

### Dealing with I/O performance issues within the Unified Model SIG-IO-UK - Wed June 6th, 2018

Adam Voysey, HPC Optimisation Team

www.metoffice.gov.uk





### A Quick introduction to I/O in the Unified Model



# The UM

The Unified Model (UM) is a numerical model of the atmosphere used for both weather and climate applications. It is in continuous development by the Met Office and its partners, adding state of the art understanding of atmospheric processes to new releases.

For more information see: https://www.metoffice.gov.uk/research/modelling-systems/unified-model





# Two Key I/O Components I/O Server (Remote I/O)

- Written in Fortran
- OpenMP + MPI Parallelism
- Multiple "Roles"



www.metoffice.gov.uk

### Common Interface allows these to be tied together

### Portio (Local I/O) • Written in C Serial Code\* Multiple "Layers"





# Two Key I/O Components

- I/O Server (Remote I/O)
- Written in Fortran
- OpenMP + MPI Parallelism
- Multiple "Roles"

www.metoffice.gov.uk





### Portio (Local I/O) • Written in C Serial Code\* Multiple "Layers"

### MPI-IO



# Fieldsfiles

- Common "Fixed Length Header"
- table, lookup table...
- Fields can be compressed e.g. WGDOS Packing

### Bespoke format (inspired by early GRIB) Secondary records include constants





# Fieldsfiles

### Row Header ! Row Data



### WGDOS Packing

www.metoffice.gov.uk











# NetCDF

### Can also output in NetCDF Not as commonly used • Not start dumps

www.metoffice.gov.uk

Some ancillary input uses NetCDF files



# Why is I/O a Problem?



### Have to deal with:

### • Performance

- I/O can be serial bottleneck
- Significant chunk of runtime
- Cost • Archive tapes are expensive!

# • Even where we can parallelise, scaling may be poor/limited

Increasing resolution is making this worse







www.metoffice.gov.uk



### What Have We Done So Far?



# Portio Layers

- Internal Buffering
- Lustre integration
- 'Helper' threads
- Add debugging layers

www.metoffice.gov.uk



# Read Buffering. (+ pre-loading) Support for e.g. striping Accelerate e.g. buffering

### Trace, Time, Throttle, Blackhole



# All-to-All Dump Reading

- Read dump fields in parallel; Use MPI Alltoall to distribute
- Based on work done by Cray
- of an operational-style n1280 global model.

www.metoffice.gov.uk

Example timings shown for three runs

estep / sec

Time to First Til



### © Crown Copyright 2018, Met Office

### With All-to-All



# MPI-IO in IO Server

- Write output in parallel
- Based on work done by Dale Roberts (Australian National University) Complications from packed fields
- Gather fields onto IO Server ranks; write in parallel (mpi\_file\_write\_at\_all)



www.metoffice.gov.uk





# Other Optimisations

- Namelists read by only one rank
- Reduction of enquiry functions

• Ensure only one rank opens files in serial case Ensure correct compiler flags. (No debug symbols with Cray Compiler)



### Current Challenges



# I/O Stalls On Lustre

- Usually get good performance

- - Episodic

• However, occasionally see a long stall (maybe ~10x I/O time) • Stalls are infrequent, but not negligibly rare • Mechanism not properly understood, but... • Hypothesis that stalls are related to Lustre response • When they occur, stall time is consistent. • More probable when higher frequency of transactions • ... when we have more threads • ... when we have smaller block sizes





### What Do We Still Need To Do?



### I/O Server Work on queues computation

www.metoffice.gov.uk

### • Want to ensure I/O server is responsive and does not block overlapping



### Portio

- Enable write helper thread
- • (Migrate to Shumlib)

www.metoffice.gov.uk

Add new layers for other file-system architectures?



### Other Changes Extend All-to-All method to other files – IAU, LBCs • Alternative compression?

www.metoffice.gov.uk





### Shumlib

- SHared UM LIBrary
- High level Fieldsfile manipulation API
- And low level I/O API
- 3-Clause BSD licence
- (Plus other non-IO related code)

# Provides bindings in Fortran, Python, and C





### Conclusions



### Conclusions

- • Why do we get stalls?

Taken some steps to deal with I/O performance in the UM Can now have parallel reads and writes of most important files • This is most important for operational runs • But still need to do more – e.g. can extend to other files We still need to do work to understand fully the low level details



### Thanks Any Questions?

www.metoffice.gov.uk

