



PP IMPACT and HPC overview

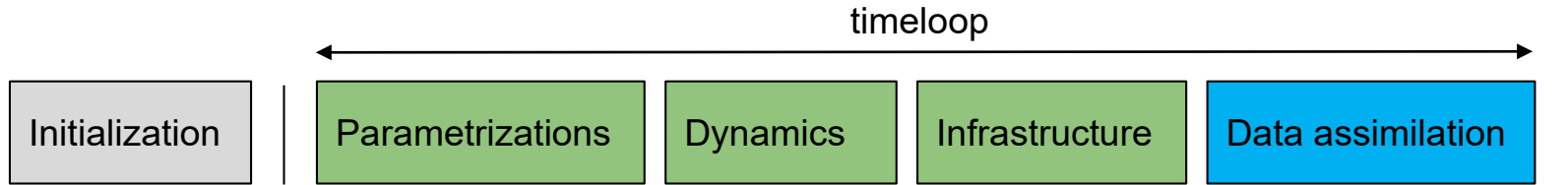
X. Lapillonne¹, C. Osuna¹, D. Hupp¹, D. Alexeev⁵, V. Cherkas¹, R. Dietlicher¹, E. Germann¹, F. Gessler¹, M. Jacob⁴, A. Jocksch³, J. Jucker², C. Müller¹, M. Röthlin¹, W. Sawyer³, U. Schättler⁴, André Walser¹

¹MeteoSwiss, ²C2SM, ³CSCS, ⁴DWD, ⁵Nvidia

06.09.2022, COSMO-GM WG6 – IMPACT parallel session



OpenACC port overview



- Most components for NWP Regional and global application ported, optimization work ongoing
- Support for both double and mixed precision
- Some components, e.g. ecRad, I need to be merged, some need to be ported
- Regular testing on builbot infrastructure



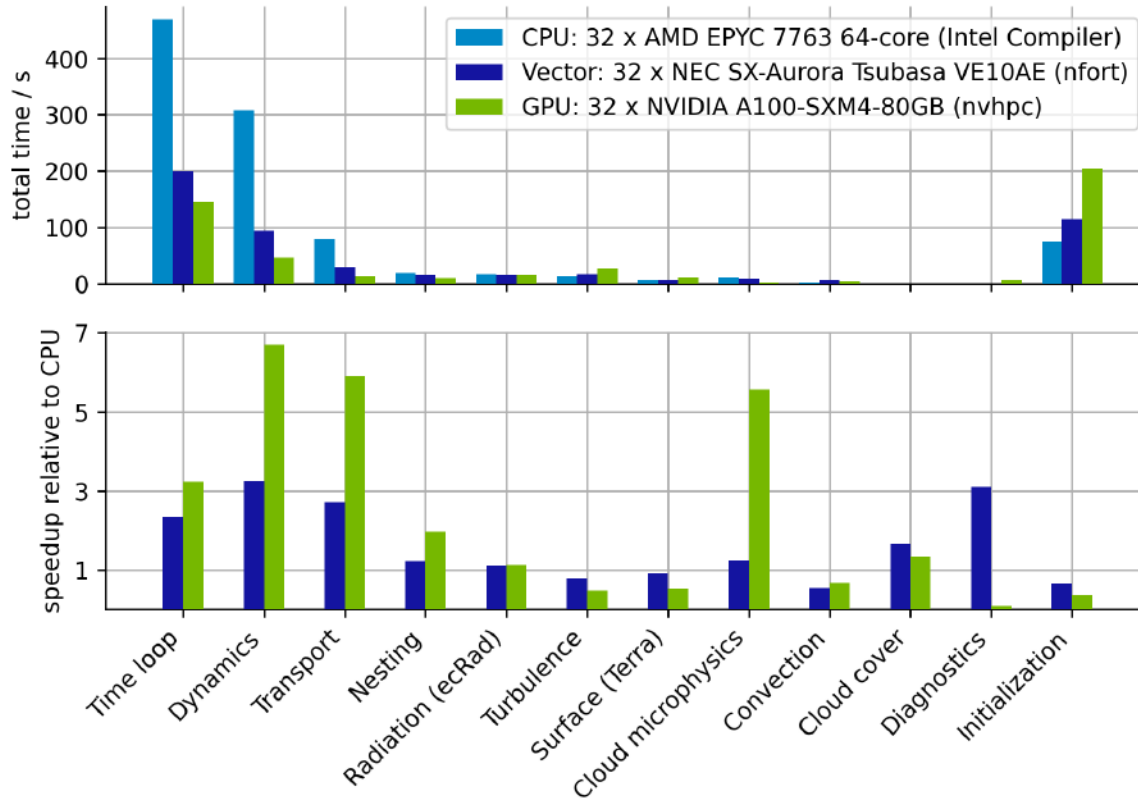
Status of the OpenACC port to GPU

	ported	merged
nh_solve	Ready	Ready
nh_hdiff	Ready	Ready
transport	Ready	Ready
2 way nesting	Ready	Ready
convection	Ready	Ready
Microphysics (graupel)	Ready	Ready
radiation	Ready	In progress
radheat	Ready	Ready
Surface (terra)	Ready	Ready
cover	Ready	Ready
turbulence	Ready	Ready
Sea-ice	Not started	Not started
sso	Ready	Ready
Non-or. Wave drag	Ready	Ready
2 mom. microphysics	In progress	Not started

	ported	merged
NWP diagnostic	In progress	In progress
DA: LHN	Ready	Ready
DA: conv. operator	On CPU + data copy and interface ported	On CPU + data copy and interface ported
DA: IAU (Incr. Anal. Update)	In progress	In progress
SPPT	In progress	Not started

Ready
In progress
On CPU + data copy and interface ported
Not started

Detailed missing features list: <https://gitlab.dkrz.de/icon/wiki/-/wikis/GPU-development/todo-list>



CPU vs GPU vs NEC

Experiment:

- ~DWD deterministic forecast Global R2B8
- 5 242 880 cells (10 km) + Nested grid over Europe 845 340 cells (5 km)
- Output disabled

Performance and energy

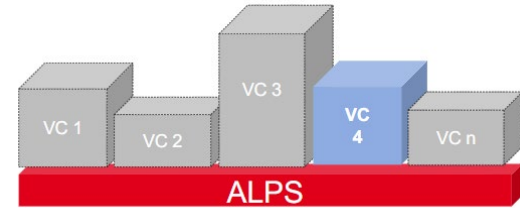
- GPU 3x faster than CPU
- 32 GPUs: 663 s/day 268 Wh/day
- 32 VEs: 914 s/day 292 Wh/day



ICON-22 Performance

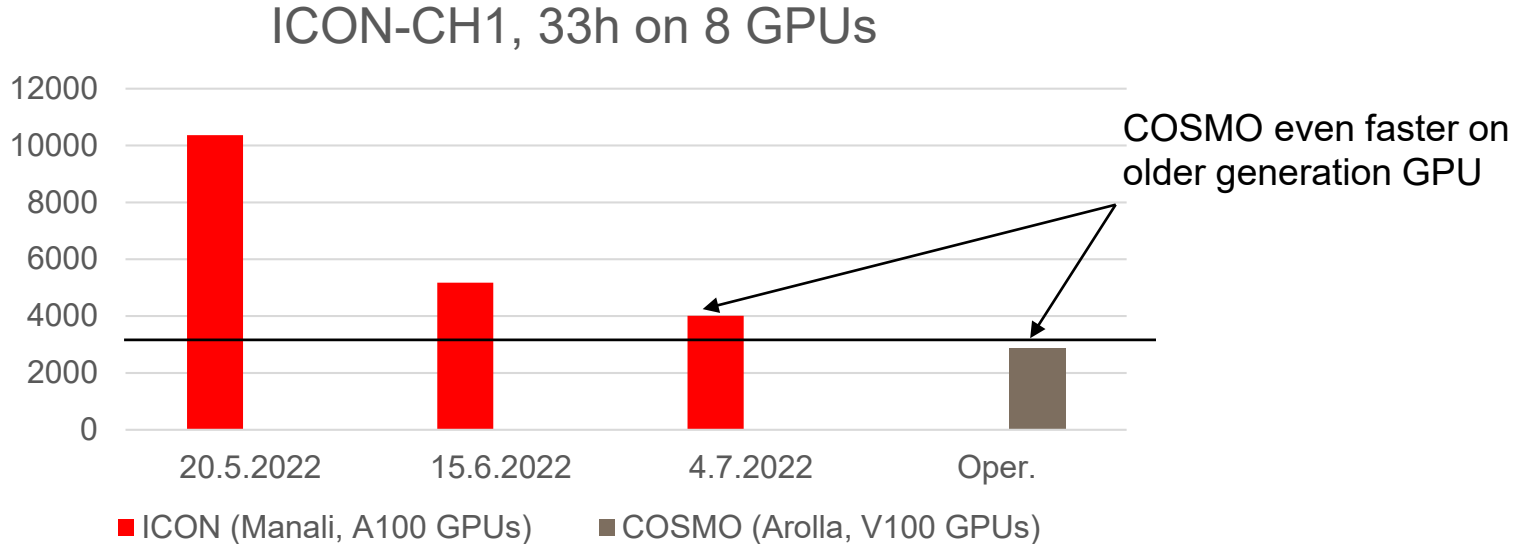
New MeteoSwiss system HPC Computing Services on Alps Plattform

- GPU nodes:
 - 4 x NVIDIA A100
 - 1 x AMD Epyc 64-cores CPU
- CPU nodes:
 - 2 x AMD Epyc 64-cores CPUs
- 2 Virtual Clusters (VC)
 - Production: 42 GPU / 15 CPU Nodes
 - R&D: 30-50 GPU / ~15 CPU Nodes (elastic)





ICON-CH1 on Alps (MeteoSwiss)



- Issues with the new system + slower GPU performance as compared to COSMO – required time to solution < 3000 s
- First optimization brought some improvement.
- ICON ca 2-3x slower than COSMO for same configuration and hardware



Domain Specific Language (DSL) in weather and climate – really ?

- DSL : computer language restricted to a particular domain
- We need performance to reach time to solution
- Separation of concern between domain and computer scientist
- Single source code for multiple target architectures
- Possible to write a new backend when a new technology emerged
- Allow aggressive optimization without degrading readability of user code
- Allow optimization across components – data centric optimization



High level DSL for ICON

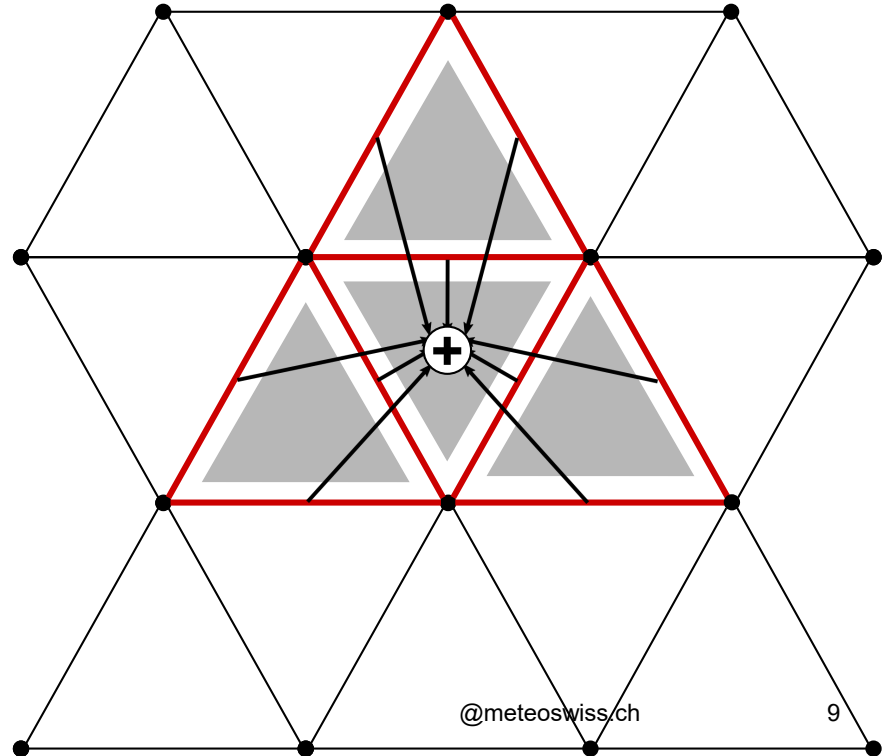
- Need to support unstructured grid, such as ICON grid
- New abstraction (e.g. neighbors operations)
- Focus on usability, productivity. Should be usable for domain scientist
- High level python dsl (gt4py)
- Development work in the EXCALIM project
 - High resolution use cases
 - Re-write code components using python DSL framework





Python DSL notation example (gt4py) : Neighbor Chains

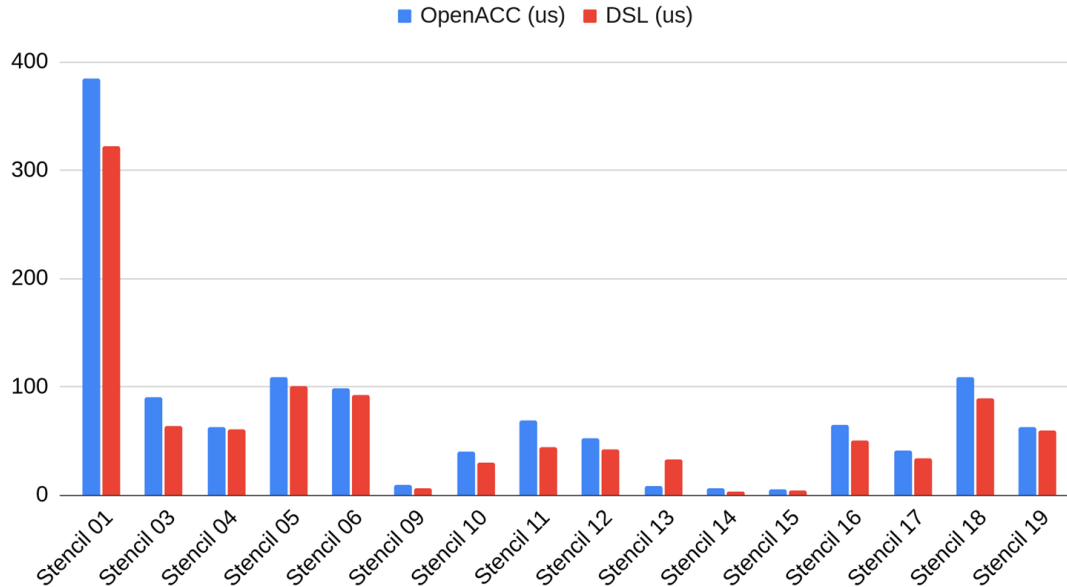
```
@field_operator
def intp(
    fe: Field[[EdgeDim], float],
    w: Field[[CellDim, C2E2C2EDim], float],
) -> Field[[CellDim], float]:
    f_c = neighbor_sum(w * f_e(C2E2C2E), axis=C2E2C2EDim)
    return f_c
```





Performance of ICON dycore (DSL) prototype

Open ACC vs DSL



DSL :
Dusk/
Dawn

- Stencil by stencil translation 10-20% faster DSL compared to OpenACC.
- Prototype DSL dycore about 40% (1.4x) faster then OpenACC - not fully optimized. Dry dycore only.



Current state gt4py re-write

- Dry dycore almost completely translated to gt4py (7 of 110 stencils still in progress)
 - Performance the same as dusk/dawn (tested for 20 stencils)
- Tracer advection partially translated (20 stencils), continuing work
- Dry dycore + tracer advection are 60%-70% of the runtime of a full run
- Next focus for dycore: optimization and robustness
- Also ongoing: microphysics using gt4py; focus here is good language support (example: if-else)



Conclusions

- First version of ICON model ported to GPU using OpenACC compiler directives
- Most components for global and regional NWP ported, and shall be soon available in icon-nwp master
 - basic version: Q4 2022
 - complete version: std NWP configurations tested with buildbot: Q4 2023
- Reasonable first performance, 3x faster than CPU, but still potential for optimization as still 2x time slower than COSMO
- Training of the ICON developers to work with OpenACC will be organized
- NWP application can likely benefit from DSL development in EXCLAIM, in particular for the dynamical core which shall be used once available

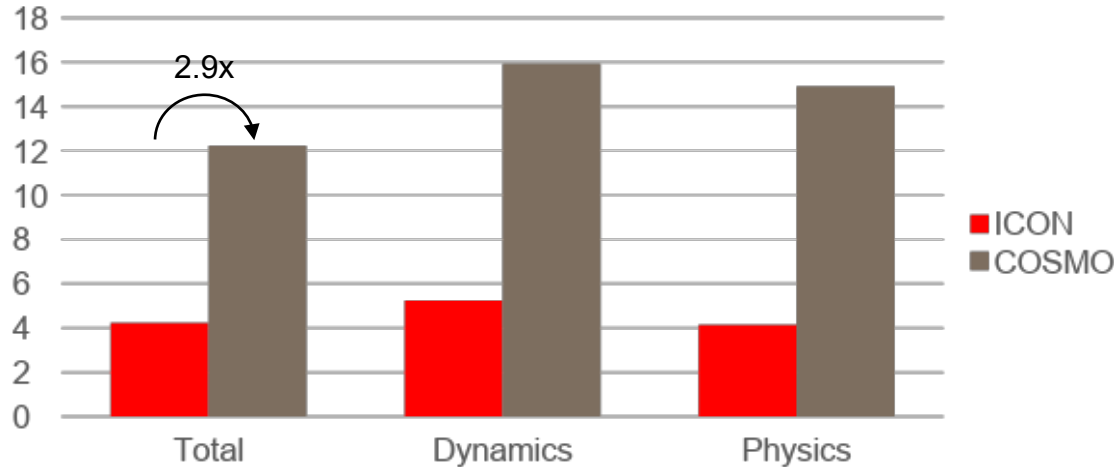


Additional slides

COSMO-ICON performance comparison on Daint

Socket to socket comparison: operational ICON-CH2 / COSMO-2E (2 km) 8 Nodes, 1h, P100 GPU vs Intel Xeon E5 12 cores (Piz Daint, CSCS, GPU node) - timeloop only

Speed up relative to ICON-CPU



- ICON CPU vs GPU : 4.2x speedup
- ICON is 2.9x slower than COSMO on GPU for an equivalent setup.

Porting and optimization challenges

OpenACC optimizations

- GPU and CPU working asynchronously
 - Reduces launch overhead
- Bundling similar loop constructs into single GPU kernels
 - Improves cache reuse
 - Reduces launch overhead
- Compiler assisted / manual inlining of function calls
 - Required for complex (deep call-trees) GPU kernels
 - Enables optimizations above

Conceptual challenges

- Tiling for surface and turbulence
 - Implicitly introduces sub-blocking which leads to underutilized GPUs
- Physics initialization on CPU
 - Prohibitively slow because of unsuitable nprma and MPI settings for CPU
- Radiation sub-blocking
 - Radiation (ec-rad) has an additional dimension which can be parallelized Sub-blocking as a memory optimization
- Code management
 - Disruptive code changes are challenging
 - ecrad: juggling diverse Institutes