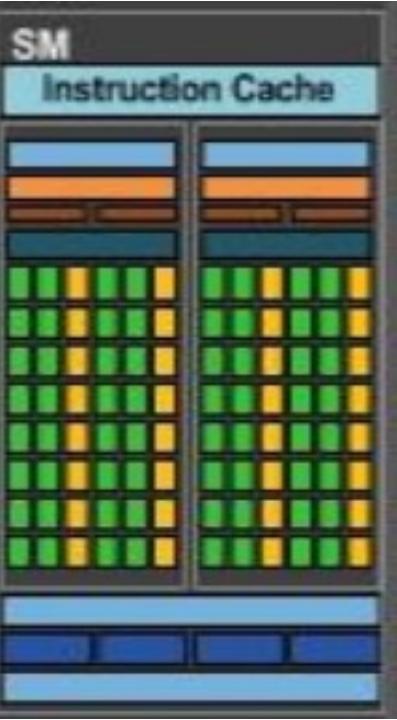
(many possibilities)



Why is this hard?



Massively parallel!



60 SMs
64 cores each
64KB shared memory

Why is this hard?



Matrix Multiply

GPU architecture is visible in the code:

Work is divided into blocks each containing a number of threads that can share memory:

- 2D block index:
blockIdx.x, blockIdx.y, block_size
- 2D thread index (within 1 block):
threadIdx.x, threadIdx.y
- Shared memory (within 1 block)
`__shared__`

```
__global__ void matmul_kernel(float *C, float *A, float *B) {  
  
    __shared__ float sA[block_size][block_size];  
    __shared__ float sB[block_size][block_size];  
  
    int tx = threadIdx.x;  
    int ty = threadIdx.y;  
    int x = blockIdx.x * block_size + tx;  
    int y = blockIdx.y * block_size + ty;  
  
    float sum = 0.0;  
    int k, kb;  
  
    for (k=0; k<WIDTH; k+=block_size) {  
        __syncthreads();  
        sA[ty][tx] = A[y*WIDTH+k+tx];  
        sB[ty][tx] = B[(k+ty)*WIDTH+x];  
        __syncthreads();  
  
        for (kb=0; kb<block_size; kb++) {  
            sum += sA[ty][kb] * sB[kb][tx];  
        }  
    }  
    C[y*WIDTH+x] = sum;  
}
```

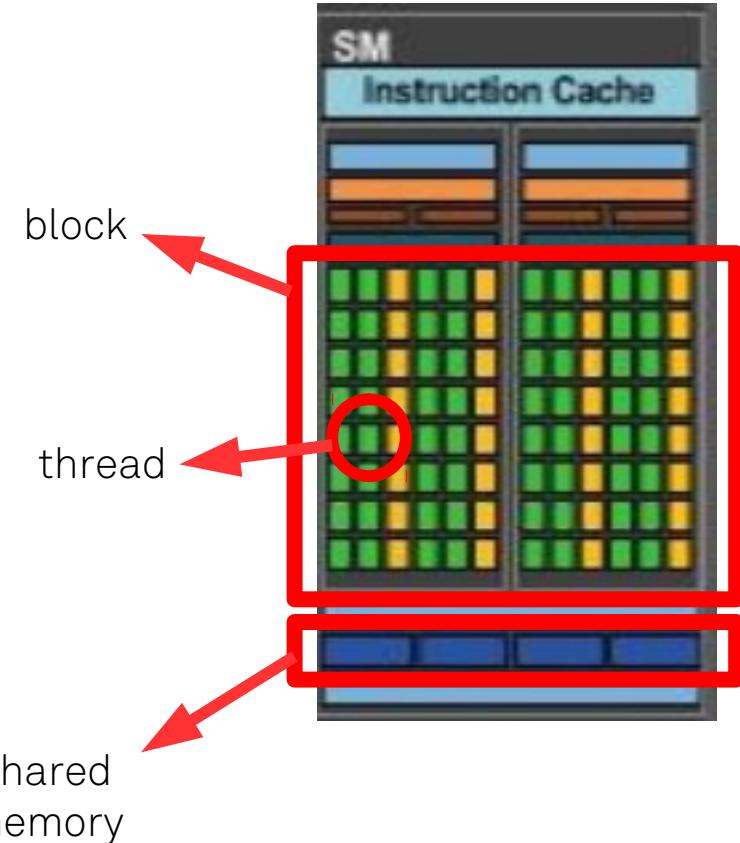
Matrix Multiply

It is very important to choose the best **block_size, thread dimesions, etc.**

When chosen correctly, the algorithm **maps directly** onto the SMs, cores and shared memory of the GPU, resulting in good performance.

When chosen incorrectly the algorithm does not map nicely, and the performance **drops fast!**

Many algorithms have tens of such parameters!



Kernel Tuner

Branch: master · New pull request · Create new file · Upload files · Find file · Clone or download ·

benvanerkhoven committed on GitHub · Update roadmap.md · Latest commit `c000e13` 6 days ago

File	Description	Last Commit
doc	added requirements and preparations for version 0.1.0	28 days ago
examples	fixed a bug for the OpenCL backend under Python 3	a month ago
kernel_tuner	option to set compiler options	12 days ago
test	option to set compiler options	12 days ago
.gitignore	updating documentation	12 days ago
.travis.yml	small fix	6 days ago
CHANGELOG.md	option to set compiler options	12 days ago
LICENSE	Initial commit.	8 months ago
MANIFEST.in	improved error reporting for when pycuda or pyopencl is not installed	25 days ago
README.md	Update README.md	22 days ago
requirements.txt	added requirements and preparations for version 0.1.0	28 days ago
roadmap.md	Update roadmap.md	6 days ago
setup.cfg	added kernel_tuner to pypi	28 days ago
setup.py	improved error reporting for when pycuda or pyopencl is not installed	25 days ago

README.md

A simple CUDA/OpenCL kernel tuner in Python

build passing · codacy A · coverage 87%

The goal of this project is to provide a - as simple as possible - tool for tuning CUDA and OpenCL kernels. This implies that any CUDA or OpenCL kernel can be tuned without requiring extensive changes to the original kernel code.

A very common problem in GPU programming is that some combination of thread block dimensions and other kernel parameters, like tiling or unrolling factors, results in dramatically better performance than other kernel configurations. The goal of auto-tuning is to automate the process of finding the best performing configuration for a given device.

This kernel tuner aims that you can directly use the tuned kernel without introducing any new dependencies. The tuned kernels can afterwards be used independently of the programming environment, whether that is using C/C++/Java/Fortran or Python doesn't matter.

The kernel_tuner module currently only contains main one function which is called `tune_kernel` to which you pass at least the kernel name, a string containing the kernel code, the problem size, a list of kernel function arguments, and a dictionary of tunable parameters. There are also a lot of optional parameters, for a complete list see the full documentation of `tune_kernel`.

netherlands
eScience center

Optimizing the choice of these parameters can be partly automated.

The **kernel tuner** is a tool that automatically searches for the optimal configuration options for a piece of code provided by the user, thereby saving a lot of a programmers time.

Porting the application to another machine is reduced to simply re-tuning it.

Kernel Tuner Example

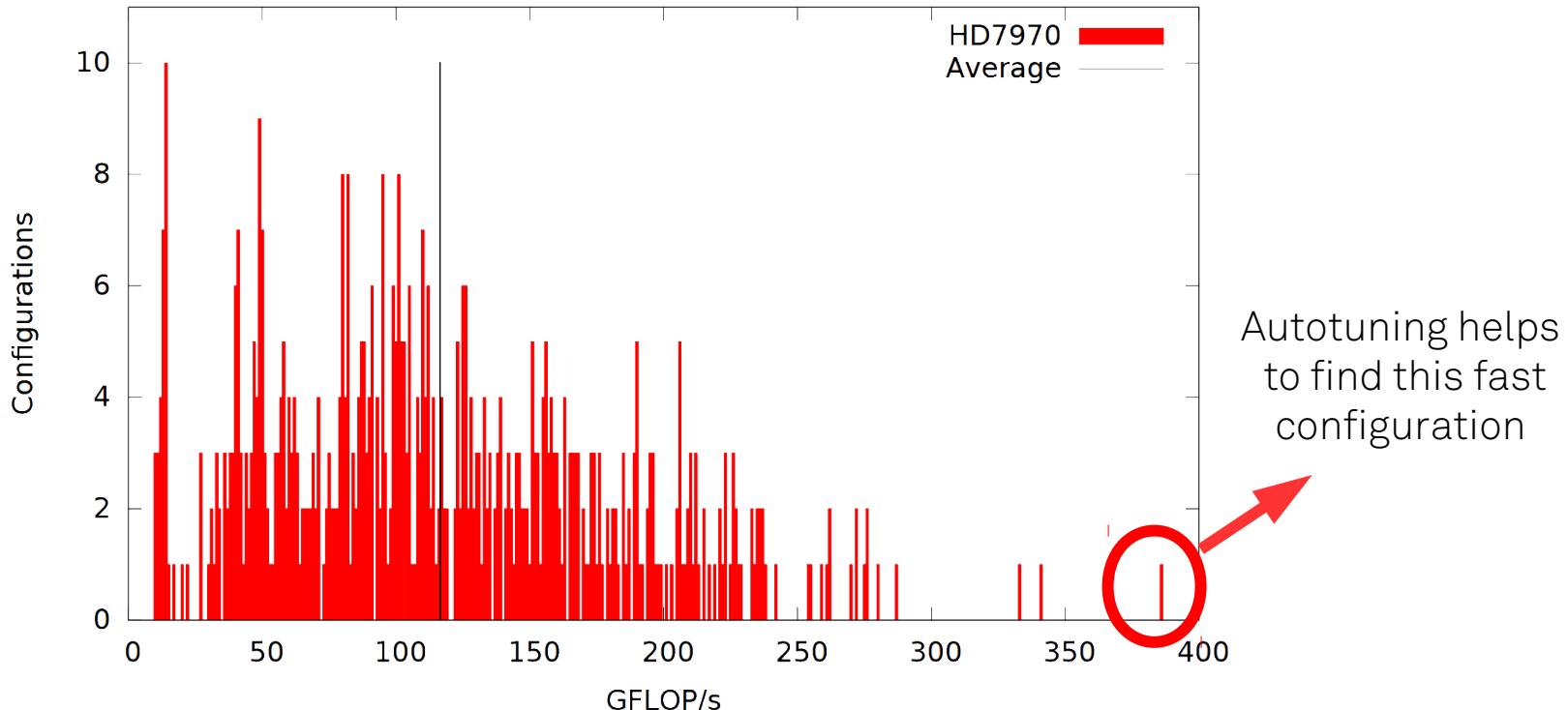
```
kernel_string = """
__global__ void vector_add(float *c, float *a, float *b, int n) {
    int i = blockIdx.x * block_size_x + threadIdx.x;
    if (i<n) {
        c[i] = a[i] + b[i];
    }
}"""

n = numpy.int32(1e7)
a = numpy.random.randn(n).astype(numpy.float32)
b = numpy.random.randn(n).astype(numpy.float32)
c = numpy.zeros_like(b)
args = [c, a, b, n]
tune_params = {"block_size_x" : [128+64*i for i in range(15)] }

kernel_tuner.tune_kernel("vector_add", kernel_string, (n,1), args, tune_params)
```

Define all possible values
for block_size_x and try
them all!

Tuning goal



Streaming Sensor data overview

Many sensors generate enormous amount of data and need real-time processing

GPUs are a good choice here, but GPU codes are hard to performance tune

Fortunately, this can be automated using tools like the **Kernel Tuner**

In conclusion ...

Real world applications often use a combination of parallelization types (MPI, OpenMP, CUDA, hierarchical simulations, ensembles, etc.)

Incorporating GPUs is often non-trivial, but potential gains are high.

Solutions are often application specific, but the knowledge and tools surrounding them are not (Xenon, Kernel Tuner, eSalsa-MPI, Marver, etc).

The goal of the eScience Center is to extract these generic knowledge and tools and make them available to all of Dutch academia.

Get in touch

Netherlands eScience Center
Science Park 140
1098 XG Amsterdam
The Netherlands

+31 (0)20 4604770
info@eScienceCenter.nl

(we also do MSc projects!)



www.eScienceCenter.nl



linkd.in/1j2uS8S



vimeo.com/eScienceCenter



twitter.com/eScienceCenter



j.maassen@esciencecenter.nl

