# MPAS Extreme Scaling Experiment

**Towards convection-resolving, global climate simulations**

[1] KARLSRUHE INSTITUTE OF TECHNOLOGY, INSTITUTE OF METEOROLOGY AND CLIMATE RESEARCH
[2] NATIONAL CENTER FOR ATMOSPHERIC RESEARCH, MESOSCALE AND MICROSCALE METEOROLOGY LABORATORY
[2] AUGSBURG UNIVERSITY, DEPARTMENT OF GEOGRAPHY, CHAIR OF REGIONAL CLIMATE AND HYDROLOGY

**16th Annual WRF Users' Workshop, 17/06/2015**

**Dominikus Heinzeller[1], Michael Duda[2] and Harald Kunstmann[1,3]**

**KIT-Campus Alpin**
IMK-IFU: Atmospheric Environmental Research

Karlsruhe Institute of Technology

NCAR
NATIONAL CENTER FOR ATMOSPHERIC RESEARCH

# MPAS variable-resolution mesh with 535554 cells

**Approximate mesh resolution (km)**



CONTOUR FROM 16 TO 52 BY 4

16   20   24   28   32   36   40   44   48   52

Model consistency:

land use / soil type classification

flow distortions of jets and waves

Research question: teleconnection between African monsoon systems and Atlantic/Indian ocean

# MPAS variable-resolution mesh with 535554 cells



Approximate mesh resolution (km)

60-12km mesh cell distribution

Research question: teleconnection between African monsoon systems and Atlantic/Indian ocean

# Variable-resolution mesh with 535554 grid cells



**Tested systems**

TGCC Curie

FZJ Juropatest

FZJ Juqueen

SCC ForHLR1

IBM Bluegene /Q: 28 racks, 1024 nodes per rack, 16 cores and 16GB memory per node.

West African Monsoon in MPAS, WRF, CCLM

Total monthly precipitation [mm] - July 1982
(initialisation: CCLM/WRF 1979-01-01, MPAS 1981-09-01)

# West African Monsoon in MPAS, WRF, CCLM



**Total monthly precipitation [mm] - July 1982**

(initialisation: CCLM/WRF 1979-01-01, MPAS 1981-09-01)

CCLM/WRF from the WASCAL regional climate simulations for West Africa using an optimised setup for the region (see poster 63).

# West African Monsoon in MPAS, WRF, CCLM

# MPAS Extreme Scaling Experiment as part of the 3rd Juqueen Porting & Tuning Workshop

# Extreme scaling on Juqueen: Pandora's box

| regular 3km mesh x1.65536002<br>ncell = 65536002, nvert = 41 | expect ≥70% parallel efficiency<br>up to 400000 threads (150 cells per task) |
|---|---|
| double precision, physics+dynamics:<br>memory requirement 0.175MB per cell | 11.5TB of memory required for 3km mesh,<br>minimum number of nodes is approx. 750 |
| I/O sizes: initial condition 1.1TB, restart 2.1TB, diagnostics 15GB, history 250GB | |

# Extreme scaling on Juqueen: Pandora's box

MPAS
Model for Prediction Across Scales

| | |
|---|---|
| regular 3km mesh x1.65536002<br>ncell = 65536002, nvert = 41 | expect ≥70% parallel efficiency<br>up to 400000 threads (150 cells per task) |
| double precision, physics+dynamics:<br>memory requirement 0.175MB per cell | 11.5TB of memory required for 3km mesh,<br>minimum number of nodes is approx. 750 |

I/O sizes: initial condition 1.1TB, restart 2.1TB, diagnostics 15GB, history 250GB

**Pitfalls:**

Test runs on 1 and 2 racks (1024, 2048 nodes) failed - memory errors. Defragmentation issue?

Model initialisation took way to long (1hr+), tuning of I/O and other parameters required

# A moment of excitement (for geeks)

# Extreme scaling results

## Key facts

- 1hr integration, cold start (init.nc)
- 12s integration time step ($\Delta x*4$) to avoid NaNs
- no output except one diagnostic file (15Gb)
- 1 I/O task per 128 tasks, i.e., 128 per rack
- hash table size increased from 27183 to 6 Mio

# Extreme scaling results

| BG size (# nodes) | Parallel efficiency integration only | Integration in 24h walltime (no output) | Estimated integration in 24h walltime (with output) |
|---|---|---|---|
| 4096 | 100.0% | 48h | 29h |
| 8192 | 91.2% | 87h | 53h |
| 16384 | 90.1% | 172h | 104h |
| 24576 | 87.7% | 250h | 152h |
| 28672 | 69.5% | 231h | 141h |

| BG size (# nodes) | Number of tasks | CPUh per 24h model run time (with output) | Speedup wrt. real time (with output to disk) |
|---|---|---|---|
| 4096 | 65536 | 1.30 Mio | 1.21 |
| 8192 | 131072 | 1.42 Mio | 2.21 |
| 16384 | 262144 | 1.45 Mio | 4.33 |
| 24576 | 393216 | 1.49 Mio | 6.33 |
| 28672 | 458752 | 1.87 Mio | 5.88 |

# Take home messages

MPAS can reproduce the dynamics of the West African Summer Monsoon using an out-of-the box setup and is a promising tool for climate modelling.

Global, convection-resolving atmospheric simulations with MPAS are within reach of current/next generation HPC facilities. Open problems are model initialisation, disk I/O and post-processing of the data.

# Take home messages

MPAS can reproduce the dynamics of the West African Summer Monsoon using an out-of-the box setup and is a promising tool for climate modelling.

Global, convection-resolving atmospheric simulations with MPAS are within reach of current/next generation HPC facilities. Open problems are model initialisation, disk I/O and post-processing of the data.

Model developers/users have to tailor and optimise models to run on modern massively parallel systems, which requires in-depth knowledge and time.

# Bonus material

# Limitations of limited area modelling

How to include new land surface and soil data or dynamic changes thereof in a consistent way?

GCM standard, Volta region (West Africa)

# Limitations of limited area modelling

How to include new land surface and soil data or dynamic changes thereof in a consistent way?



GCM standard, Volta region (West Africa)

RCM standard land cover (MODIS)

# Limitations of limited area modelling

How to include new land surface and soil data or dynamic changes thereof in a consistent way?



GCM standard, Volta region (West Africa)

RCM standard land cover (MODIS)

RCM high-resolution land cover
(DLR: MODIS, ASAR, TanDEM-X)

Credits: Ursula Gessner, DLR

# Limitations of limited area modelling

How to include new land surface and soil data or dynamic changes thereof in a consistent way?



Credits: Ursula Gessner, DLR

GCM standard, Volta region (West Africa)

RCM standard land cover (MODIS)

RCM high-resolution land cover
(DLR: MODIS, ASAR, TanDEM-X)

Consistency? Feedback processes?

Dynamic changes of land use?

CMIP5 vs CORDEX issue:
land use classification
and land use change in
ESMs not reflected in RCMs

# MPAS, WRF, CCLM vs. observations - any good?

**Near surface temperature [°C] - July 1982 monthly mean**

(initialisation: WRF 1979-01-01, MPAS 1981-09-01)

# At least one monsoon period needed for spinup

## Soil temperature [°C] (top) and relative soil moisture [%] (bottom) - July 1982 mean

# At least one monsoon period needed for spinup



Mean sea level pressure [hPa] - July 1982 monthly mean

# Take home messages

MPAS can reproduce the dynamics of the West African Summer Monsoon using an out-of-the box setup and is a promising tool for climate modelling.

Global, convection-resolving atmospheric simulations with MPAS are within reach of current/next generation HPC facilities. Open problems are model initialisation, disk I/O and post-processing of the data.

# Take home messages

MPAS can reproduce the dynamics of the West African Summer Monsoon using an out-of-the box setup and is a promising tool for climate modelling.

Global, convection-resolving atmospheric simulations with MPAS are within reach of current/next generation HPC facilities. Open problems are model initialisation, disk I/O and post-processing of the data.

Model developers/users have to tailor and optimise models to run on modern massively parallel systems, which requires in-depth knowledge and time.