

MPAS Extreme Scaling Experiment

Towards convection-resolving, global climate simulations

 ¹ INSTITUTE OF METEOROLOGY AND CLIMATE RESEARCH, ATMOSPHERIC ENVIRONMENTAL RESEARCH (IMK-IFU) Regional Climate Systems/Regional Climate and Hydrology
² NATIONAL CENTER FOR ATMOSPHERIC RESEARCH, EARTH SYSTEM LABORATORY (NESL) Mesoscale Prediction

Third JUQUEEN Porting and Tuning Workshop, 2th-6th February 2015

Dominikus Heinzeller¹ and Michael Duda²







KIT – University of the State of Baden-Wuerttemberg and National Research Center of the Helmholtz Association

www.imk-ifu.kit.edu

MPAS - Model for Prediction Across Scales



Collaborative project (NCAR/LANL) for developing earth-system simulation components for use in climate, regional climate and weather studies.





MPAS - Model for Prediction Across Scales



Collaborative project (NCAR/LANL) for developing earth-system simulation components for use in climate, regional climate and weather studies.



MPAS Website: http://mpas-dev.github.io

Skamarock et al.: A Multi-scale Nonhydrostatic Atmospheric Model Using Centroidal Voronoi Tesselations and C-Grid Staggering. 2012 MWR, 240, 3090–3105

Ringler et al: Ocean Modelling. 2013 Ocean Modelling, 69(C), 211–232





Key features of MPAS



Variable resolution meshes, smooth transitions





Key features of MPAS



Variable resolution meshes, smooth transitions

Key features of MPAS-A

Split-explicit time integration scheme Numerical schemes ported from WRF,

work underway to make scale-aware







Key features of MPAS-A

Key features of MPAS







Key features of MPAS-A

Key features of MPAS





Variable resolution meshes and filtering





WRF Characteristics

- Lat-Lon global grid
 - Anisotropic grid cells
 - Polar filtering required



- Grid refinement through domain nesting
 - Flow distortions at nest boundaries



Variable resolution meshes and filtering



WRF Characteristics

- Lat-Lon global grid
 - Anisotropic grid cells
 - Polar filtering required

MPAS Characteristics

- Unstructured Voronoi (hexagonal) grid
 - Good scaling on massively parallel computers





- Grid refinement through domain nesting
 - Flow distortions at nest boundaries
- Smooth grid refinement on a conformal mesh
 - Increased accuracy and flexibility in varying resolution





Variable resolution meshes and filtering



WRF Characteristics

- Lat-Lon global grid
 - Anisotropic grid cells
 - Polar filtering required

MPAS Characteristics

- Unstructured Voronoi (hexagonal) grid
 - Good scaling on massively parallel computers







- Grid refinement through domain nesting
 - Flow distortions at nest boundaries
- Smooth grid refinement on a conformal mesh
 - Increased accuracy and flexibility in varying resolution

- Short-wavelength modes will be reflected in a fine-coarse mesh transition if not filtered.
- Smooth transitions result in reflection of only the very-shortest wavelength modes (locally), where filters are efficient.
- Abrupt transitions result in reflection of longer waves, where filters are less efficient.
- In MPAS: Smagorinski, viscosity, hyperviscosity filters





Credits: Bill Skamarock

Parallel decomposition





The dual mesh of a Voronoi tessellation is a Delaunay triangulation – the connectivity graph of the cells.

Parallel decomposition of an MPAS mesh becomes a graph partitioning problem: equally distribute nodes among partitions (each process gets equal work) while minimizing the edge cut (minimizing communication).

Graph partitioning

Parallel graph decomposition is done in a pre-processing step using the Metis package (can also handle weighted graph partitioning).



Parallel decomposition (continued)



Given an assignment of cells to a process, any number of layers of halo (ghost) cells may be added.

With a list of cells stored in a block, adjacent edge and vertex locations can be found. A simple rule is applied to determine the ownership of edges and vertices adjacent to real cells in different blocks.

nEdges



Block of cells owned by a process



Block plus two layers of halo/ghost cells





nEdgesSolve

Cells are stored in a 1D array (2D with vertical dimension...), with halo cells at the end.

Overall MPAS performance



Consider performance in terms of number of model time steps completed in one wallclock hour (allows for easier comparison of performance and scaling).



Simulation rate of double-precision MPAS-A on NCAR's "Yellowstone" computer; 45 vertical layers (41 layers for 7.5, 5, and 3 km), full suite of physics, no I/O.



Overall MPAS-A parallel efficiency



The full MPAS-A solver (30-km, double-precision, physics+dynamics, no I/O) achieves 68.69% efficiency down to about 160 owned cells per MPI task.



MPAS-A parallel efficiency relative to 10240 owned cells per MPI task, Yellowstone MPAS-A scaling – 30-km mesh, Yellowstone (ncell = 655362)





Test case 1: regular 120km mesh x1.40962







Test case 1: regular 120km mesh x1.40962





1PAS

Test case 2: variable 120/25km mesh x4.163842 MPAS









Test case 2: variable 120/25km mesh x4.163842 MPAS





Test case 2: variable 120/25km mesh x4.163842 MPAS 10^{2} JUQUEEN, ideal **FZJ JUQUEEN** JUQUEEN, real 32 to 1024 nodes CURIE, ideal CURIE, real 512 to 16384 threads JUROPATEST, ideal JUROPATEST, real **TGCC CURIE** 1 to 192 nodes Scaling 16 to 3072 threads 10¹ **FZJ JUROPATEST** ncell/160 1 to 60 nodes (max) 14 to 840 threads ncell = 163842, nvert = 41, double-precision, physics+dynamics, 10^{0} 10² 10^{3} 10⁴ 10¹ parallel I/O, standard opt. Number of threads



















MPAS





MPAS











Extreme scaling: regular 3km mesh x1.65536002 MPAS



Notos

Pure MPI parallelization in MPAS

Expect breakdown of parallel • efficiency at about 400000 threads

Notes

Trial OpenMP hybrid mode during tuning and porting workshop (OpenMP test code in place)

Exploit limiting factors in the model when the scaling breaks for extreme numbers of threads



Extreme scaling: regular 3km mesh x1.65536002 MPA



Key data

I/O sizes: initial condition file 1.1Tb, restart file 2.1Tb, history file 250Gb

Notes

Pure MPI parallelization in MPAS

Expect breakdown of parallel • efficiency at about 400000 threads

Trial OpenMP hybrid mode d Probably Not! porting workshop (OpenMP test code in place)

Exploit limiting factors in the model when the scaling breaks for extreme numbers of threads

