

DESTINATION EARTH

Containerization and deployment of weather models on EuroHPC JU infrastructures

*Massimo Gisonni, Giuseppa Muscianisi (CINECA),
Francesco Cola (CINECA/JSC), Thomas Geenen (ECMWF)*

21st ECMWF workshop on high performance computing in meteorology, Bologna (Italy), 15-19 September 2025



Funded by
the European Union

Destination Earth

implemented by



CINECA

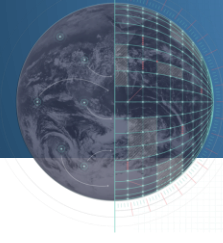




Funded by
the European Union

Destination Earth

implemented by



GOAL OF THE PROJECT

PoC: design of a workflow to containerize weather models & run them at performance in the clusters of the EuroHPC JU

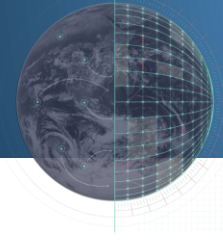
- A single immutable container image (.sif) holding the application's binary, the required libraries, interfacing at runtime to MPI



Funded by
the European Union

Destination Earth

implemented by



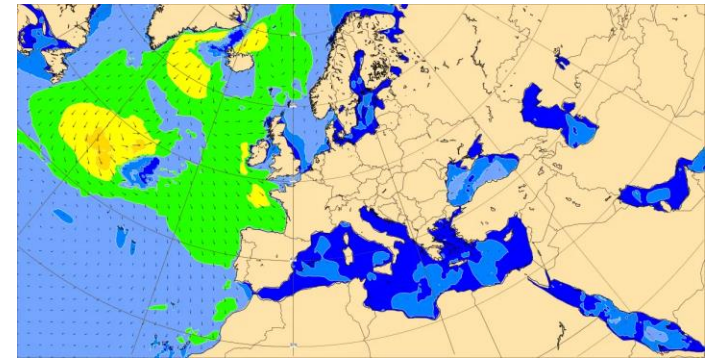
GOAL OF THE PROJECT

PoC: design of a workflow to containerize weather models & run them at performance in the clusters of the EuroHPC JU

- A single immutable container image (.sif) holding the application's binary, the required libraries, interfacing at runtime to MPI

- **Target applications:**

- *ecWAM* (Ocean Wave model)
- *ecTrans* (Spectral transforms)
- *RAPS20* (GPU support)
- *RAPS21* (CPU only)

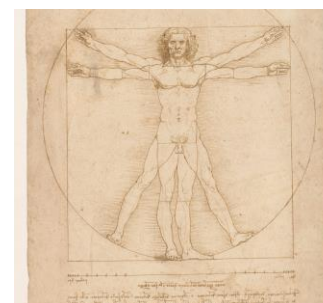
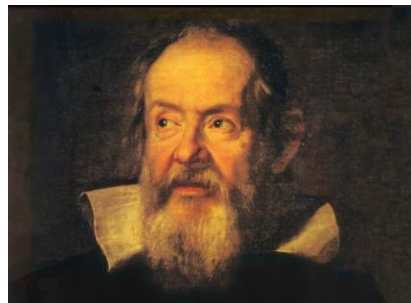


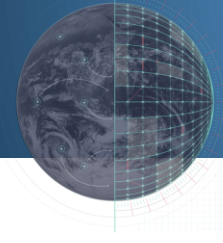
- **Target clusters:**



EuroHPC
Joint Undertaking

- *GALILEO100* (CINECA)
- *LEONARDO* (CINECA)
- *MeluXina* (LuxProvide)
- *LUMI* (CSC)





THE ASSETS OF CONTAINERS

Portability:

- No need to rebuild the same software for each system
- Different containers for different **architectures** (x86 vs CRAY)

Reproducibility:

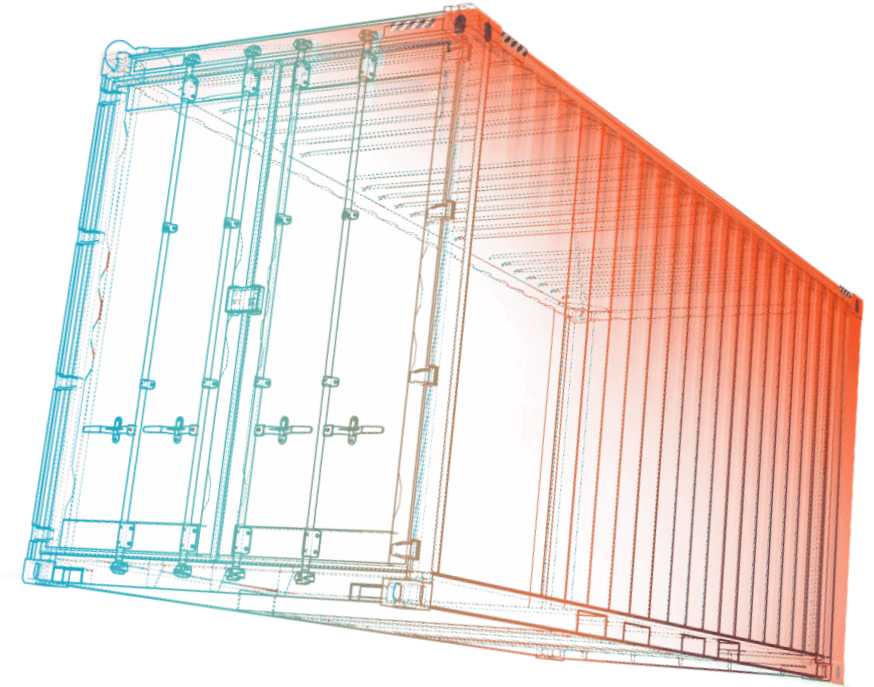
- The same container can be built in different sites from the **same recipe**
- Software stack inside the container is fixed

Legacy:

- Testing old/new software stack against old/new host software stack
- Quick deployment in case of issues in the host

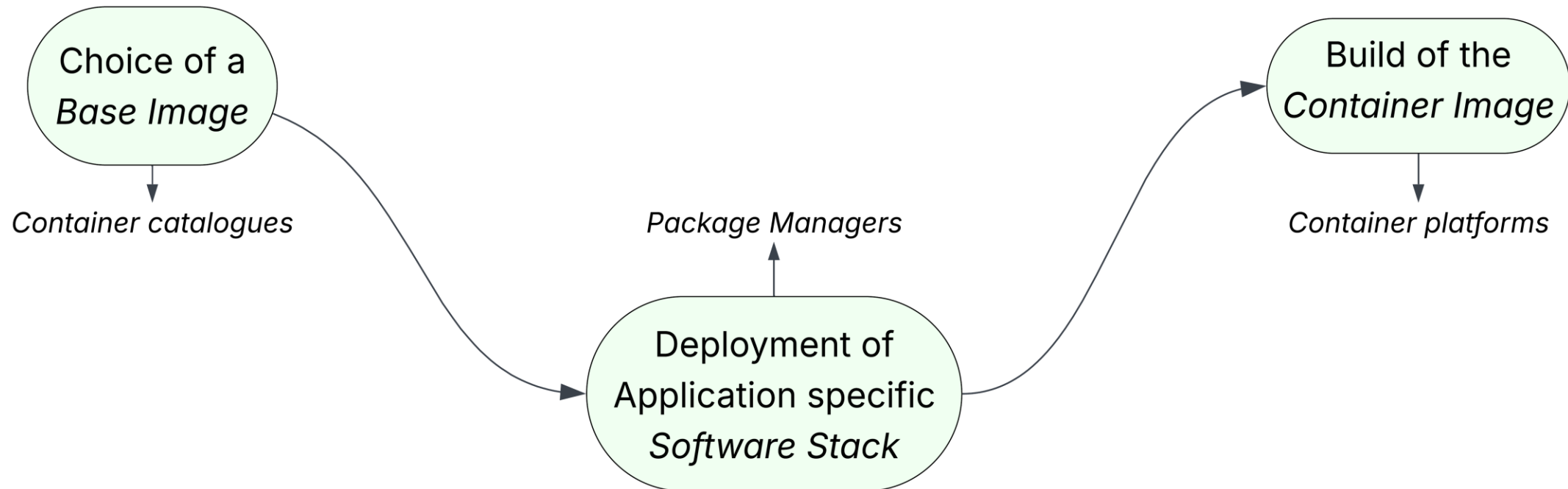
Distribution:

- Containers are **lightweight**, easily transferable
- Integration on CI/CD pipeline allows for controlled deployment





THE TOOLS OF THE PROJECT





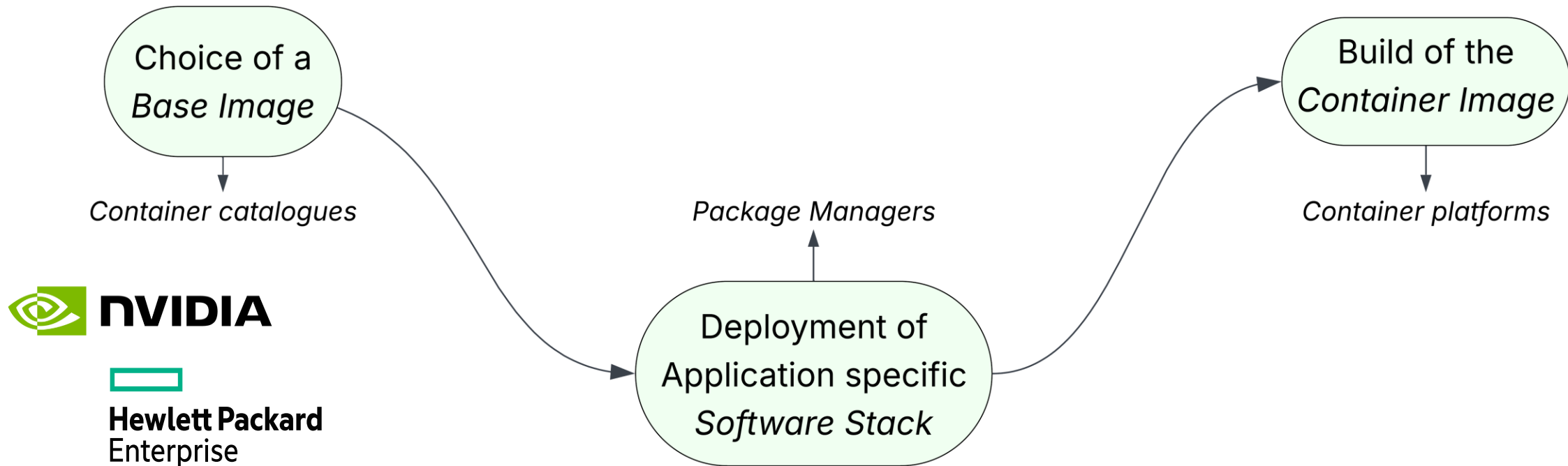
Funded by
the European Union

Destination Earth

implemented by



THE TOOLS OF THE PROJECT

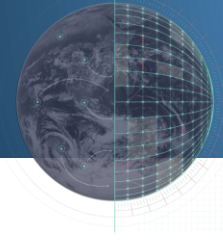




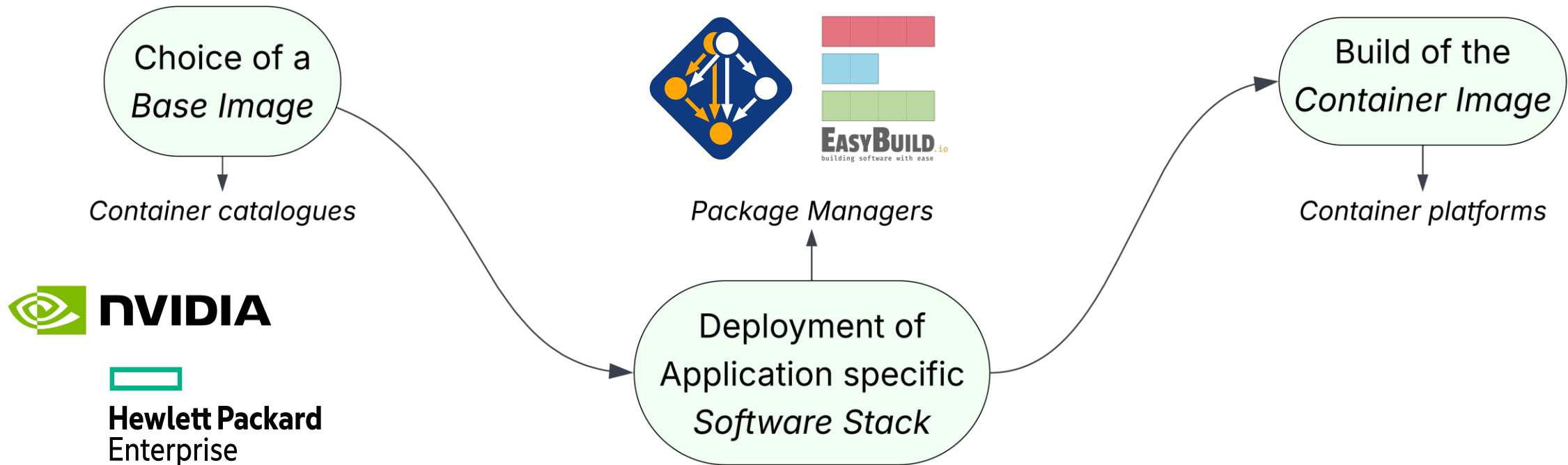
Funded by
the European Union

Destination Earth

implemented by



THE TOOLS OF THE PROJECT





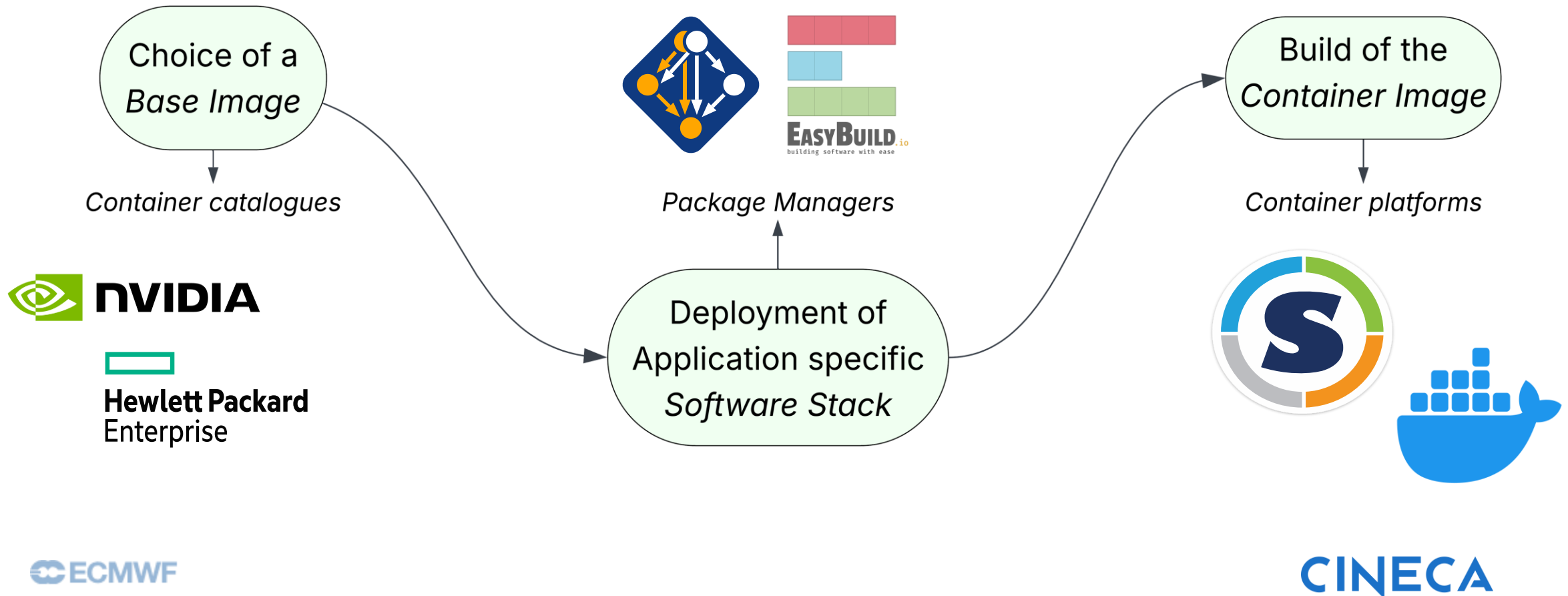
Funded by
the European Union

Destination Earth

implemented by



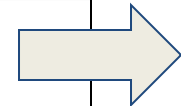
THE TOOLS OF THE PROJECT





THE TOOLS OF THE PROJECT: SINGULARITY

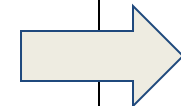
```
Bootstrap: docker
From: nvcrl.io/nvidia/nvhpc:23.1-devel-cuda_multi-ubuntu22.04
Stage: build
```



HEADER

```
%files
  /model.tar.gz      /tmp
  /spack.yaml        /spack_env/spack.yaml

%post
  apt-get update && apt-get -y install python3 git wget
  tar -xzf model.tar.gz
  ./install_model.sh --builddir=/build
```

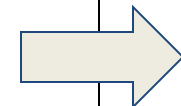


SECTIONS

```
%environment # allows to define environment variables
  export VARIABLE=MEATBALLVALUE
```

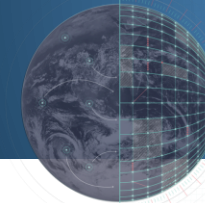
```
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
```

```
Bootstrap: docker
From: nvcrl.io/nvidia/nvhpc:23.1-runtime-cuda_multi-ubuntu22.04
Stage: deploy
```



MULTISTAGE

```
%files from build
  /build/bin/model    /deploy/model
```



THE TOOLS OF THE PROJECT: SPACK

Host compilers and pre-installed packages

config:

```
source_cache: $env/source_cache  
install_tree: $env/install
```

build_stage:

- \$env/build/stage
- \$env/spack-stage

compilers:

```
- compiler:  
  spec: gcc@=11.3.0  
  paths:  
    cc:  /bin/gcc  
    cxx: /bin/g++  
    f77: /bin/gfortran  
    fc:  /bin/gfortran  
  flags: {}  
  operating_system: ubuntu22.04  
  target: x86_64
```

packages:

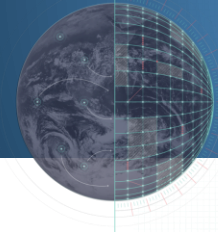
```
cuda:  
  externals:  
    - spec: cuda@12.0.76  
      prefix: /opt/cuda  
  
openmpi:  
  externals:  
    - spec: openmpi@4.1.5 %gcc@=11.3.0  
      prefix: /opt/mpi
```

specs:

```
- hdf5@1.12.2      %nvhpc      -cxx +fortran +mpi  
- zlib-ng          %gcc  
- netcdf-c@4.9.0   %nvhpc      +mpi +optimize  
- netcdf-fortran@4.6.0 %nvhpc    build_system=autotools  
- cmake@3.24
```

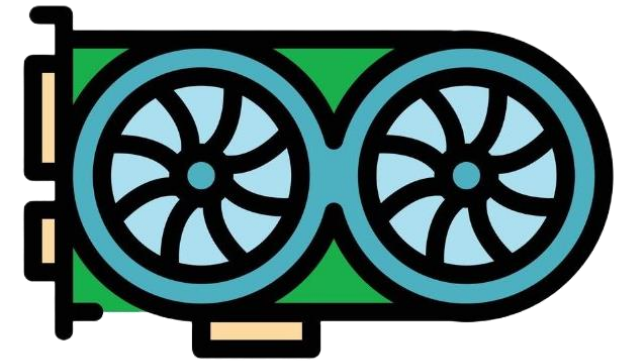
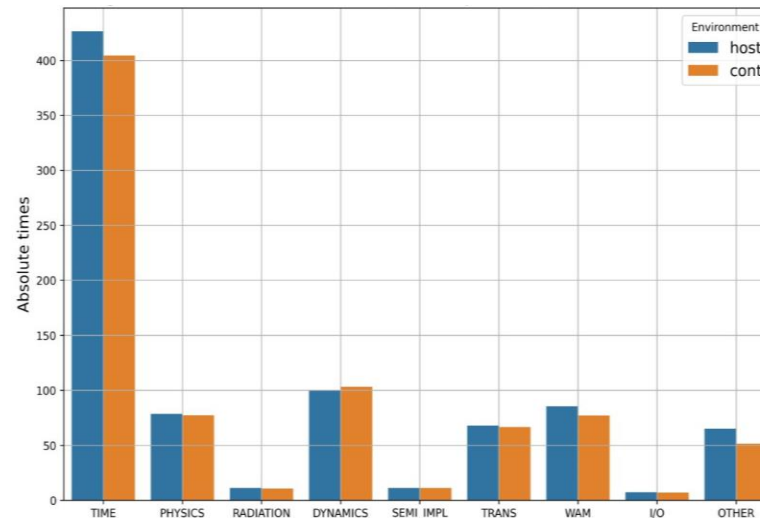
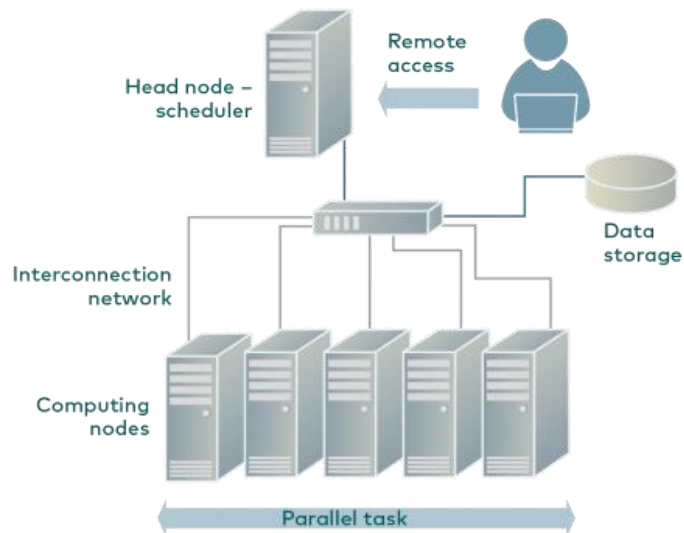
```
- package@version    %compiler  +options
```

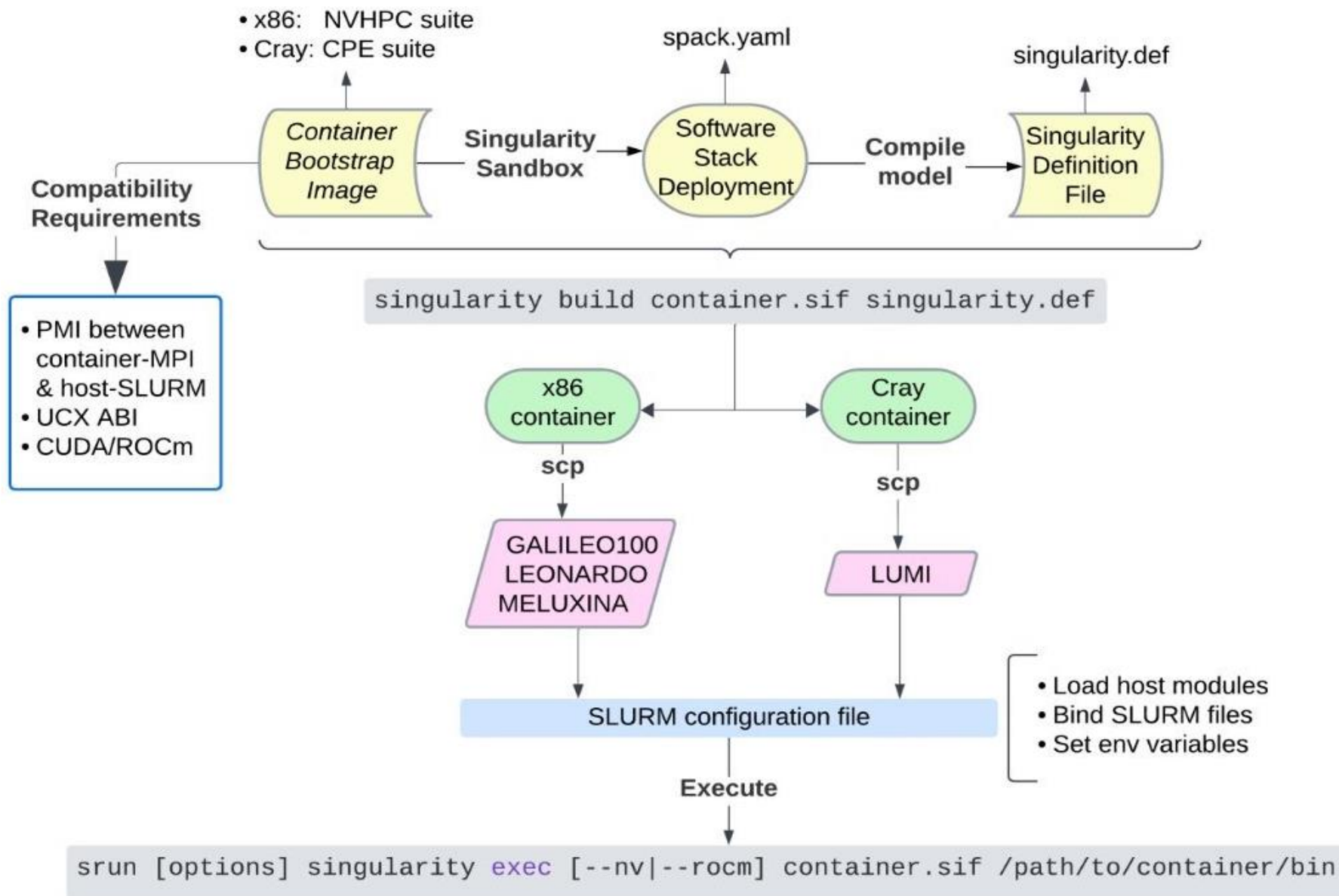
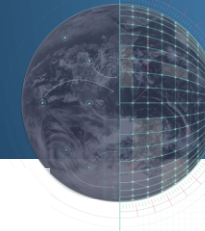
Packages to be installed

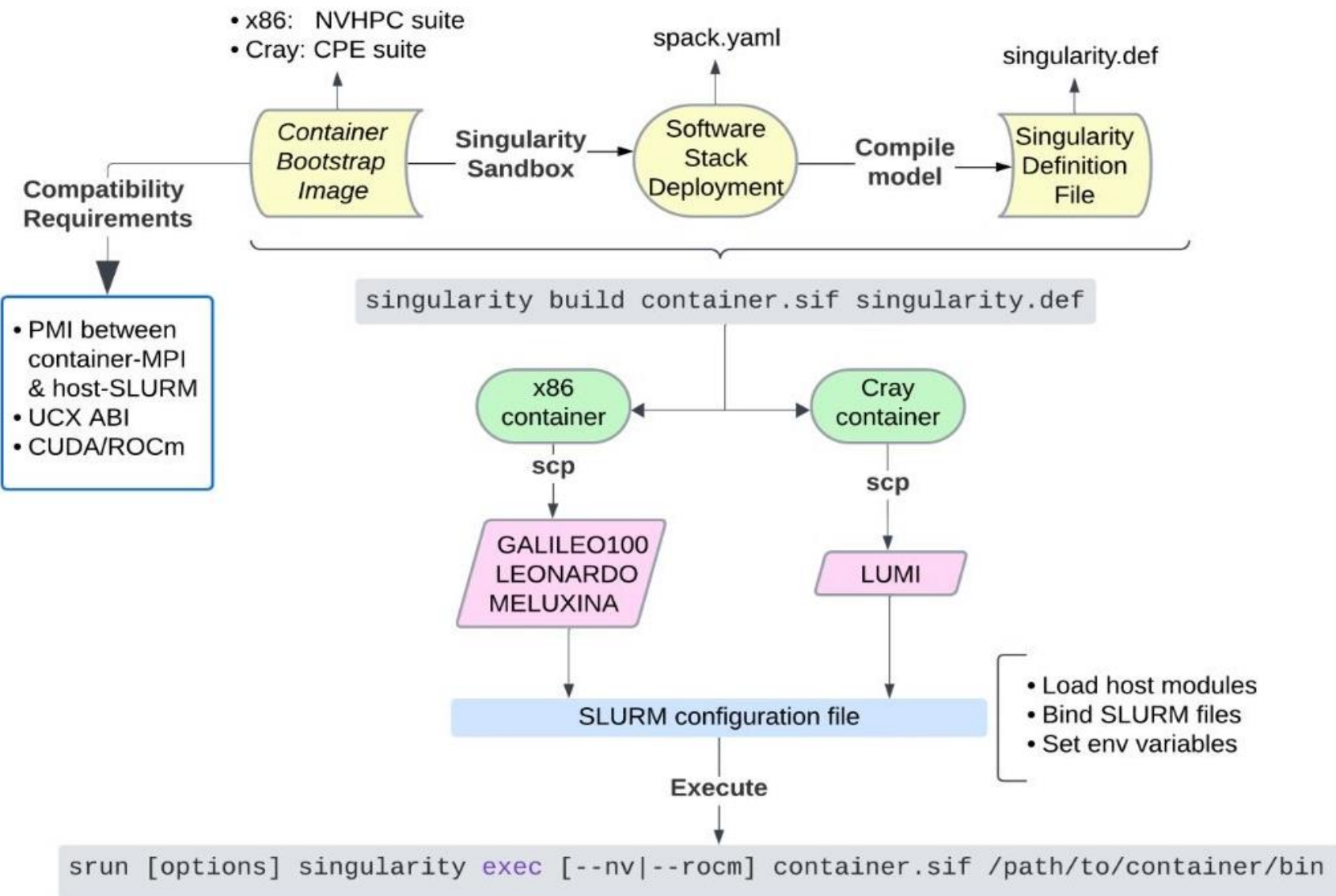
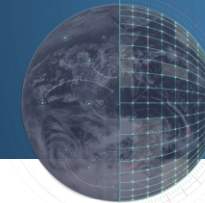


THE REQUIREMENTS

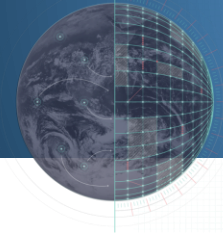
- *Integration in a HPC environment*
 - Integration with resource manager & host software stack
- *Run at performance wrt host*
 - Usage of bare metal capabilities (parallelism, infiniband)
- *Ensure usage of accelerators*
 - Correctly offloading to NVIDIA/AMD GPUs





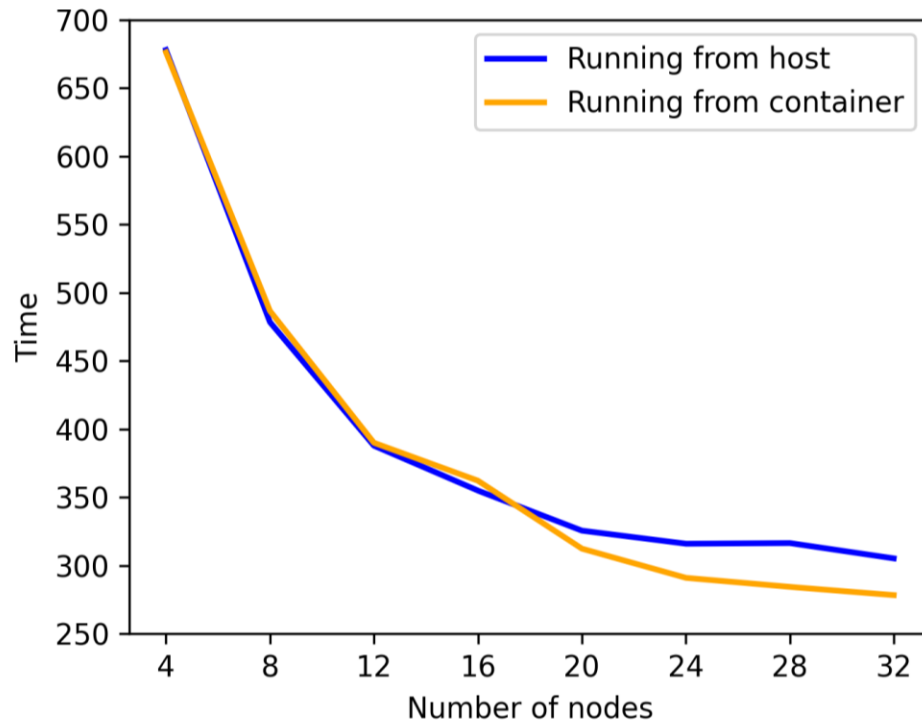


- Build on dedicated VMs/clusters allowing
- Singularity to Docker and viceversa
- Leveraging base container images:
 - NVIDIA images holds NVHPC suite (compiler, OpenMPI, CUDA)
 - HPE images holds CPE suite + help from CSC for ROCm integration
- Reproducing host software stack:
 - as faithfully as possible via Spack
- Multistage Singularity builds:
 - lightweight final container images
- Singularity hybrid model approach:
 - Weak compatibility requirements



Performance Assessment: ecWAM on GALILEO100

- Running on the O640 grid, CPU only, 48 threads per MPI task

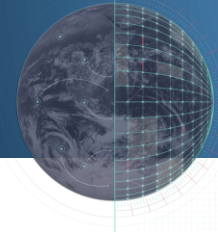


Raw timings

	host	container
4	677.98	676.26
8	478.42	486.64
12	387.85	389.98
16	354.93	362.21
20	325.57	312.37
24	315.97	290.96
28	316.45	284.32
32	305.24	278.30

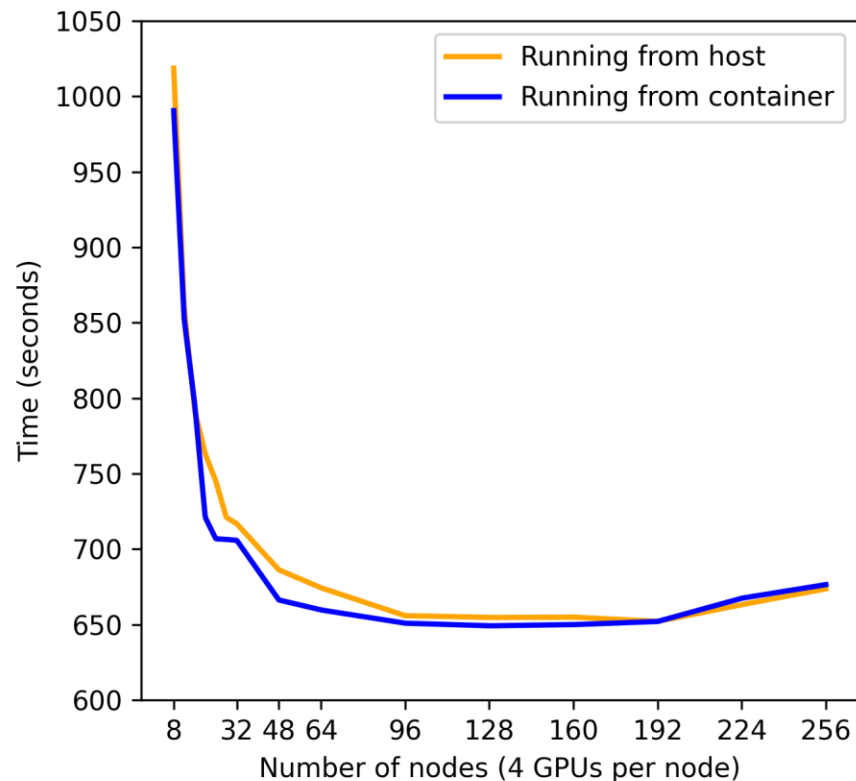
Percentage Difference

-0.25
1.72
0.55
2.05
-4.05
-7.92
-10.15
-8.83

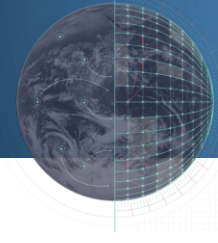


Performance Assessment: ecWAM on LEONARDO

- Running on the O1280 grid, GPU enabled, 4 MPI tasks per node, 8 threads per MPI task, 4 GPUs per node

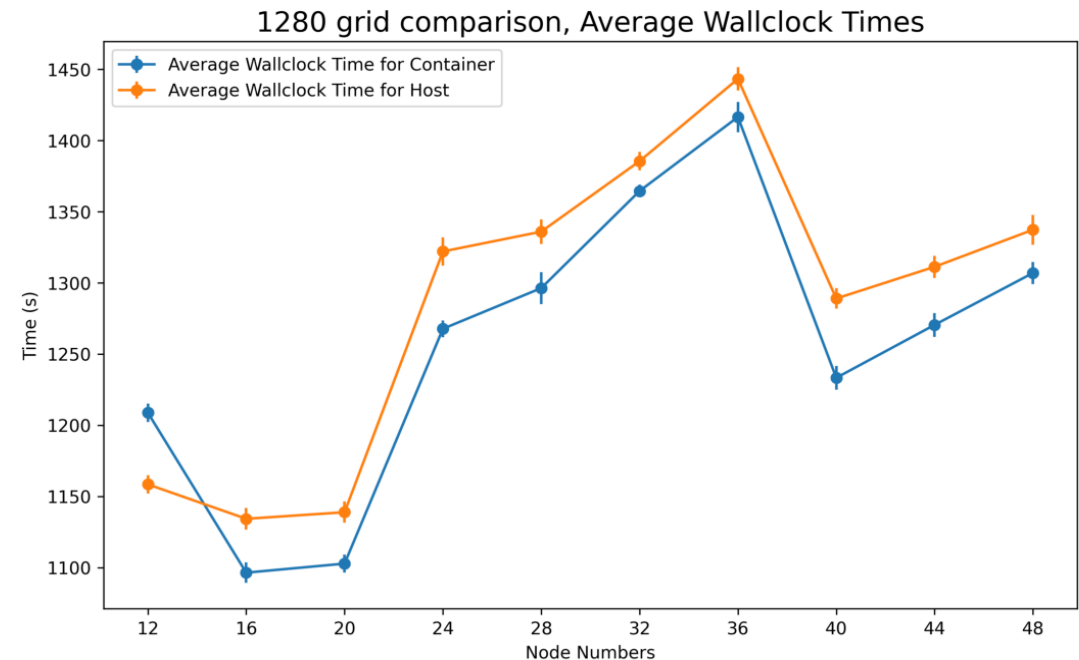
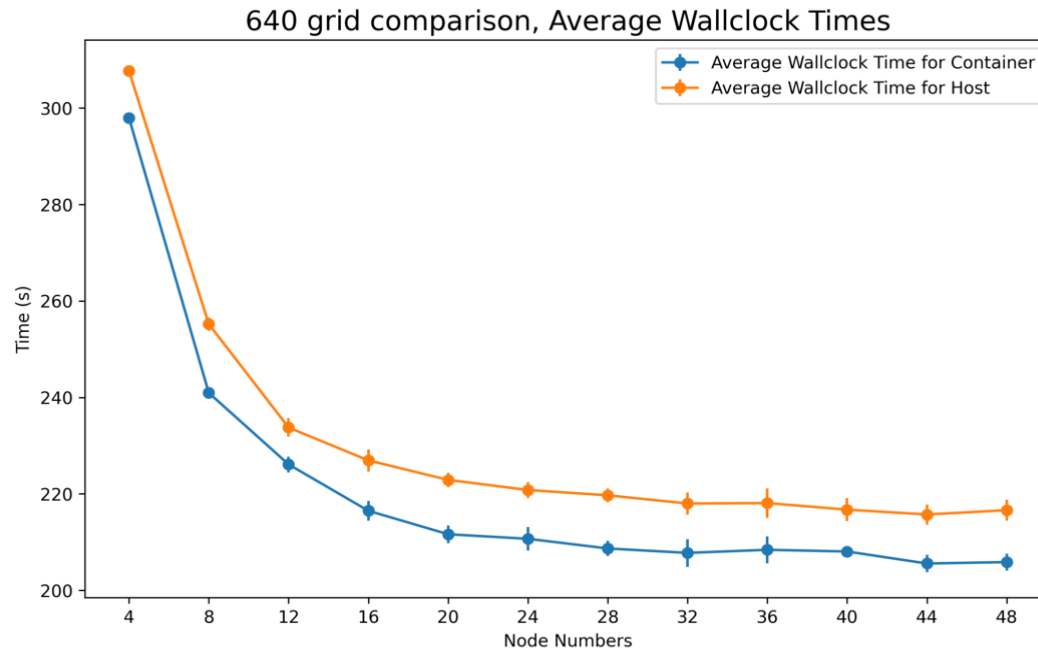


Raw timings			Percentage Difference		
8	1018.68	990.56	-2.76	8	
12	859.82	852.02	-0.91	12	
16	791.95	795.85	0.49	16	
20	763.04	721.04	-5.50	20	
24	745.15	706.89	-5.13	24	
28	721.10	706.48	-2.03	28	
32	716.81	705.83	-1.53	32	
48	686.29	666.32	-2.91	48	
64	674.42	659.73	-2.18	64	
96	655.95	650.94	-0.76	96	
128	654.75	649.21	-0.84	128	
160	655.00	649.98	-0.77	160	
192	652.13	652.10	-0.01	192	
224	663.24	667.50	0.64	224	
256	673.73	676.51	0.41	256	
host		container			



Performance Assessment: ecWAM on MeluXina

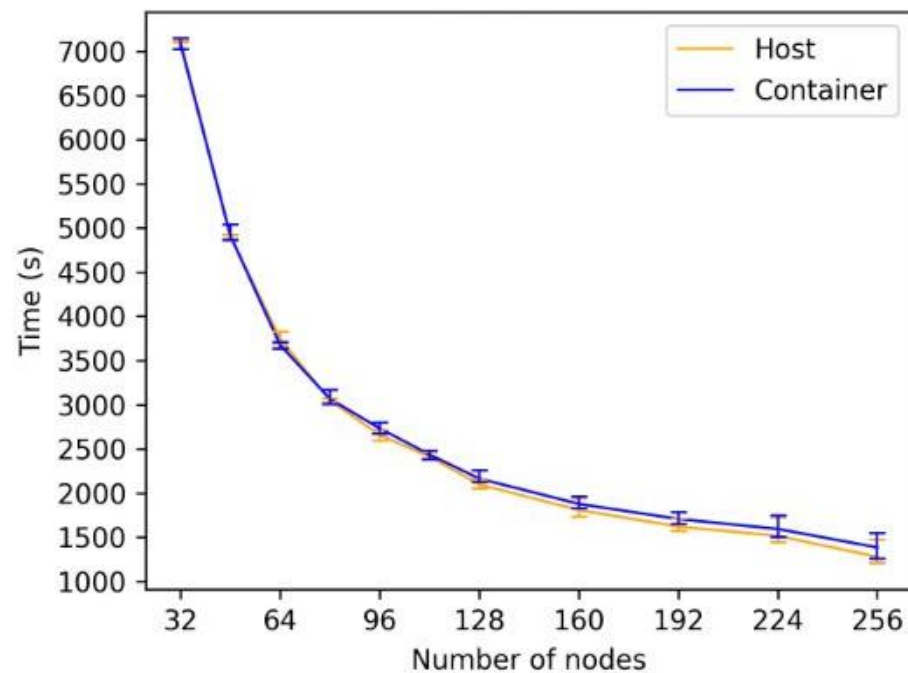
- Running on the 0640 & O1280 grid, GPU enabled, 4 MPI tasks per node, 8 threads per MPI task, 4 GPUs per node
- The **very same** container image deployed on LEONARDO was used





Performance Assessment: RAPS20 on LEONARDO

- Running *RAPS20* on tco2559l137, ecWAM on GPU, full nodes of LEONARDO booster partition

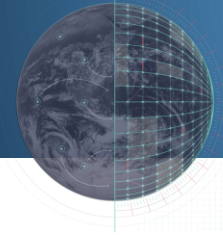


Raw timings

Nodes	Host	Container
32	7106.16	7081.06
48	4908.57	4916.72
64	3745.61	3676.47
80	3051.12	3067.47
96	2656.69	2732.33
112	2414.40	2437.75
128	2098.33	2164.45
160	1809.10	1878.06
192	1621.41	1707.23
224	1516.67	1594.50
256	1279.18	1385.41

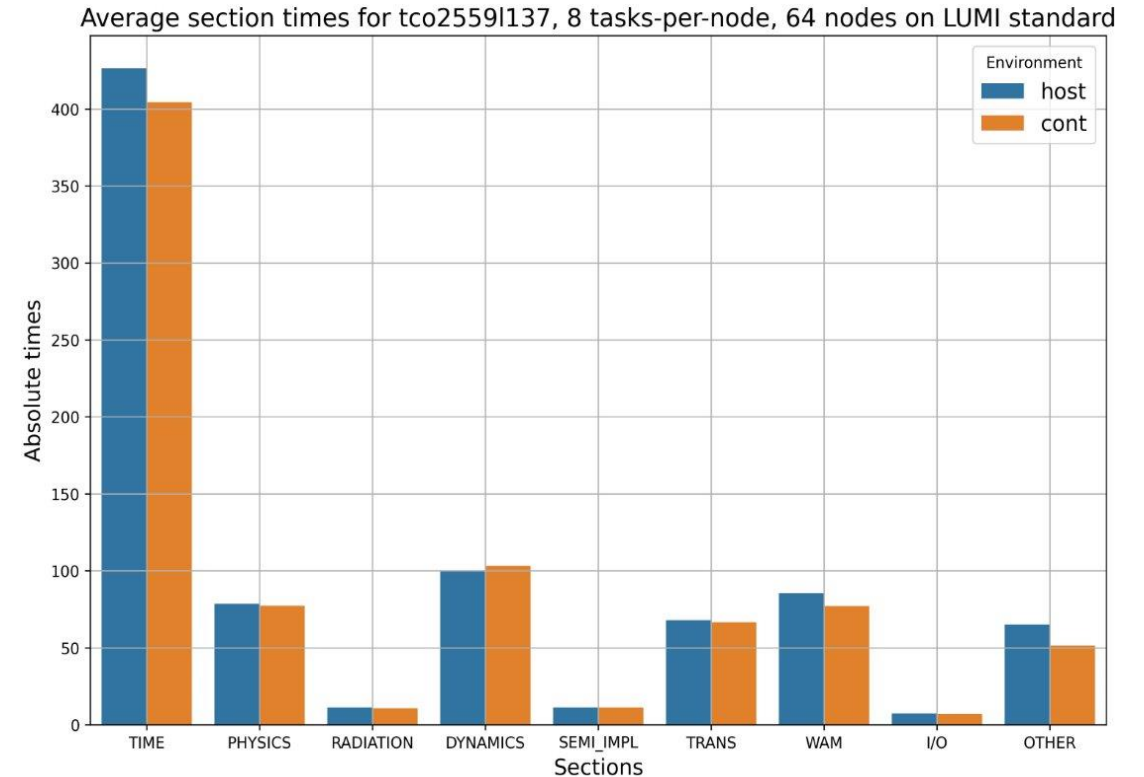
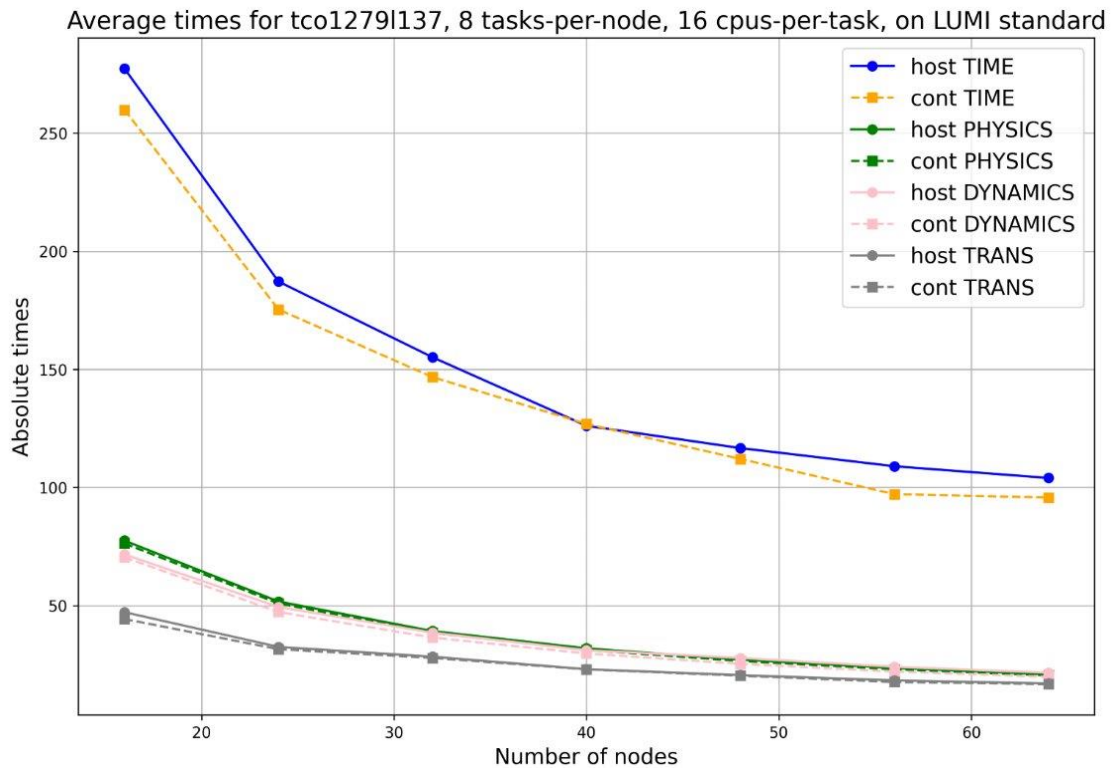
Percentage Difference

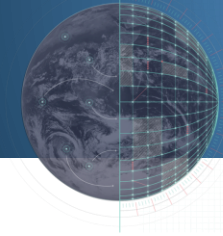
Nodes	Percentage Difference
32	-0.35
48	0.17
64	-1.85
80	0.54
96	2.85
112	0.97
128	3.15
160	3.81
192	5.29
224	5.13
256	8.30



Performance Assessment: RAPS21 on LUMI

- Running *RAPS21* on full nodes of LUMI-C, tco1279l137 and tco2559l137 resolution



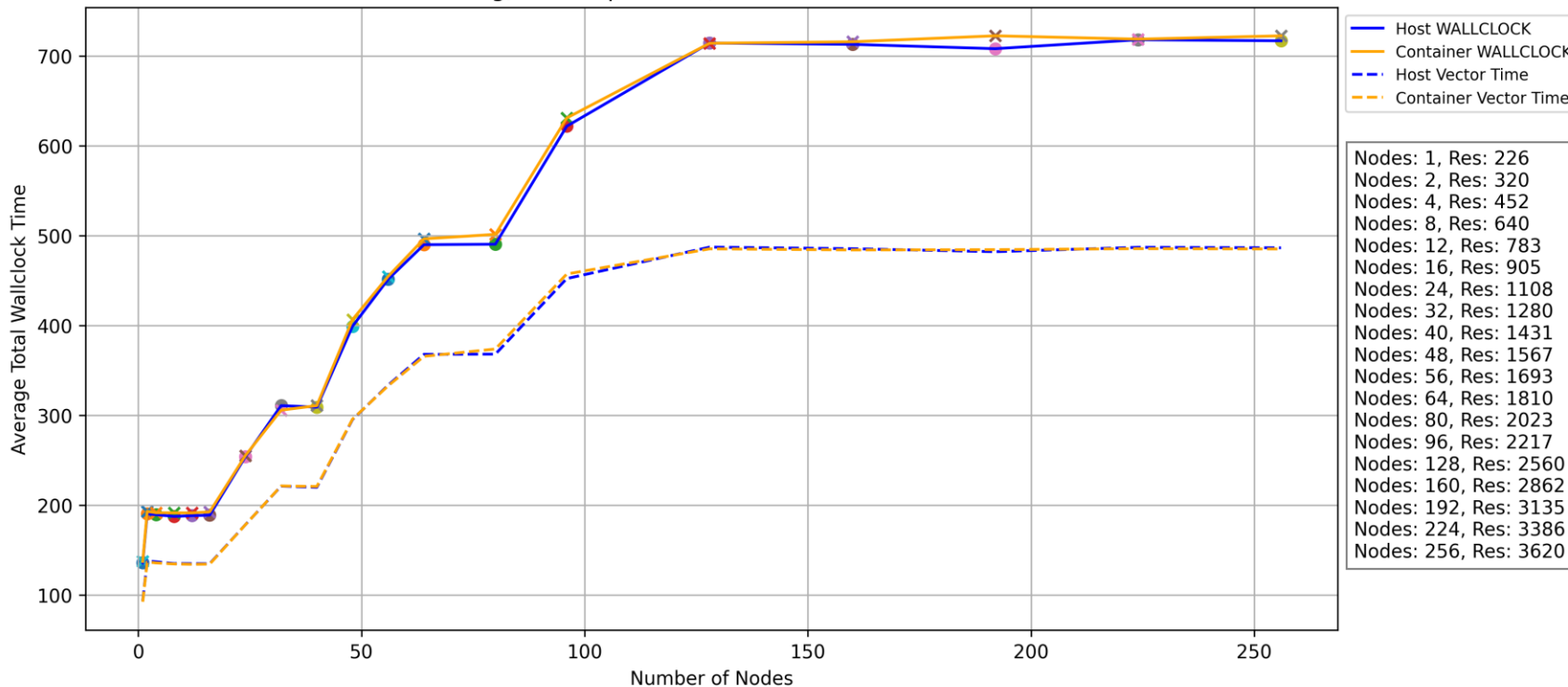


Performance Assessment: ecTrans on LUMI

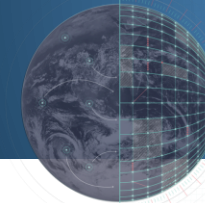
- Running *ecTrans* (redgreengpu branch) on full nodes of LUMI-G, tco1279|137 and tco2559|137 resolutions

ecTrans-redgreengpu tests on LUMI-g, Host vs Container, CPE/24.03, cce/17.0.1

Weak scaling, 8 GPUs per node, resolution is $\sqrt{\text{Nodes} \times 8 \times 80^2}$



135.84	137.28	1.06
190.28	193.04	1.45
189.39	191.57	1.15
187.48	191.52	2.16
188.29	191.35	1.63
188.92	192.57	1.93
253.93	255.21	0.50
310.82	305.76	-1.63
309.14	310.84	0.55
399.19	406.65	1.87
451.61	454.83	0.71
489.98	496.48	1.33
490.52	501.49	2.24
621.93	631.09	1.47
714.50	714.29	-0.03
713.02	715.91	0.41
708.21	722.55	2.03
718.14	718.84	0.10
717.00	722.53	0.77
host_wt	cont_wt	%_diff_wt



Sanity Check: IFS_TEST_BITIDENTICAL=check

```
#-----#
#####
### CONTAINER INTEGRATION ###
#####

if [[ -n "${CONTAINER_ENV}" ]]; then
  if [[ ! -f "${CONTAINER_ENV}" ]]; then
    echo "Error: CONTAINER_ENV is set but file '${CONTAINER_ENV}' does not exist."
    exit 1
  fi
  echo "CONTAINER_ENV is set to '${CONTAINER_ENV}'. Reading variables..."
  source "${CONTAINER_ENV}"
  # Define the required variables (they must be declared in CONTAINER_ENV)
  required_vars=("CONTAINER_SW" "CONTAINER_FLAGS" "CONTAINER_GPU_FLAGS" "CONTAINER_IMAGE" "CONTAINER_BINARIES")
  # Check that all required variables exist (even if empty)
  for var in "${required_vars[@]"; do
    if [[ ! -v $var ]]; then
      echo "Error: Required variable '$var' is not set in '${CONTAINER_ENV}'."
      exit 1
    fi
  done

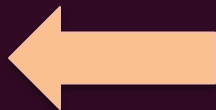
  # Construct the LAUNCH command (as above without IFS_SOURCE_DIRS, LAUNCHER_OTHER_FLAGS to be appended at the end?)
  LAUNCH="${LAUNCHER} ${LAUNCHER_NPROC_FLAG} ${NPROC} ${LAUNCHER_NTHREAD_FLAG} ${NTHREAD} \
    ${LAUNCHER_MEM_FLAG} "

  # Concatenate container setup
  LAUNCH+="${CONTAINER_SW} ${CONTAINER_FLAGS} ${CONTAINER_GPU_FLAGS} \
    ${CONTAINER_IMAGE}"

  # Set the IFS executable to the container one
  # Extract basename (could be either ifsMASTER.SP or osmMASTER)
  IFS_EXECUTABLE=$(basename "${IFS_EXECUTABLE}")
  # Prepend RAPS bin folder of container
  IFS_EXECUTABLE=${CONTAINER_BINARIES}${IFS_EXECUTABLE}
else
  echo "Launching Host installation"
  echo "Set variable CONTAINER_ENV to a file holding container variables to launch a containerized instance"
fi
#-----#
```

Checked that the containerized RAPS21 image not only runs at performance, but runs meaningfully

- Ran the ctests of the *ifs-test suite* with a host-based installation
- Compared bitwise against a containerized version, with simple changes on the ifs-run script
- *Passed successfully*
- Failed using different compilers on host and container (as expected)
- Tested on LUMI, ECMWF ATOS, LEONARDO

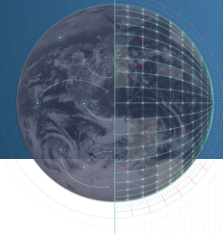




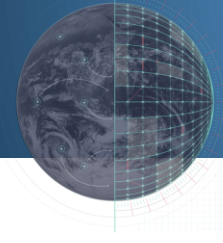
Funded by
the European Union

Destination Earth

implemented by

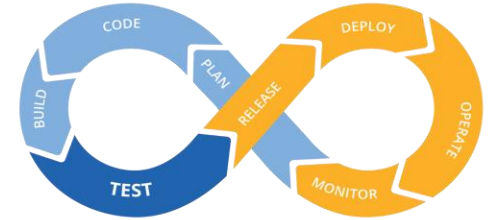


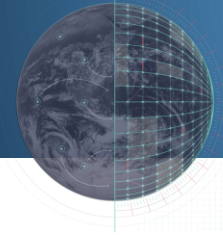
NEXT ON THE LIST



NEXT ON THE LIST

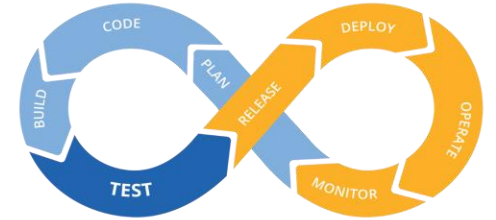
- **CI/CD integration:** container build and test workflows integrated on the CI/CD pipeline
- **CW ---> Continuous Work!**
 - Containerize a new application takes time.
 - Containerize the same application for different versions/clusters does not take that much time





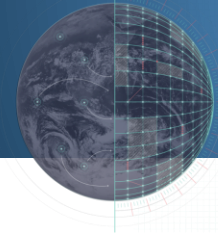
NEXT ON THE LIST

- **CI/CD integration:** container build and test workflows integrated on the CI/CD pipeline
- **CW ---> Continuous Work!**
 - Containerize a new application takes time.
 - Containerize the same application for different versions/clusters does not take that much time



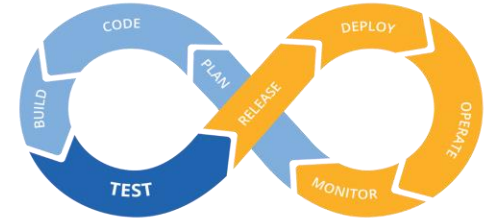
Ideally: one container to rule them all

- *Use the same container across all EuroHPC supercomputers*



NEXT ON THE LIST

- **CI/CD integration:** container build and test workflows integrated on the CI/CD pipeline
- **CW ---> Continuous Work!**
 - Containerize a new application takes time.
 - Containerize the same application for different versions/clusters does not take that much time

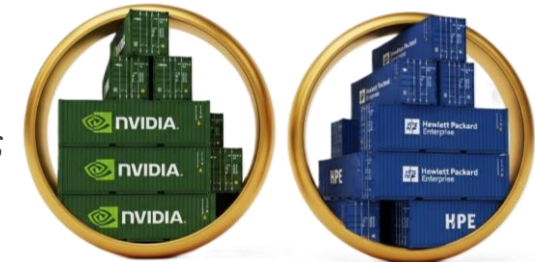


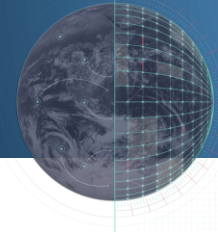
Ideally: one container to rule them all

- Use the same container across all EuroHPC supercomputers

This project: two sets of containers to rule some of them

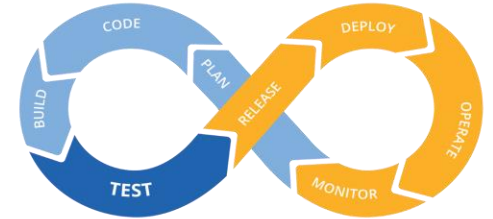
- The same container is used across similar architectures





NEXT ON THE LIST

- **CI/CD integration:** container build and test workflows integrated on the CI/CD pipeline
- **CW ---> Continuous Work!**
 - Containerize a new application takes time.
 - Containerize the same application for different versions/clusters does not take that much time



Ideally: one container to rule them all

- Use the same container across all EuroHPC supercomputers

This project: two sets of containers to rule some of them

- The same container is used across similar architectures



Practical?: a set of container for each cluster

- Registry holding snapshots of applications and/or clusters software stack

Thanks for the attention!

QUESTIONS?

