rid id cells required

through

ns at nest

- Unstructured Voronoi (hexagonal) grid
 - Good scaling on massively parallel computers
- Smooth grid refinement on a conformal mesh
 - Increased accuracy and flexibility in varying resolution





oped for global applications on uniform and variable-mesparing for Exascale: Convection-permitting, global atmospheric simulations replace IDR Flein Hein MRAS Duble, drie Kunstmann, 2016: Geosci. Model Dev., 9, 77-110, http://www.geosci-model-dev.net/9/77/2016

complementary models!

3

3

2

bootstrapping

stream input

time integration

stream output

| | Convection-permitting the next grand challer | nvection-permitting global model applications are next grand challenge in NWP and on the horizon next-generation massively parallel HPC systems | | |
|--|---|---|----------------|--|
| | Extreme scaling experiment with MPAS on E7.L.III.QUEEN (IBM Bluegene (Q) in 2015 | | | |
| | Uniform 3km mosh 65 536 002 colls | | | |
| | 41 vertical levels, deuble precision | | | |
| | • 41 vertical levels, double precision | | | |
| | • Inr model integration, no file output | | | |
| | Initial conditions: 1.11B pnetCDF CDF5 | | | |
| | Min. 4096 nodes, 65TB memory; max: 28,672 node | | | |
| | Fastest run: 6.3 x real-time, 1.6 Mio CPUh/24h | | | |
| | The dynamical solver of MPAS-Atmosphere scales up | | | |
| | 400,000 MPI tasks (160 owned cells per task). Show-s | | | |
| | at extreme scale are file I/O and model setup (bootstra | | | |
| | Step 1. Addressing the file I/O performance | | | |
| | SIONIIb I/O layer (http://www.fz-juelich.de/jsc/sion | | | |
| | massively parallel I/O in addition to existing I/O fo | | | |
| | Post-processor of | Post-processor core for converting to netCDF, rec | | |
| | to lat-lon grids a | to lat-lon grids and interpolation to station location | | |
| | Reading/writing in SIONlib format requires to use | | | |
| | same number of MPI tasks and the same graph p | | | |
| | Information encoded in SIONlib data can be used | | | |
| | skip parts of the bootstrapping at model startup | | | |
| | The SIONIib I/O layer addresses file I/O and model | | | |
| | Timer name | pnetCDF, CDF5 | SIONIib | |
| | 1 total time | 3585 | 211 | |
| | 2 initialiaa | 1176 | 2 · · 2 0 / | |
| | | 11/0 | 24 | |

540

612

1580

818



Step 2. Reducing MPI communication overhead

Hybrid MPI+OpenMP parallelisation to speed up bootstrapping and decrease halo exchange time in dynamical solver of MPAS-A

Optimisation for latest many-core architectures Combine with SIONIIb I/O layer improvements for maximum performance at extreme scale

Threading of one additional routine in solver

Avoid repeated creation and destruction of threads

Optimised code

Original code

call non-threaded function !\$OMP parallel do do thread=1,nThreads ! call threaded function

!\$OMP end parallel do

!\$OMP parallel do do thread=1,nThreads ! call threaded function

!\$OMP end parallel do

!\$OMP parallel !\$OMP do schedule(static,1) do thread=1,nThreads ! call threaded function end do !\$OMP end do !\$OMP master ! call non-threaded function !\$OMP end master !\$OMP end parallel

4096 x 16 x 1 4096 x 8 x 2 4096 x 2 x 8 4096 x 1 x 16

Original vs. optimised hybrid code on Intel Xeon Phi KNL for a uniform 240km mesh (10242 cells).

Note: Using optimised 🚊 1000 compiler flags for the Intel KNL, performance gains are smaller (≈20%)

2.5

Putting it to the test: extreme scaling experiment on FZJ JUQUEEN in 2017

- Jablonowski & Williamson (2006) baroclinic wave 2km mesh, 147,456,002 cells, single precision
- 26 vertical levels, 10min model integration, $\Delta t=5s$
- File input 1.8TB, output 4TB (SIONlib)
- Dynamics and file I/O only, no physics: more stringent test of dynamical solver
- MPI only: threading model not
- supported on JUQUEEN
- Time integration speedup (dynamics and file I/O) (100000



- Convection-permitting global simulations are within reach of current and next-generation HPC systems
- Efficient parallel I/O, code preparation for novel many-core architectures and HPC-specific adaptation are key to success

