

Earth Sciences Department





BSC Barcelona Supercomputing Center
Centro Nacional de Supercomputación

SBC PRESENTS
SUSTAINABILIT


Destination Earth




Pioneering the Future with Digital Twins

Okke van Eck
GPU Research Engineer
Barcelona Supercomputing Center

okke.vaneck@bsc.es





27th of November 2024, Bussum

 Funded by the European Union

Destination Earth implemented by   

Short:

- Talk about my team's contribution to EC's Destination Earth project
- GPU Research Engineer at BSC
- Performance team at BSC under earth sciences department
- Links to social media

Long

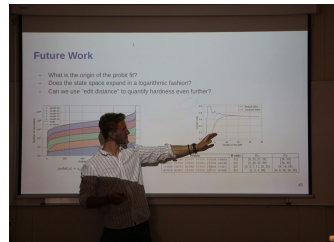
Welcome to this talk where I'll be talking about my team's contributions to the European Commission's Destination Earth project.

My name is Okke van Eck, I'm a research engineer at the Barcelona Supercomputing Center, short BSC.

I'm specialized in GPU programming, but do performance engineering in general.

Here you can find my email address and links to my social media in case you want to ask some questions afterwards.

Who's Okke?



Funded by
the European Union

Destination Earth implemented by



2

Short

- Bachelor CS at UvA and board 19 at VIA.
 - SNIC advisory board in 2019-2020
- Master's PDCS at VU
 - Awarded the Unilever Research Prize 2023
- Moved to Barcelona for BSC
 - Here small part of the team in the beginning
- Now we have a much bigger team, mostly consisting of juniors
 - Having a lot of fun, for example dinners together and go rock climbing on the weekends

Climate is changing..



Funded by
the European Union

Destination Earth

implemented by



3

The core problem we analyze is climate change.

I want to show you how bad it really was last year in 2023.

This video from the Guardian (news Australia) highlights some of the effects of climate change.

Climate is changing..



Funded by
the European Union

Destination Earth

implemented by

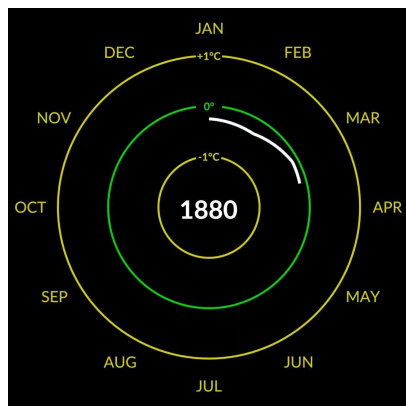


4

VIDEO PLAYS AUTOMATICALLY

https://www.youtube.com/watch?v=pf4aOIE_3bc

Climate is changing..

Funded by
the European Union**Destination Earth** implemented by

5

VIDEO PLAYS AUTOMATICALLY

So as you saw, the effects of climate change are very real and very destructive.
To make it even worse, this is just the beginning and the problems will get way worse.

In this circle, you see the temperature difference per month relative to the average temperature of 1900-2000 (20th century).

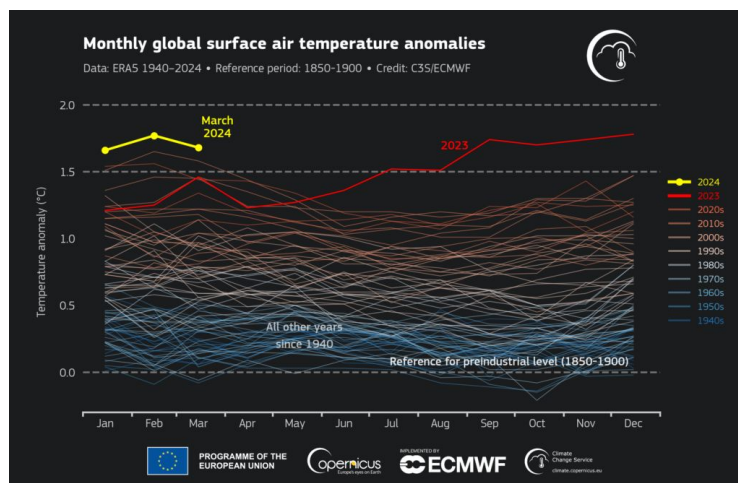
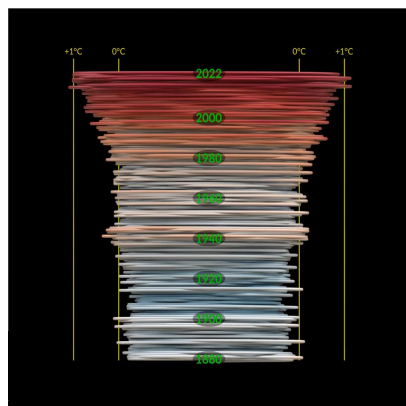
As you can see we have first colder moments.

Then there is a small increase at 1945 due to the 2nd world war.

But from 1980 the temperature really starts to rise.

And from 2019 it got even worse.

Climate is changing..

Funded by
the European Union**Destination Earth** implemented by

6

Here you can see the spiral in full, and it's clear that it's getting incrementally worse from the 2000's onwards.

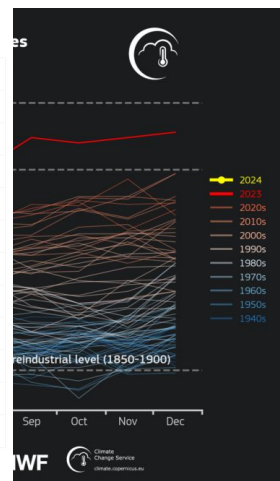
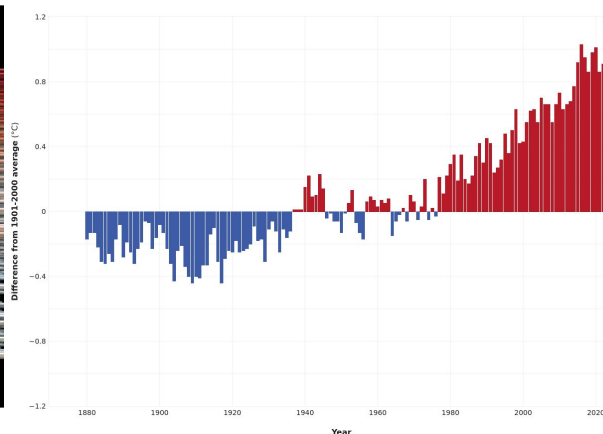
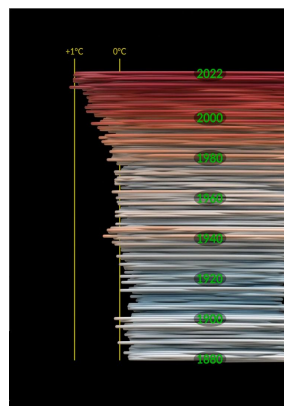
On the right side, you can see it in another view from last march where 2023 and 2024 were clearly breaking previous records.

Look at the gap in the summer between 2023 and the rest, we are really accelerating now.

And unfortunately, 2024 broke the records of 2023 afterwards.

Climate is changing..

GLOBAL AVERAGE SURFACE TEMPERATURE

Funded by
the European Union**Destination Earth** implemented by

7

So things are getting a lot worse lately.

Here you have a bar-chart version of the same data.

There is one bar missing for 2024, which will be higher than 2023.

Destination Earth: Digital Twins



Funded by
the European Union

Destination Earth

implemented by



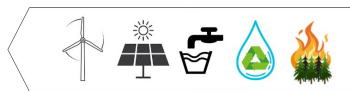
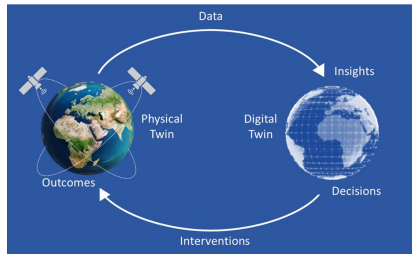
8

And so, clearly, something needs to be done.

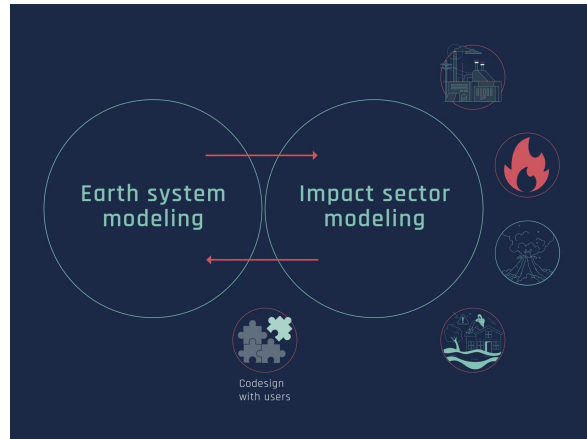
I hope you are not feeling too hopeless, but this is reality, and we achieve nothing by sitting on our hands.

Thus the European Commission has started the Destination Earth project in 2022, where they aim to build a digital twin of the earth.

Why digital twins?



Impact sector applications

Funded by
the European Union**Destination Earth**

implemented by



Let me explain what a digital twin exactly is.

Here you can see the basic premise.

We have a physical twin, the earth, where we can measure current activity with the help of satellites.

Then we want to use that data to simulate what our climate will do in the future.

This is done in the digital twin, which is a model of the earth.

Based on the predictions of this model, we can intervene and change policies to better our climate.

These policies have an effect, which we can measure again through satellites, feed to the model, and repeat.

This idea of using a digital twin consists of 2 sectors.

We have Earth System Modelling, which models the physics of the earth.

Then you have Impact Sector Modelling, which models the effects of the climate on specific sectors.

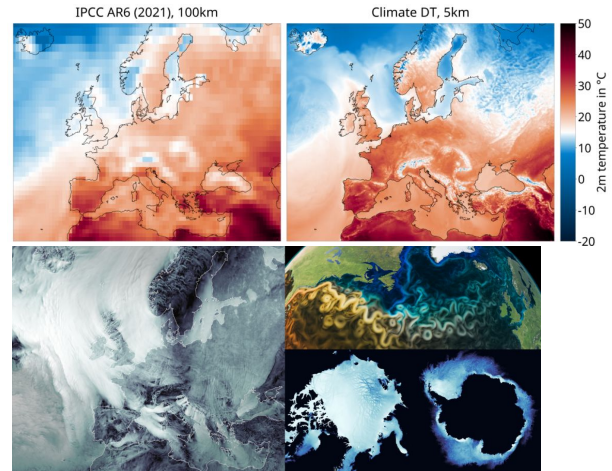
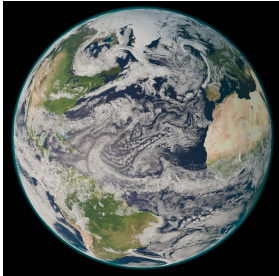
For example the effect on forest fires or floods.

These impact sector applications are viewed by ESA to advice policy makers on potential interventions.

We will see this later when I'll explain how we run a simulation.

Destination Earth: Digital Twins

- Phase 1 ran from 2022-2024
- 3 different models
- One simulation takes months
- 300+ supercomputer nodes
- Multiple EuroHPC supercomputers



Funded by
the European Union

Destination Earth

implemented by



10

Thus we have the Destination Earth project, where the Digital Twins are a subproject of.

The goal is to create a highlight accurate digital twin at global scale, which monitors, simulates, and predicts the intersection between natural phenomena and human activities

Phase 1 finished last april and took 2 years.

There are 3 or 4 other phases of 2 years coming up, so in total it will be a project till +-2030.

In this first phase, we were able to achieve a 5km resolution as you can see in the left upper corner.

This is much more detailed than the 100km standard from before the project. Russia is actually cold now for example.

One important thing is that we do not use 1 model, but 3 different ones.

Each model has a bit of a bias towards some part of the physics.

It is super complex to get it exactly right, because the laws of physics are incredibly complex.

Instead of optimizing 1 model, we have 3 and take the average of them for the predictions.

One of these simulations takes multiple months for a 50 year prediction.

It requires 300+ supercomputing nodes to achieve this simulation speed.

As you can imagine, this takes a lot of power as well.

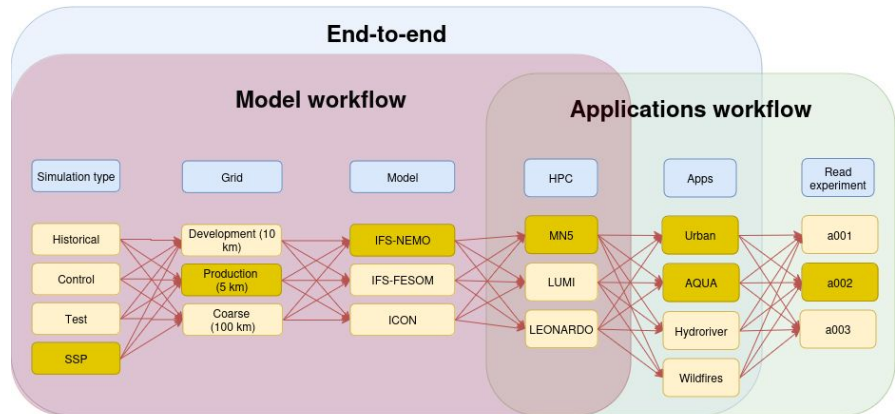
Lastly, we have to develop the models to run on 3 different EuroHPC supercomputers.

Each have different architectures so the portability becomes a bit of a problem, but about that later more.

But performance is critical, so we tune the models to work as best as possible.

What does it take to run a model?

- **AUTOSUBMIT** 



Funded by
the European Union

Destination Earth implemented by



11

And so you might wonder, what does it take to run a model?
Well, for sure you will need a workflow manager.

AutoSubmit is our in-house open-source workflow manager, that is specially created for running climate simulations.

It gives scientists options for preparing, launching, and analyzing climate simulations. As you can see here, it has different configurations we can choose from, depending on the simulation we want to run.

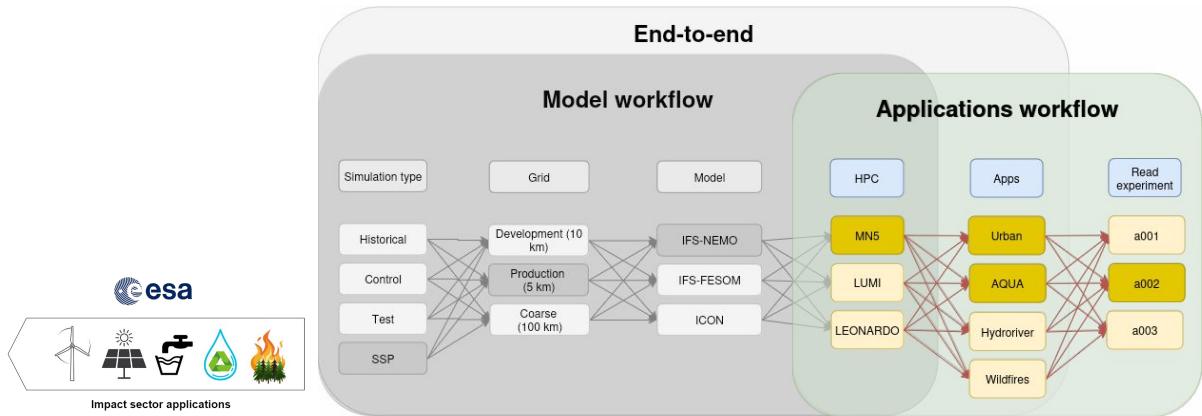
There is the "Model Workflow" and the "Applications Workflow".

SSP = Shared Socioeconomic Pathways

<https://www.carbonbrief.org/explainer-how-shared-socioeconomic-pathways-explore-future-climate-change/>

What does it take to run a model?

- 



Funded by
the European Union

Destination Earth

implemented by



12

Let's focus on the application workflow first.

Here you have three parts, of which the first one is to select a high-performance computer (HPC).

Then on the 2nd column, we have different applications that analyze the experiments.

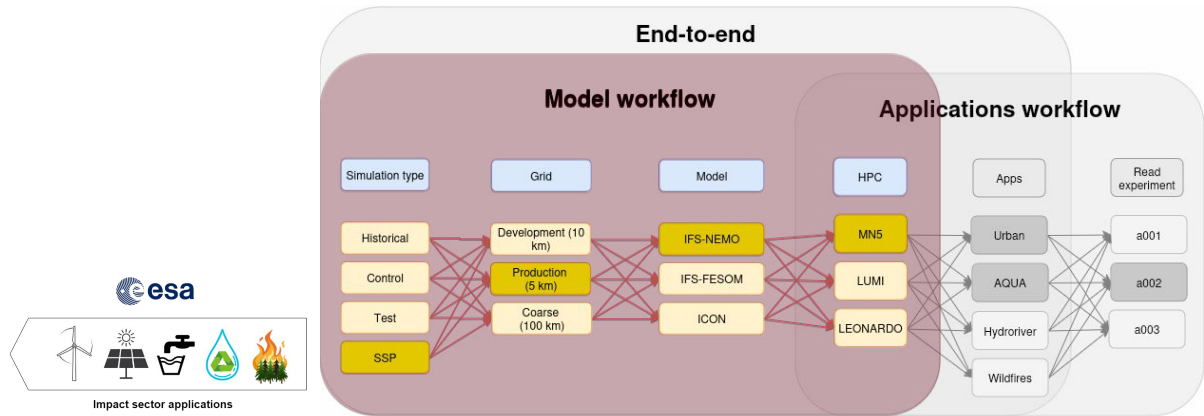
These correspond to the impact sector applications that we have seen before.

All of these need to read data from experiments, which is the last column.

So this is really easy. Load data, run the applications on HPC, and have a climate researcher look at the data.

What does it take to run a model?

- 



Funded by
the European Union

Destination Earth implemented by



13

Then we can also look at the Model Workflow.

This is the part where we actually start up simulations.

In the first column we select the simulation type that we want to do.

The meaning of SSP doesn't matter much for this talk, but it's basically some socio-economic situations.

Then in the 2nd column we select the grid size of the model, like we saw before.

We typically run with a 10km resolution for development and official 5km during simulations

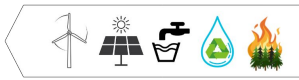
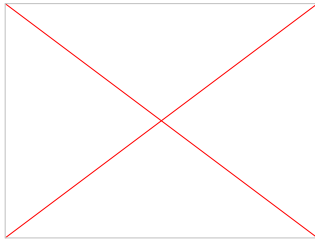
NEXT SLIDE

SSP = Shared Socioeconomic Pathways

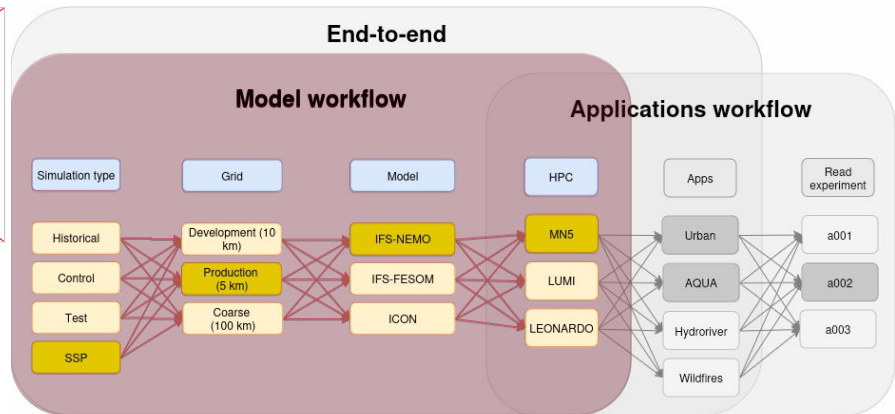
<https://www.carbonbrief.org/explainer-how-shared-socioeconomic-pathways-explore-future-climate-change/>

What does it take to run a model?

• **AUTOSUBMIT**



Impact sector applications



Funded by
the European Union

Destination Earth implemented by



14

VIDEO PLAYS AUTOMATICALLY

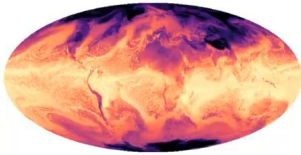
In this video you can see how the different grids are relative to each other.

This is managed through something called HEALPix.

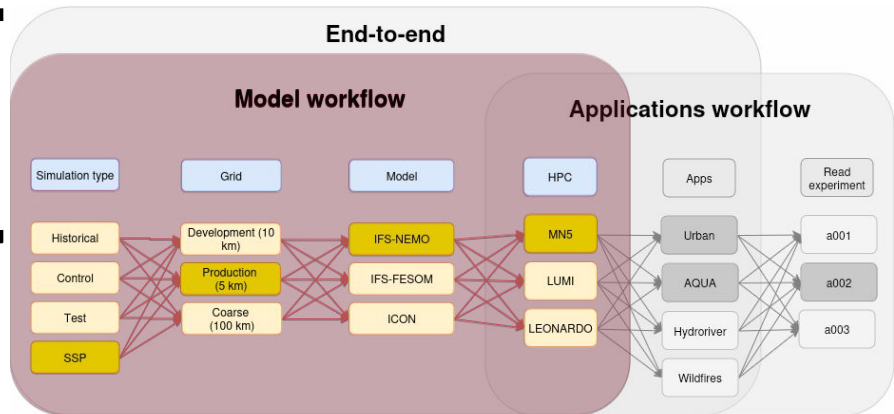
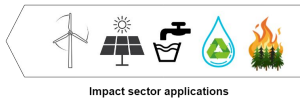
It's a pixelation algorithm for sphere-like geometry, where each pixel always covers the same amount of area.

What does it take to run a model?

• **AUTO****SUBMIT**



 **esa**



Funded by
the European Union

Destination Earth implemented by



15

Then in the 3rd column we select what model to run.
Like I said before, we have 3 models that we develop for averaging out biases.

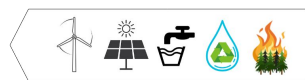
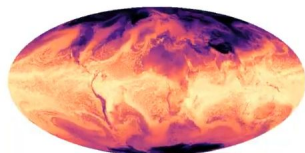
IFS is an atmosphere model, which is coupled with a sea model.
In case of the first 2 options, the sea models are called NEMO and FESOM.
Hence you can imagine how difficult the sea is to model, especially the sea ice.

Lastly we have ICON, which is a singular model which combines an atmosphere and a sea part.

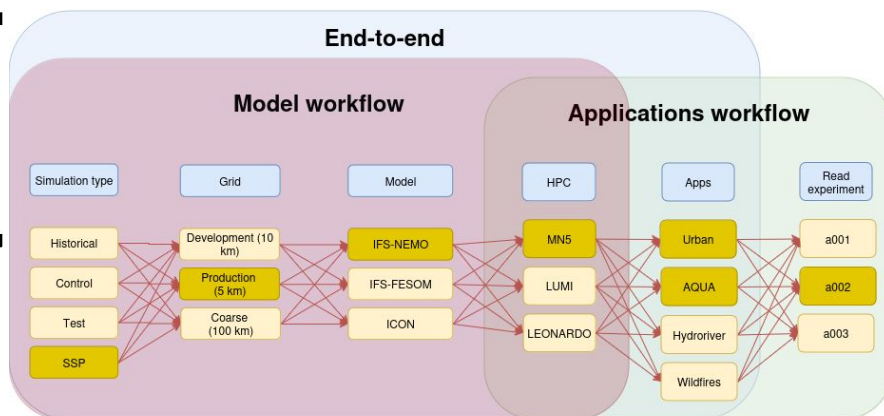
Then in the 4th column we have the different supercomputers that we have to develop for, which I'll cover later.

What does it take to run a model?

• **AUTO** **SUBMIT**



Impact sector applications



Funded by
the European Union

Destination Earth implemented by



16

Going back to the full picture, this is basically what is required for running a climate model.

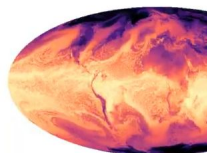
Note the "End-to-End" workflow in blue, which basically forwards the output data straight to the apps, which eliminates the storage.

What does it take to run a model?

• **AUTOS**

esiwace
Centre of Excellence in Simulation of Weather
and Climate in Europe

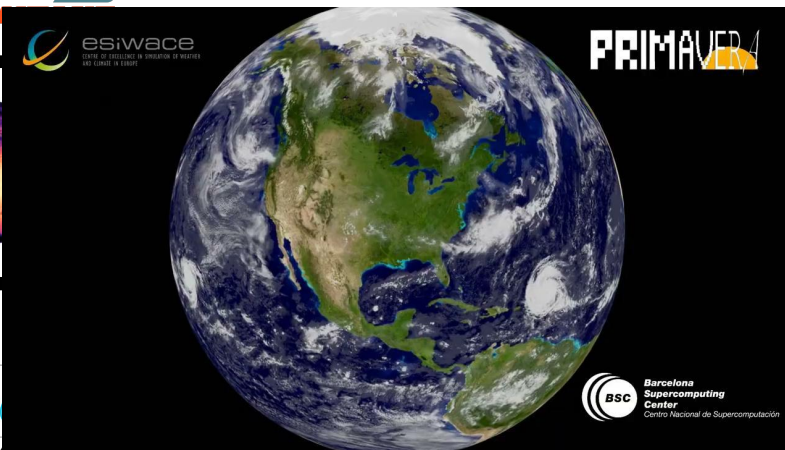
PRIMAVERA



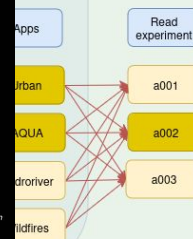
esa



Impact sector applications



Workflow



Funded by
the European Union

Destination Earth implemented by



17

And thus if you combine everything together, you will be able to simulate the earth on different resolutions.

Here you can see the different resolutions for the sea part.

Performance Engineering



Funded by
the European Union

Destination Earth

implemented by



18

That sums up my introduction to climate modelling, and we can talk about what my team does.

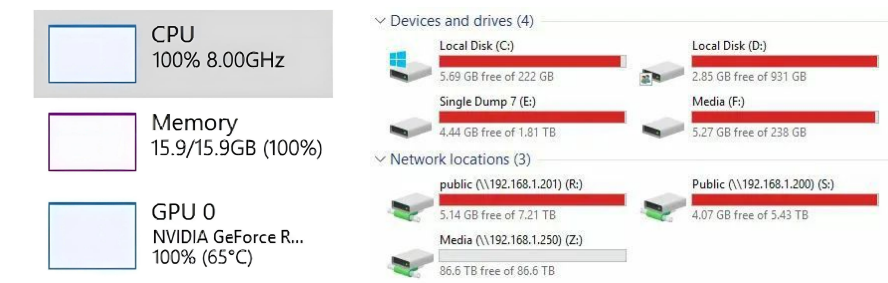
As the Performance team, we look at models and try to understand inefficiencies that we can improve upon.

Climate models are generally written by climate scientists, and thus their performance is not really optimal for the modern day systems.

Let me explain a bit what performance engineering is about.

Performance Engineering

"I paid for the whole PC,  gonna use the whole PC"



Funded by
the European Union

Destination Earth implemented by



19

So, what is performance engineering?

One might think that performance engineering is about fully utilizing a system.

However, that's not really the case.

More optimally utilizing the system is our solution to some problems organizations face.

Performance Engineering

"I paid for the whole PC,  gonna use the whole PC"



Funded by
the European Union

Destination Earth implemented by



21

Which are for example to reduce expenses.

Hardware is expensive, and thus buying extra machines to scale up is not always the best choice.

Maybe you can compute more with the hardware you already have.

Next, it's about sustainability.

By optimizing models, we reduce their compute time and thus we need less hardware to begin with.

Lastly, it's also about cutting CO2 emissions.

Reducing runtime makes supercomputers consume less energy and thus the emission is reduced.

Credits for the 3 icons:

https://www.freepik.com/free-vector/cut-price-bargain-offering-reduced-cost-discount-low-rate-special-promo-scissors-dividing-banknote-crisis-bankruptcy-cheapness-market-vector-isolated-concept-metaphor-illustration_12470246.htm#fromView=search&page=1&position=0&uuid=2212521e-e19b-4c7c-8e37-8f20c0c2fb4d

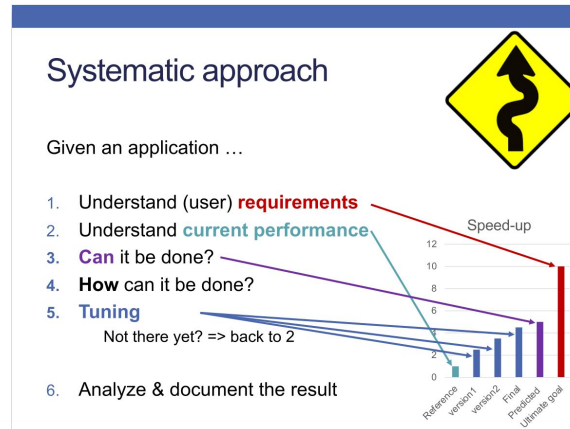
https://www.freepik.com/free-vector/resources-protection-abstract-concept-illustration-protection-natural-resources-land-conservation-safeguarding-nature-smart-water-use-environment-preservation_12145555.htm#fromView=search&page=1&position=0&uui

[d=2d973d43-fda8-453f-8df7-9bd3a2177589](#)

https://www.freepik.com/free-vector/stop-air-pollution-carbon-dioxide-reduction-environmental-damage-atmosphere-protection-toxic-emission-problem-ecology-volunteer-cartoon-character_12145539.htm#fromView=search&page=1&position=2&uuid=18f7ce33-9c1b-4904-a098-a0a807c4d5fb

Performance Engineering

Follow these simple steps:



Funded by
the European Union

Destination Earth implemented by



21

So how does the process of performance engineering look like?
Follow these simple steps:

1. Understand the requirements that the user is asking for.
Usually this is too ambitious.
2. Understand what we currently can do.
3. Predict what performance is possible.
4. Think of techniques to achieve this predicted performance.
5. Improve bit by bit and compare your versions!
Can we actually achieve the prediction?
6. Analyze your final version and document why what was possible

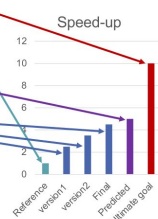
Performance Engineering

Follow these simple steps:

Systematic approach

Given an application ...

1. Understand (user) **requirements**
2. Understand **current performance**
3. **Can** it be done?
4. **How** can it be done?
5. **Tuning**
Not there yet? => back to 2
6. Analyze & document the result



Thank you Ana Varbanescu!



Funded by
the European Union

Destination Earth implemented by



22

If these slides looked a bit familiar, I stole them from Ana Varbanescu. She gave me my performance engineering course at the UvA and VU, and this slide is from her slide deck.
Thank you Ana for the great course (and slides)!
Now she mainly works at the TU Twente I believe, highly recommend following her courses.

Performance Engineering

Always optimize your application for a specific system!



Funded by
the European Union

Destination Earth

implemented by



23

One thing Ana taught us, was to always optimize your application for a specific system.

This is due to caching behavior, network behavior, etc.

Performance Engineering

Always optimize your application for a specific system!



Funded by
the European Union

Destination Earth

implemented by

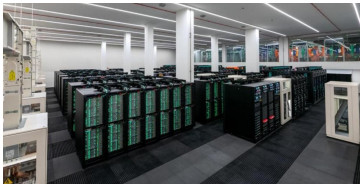


24

Well, then you get a European project, which doesn't ask you to optimize for one system..

Performance Engineering

Always optimize your application for a specific system!



Funded by
the European Union

Destination Earth implemented by

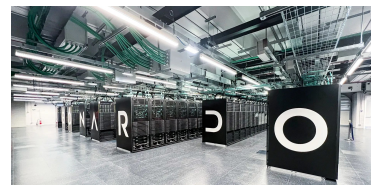


25

Not two systems..

Performance Engineering

Always optimize your application for a specific system!



Funded by
the European Union

Destination Earth implemented by



26

But three different systems!

And the model should be optimized for all.

On the left you see MareNostrum 5, our supercomputer in Barcelona, Spain.

Managed by Barcelona Supercomputing Center

Then in the middle you see LUMI, which is the supercomputer in Kajaani, Finland.

Managed by CSC.

And on the right Leonardo, which is the supercomputer of Bologna, Italy.

Managed by Cineca.

Performance Engineering



Rank	TOP500 Rank	System	Cores	Rmax (PFlop/s)	Power (kW)	Energy Efficiency (GFlops/watts)
12	5	LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland	2,752,704	379.70	7,107	53.428
15	8	MareNostrum 5 ACC - BullSequana XH3000, Xeon Platinum 8460Y+ 32C 2.36GHz, NVIDIA H100 64GB, Infiniband NDR, EVIDEN EuroHPC/BSC Spain	663,040	175.30	4,159	48.320
28	7	Leonardo - BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, NVIDIA A100 SXM4 64 GB, Quad-rail NVIDIA HDR100 Infiniband, EVIDEN EuroHPC/CINECA Italy	1,824,768	241.20	7,494	32.187

Funded by
the European Union

Destination Earth

implemented by



27



These 3 computers are the most powerful in Europe, and also very green!
Here are their ranks according to the Green500 and Top500, which are world wide rankings.

You can see that they are the numbers 5, 8, and 7 world wide.

These are all partitions which have GPUs available, which make them more energy efficient than the CPU partitions.

Performance Engineering

Performed optimizations:

- CPU 
- GPU 



Funded by
the European Union

Destination Earth implemented by



28

So we used these machines to perform some optimizations on the IFS-NEMO climate model, and I'd like to share some outcomes of those.

The problem with the optimizations is that they are quite technical and thus I don't think I can fully explain them within 45 minutes.

Hence, I will briefly mention them and continue explaining how we detect inefficiencies within code.

CPU Icon:

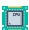

https://www.flaticon.com/free-icon/cpu_984391

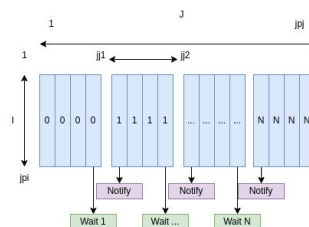
GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

Performance Engineering

Performed optimizations:

- CPU 
 - OpenMP regions (1.2x - 2.8x depending on the region)
- GPU 



Funded by
the European Union

Destination Earth implemented by



29

First for the CPU optimizations, we have OpenMP regions that are optimized. Basically, OpenMP is a multiprocessing paradigm which allows for shared memory programming.

Sometimes, these processes have to wait for each other, creating time in which nothing is computed.

We optimized some of these regions and got a 1.2x to 2.8x improvement depending on the region.

Note that this is a REGIONAL improvement, not a GLOBAL one!

The global improvement is much less, typically we are happy with a 5% global improvement, which is 1.05x.

CPU Icon:

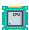

https://www.flaticon.com/free-icon/cpu_984391

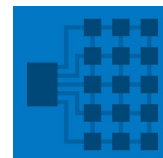
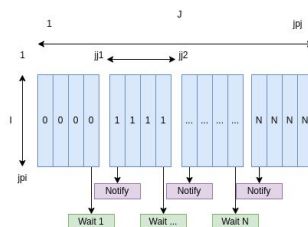
GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

Performance Engineering

Performed optimizations:

- CPU 
 - OpenMP regions (1.2x - 2.8x depending on the region)
 - MultIO improvement (1.03x)
- GPU 



Funded by
the European Union

Destination Earth implemented by



30

Next, we did an improvement within the MultIO part, which resulted in a 1.03x GLOBAL improvement.

Might seem little, but is actually significant.

MultIO handles the parallel input/output of the model, which allows for writing results in parallel to computation.

CPU Icon:

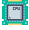

https://www.flaticon.com/free-icon/cpu_984391

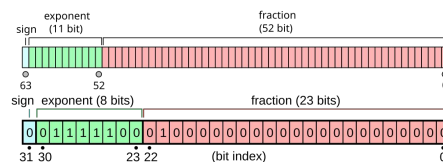
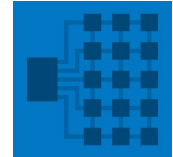
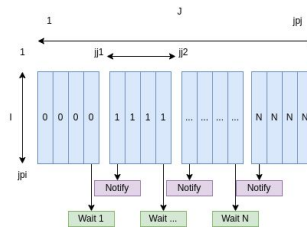
GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

Performance Engineering

Performed optimizations:

- CPU 
 - OpenMP regions (1.2x - 2.8x depending on the region)
 - MultIO improvement (1.03x)
 - Mixed Precision (1.2x)
- GPU 



Funded by
the European Union

Destination Earth

implemented by



31

Then we also introduced mixed precision to the model, which resulted in a 1.2x GLOBAL speedup.

Within computing you have single and double precision, which are 32 and 64 bit variables respectively.

Like you can see in the figures.

As you can imagine, single precision with 32-bits allow you to store twice the amount of variables in the same space.

This means that you can load twice the amount of variables in the cache, and thus are more efficient when computing the data due to fewer memory loads.

CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

GPU Icon:

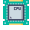

https://www.flaticon.com/free-icon/gpu_4617742

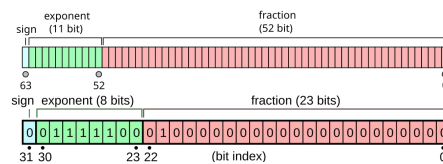
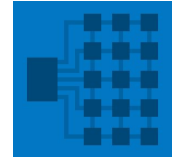
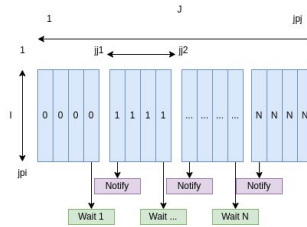
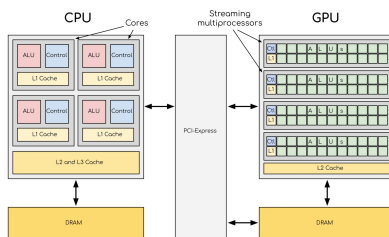
Server icon:

<https://dribbble.com/shots/4189516-Parallel-Computing-Icon>

Performance Engineering

Performed optimizations:

- CPU 
 - OpenMP regions (1.2x - 2.8x depending on the region)
 - MultIO improvement (1.03x)
 - Mixed Precision (1.2x)
- GPU 
 - Porting pilot-regions of Sea-Ice module to GPUs (5.9x)



Funded by
the European Union

Destination Earth implemented by



32

And lastly, my sub-team has been porting parts of the model to GPU to see if we can speedup the computation.

This resulted in a 5.9x times speedup for the first module that we ported. GPU computing requires memory to be transferred from the CPU to the GPU, which takes time.

Hence, we hope to achieve higher global speedup when we include other modules as well that reuse the same memory.

CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

Server icon:

<https://dribbble.com/shots/4189516-Parallel-Computing-Icon>

CPU-GPU figure:

https://enccs.github.io/gpu-programming/_images/CPUAndGPU.png

Performance Profiling



Funded by
the European Union

Destination Earth

implemented by



33

Alright, so let's take a look at how we detect inefficiencies within the models.
I'll share some of our tools and techniques we use on a daily basis.

HPC 101



Funded by
the European Union

Destination Earth

implemented by

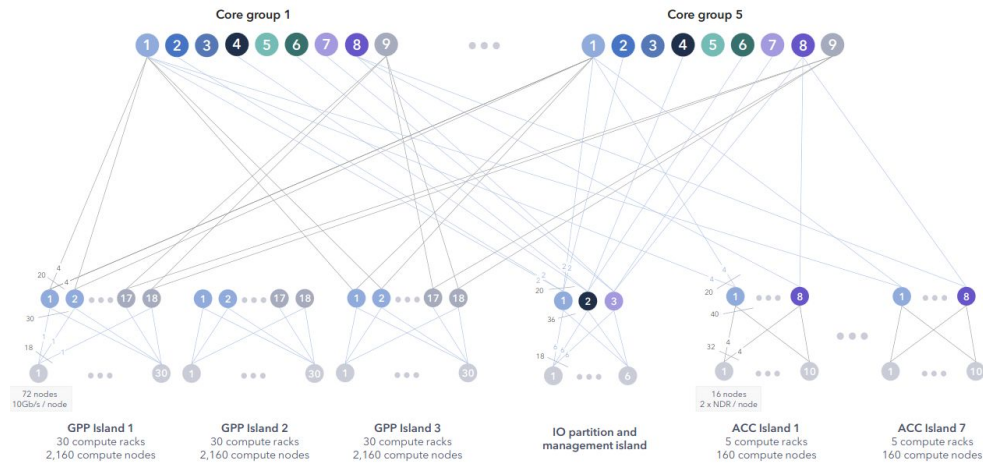


34

However, in order to do so I'll give a super brief introduction to high performance computing,

As you will need to know how supercomputers work to understand the tools.

HPC 101

Funded by
the European Union**Destination Earth**

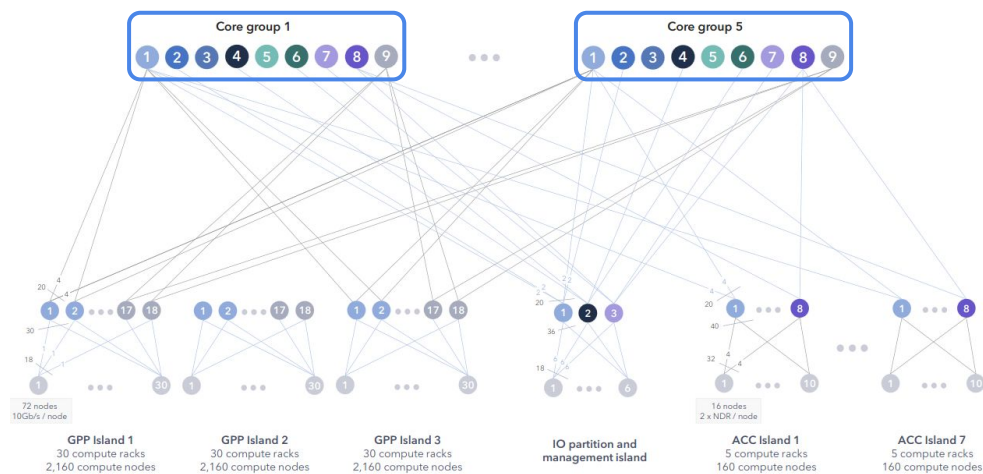
implemented by



35

This is an overview of the MareNostrum 5 system topology.
It might look a bit daunting, but don't worry it's secretly quite simple.

HPC 101

Funded by
the European Union**Destination Earth**

implemented by



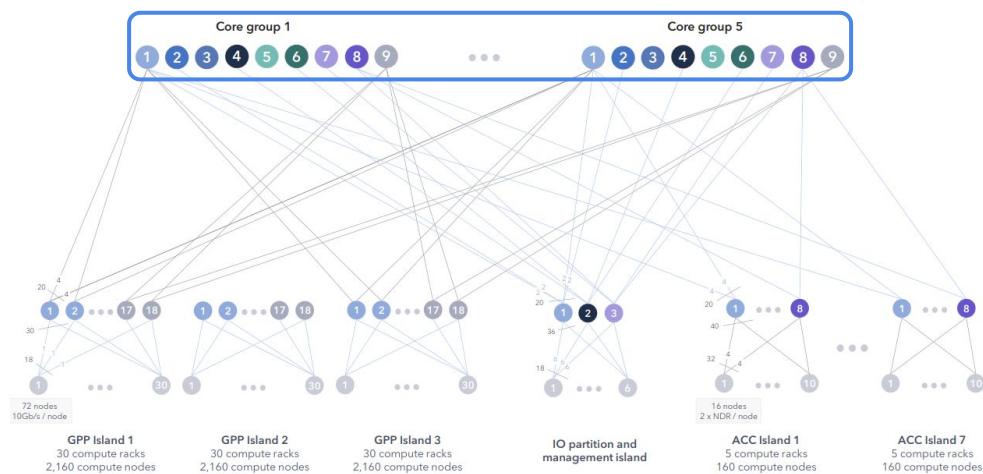
36

These core groups are basically groupings of different types of nodes that you can use.

A node is basically the smallest building block of a supercomputer, which consists of a CPU, Memory, and a Network card.

Some also come with GPUs.

HPC 101

Funded by
the European Union**Destination Earth**

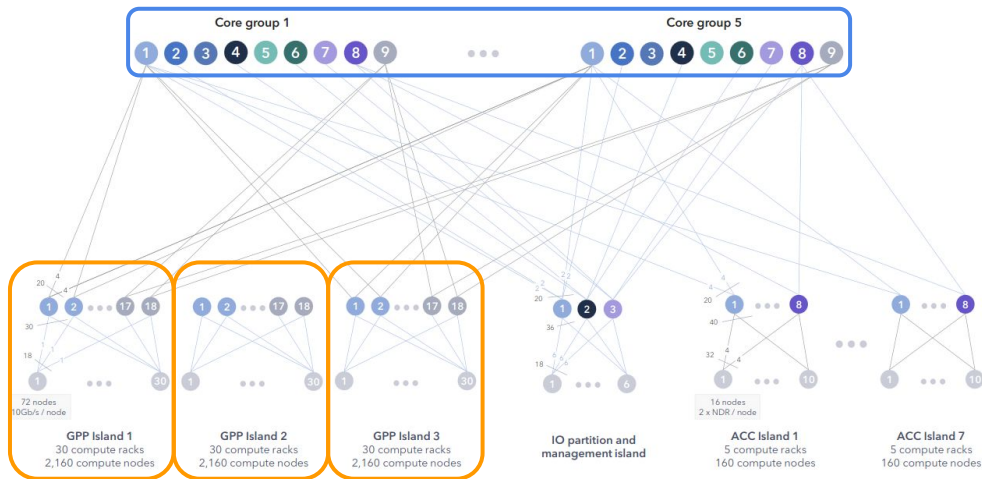
implemented by



37

In this example we will use the core groups as the main controller of the machine which instructs the nodes that we see below.
So we say that this blue rectangle is one entity.

HPC 101

Funded by
the European Union**Destination Earth**

implemented by

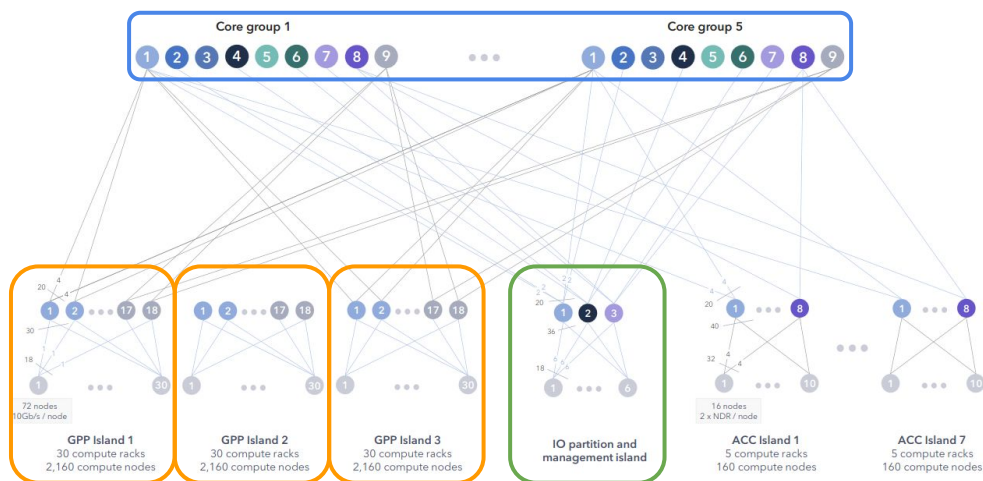


38

Then on the bottom left we have the GPP or so called "General Purpose Partition" islands.

These are groups of nodes that only have CPUs and no GPUs.

HPC 101

Funded by
the European Union**Destination Earth**

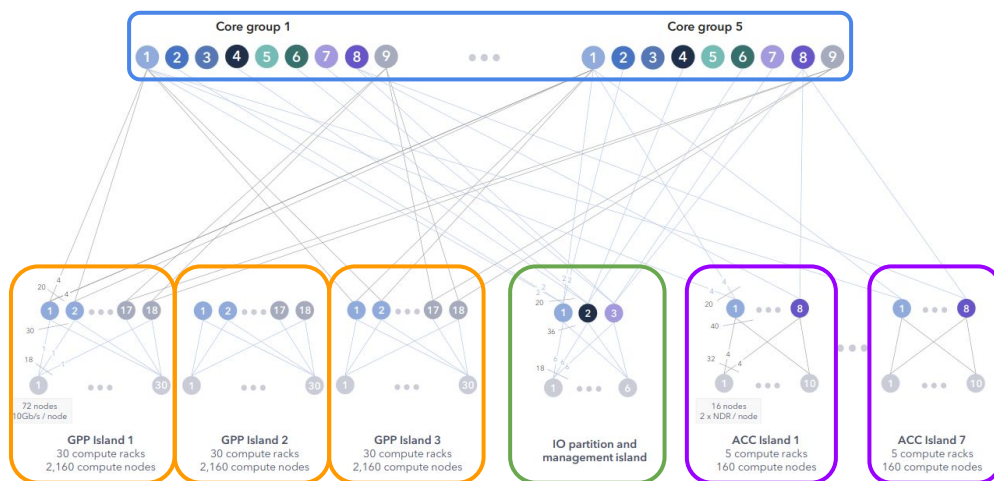
implemented by



39

In the middle there is the IO partition, which manages reading and writing to files within storage.

HPC 101

Funded by
the European Union**Destination Earth**

implemented by

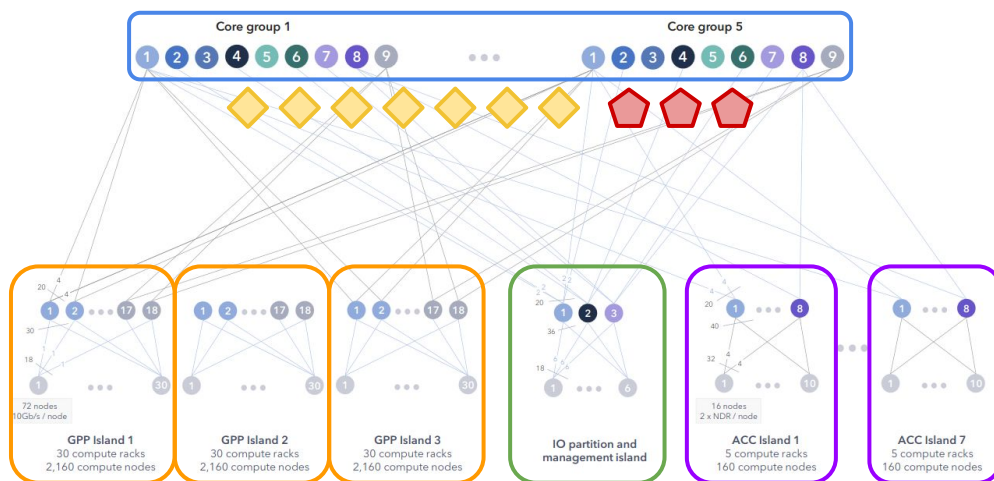


40

And on the right side we have the ACC or so called "Accelerated Partition" islands. These are groups of nodes which have both CPUs and GPUs available. Thence, when using the supercomputer, it is these nodes that you use for GPU computing.

For each of these islands, we simply say that each rectangle is one node, just as we did for the core groups.

HPC 101

Funded by
the European Union

Destination Earth

implemented by



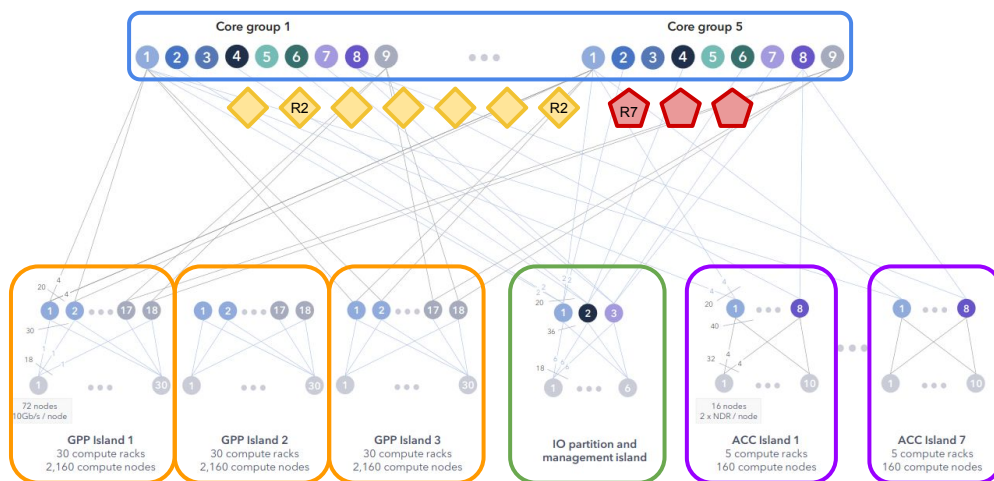
41

Let's add some work to the picture.

The yellow diamonds are CPU work, and the red hexagons are GPU work.

So these are parts of a problem that we want to compute.

HPC 101

Funded by
the European Union**Destination Earth**

implemented by

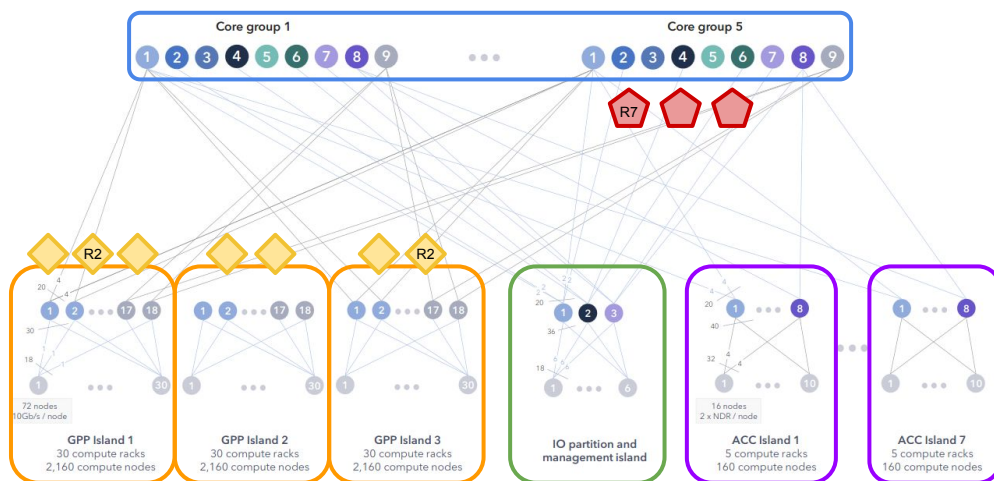


42

Sometimes, a piece of work needs the results from another piece of work before it can start.

I've written that as R followed with a number, where the number represents the island it needs the data from.

HPC 101

Funded by
the European Union**Destination Earth**

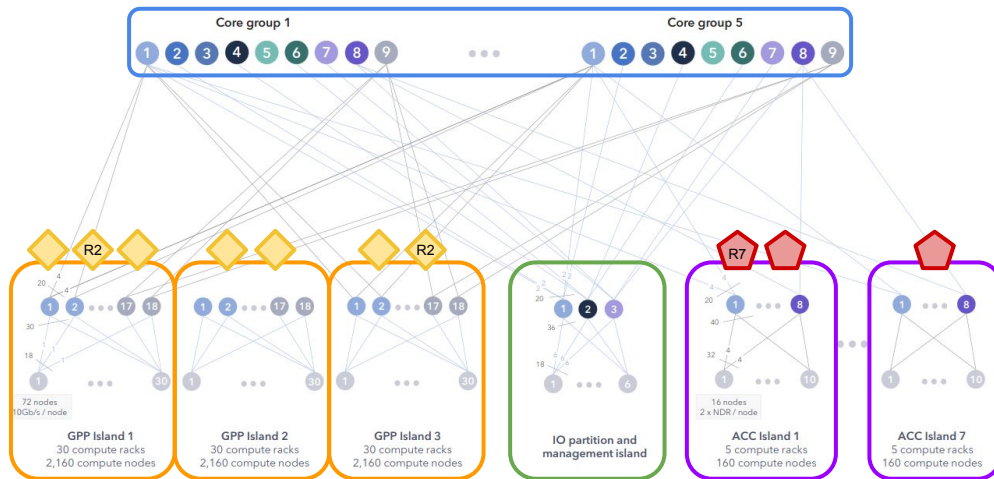
implemented by



43

Let's distribute the CPU work among our nodes.

HPC 101

Funded by
the European Union**Destination Earth**

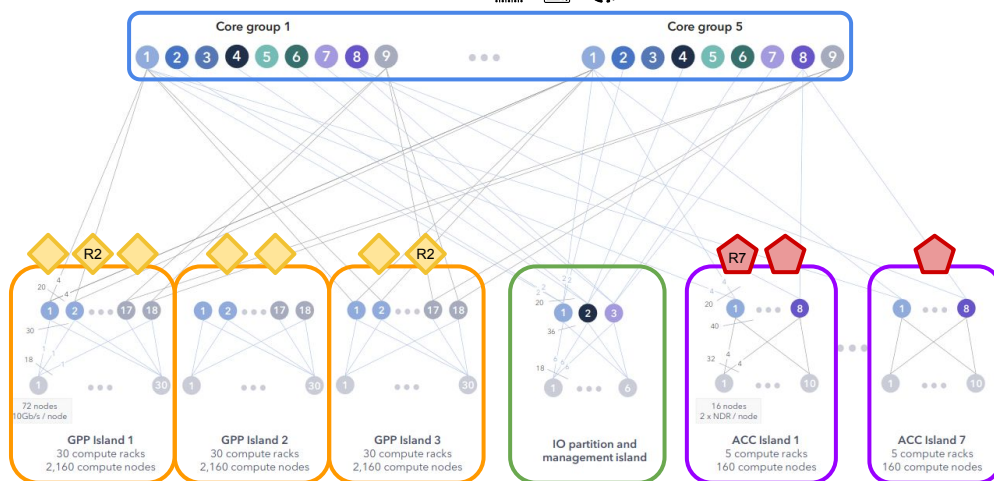

implemented by



44

And do the same for our GPU work.

HPC 101

Operations:   Funded by
the European Union

Destination Earth

implemented by



45

Within a supercomputer, nodes can do a couple different operations.
For this simplified example, we limit us to:

Computing work



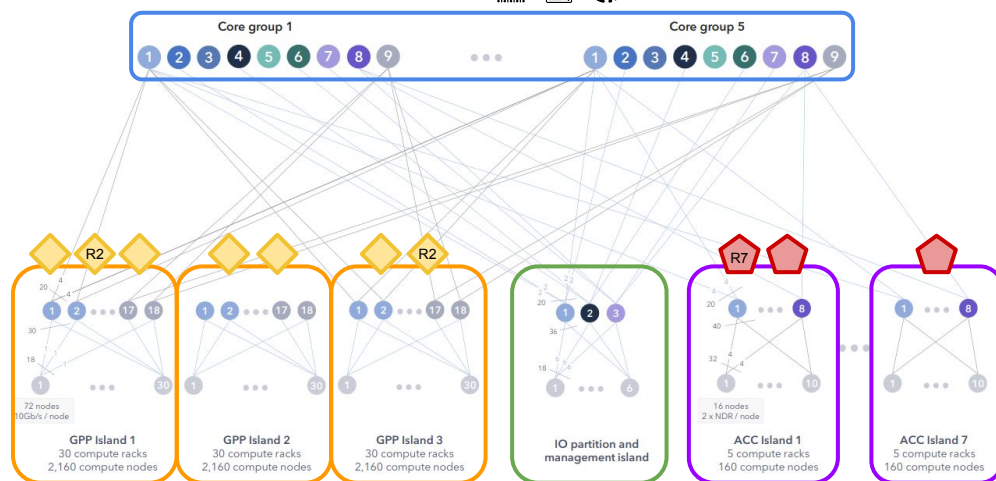
Transferring data

And waiting

Icons:

- Time left: https://www.flaticon.com/free-icon/time-left_66163
- Microprocessor: https://www.flaticon.com/free-icon/microprocessor_5905354
- Cloud: https://www.flaticon.com/free-icon/cloud_860276

HPC 101

Operations:   Funded by
the European Union

Destination Earth

implemented by



46

To help you understand where inefficiencies within supercomputers come from, we'll play out this little example.

On the right top side we keep track of what cycle we are currently working in.

A node can only do one thing in a cycle, which is compute work, transfer data, or wait.

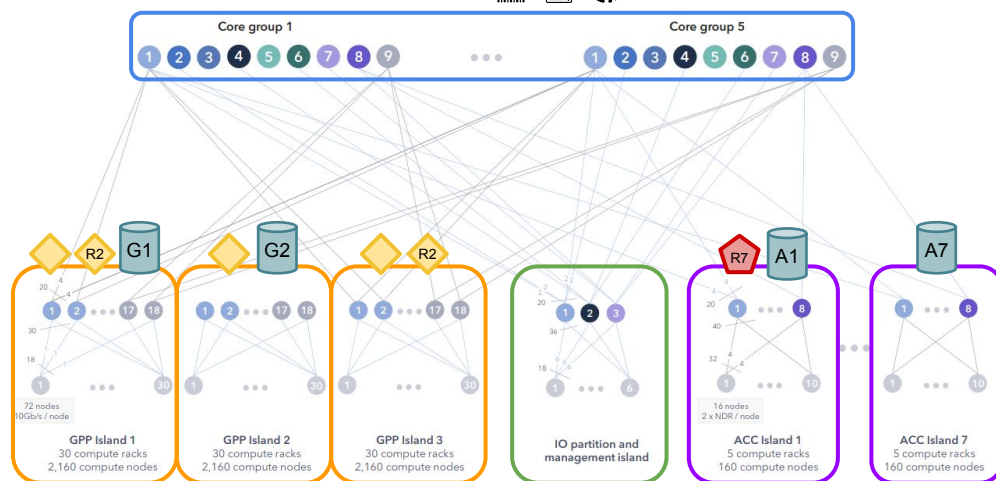
So we get into cycle one and the nodes get busy.

Watch what happens when we go into cycle 1.

HPC 101

Operations:   

Cycle 1

Funded by
the European Union

Destination Earth

implemented by




47

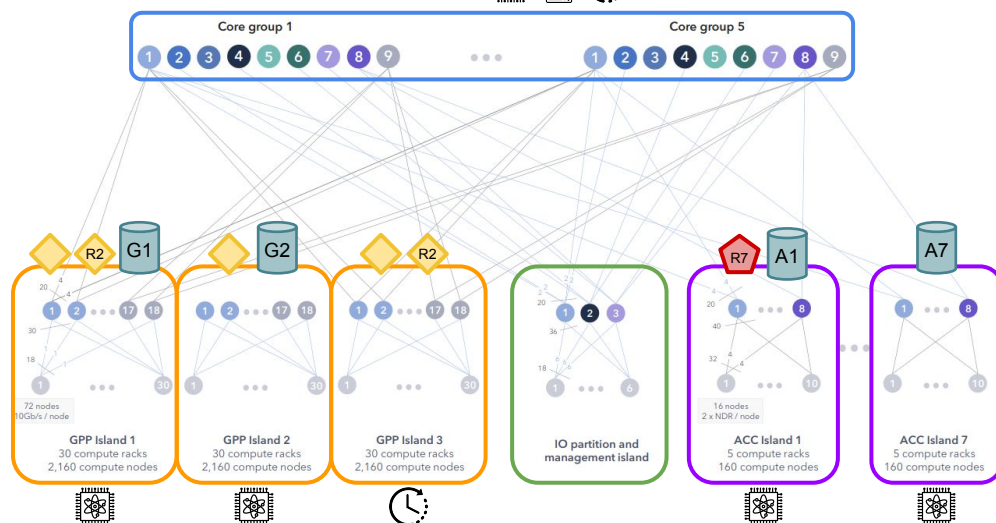
As you can see, some of the icons changed from shape into a cylinder.
We'll use this to represent solutions to the processed work.
Within the cylinder we write where the work was performed,
So for the GPP partitions we have a G,
And for the ACC partitions we have an A.

So what operations did we perform?

HPC 101

Operations:   

Cycle 1



Well, the GPP 1 computed their problem into a solution.

The same applies for GPP 2.

For GPP 3, nothing happened.

It had to wait. Why?

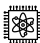


Well, because it needs the results from GPP 2 before it can start computing.

So you can already see that waiting for each other is inefficient.

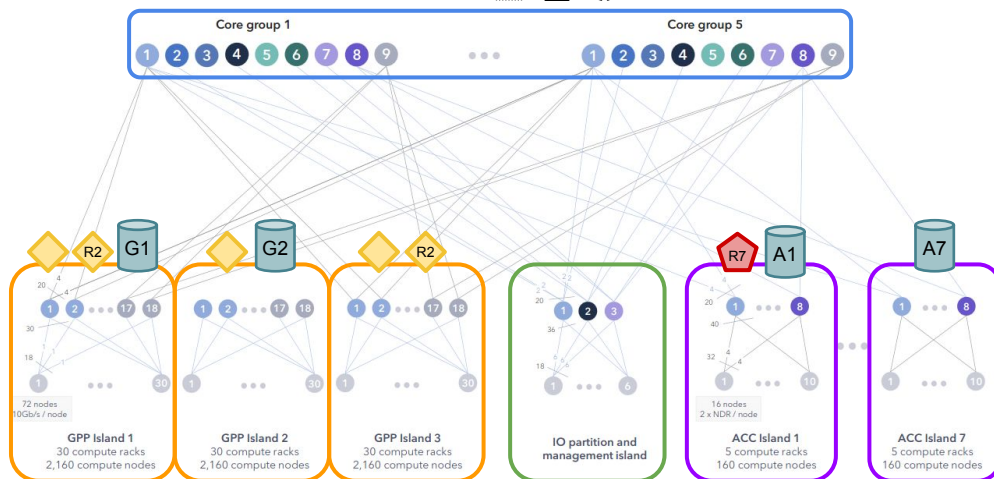
Then on the right side we have ACC 1 which could compute.

And the same is the case for ACC 2.

HPC 101

Operations:   

Cycle 2

Funded by
the European Union

Destination Earth

implemented by



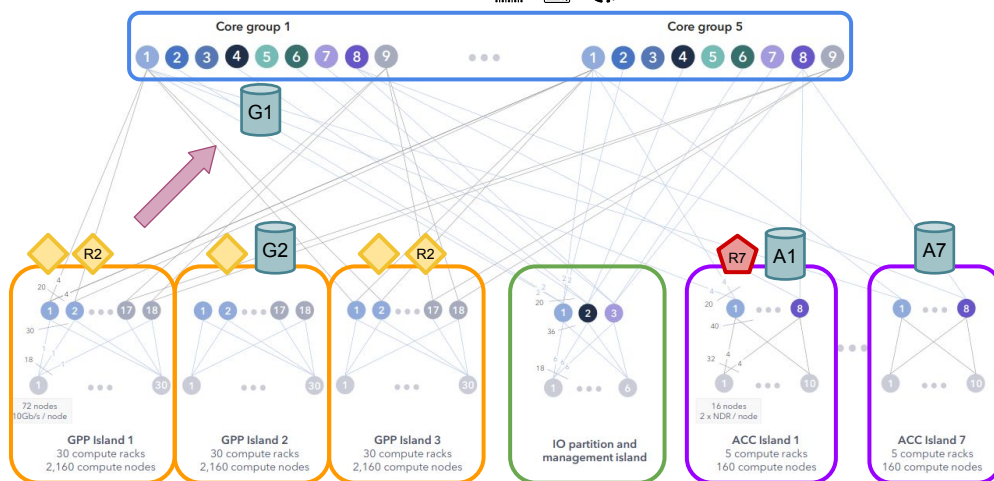
49

Then we go to cycle 2.

HPC 101


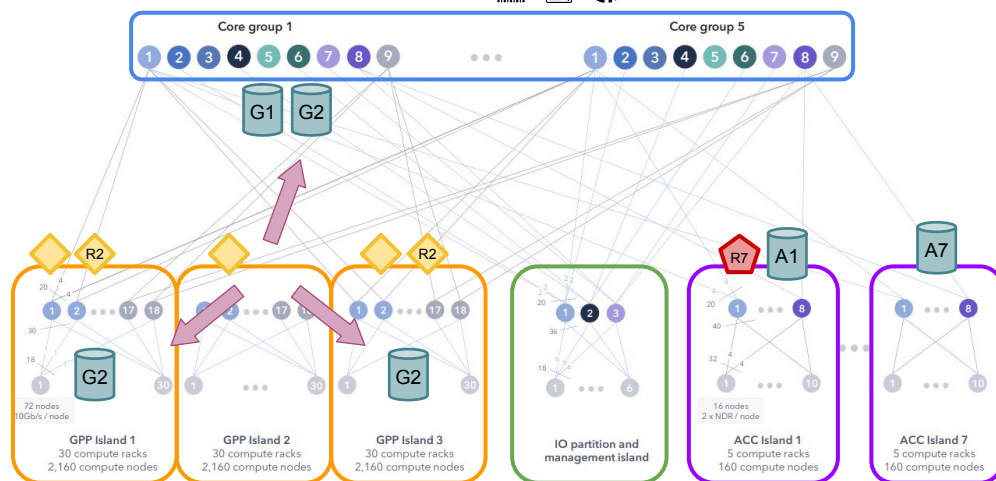
Operations:   

Cycle 2



GPP 1 node has data that can be returned to the controller.
Thus it's performing a data transfer operation.

HPC 101

Operations:   Funded by
the European Union

Destination Earth

implemented by



51

The same applies for GPP 2.

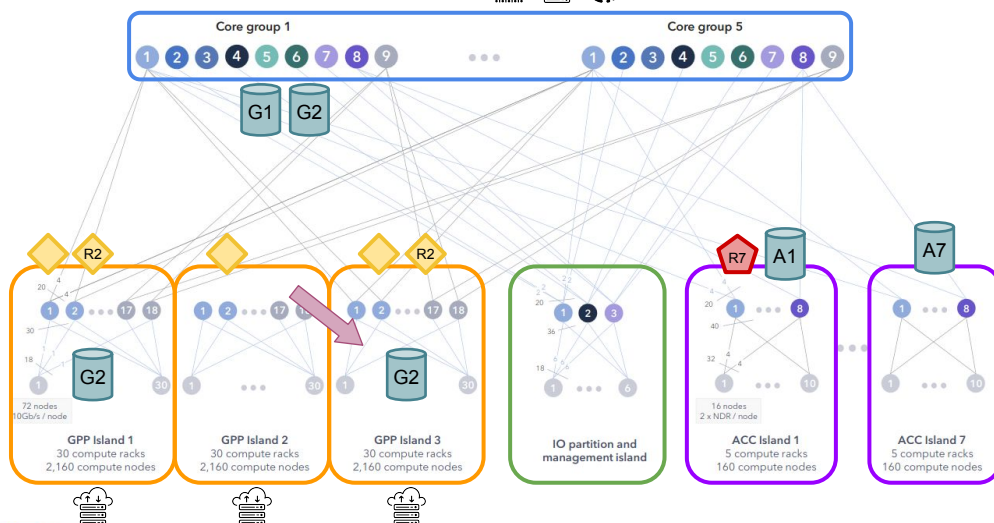
However, now we see that GPP 1 and GPP 3 need input data from GPP 2.

Thus, we not only transfer the results back, but we also send them to the other nodes.

HPC 101

Operations:   

Cycle 2

Funded by
the European Union

Destination Earth

implemented by



52

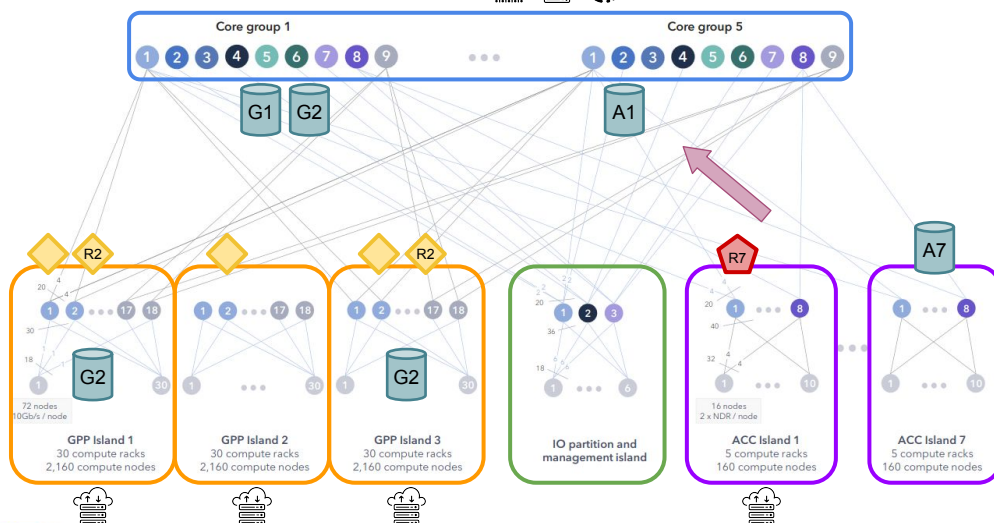
Then for GPP 3, you might think it waits, but now it spends its time on data transfer as well.

But then the transfer is receiving data.

HPC 101

Operations:   

Cycle 2

Funded by
the European Union

Destination Earth

implemented by



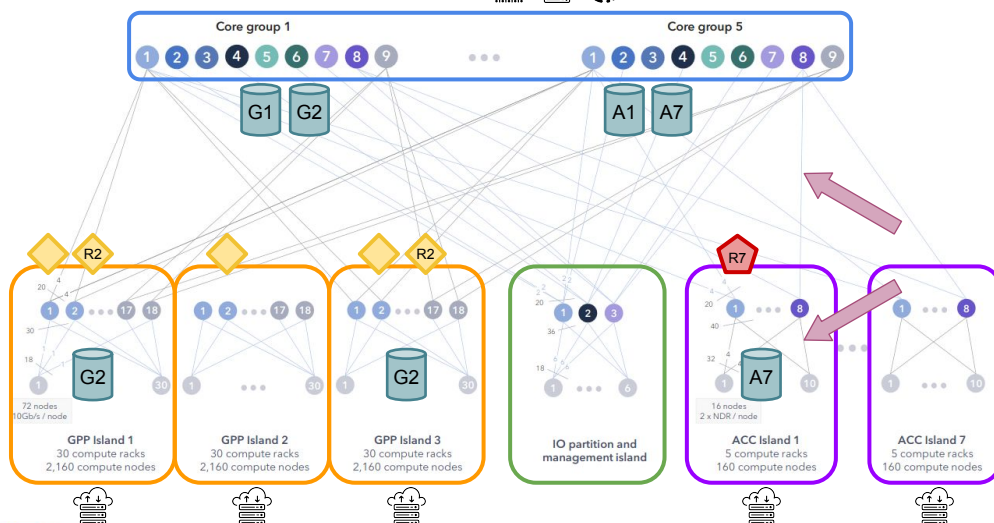
53

For ACC 1 we have results, which can also be send back to the controller.

HPC 101

Operations:   

Cycle 2

Funded by
the European Union

Destination Earth

implemented by



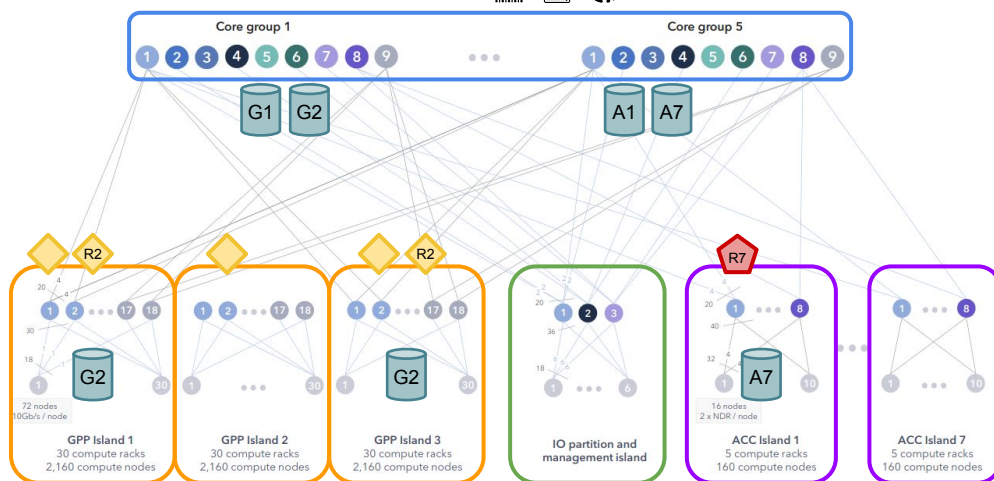
54

And ACC 7 has results that ACC 1 needs, thus we see again a data transfer that goes to both the controller and another node.

HPC 101

Operations:   

Cycle 3

Funded by
the European Union

Destination Earth

implemented by



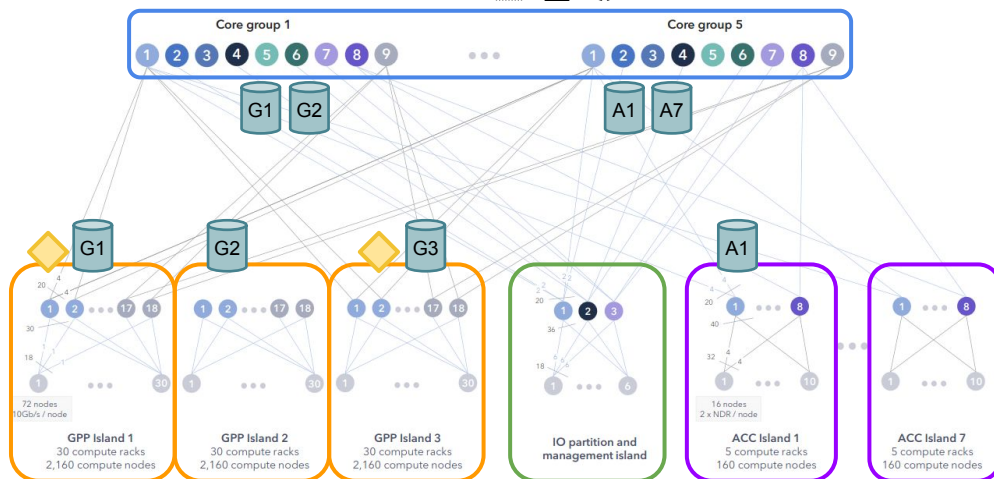
55

And we continue to cycle 3.

HPC 101

Operations:   

Cycle 3

Funded by
the European Union

Destination Earth

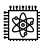
implemented by



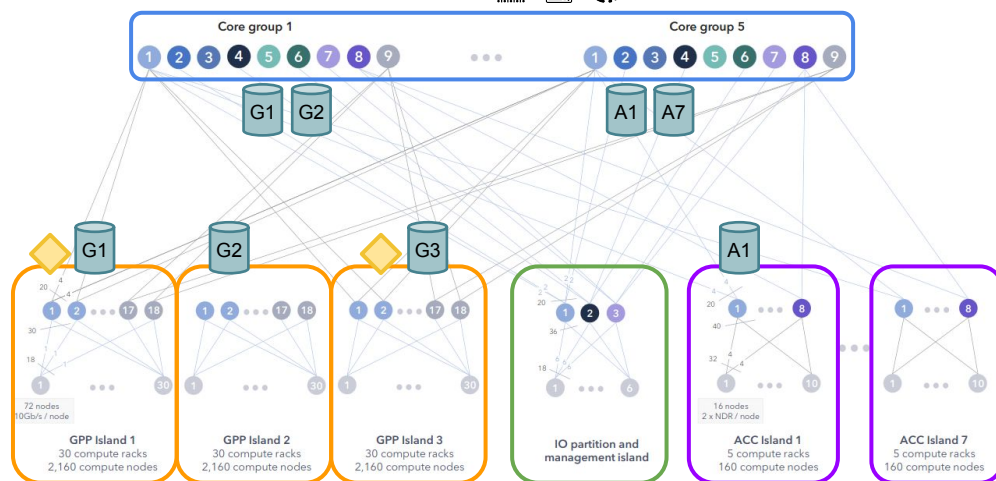
56

We again perform the operations, nothing needed to be transferred.

HPC 101

Operations:   

Cycle 3

Funded by
the European Union

Destination Earth

implemented by



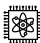


57

And thus we see, again that all nodes could spend time computing, except for ACC 7. Why? Well, because ACC 7 has no work anymore.

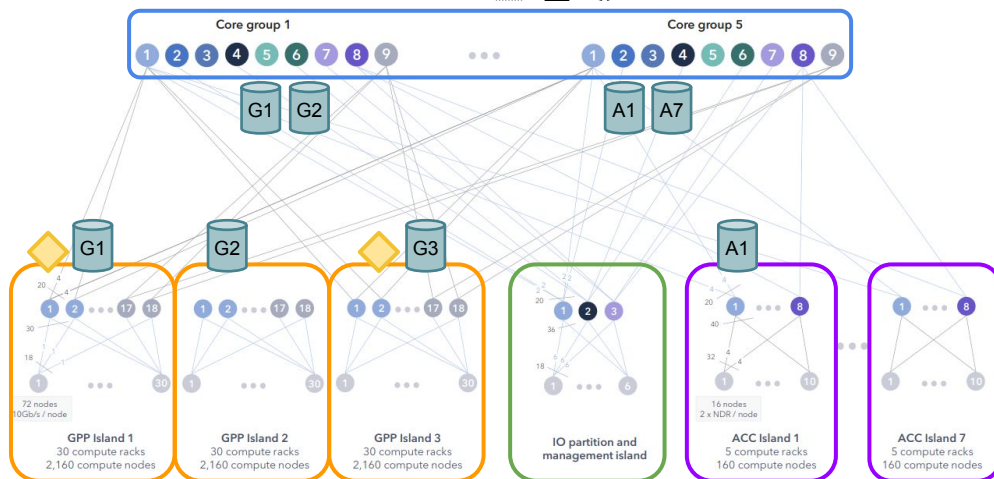
So this is another inefficiency that we have to watch out for.

We could this the load balancing, where we have an imbalance if some nodes have more work to do than other nodes.

HPC 101

Operations:   

Cycle 4

Funded by
the European Union

Destination Earth


implemented by



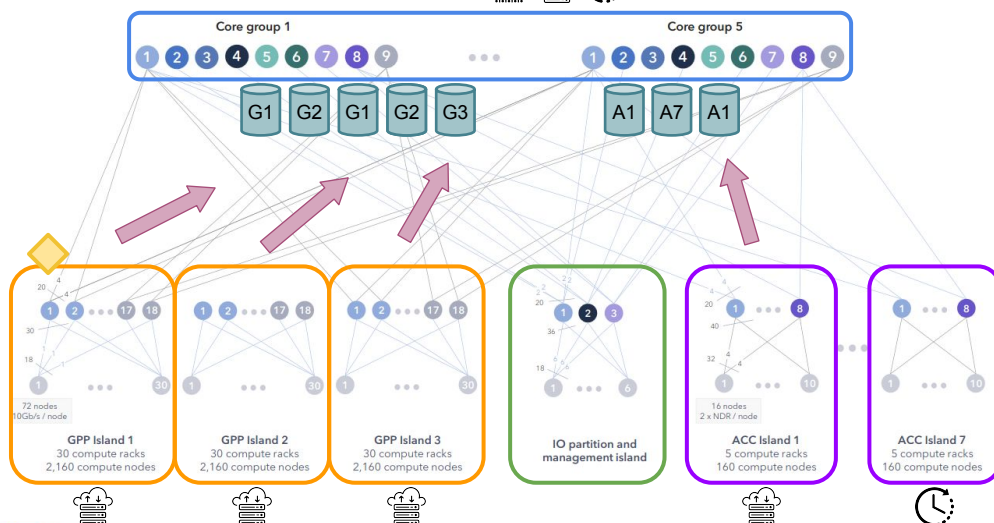
58

Going to cycle 4.

HPC 101

Operations:   

Cycle 4

Funded by
the European Union

Destination Earth

implemented by



59

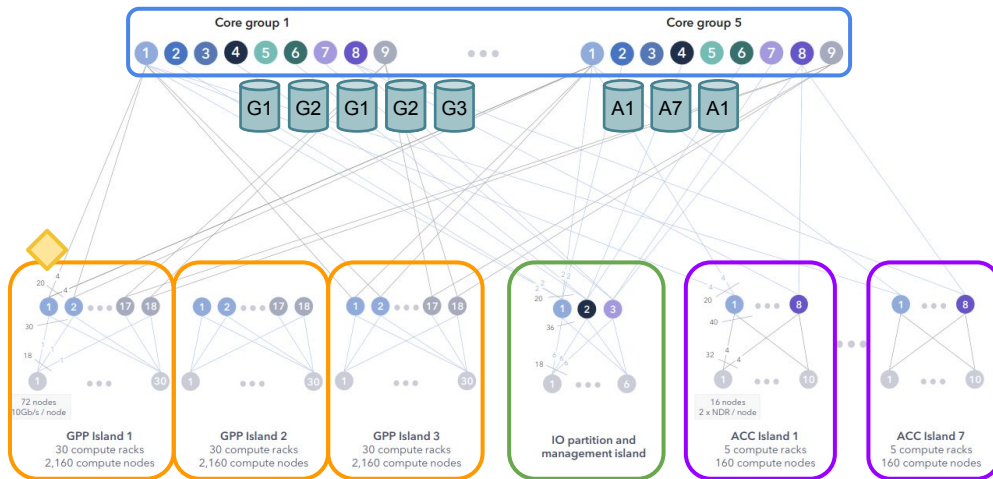
All of the results that we computed in cycle 3 need to be transferred to the controller.
There is no node that requests any input from another node.

ACC 7 again waits.

HPC 101

Operations:   

Cycle 5

Funded by
the European Union

Destination Earth


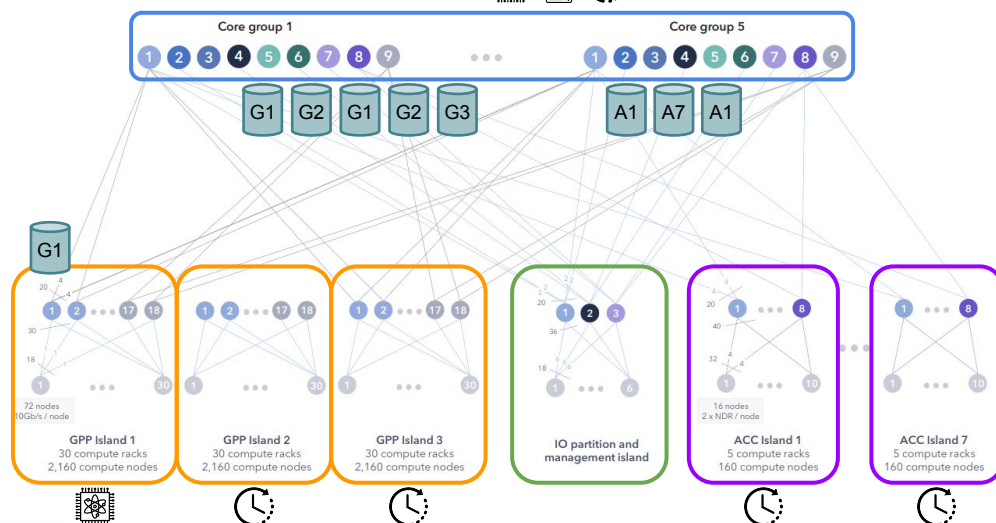
implemented by



60

Then for cycle 5, everything looks quite empty.

HPC 101

Operations:   

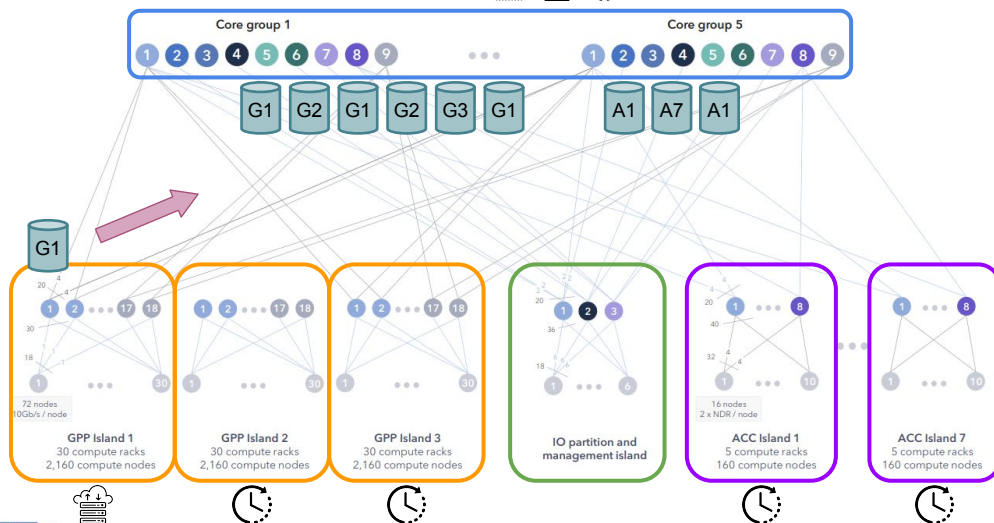
GPP 1 computes the last piece of work it has, but all the other nodes spend their cycle waiting.

And now you can see that the load imbalance can have quite a big impact on the system.

HPC 101

Operations:   

Cycle 6

Funded by
the European Union

Destination Earth


implemented by



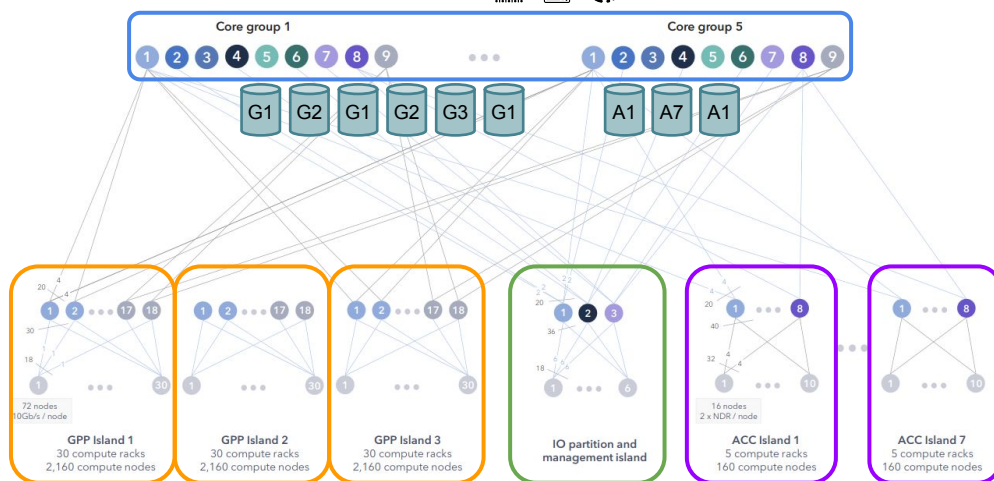
62

In cycle 6 we move the result of GPP 1 to the controller.

HPC 101

Operations:   

Cycle 6

Funded by
the European Union

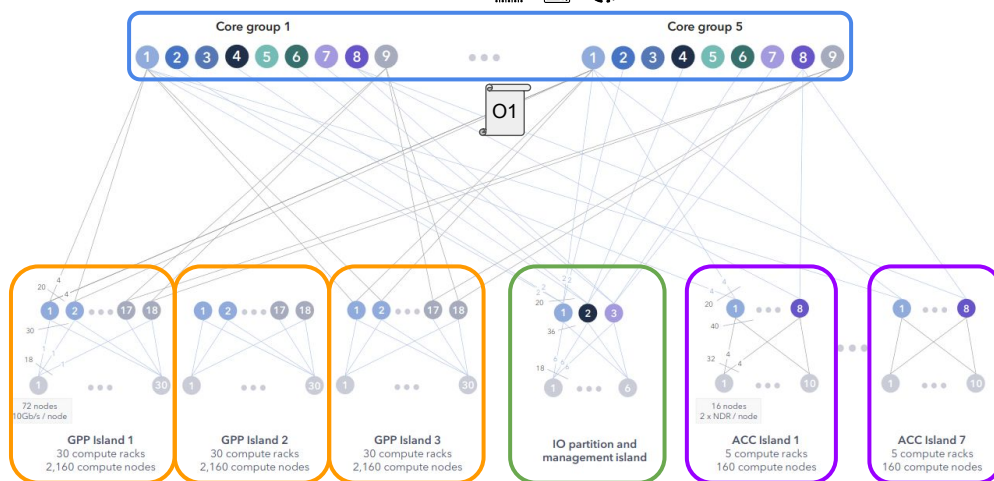
Destination Earth

implemented by



63

HPC 101

Operations:   Funded by
the European Union**Destination Earth**

implemented by



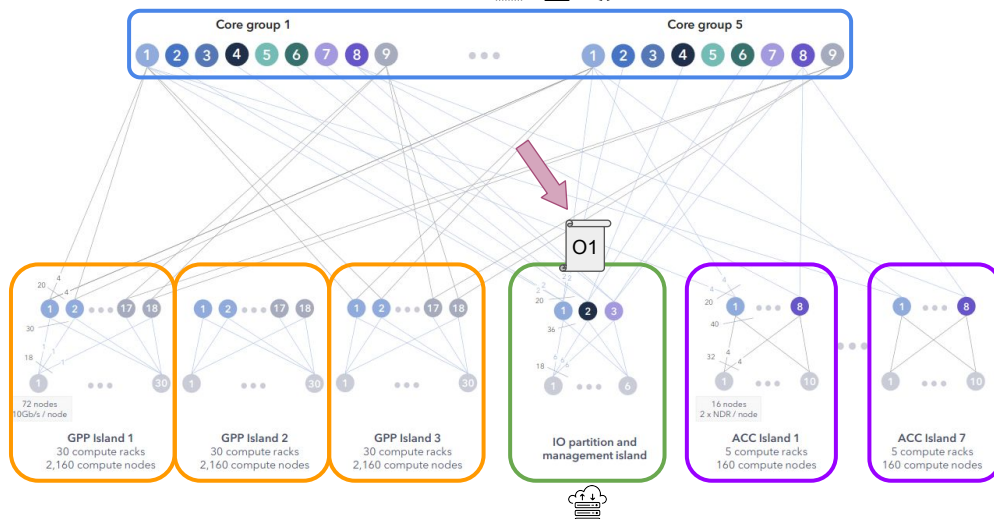
64

And then in cycle 7, we convert all the results in output that can be written to a file.

HPC 101

Operations:   

Cycle 8

Funded by
the European Union

Destination Earth

implemented by



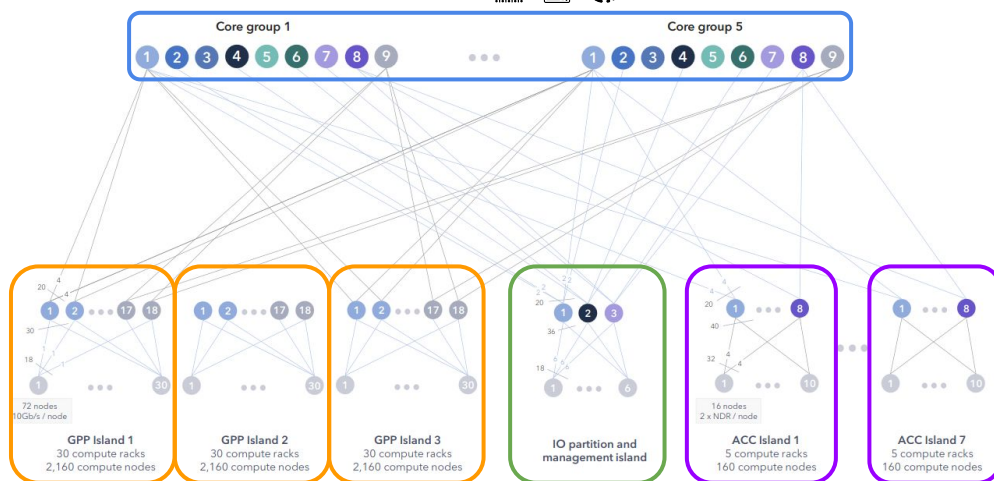
65

This is moved to the IO partition in cycle 8.

HPC 101

Operations:   

Cycle 9

Funded by
the European Union

Destination Earth

implemented by



66

And in cycle 9 we write the file to disk, finishing the execution.

As you can see, there is a lot of communication going on with the sending and receiving of data.

Performance Profiling

Funded by
the European Union**Destination Earth**

implemented by



67

Now that you sort of know how a high-performance computer works, we can talk about Performance Profiling.

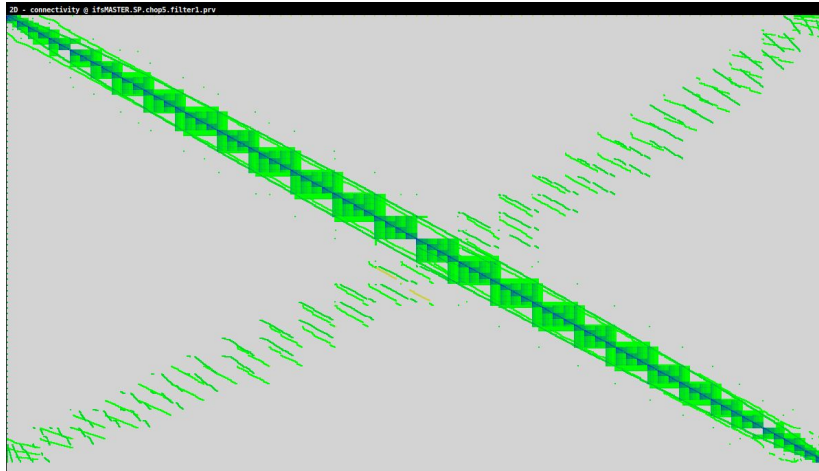
This is basically the task of analyzing the model in order to find potential inefficiencies.

Most of the times, we visualize the data we collect in order to make it easier.

Hence I will show you many visualizations and explain what you are looking at.

As a first example, the communication pattern that you saw on the previous slide is one of those visualizations.

Performance profiling



Funded by
the European Union

Destination Earth

implemented by



So in order to perform performance profiling, we need some tools.

Performance profiling

Tools to use:



Funded by
the European Union

Destination Earth

implemented by



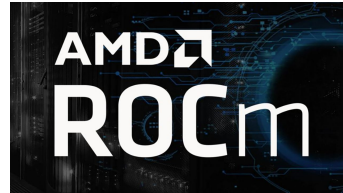
69

So in order to perform performance profiling, we need some tools.

Performance profiling

Tools to use:

- GPU → Depends on platform!



Funded by
the European Union

Destination Earth

implemented by



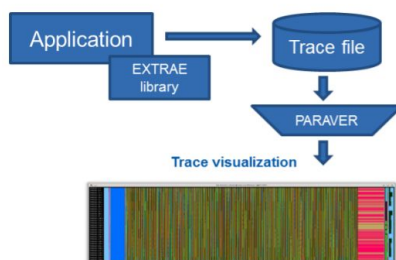
70

For the GPU, it really matters what GPU vendor you have.
NVIDIA has their Nsight platform, while AMD has their ROCm platform.

Performance profiling

Tools to use:

- GPU → Depends on platform!
- CPU → Extrae & Paraver



Funded by
the European Union

Destination Earth implemented by



71

For CPU, you usually can use a tool for multiple platforms, like Intel and AMD. We have two in-house open-source tools, called Extrae and Paraver. Extrae is used to collect information on the model. And Paraver is used to visualize that data.

In this visualization you can see how Extrae acts as a library that is loaded during the application execution.

This produces a so-called trace file of the execution.

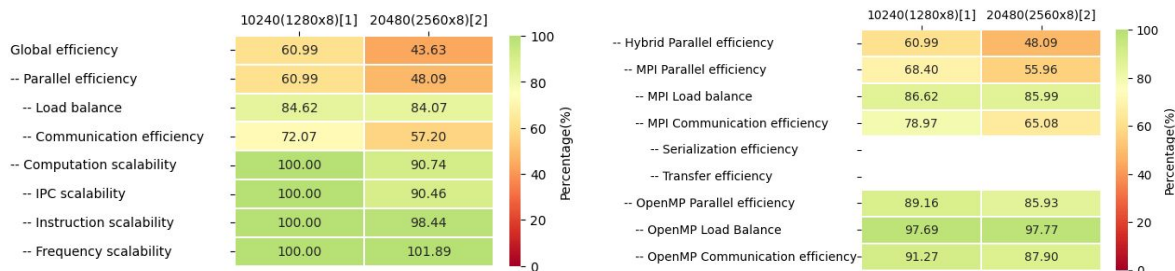
Then we feed the trace file to Paraver, which visualizes it in a timeline.

Later you'll see more examples of these trace files, so don't worry if it sounds too vague.

CPU Profiling



- CPU Model factors



Funded by
the European Union

Destination Earth implemented by



72

The first visualization I want to show you, is what we call "Model factors". Normally it is 1 big table, but due to the slides I chopped it up in 2 parts.

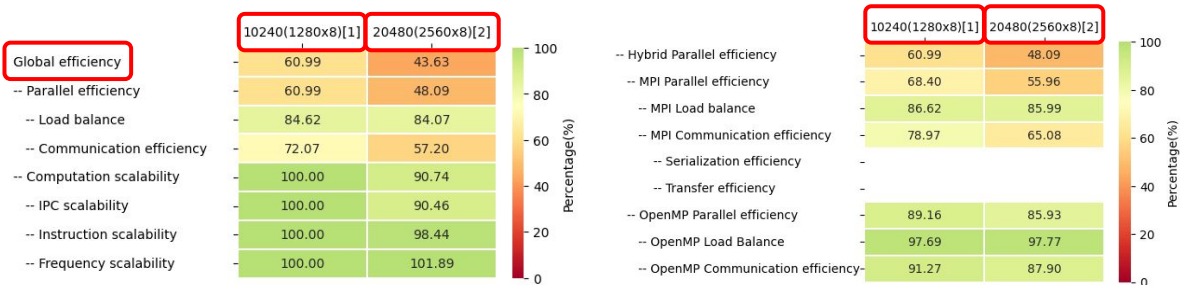
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- CPU Modelfactors



Funded by
the European Union

Destination Earth implemented by



73

The Modelfactors compute the efficiency of your code for different workloads. So in these rows you can see the global efficiency for two system setups, namely 10240 and 20480 cores.

Basically 80 and 160 nodes on HPC2020, 128 cores per node. These numbers correspond to the number of nodes that we have used, so the 2nd column distributes the work between more parts of the supercomputer than the first one.

You can see that the global efficiency drops from 60.99% to 43.63%.

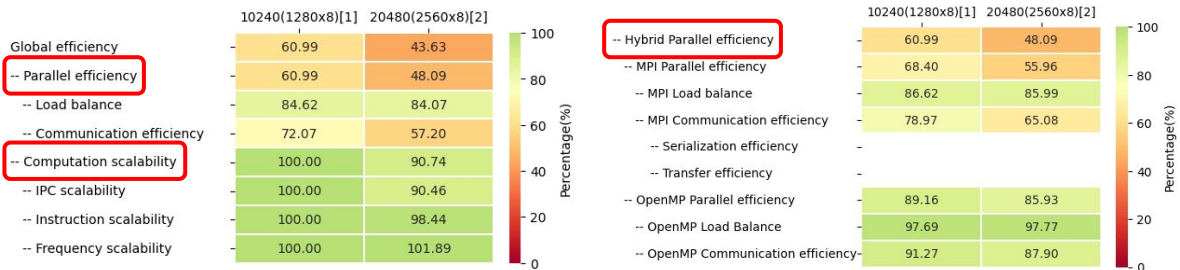
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- CPU Model factors



Funded by
the European Union

Destination Earth implemented by



74

The rows underneath are sub-sections, which each describe a different aspect of the efficiency.

Together they form the percentage of the row above,

so the 60.99% of Parallel Efficiency, with the 100% of Computation Scalability, with the 60.99% Hybrid Parallel Efficiency form the Global Efficiency.

And basically, you can go through each row to figure out where the inefficiency is coming from.

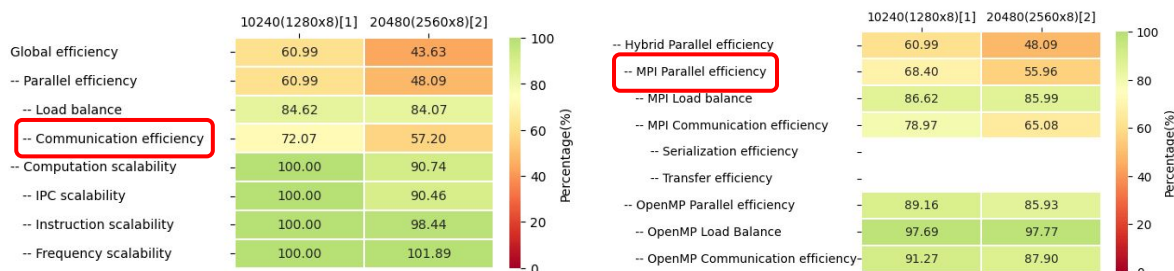
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- CPU Model factors



Funded by
the European Union

Destination Earth implemented by



75

In this case, we can see that the inefficiency mainly comes from a decrease in Communication Efficiency and MPI Parallel Efficiency. So clearly, something in the communication of this model needs to be approved.

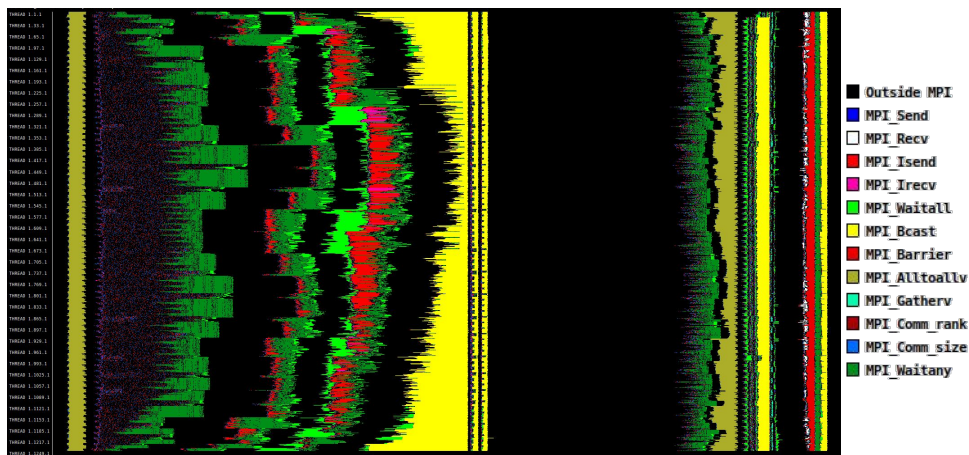
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver

Funded by
the European Union**Destination Earth** implemented by

76

So we now know that we need to focus on the communication happening between the different nodes.

We again use our trace file of Extrae to visualize the communication that's happening.

That gives us an image like this.

It might look very scary to you, so let me explain what we are looking at.

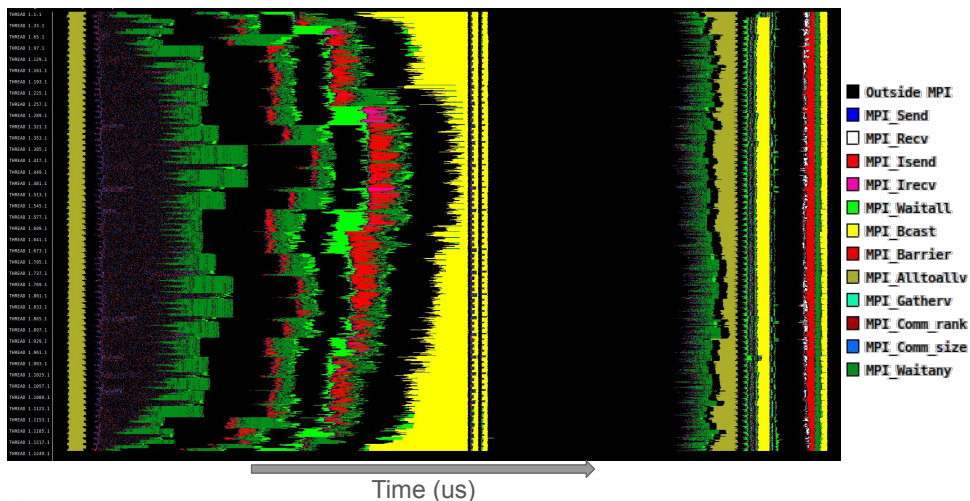
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



Over the horizontal axis, we have time.
It is basically a timeline what we are looking at.

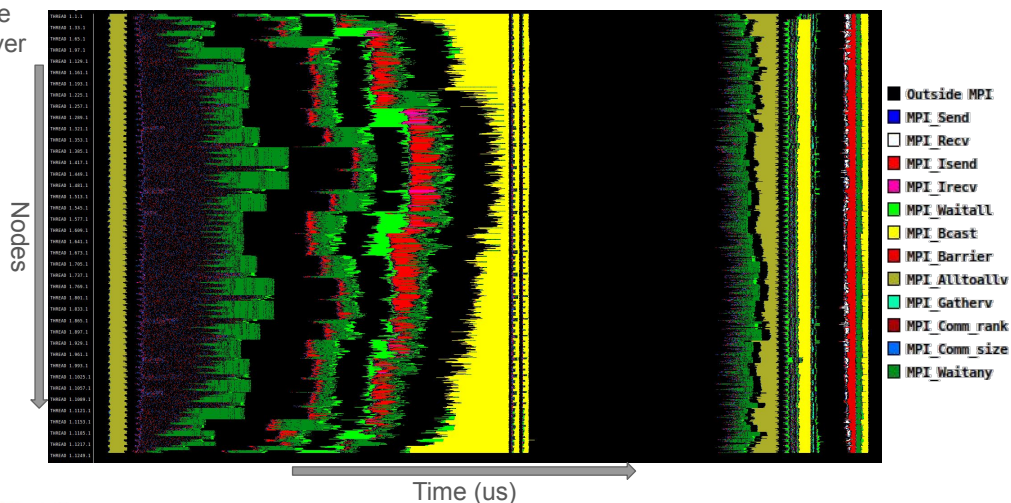
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver

Funded by
the European Union**Destination Earth** implemented by

78

Then over the vertical axis, we have the different threads.

These are basically individual workers that are grouped per node, as in our HPC101 example.

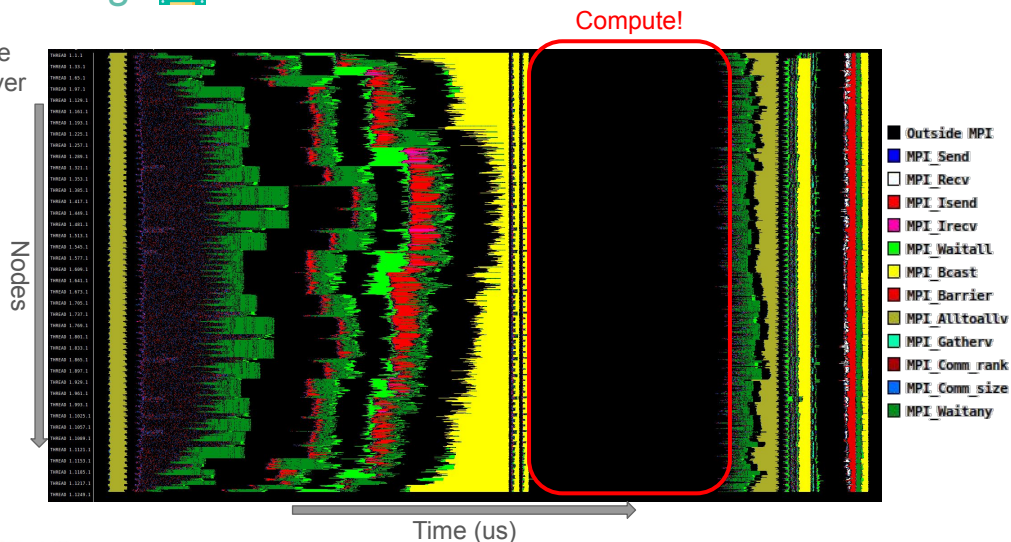
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



80

The black part is time spend outside of MPI, in other words, time where we compute or wait!

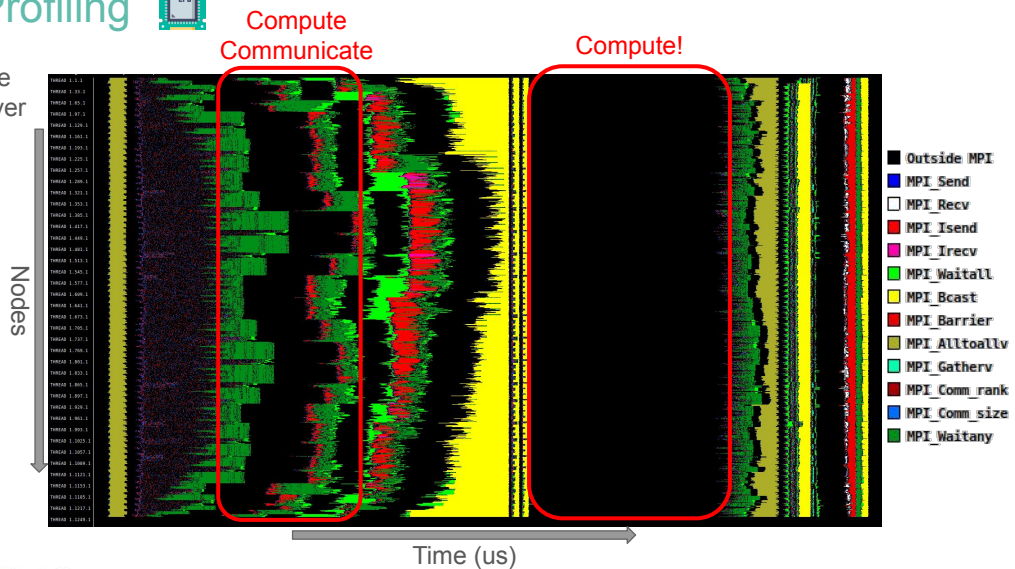
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



81

On the left we also see a very empty piece, where there is still some communication happening.

So here we have a mix of Compute, Communicate, and waiting.

This becomes even more clear if we switch our view to something called "Useful duration".

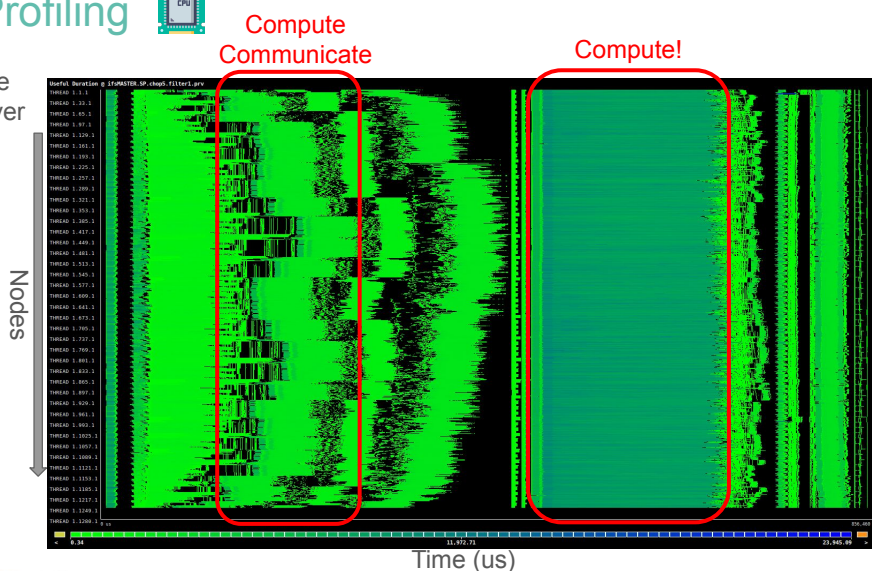
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



Here you have a scale going from green to blue, where the color corresponds to how much computation there is.

So you can see that the empty black part on the right is indeed mostly compute.

And on the right side we have an alternation between compute and communicate.

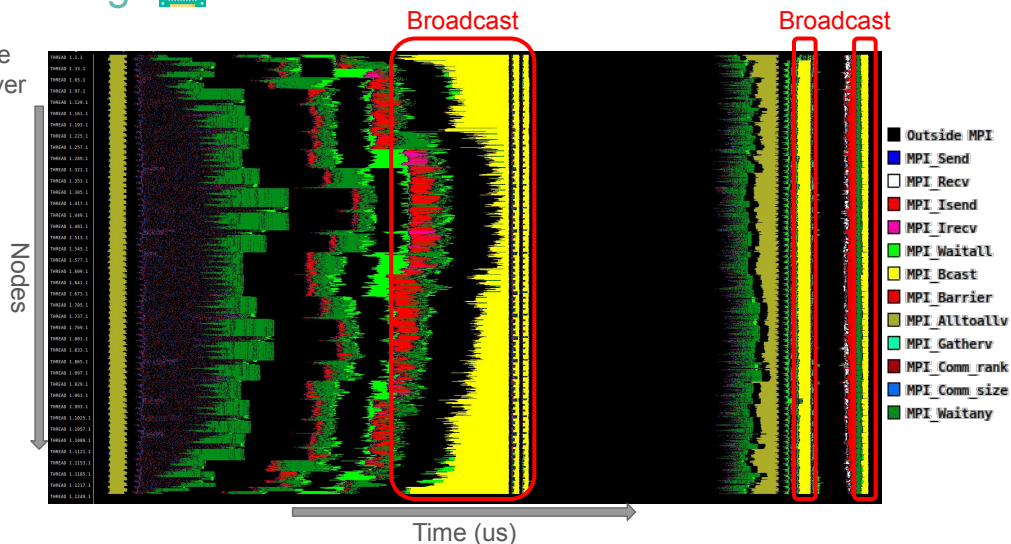
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



84

The Broadcast communication type is one of those.

In this case, some nodes send data which has to be received by all nodes before the process can continue.

Thus results in blocking communication, which makes nodes wait for each other instead of computing the problem.

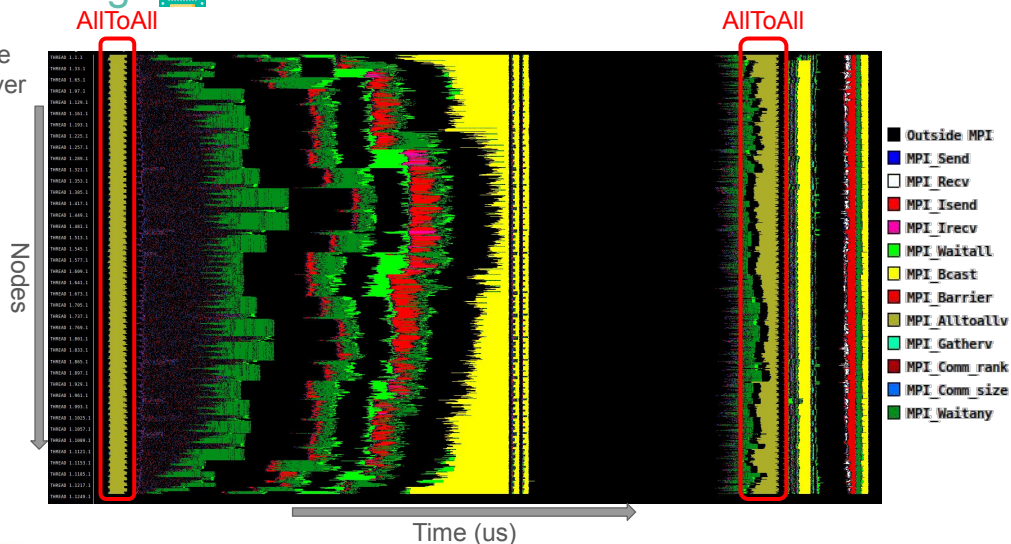
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



85

The AllToAll communication type has exactly the same problem, where the only difference is that all nodes are sending now as well. This as well cause all nodes to wait for each other.

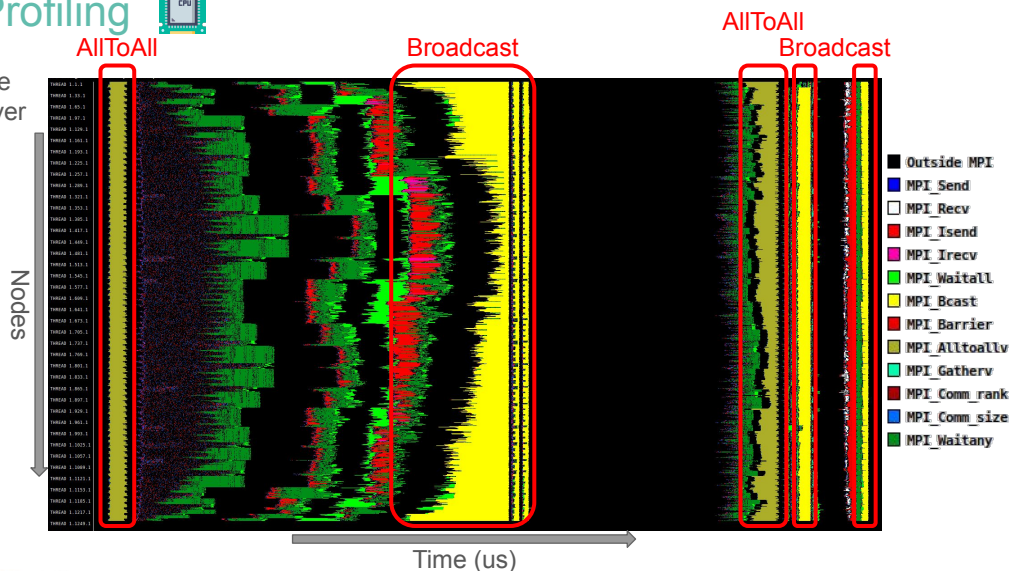
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver

Funded by
the European Union

Destination Earth implemented by



86

So take a close look at these yellow bits in the timeline.

When we switch to the useful duration view, where we see computation, you will see that nothing happens there.

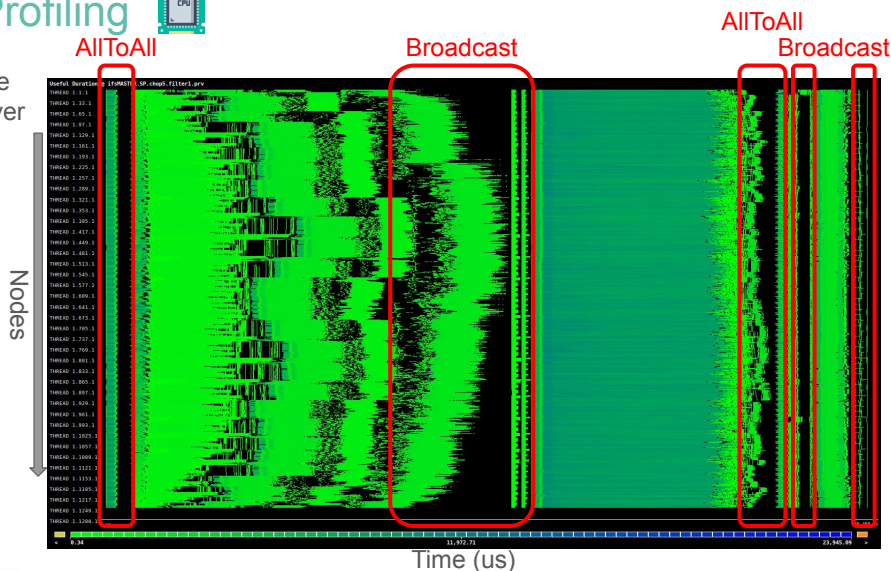
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver

Funded by
the European Union**Destination Earth** implemented by

87

So here you can see that these communication types completely kill performance.

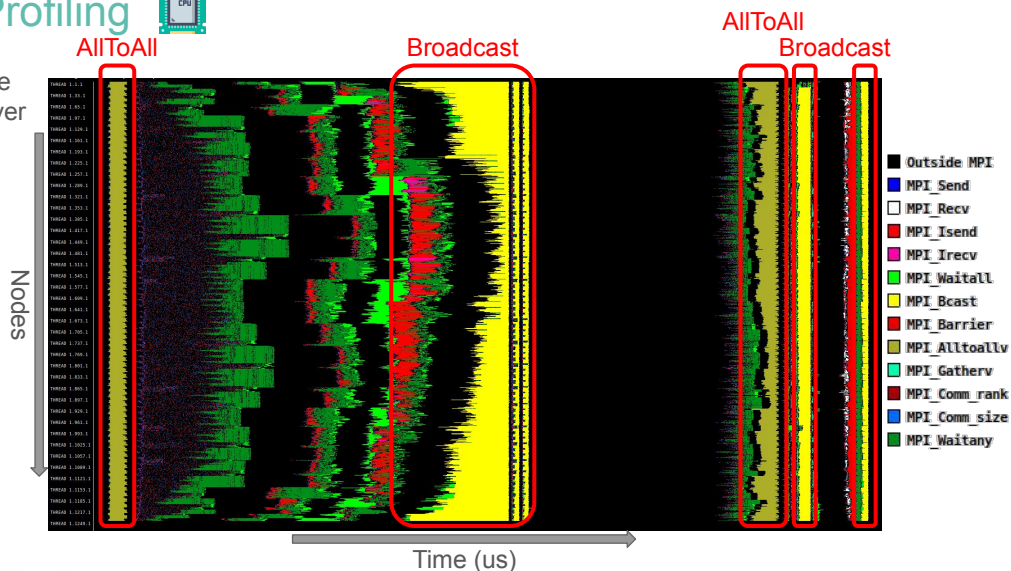
CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

CPU Profiling



- Extrae
- Paraver



Funded by
the European Union

Destination Earth implemented by



88

Thus these are patterns that we look out for when analyzing a timestep of a model.

CPU Icon:

https://www.flaticon.com/free-icon/cpu_984391

GPU Profiling

Platform matters!



Funded by
the European Union

Destination Earth

implemented by



89

Okay, so that was how we detect communication inefficiencies for CPUs.
Let's take a look at what we can do for GPUs.
First thing you need to know: platform matters!

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

Platform matters!



GPU



Funded by
the European Union

Destination Earth implemented by



90

It matters what type GPU you have.
Generally you have NVIDIA, AMD, or Intel.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

NVIDIA:

https://en.wikipedia.org/wiki/CUDA#/media/File:Nvidia_CUDA_Logo.jpg

AMD:

<https://github.com/ROCm>

GPU Profiling

Platform matters!



GPU

Compiler



Funded by
the European Union

Destination Earth implemented by



91

Then it also matters what compiler you have,
since compilers can output debug information, but all do it in a different way.
Some examples here are OpenMPI, GNU, or Cray.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

NVIDIA:

https://en.wikipedia.org/wiki/CUDA#/media/File:Nvidia_CUDA_Logo.jpg

AMD:

<https://github.com/ROCm>

Intel:

https://es.wikipedia.org/wiki/Archivo:Intel_logo_%282006-2020%29.svg

OpenMPI:

<https://www.open-mpi.org/doc/v3.1/man1/mpifort.1.php>

Cray:

<https://insidehpc.com/wp-content/uploads/sites/2/2013/11/images.jpg>

GNU:

https://es.wikipedia.org/wiki/GNU#/media/Archivo:Heckert_GNU_white.svg

GPU Profiling

Platform matters!



Funded by
the European Union

Destination Earth implemented by



94

And lastly, it also matters what programming language your source code is in. Because Fortran for example, has less debug information during compilation time than when you compile C++.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

NVIDIA:

https://en.wikipedia.org/wiki/CUDA#/media/File:Nvidia_CUDA_Logo.jpg

AMD:

<https://github.com/ROCm>

Intel:

https://es.wikipedia.org/wiki/Archivo:Intel_logo_%282006-2020%29.svg

OpenMPI:

<https://www.open-mpi.org/doc/v3.1/man1/mpifort.1.php>

Cray:

<https://insidehpc.com/wp-content/uploads/sites/2/2013/11/images.jpg>

GNU:

https://es.wikipedia.org/wiki/GNU#/media/Archivo:Heckert_GNU_white.svg

C++:

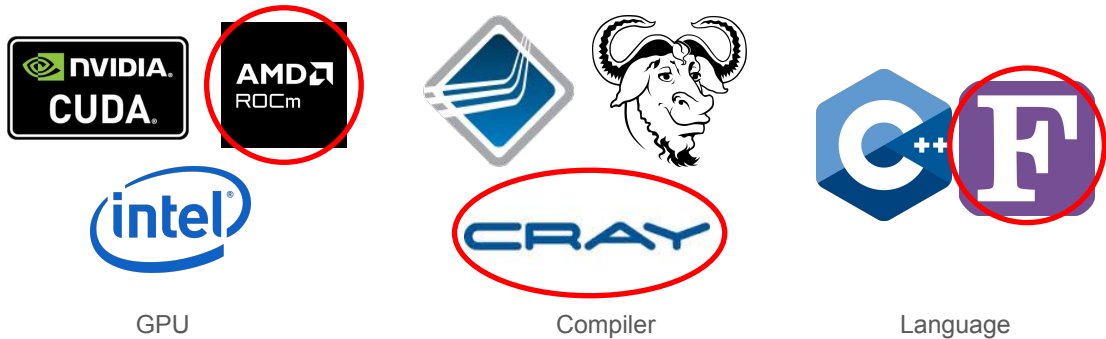
https://en.wikipedia.org/wiki/C%2B%2B#/media/File:ISO_C++_Logo.svg

Fortran:

https://commons.wikimedia.org/wiki/File:Fortran_logo.svg

GPU Profiling

Platform matters!



Funded by
the European Union

Destination Earth implemented by



96

In this example, I'll take an example case for a report I made last year on the LUMI supercomputer in Finland.

This means we have AMD GPUs, the Cray FTN compiler, for our Fortran code base.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

NVIDIA:

https://en.wikipedia.org/wiki/CUDA#/media/File:Nvidia_CUDA_Logo.jpg

AMD:

<https://github.com/ROCm>

Intel:

https://es.wikipedia.org/wiki/Archivo:Intel_logo_%282006-2020%29.svg

OpenMPI:

<https://www.open-mpi.org/doc/v3.1/man1/mpifort.1.php>

Cray:

<https://insidehpc.com/wp-content/uploads/sites/2/2013/11/images.jpg>

GNU:

https://es.wikipedia.org/wiki/GNU#/media/Archivo:Heckert_GNU_white.svg

C++:

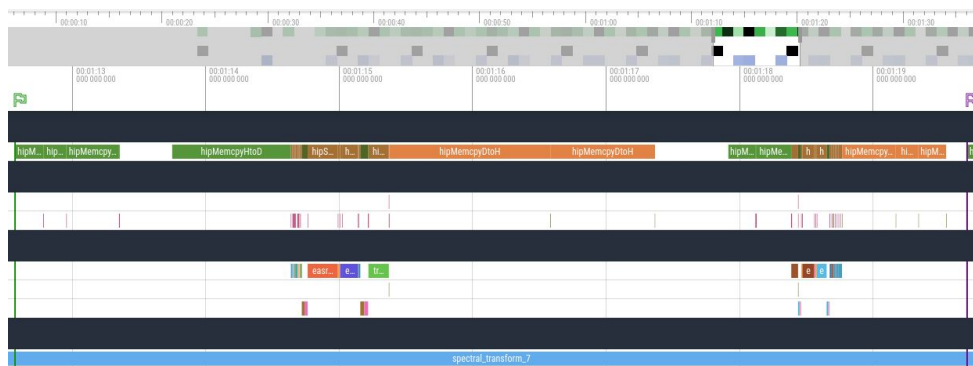
https://en.wikipedia.org/wiki/C%2B%2B#/media/File:ISO_C++_Logo.svg

Fortran:

https://commons.wikimedia.org/wiki/File:Fortran_logo.svg

GPU Profiling

```
rocprof --hip-trace --roctx-trace -o ./rocprof/rocprof.csv $BIN
```



Perfetto visualization



Funded by
the European Union

Destination Earth implemented by



94

With the AMD ROCm stack, you get a tool called Rocprof.
This is the most fundamental tool when profiling on AMD GPUs.
You can also collect a trace for a timeline, just like we did on CPU.
This is again 1 timestep of the code that we are analyzing, which is ecTrans, a spectral transformation code.

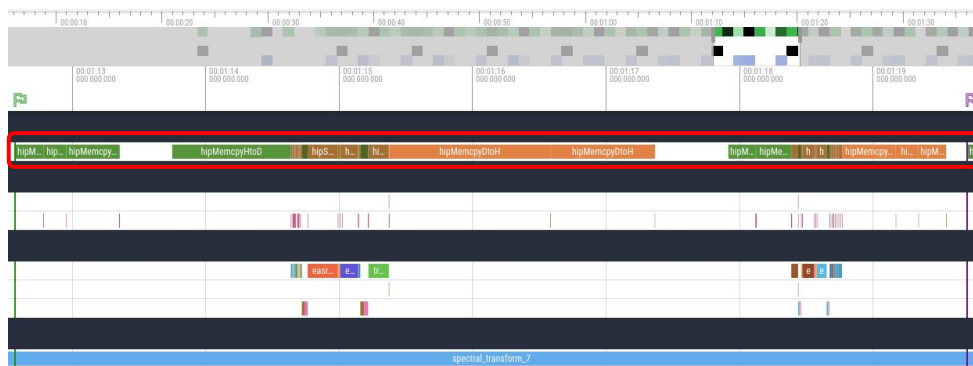
GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling



```
rocprof --hip-trace --roctx-trace -o ./rocprof/rocprof.csv $BIN
```



Perfetto visualization



Funded by
the European Union

Destination Earth implemented by



95

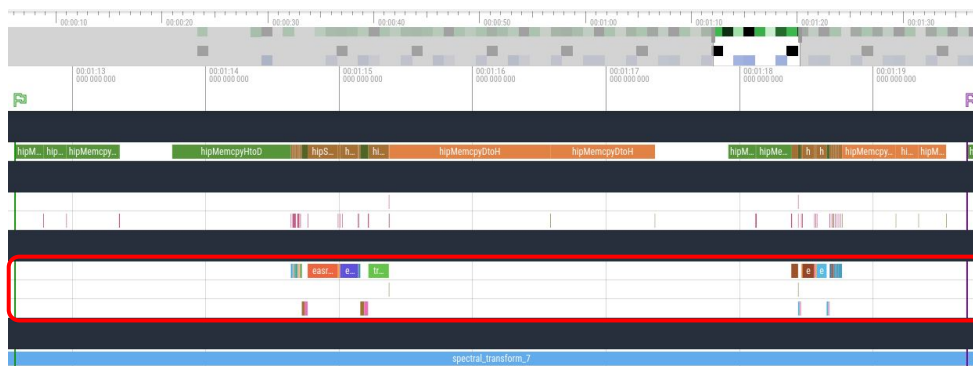
On the top you have the memory transfers that occur between the CPUs and GPUs. For those we don't know, a GPU is controlled from the CPU. So the CPU will transfer the required data and instructions to the GPU, which executes them, and then the results are copied back. These memory copies are very slow compared to the computation, and thus this is usually the main bottleneck for performance.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

```
rocprof --hip-trace --roctx-trace -o ./rocprof/rocprof.csv $BIN
```



Perfetto visualization



Funded by
the European Union

Destination Earth implemented by



96

Here you can see the GPU kernels, which are isolated programs that perform computations.

Ideally, you would like to spend most time executing kernels.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

Funded by
the European Union**Destination Earth**

implemented by



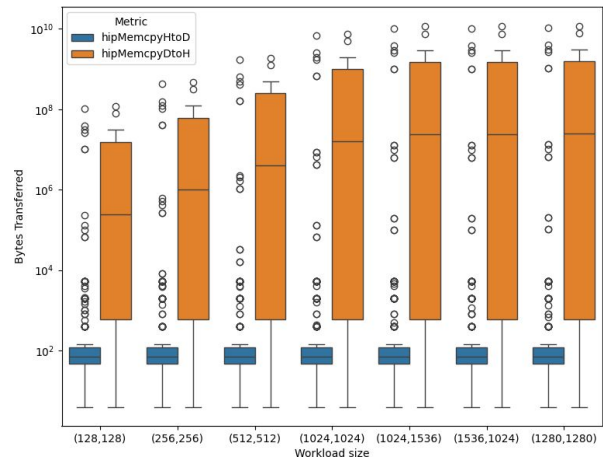
97

It might be a bit difficult to see, but here you have the two computation parts enlarged so it's easier to see the distribution between their runtimes.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

Funded by
the European Union

Destination Earth

implemented by



98

Since memory copies are usually a bottleneck, the first thing we do is parse this JSON into a boxplot as shown on the right side.

Here you see a column for each workload, which is basically how big the areas are that we are updating.

In blue you see the HostToDevice copies, while in orange you see the DeviceToHost copies.

Device in this case refers to the GPU, while the Host is the CPU.

The boxplots capture how big the transfers are that have been made.

And you can clearly see that the CPU to GPU transfers here are super small, while the GPU to CPU transfers are of okay size.

These super small transfers bring an overhead compared to large ones, since you keep calling the API.

Hence, this is causing quite a severe slowdown of the GPU performance, and thus our advice to our customer was to merge these transfers as much as possible.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling



Kernel name	Calls	TotalDurationNs	AverageNs	Percentage
usec18\$sevalh_mod.Sck.L75.1	10	2164195747	216419574	19.8255
trigtoLcudaaware\$trigtoLmod.Sck.L559.4	10	1531345585	153134558	14.0292
efourier_in\$efourier_in_mod.Sck.L61.2.ccs\$ndloop8form	10	1198170052	119817005	10.9760
efourier_out\$efourier_out_mod.Sck.L56.1	10	875120791	87512079	8.0167
epil28\$epil28_mod.Sck.L84.1.ccs\$ndloop8form	10	740841358	74084135	6.7866
expdapl\$expdapl_mod.Sck.L86.1	60	612327357	10213789	5.6139
void real_pro.process_kernel<...>(...)	20	513296858	25664842	4.7021
trigtoLcudaaware\$trigtoLmod.Sck.L563.3	10	479085073	47908507	4.3887
epil18\$epil18_mod.Sck.L90.2	60	299659753	4994329	2.7451
void ip.inverse_length768.SBRR<...>(...)	10	282099617	28209961	2.5834
void real_post_process_kernel.interleaved_1D<...>(...)	20	226225138	11311256	2.0724
void ip.inverse_length512.SBRR<...>(...)	10	216695188	21669518	1.9796
trigtoLcudaaware\$trigtoLmod.Sck.L205.1	10	178045797	17804579	1.6310
evituv\$evituv_mod.Sck.L132.2	10	149741079	14974107	1.3717
efourier_in\$efourier_in_mod.Sck.L58.3.ccs\$ndloop8form	10	148392992	14839299	1.3593
clcdlr\$clcdlr_mod.Sck.L98.1	10	139369528	13936952	1.2767
espnde\$espnde_mod.Sck.L87.1	10	115237936	11523793	1.0557
cltin\$cltin_mod.Sck.L136.2.ccs\$ndloop8form	10	113872491	11387249	1.0431
tritonLcudaaware\$tritonLmod.Sck.L218.1	10	110770720	11077072	1.0147
void ip.forward_length768.SBRR<...>(...)	10	110712645	11071264	1.0142
void ip.forward_length512.SBRR<...>(...)	10	110086079	11008607	1.0085
cltdlr\$cltdlr_mod.Sck.L89.2.ccs\$ndloop8form	10	94305395	9430539	0.8639
clsc\$clsc_mod.Sck.L97.2	10	89873345	8987334	0.8233
evitv\$evitv_mod.Sck.L108.2	10	85680570	8568057	0.7849
evituv\$evituv_mod.Sck.L89.1	10	68517640	6851764	0.6277
evituv\$evituv_mod.Sck.L113.1	10	66336511	6633651	0.6077
clsc\$clsc_mod.Sck.L75.1	10	58245610	5824561	0.5336
epil18\$epil18_mod.Sck.L85.3.ccs\$ndloop8form	60	57369265	956053	0.8258
cltdlr\$cltdlr_mod.Sck.L157.2.ccs\$ndloop8form	10	57060646	5706064	0.5227
cltdlr\$cltdlr_mod.Sck.L156.2.ccs\$ndloop8form	10	22749187	2274918	0.2084
<barrier packet>	20	753461	37673	0.0070
evituv\$evituv_mod.Sck.L170.4	10	97760	9776	0.009
evitvLcomm\$evitvLcomm_mod.Sck.L106.2.ccs\$ndloop8form	10	63081	6308	0.0006
trigtoLcudaaware\$trigtoLmod.Sck.L456.8	10	32960	3296	0.0003
trigtoLcudaaware\$trigtoLmod.Sck.L457.9	10	23040	2304	0.0002

Table 8: Global overview of the eTrans kernels for 10 iterations. The **blue** highlighted rows are the 6 user kernels which are analysed. The **green** highlighted rows have a smaller kernel granularity. The **yellow** highlighted rows are library calls.



Funded by
the European Union

Destination Earth implemented by



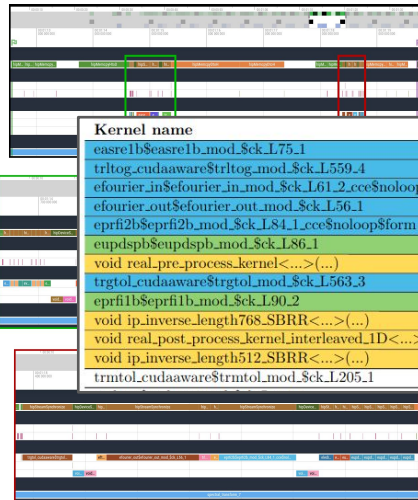
99

Next, we can take a look at these kernels where the computations are happening. The table here to the right shows all kernels in the timestep and how long they took.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling



Kernel name	Calls	TotalDurationNs	AverageNs	Percentage
easre1b\$easre1b_mod.\$ck.L75.1	10	2164195747	216419574	19.8255
trrtog.cudaaware\$trrtog_mod.\$ck.L559.4	10	1531345585	153134558	14.0282
efourier_in\$efourier_in_mod.\$ck.L61.2.cce\$noLoop\$form	10	1198170052	119817005	10.9760
efourier_out\$efourier_out_mod.\$ck.L56.1	10	875120791	87512079	8.0167
epri2b\$epri2b_mod.\$ck.L84.1.cce\$noLoop\$form	10	740841358	74084135	6.7866
eupdspb\$eupdspb_mod.\$ck.L86.1	60	612827387	10213789	5.6139
void real_pre_process_kernel<...>(...)	20	513296858	25664842	4.7021
trrtol.cudaaware\$trrtol_mod.\$ck.L563.3	10	479085073	47908507	4.3887
epri1b\$epri1b_mod.\$ck.L90.2	60	299659753	4994329	2.7451
void ip_inverse_length768_SBR<...>(...)	10	282009617	28200961	2.5834
void real_post_process_kernel_interleaved_1D<...>(...)	20	226225138	11311256	2.0724
void ip_inverse_length512_SBR<...>(...)	10	216095188	21609518	1.9796
trmtol.cudaaware\$trmtol_mod.\$ck.L205.1	10	178045797	17804579	1.6310
<barrier packet>	20	753461	37673	0.0070
evdiv\$evdiv_mod.\$ck.L170.4	10	97760	9776	0.009
evdiv\$evdiv_mod.\$ck.L170.4	10	63081	6308	0.006
trrtol.cudaaware\$trrtol_mod.\$ck.L456.8	10	32960	3296	0.0003
trrtol.cudaaware\$trrtol_mod.\$ck.L457.9	10	23040	2304	0.0002

Table 8: Global overview of the eTrans kernels for 10 iterations. The blue highlighted rows are the 6 user kernels which are analysed. The green highlighted rows have a smaller kernel granularity. The yellow highlighted rows are library calls.



Funded by
the European Union

Destination Earth implemented by



104

Let me enlarge the top rows so you can see better.

These numbers are for 10 time steps, so 10 calls means that there is 1 call per time step.

The blue kernels are the ones that we were investigating, since they looked the most promising as they took more than 3% of the computation time.

The green kernels also took significant time, but they had 60 calls, which means they are called 6 times.

This gives a bit of a skewed view, since you have 6 times the overhead of calling the function.

It could still be interesting to analyze, but with limited time we preferred to analyze the others.

And lastly we have the yellow kernels, which are system or library calls.

This is code that is called from another kernel and is not written by our client.

Hence, they cannot optimize these codes.

So let me show you how we would analyze a single kernel.

We could for example take a look at the top one, called easre1b with around 20% of computation time.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

```
rocprof -i rocprof_counters.txt -o ./rocprof/rocprof.csv $BIN
```

```
Index,KernelName,gpu-id,queue-id,queue-index,pid,tid,grd,wgr,lds,scr,vgpr,sgpr,fbar,sig,obj,GPUBusy,Wavefronts,VALUInsts,VFetchInsts,VWriteInsts,VALUUtilization,VALUBusy,WriteSize,L2CacheHit,MemUnitBusy,MemUnitStalled,LDSBankConflct
```

```
0,"eltinv$eltinv_mod_$ck_L136_2_cce$nooop$form.kd"0,0,0,26071,26071,1579087872,256,0,0,4,48,33664,0x0,0x14b05013c040,100.0000000000,24673248.0000000000,17.0000000000,0.0000000000,0.0000000000,100.0000000000,20.7725185968,12332528.0000000000,49.9978626693,77.1977044758,55.9370012636,0.0000000000
```

```
1,"eprfilb$eprfilb_mod_$ck_L85_3_cce$nooop$form.kd"0,0,2,26071,26071,157593600,256,0,0,4,72,38592,0x0,0x14b05013c100,100.0000000000,2462400.0000000000,20.0000000000,0.0000000000,0.0000000000,100.0000000000,24.2834884261,1227104.0000000000,49.9996649935,71.7003018579,49.2693704348,0.0000000000
```

```
2,"eprfilb$eprfilb_mod_$ck_L90_2.kd"0,0,4,26071,26071,56320,256,0,0,24,80,37696,0x0,0x14b05013c080,100.0000000000,880.0000000000,111907.6136363636,0.0000000000,0.0000000000,96.8019501699,9.6155911525,1008455.6562500000,81.3899031130,60.5181117611,43.6472872603,0.0000000000
```

```
3,"eprfilb$eprfilb_mod_$ck_L85_3_cce$nooop$form.kd"0,0,6,26071,26071,157593600,256,0,0,4,72,38592,0x0,0x14b0...
```



Funded by
the European Union

Destination Earth implemented by



101

To get performance metrics on a kernel, we have to collect so called hardware counters.

These are shown in the back of the header.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

```
rocprof -i rocprof_counters.txt -o ./rocprof/rocprof.csv $BIN
```

```
Index,KernelName,gpu-id,queue-id,queue-index,pid,tid,grd,wgr,lds,scr,vgpr,sgpr,fbar,sig,obj,GF,GBusy,Wavefronts,VA
LUInsts,VFetchInsts,VWriteInsts,VALUUtilization,VALUBusy,WriteSize,L2CacheHit,MemUnitBusy,MemUnitStalled,LDSBankC
onflict
```

```
0,"eltinv$eltinv_mod_$ck_L136_2_cce$no-loop$form.kd"0,0,0,26071,26071,1579087872,256,0,0,4,48,33664,0x0,0x14b0501
3c040,100.0000000000,24673248.0000000000,17.0000000000,0.0000000000,0.0000000000,100.0000000000,20.7725185968,123
32528.0000000000,49.9978626693,77.1977044758,55.9370012636,0.0000000000
1,"eprfilb$eprfilb_mod_$ck_L85_3_cce$no-loop$form.kd"0,0,2,26071,26071,157593600,256,0,0,4,72,38592,0x0,0x14b0501
3c100,100.0000000000,2462400.0000000000,20.0000000000,0.0000000000,0.0000000000,100.0000000000,24.2834884261,1227
104.0000000000,49.9996649935,71.7003018579,49.2693704348,0.0000000000
2,"eprfilb$eprfilb_mod_$ck_L90_2.kd"0,0,4,26071,26071,56320,256,0,0,24,80,37696,0x0,0x14b05013c080,100.0000000000
0,880.0000000000,111907.6136363636,0.0000000000,0.0000000000,96.8019501699,9.6155911525,1008455.6562500000,81.38
99031130,60.5181117611,43.6472872603,0.0000000000
3,"eprfilb$eprfilb_mod_$ck_L85_3_cce$no-loop$form.kd"0,0,6,26071,26071,157593600,256,0,0,4,72,38592,0x0,0x14b0
...
```



Funded by
the European Union

Destination Earth implemented by



102

To get performance metrics on a kernel, we have to collect so called hardware counters.

You write these in a file called rocprof_counters.txt, but they are also shown in the back of the header.

And so for each kernel, we collect the value of these counters, for each time step, and for different workloads.

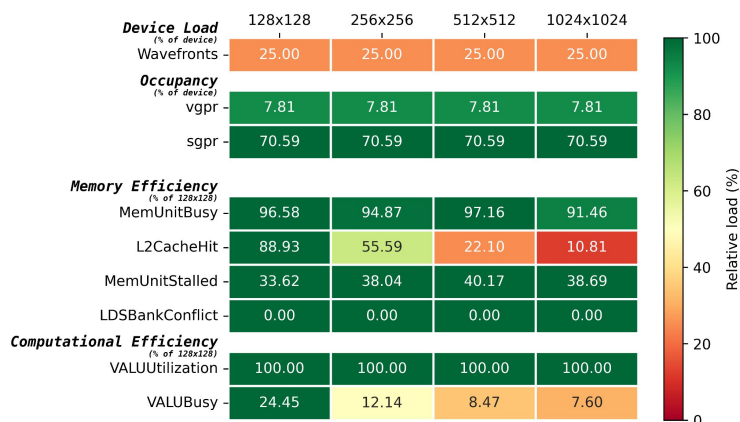
There is no default visualization, and thus I have written my own matrix visualization of these counters for different workloads.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

GPU Profiling

GPU Model factors for easre1b\$easre1b_mod_\$ck_L75_1.kd.



Funded by
the European Union

Destination Earth implemented by



103

Here you can see that visualization, which is very much like the CPU Model factors of Paraver.

The columns are different workloads, or grid sizes, which were also used for the memory plot.

The rows are the hardware counters that we test for.

I won't get into detail what each hardware counter resembles, but I'd say it's easy to see what the problem is with this kernel.

The L2CacheHit rate goes down from 88.93% to 10.81% when increasing the workload.

At the same time we see the VALUBusy metric go down from 24.45% to 7.60%.

So clearly, we have a lot of memory loads that make the computational ALU units wait.

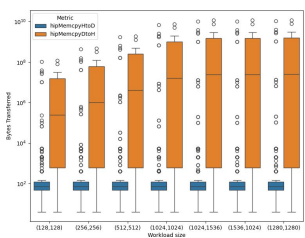
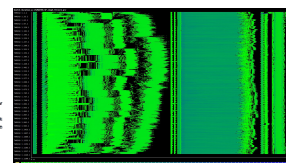
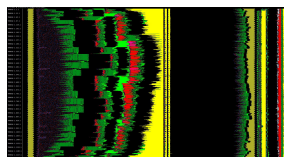
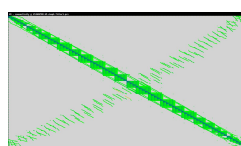
Hence, our advice is to optimize the memory layout for this kernel.

We expect quite an improvement, since it took close to 20% of the computational time.

GPU Icon:

https://www.flaticon.com/free-icon/gpu_4617742

Performance profiling

[illegible]

GPU Model factors for aresb1baresb.1k_sml_175_1.kd.				
Device Load	128x128	256x256	512x512	1024x1024
Occupancy	~35.00	~35.00	~35.00	~35.00
sgpr	7.81	7.81	7.81	7.81
sgpr	70.59	70.59	70.59	70.59
Renory Efficiency				
MemHwBw	96.58	94.87	97.16	91.46
L2CacheHit	83.93	55.59	72.19	38.81
MemHwStall	11.62	38.04	44.07	10.49
LDSBarConflict	0.00	0.00	0.00	0.00
Utilational Efficiency				
VALUJIDBusy	100.00	100.00	100.00	100.00
VALUBusy	24.45	12.14	8.47	7.60

And thus you can see, we have many different ways to visualize the performance of a model.

In this presentation we only covered the tip of the iceberg, but you can go much deeper into questions like why a model is computing more/less in certain regions.

The earth is still alive..



Funded by
the European Union

Destination Earth implemented by



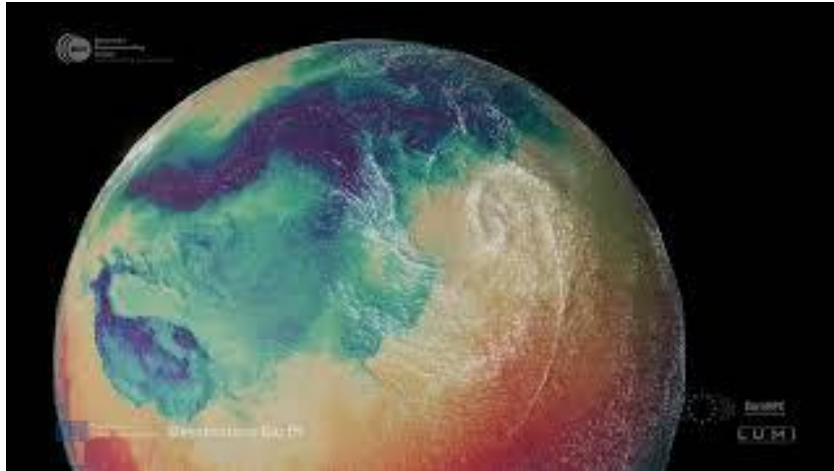
105

That sums up my quick overview of how to detect
Communication inefficiencies on the CPU
And memory issues on the GPU

Of course, there is much much more to explain, but I hope this got you a bit triggered
to think more about the performance of you code when you program.

To close my talk, I'd like to end on a positive note.
Climate change is getting worse and worse, but the Earth is still trying to adapt.
And many of the important climate systems still live.
So I want to show you our latest visualizations of the earth, which shows many
climate related phenomena.

The earth is still alive..



Funded by
the European Union

Destination Earth

implemented by



106

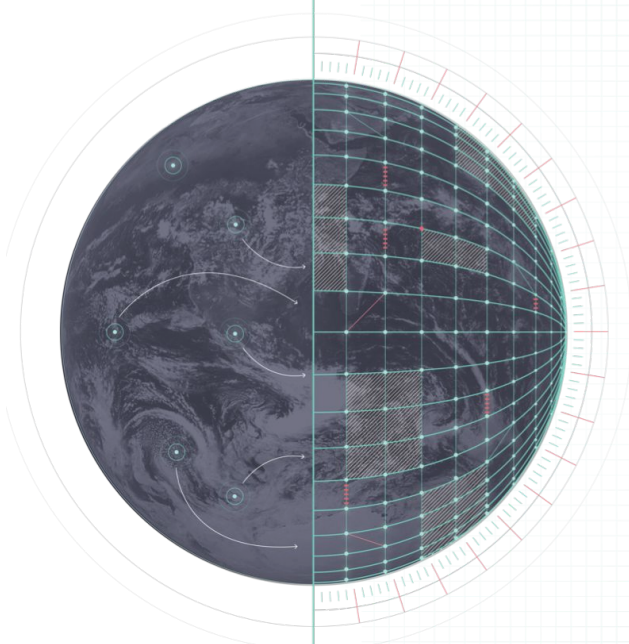
https://www.youtube.com/watch?v=hhJW0RBq0hA&list=PLbABsNMD2jhy5zsINz9VluOzVsrM0_vwu

0:00-0:45: clouds forming and moving + hurricanes

1:12-3:00: earth heightmap, sea temperature streams,

5:10-6:30: wind map, to pressure + temperature map, heart beat from day night

Total: $45+108+80=233s = 3:53m$



Earth Sciences
Department

**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación





SBC PRESENTS
SUSTAINABILIT

Destination Earth

Pioneering the Future with Digital Twins

Okke van Eck
GPU Research Engineer
Barcelona Supercomputing Center



okke.vaneck@bsc.es

27th of November 2024, Bussum

**Funded by
the European Union**

Destination Earth implemented by

That was all from my side.

I hope you learned something about climate modellen or performance engineering.

Here you have my contact information in case you want to contact me later, but please ask any questions you have now as well.

Thank you!