

Porting and Optimisation of UM on ARCHER

Karthee Sivalingam, NCAS-CMS

HPC Workshop ECMWF



**National Centre for
Atmospheric Science**

NATURAL ENVIRONMENT RESEARCH COUNCIL



Acknowledgements

➤ NCAS-CMS

➤ Bryan Lawrence

➤ Jeffrey Cole

➤ Rosalyn Hatcher

➤ Andrew Heaps

➤ David Hassell

➤ Grenville Lister

➤ William McGinty

➤ Simon Wilson

➤ Annette Osprey

➤ Met Office

➤ Mick Carter

➤ Jean-Christophe Rioual

➤ Matthew Hambley

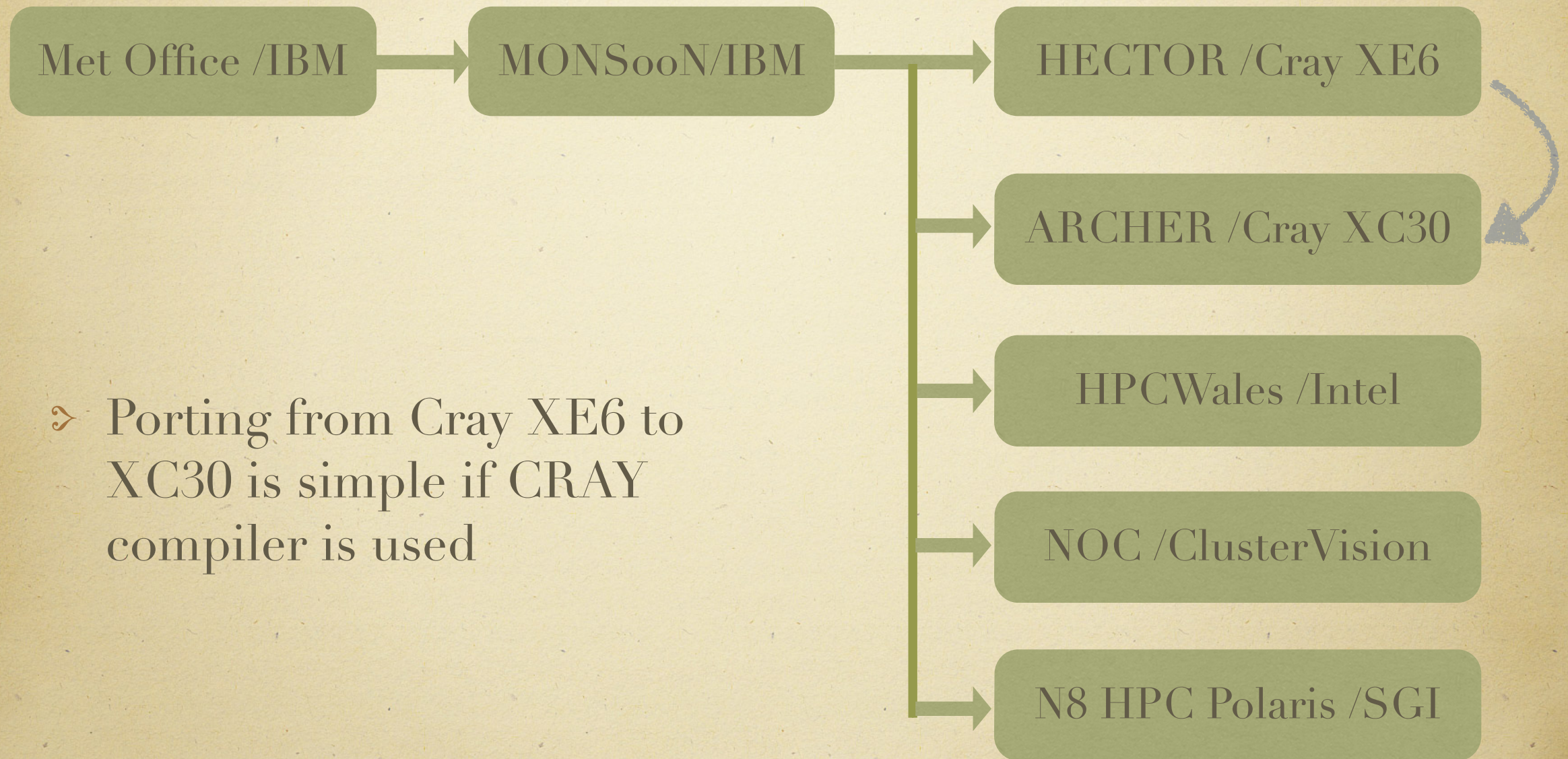
➤ Jamie Kettleborough

➤ CRAY & EPCC

Summary

- Porting UM to ARCHER (XC30)
- Initial Performance on ARCHER
- Performance analysis
- Tuning MPI ~~and IO~~
- Optimising thread performance

Porting



➤ Porting from Cray XE6 to XC30 is simple if CRAY compiler is used

HECTOR and ARCHER

	AMD Opteron	Intel Xeon
Clock speed	2.3	2.7
Cores per node	32	24
threads per core	1	2
vector	sse	avx
bits wide	128	256
L2 cache	512 KB	256 KB
L3 cache	6 MB	30 MB
Peak DP flops per node	294 GFlops	518 GFlops

Memory bandwidth , MPI bandwidth and IO bandwidth ?

ARCHER ...

- 3008 compute nodes
- Dual socket , 12 cores - forming 2 NUMA nodes
- 64 (or 128) GB memory per node
- Cray Aries network
- Cray XC30 Dragonfly Topology
- Copper and Optical cabling
- Lustre parallel file system

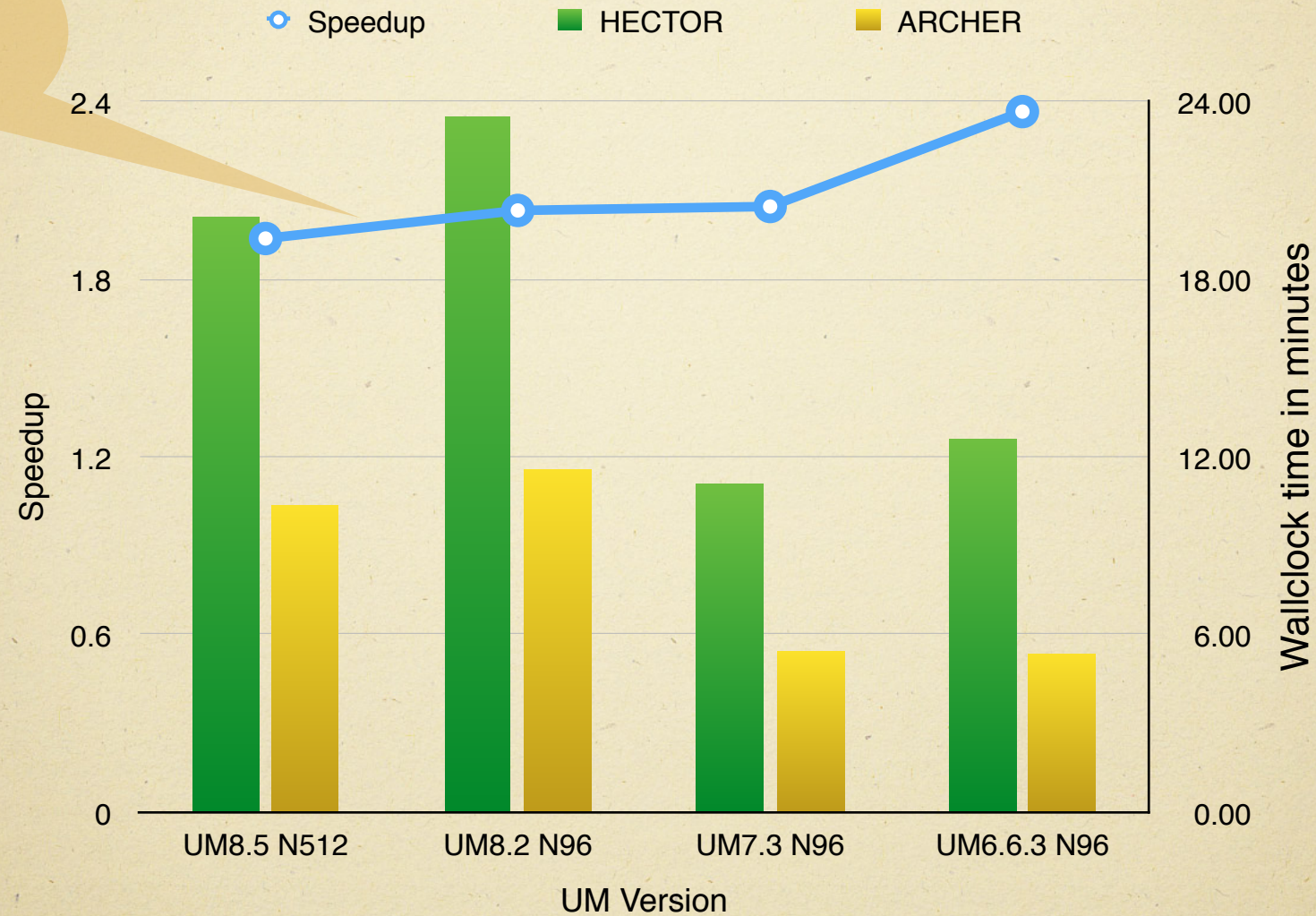
UM jobs ...

Job name	Columns	Rows	Land points	Vertical levels
N96	192	144	11271	85
N216	432	324	52614	85
N512	1024	768	280592	85

All the jobs are standard jobs ported from MetOffice to ARCHER. All the jobs use Cray compiler and strict bit comparability is enforced.

Initial Performance Comparison

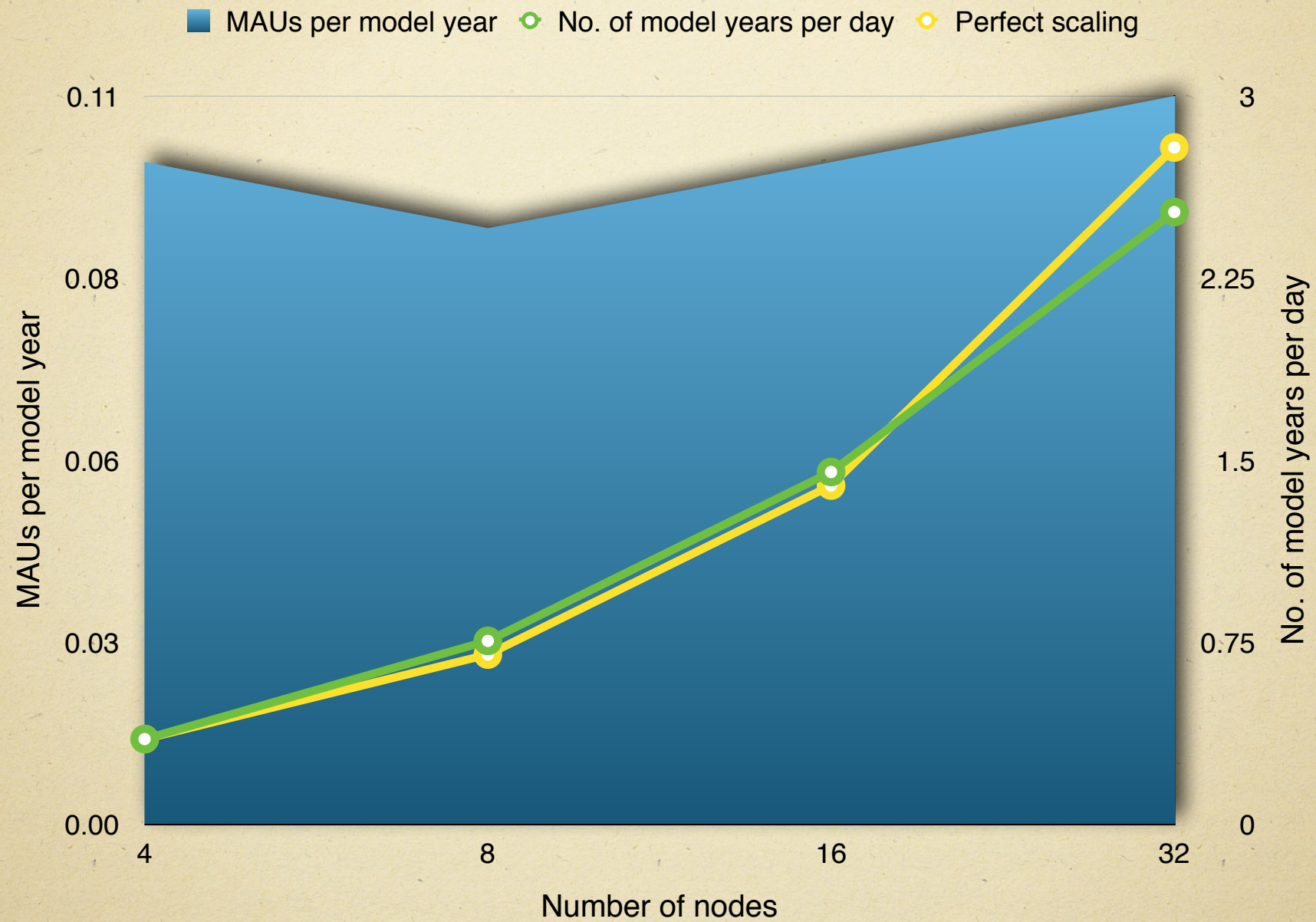
4 times
throughput?



Above jobs in HECTOR and ARCHER use similar MPI decomposition and same number of threads

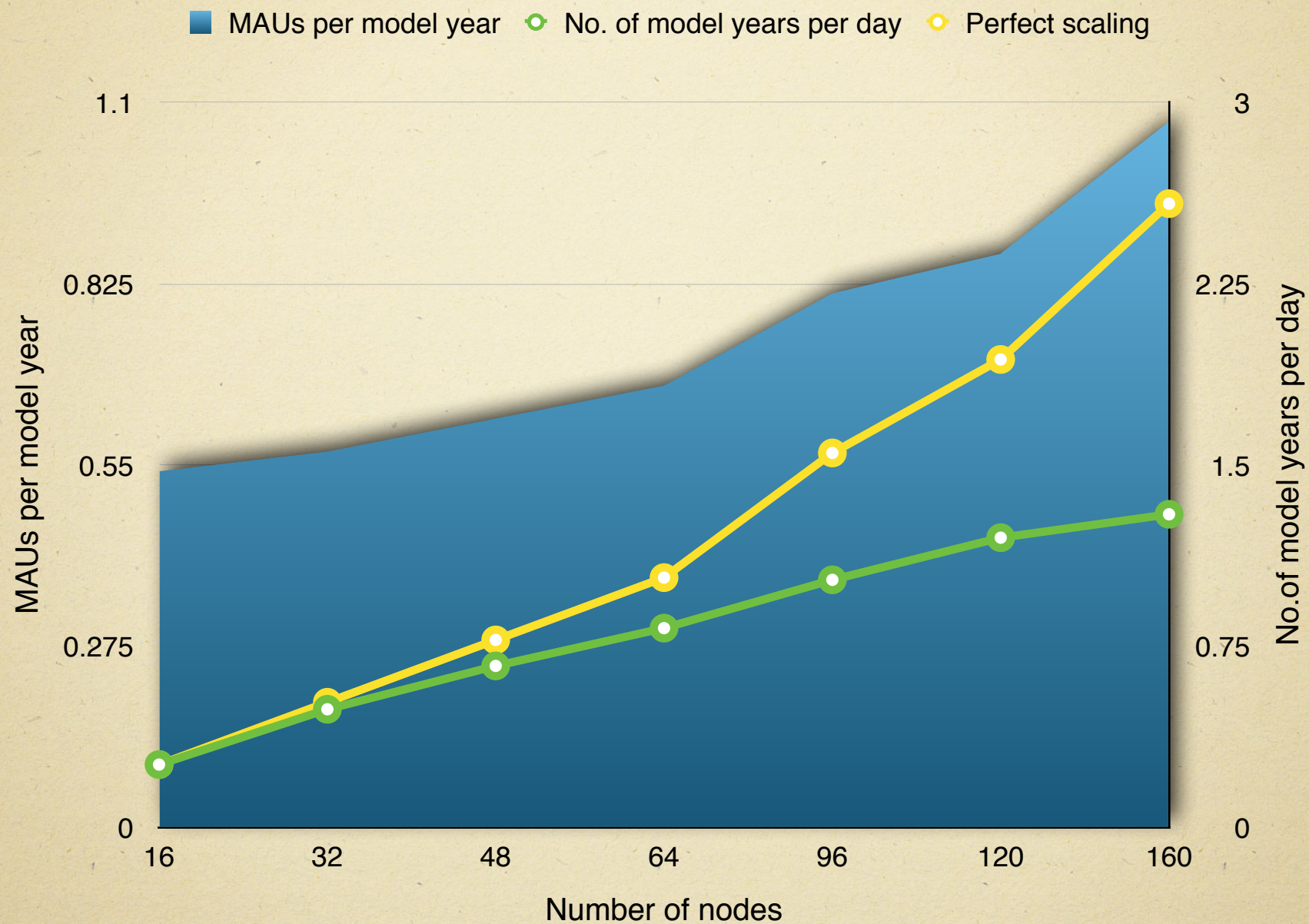
Wallclock time refers to the time taken to complete simulation of 1 model day

N96 job on ARCHER



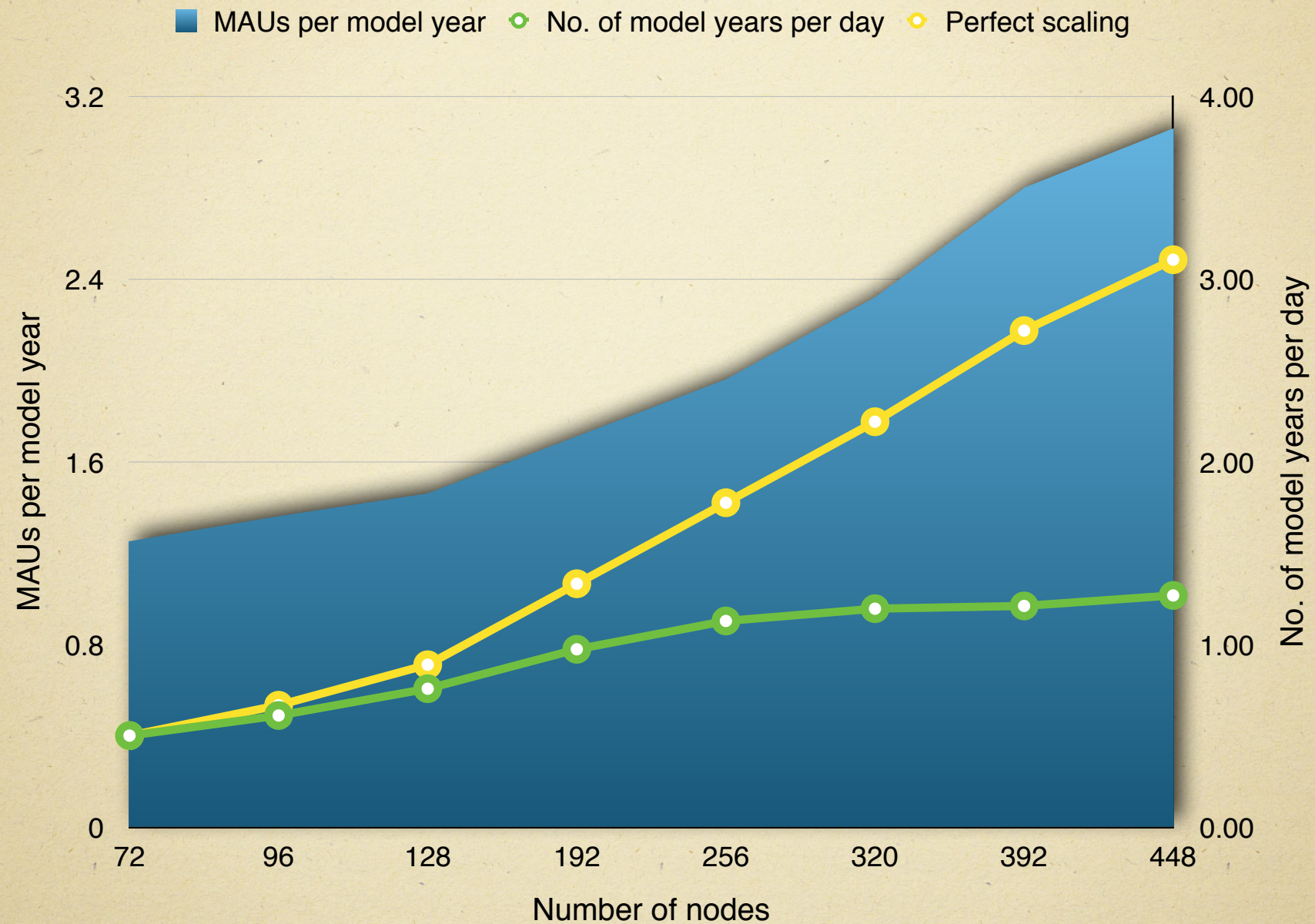
All the above runs uses 12 MPI PEs and 2 threads per node. Cost per MAU of ARCHER is £1640

N216 job on ARCHER



All the above runs uses 12 MPI PEs and 2 threads per node. Cost per MAU of ARCHER is £1640.

N512 job on ARCHER



All the above runs uses 12 MPI PEs and 2 threads per node. Cost per MAU of ARCHER is £1640.

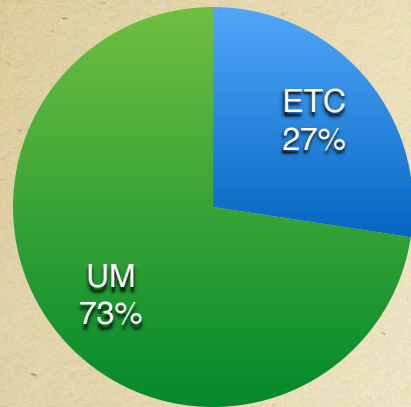
When number of PEs in EW or NS goes more that 120 or 60 respectively, the model does not bit compare.

What to do?

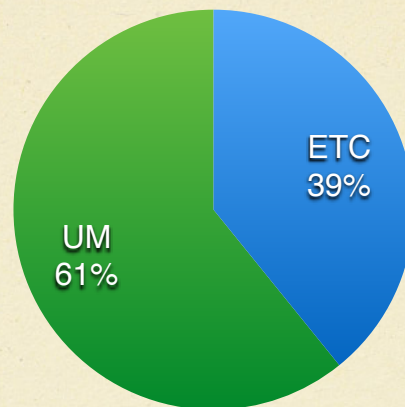
- N512 job does not scale beyond 250 nodes.
- N216 job does not scale beyond 96 nodes.
- Is 12 MPI PEs per node ideal? Is 2 threads per PE ideal ?
Can we use hardware threads ?
- Need to analyse the performance.
- CRAY performance analysis tools readily available.
 - CrayPAT
 - Apprentice

Profile analysis

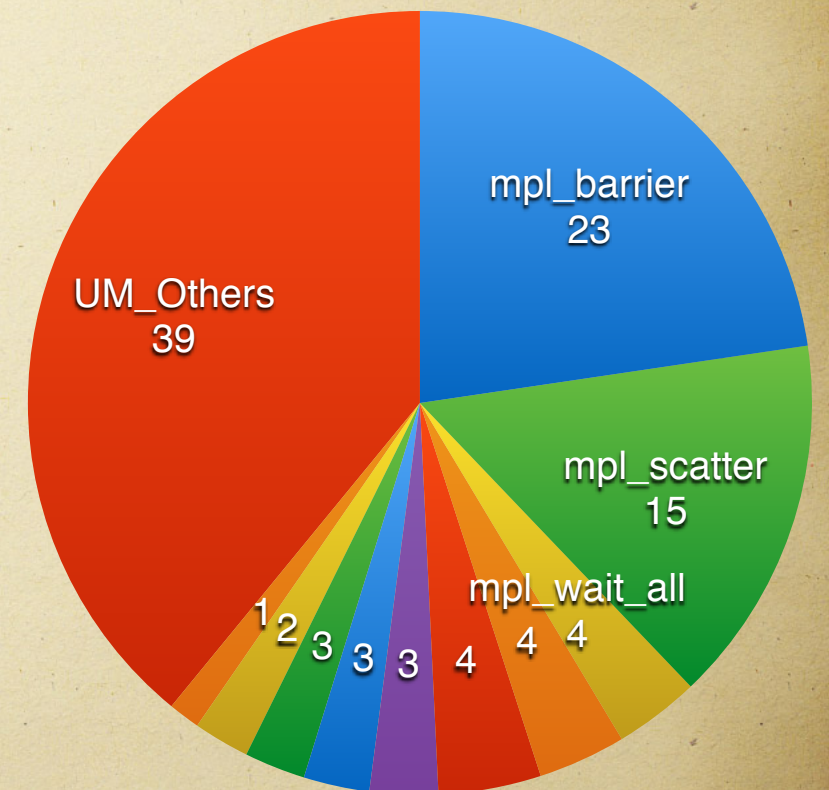
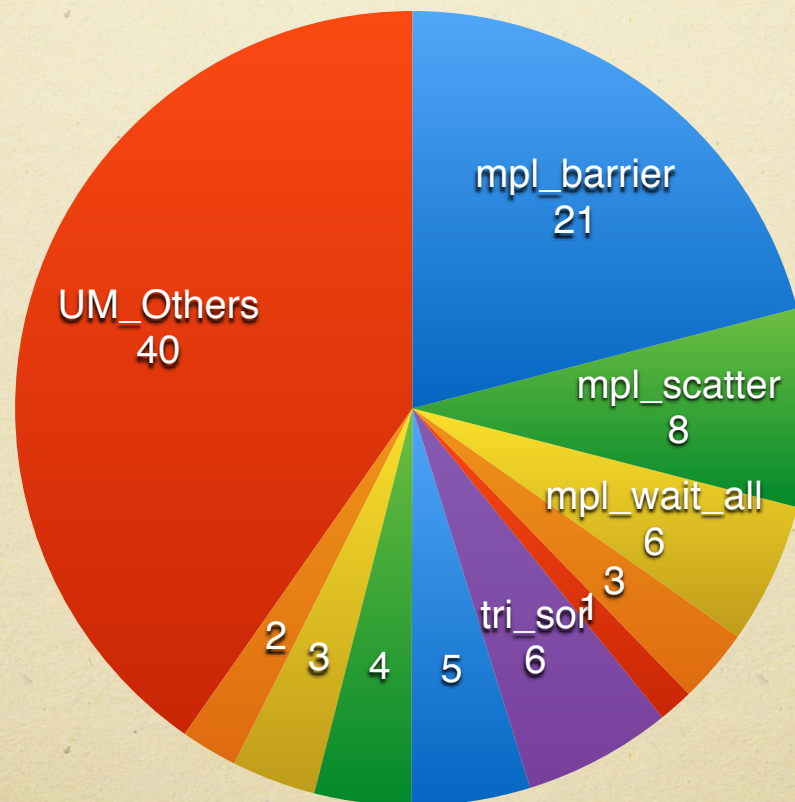
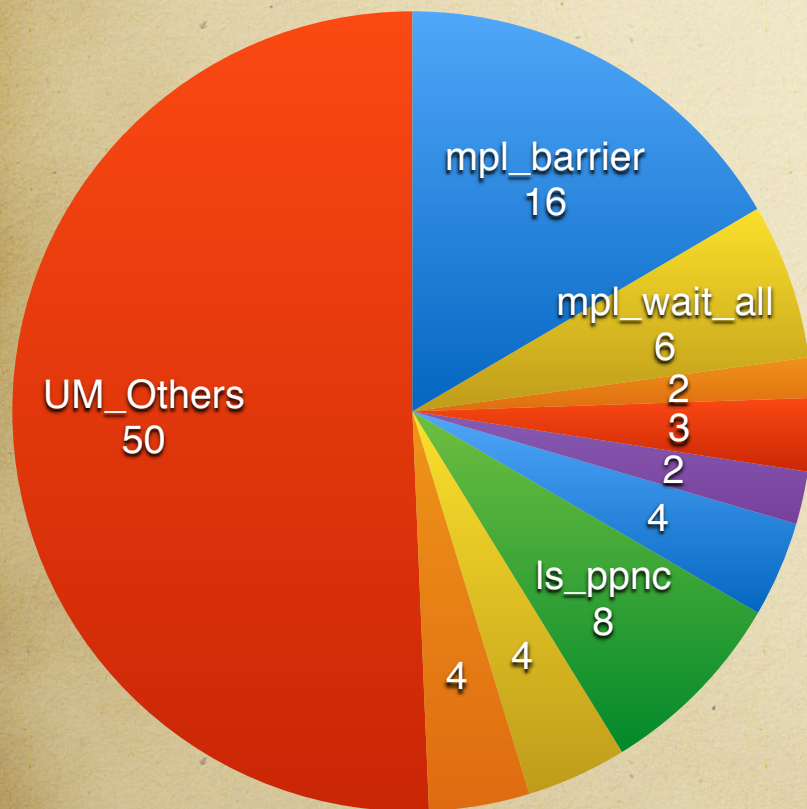
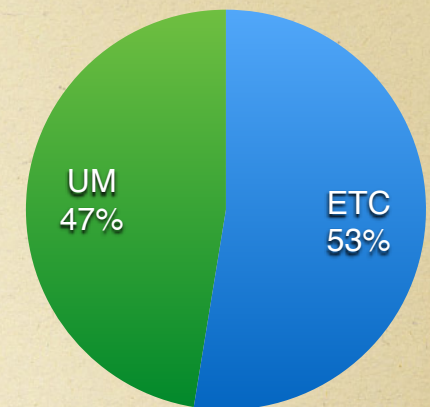
N96 - 8 nodes



N512 - 72 nodes



N512 - 192 nodes



- mpl_barrier
- mpl_scatter
- mpl_wait_all
- mpl_allreduce
- ETC_Others
- tri_sor
- eg_cubic_lagrange
- ls_ppnc
- ni_conv_ctl
- glue_rad
- UM_Others

MPI Bandwidth

- Message passing in UM is done only by thread 0.
- Message passing overhead increases by 14% as the number of nodes is increased from 72 to 192.
- Scaling of MPI bandwidth improved by using
 - MPI Asynchronous progress
 - Use unassigned CPUs for MPI
 - MPI rank placement (Grid order)

PE decomposition 24 x 36

0	1	2	3	4	...	35
36	37	38	39	40	...	71
72	73	74	75	76	...	107
108	109	110	111	112	...	143
...
828	829	829	829	829	...	863

12 PEs per node

Nearest neighbour communications

SMP - MPI Rank Order

GRID - MPI Rank Order

Node 0

Node 4

0	7	8
1	6	9
2	5	10
3	4	11

36	43	44
37	42	45
38	41	46
39	40	47

Node 10

Node 6

108	115	116
109	114	117
110	113	118
111	112	119

72	79	80
73	78	81
74	77	82
75	76	83

Node 0

Node 4

0	1	2
36	37	38
72	73	74
108	109	110

111	112	113
75	76	77
39	40	41
3	4	5

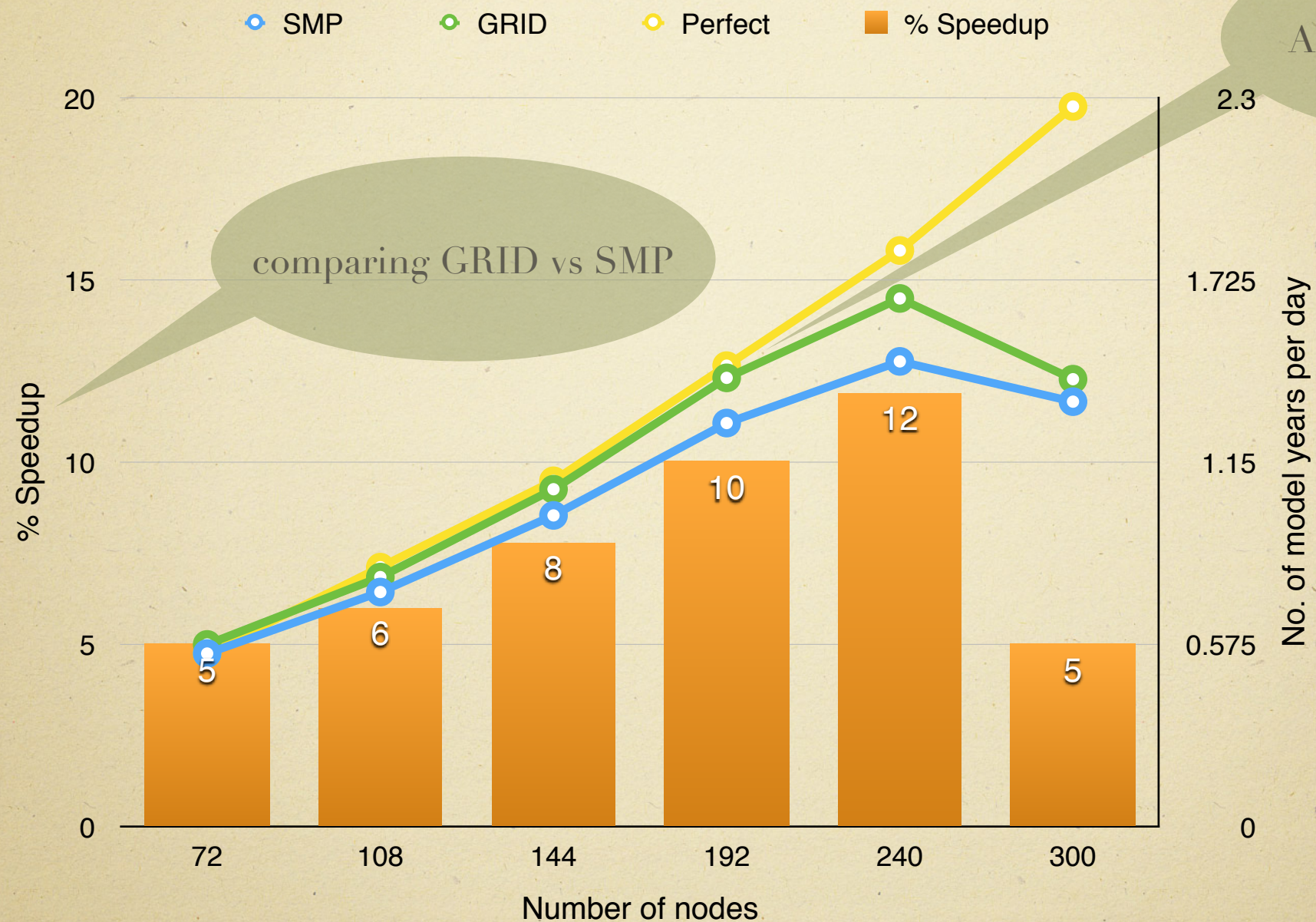
Node 2

Node 3

144	145	146
180	181	182
216	217	218
252	253	254

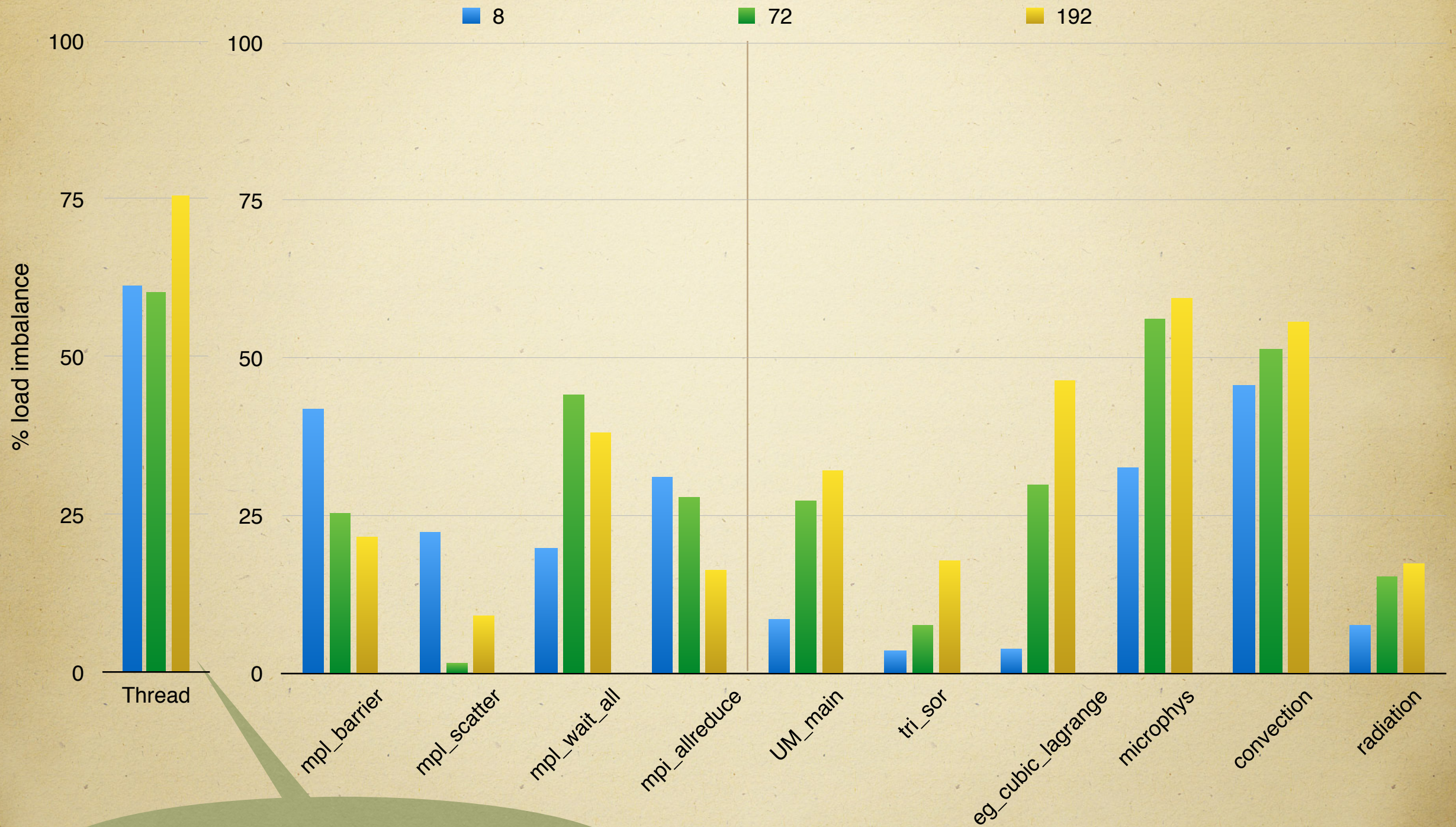
255	256	257
219	220	221
183	184	185
147	148	149

GRID vs SMP MPI rank order



All the above results are obtained using 12 PEs per node and two OPENMP threads. IO is not included in the measurement as the IO resources are shared and accurate measurement is difficult.

Load imbalance



Can expect poor thread scaling

OPENMP directives

- 60% to 75% thread imbalance in UM jobs.
- coverage of OPENMP regions is less than 50%.
- Cray Reveal - integrated performance analysis and code optimisation tool.
 - provides loop analysis and scoping of serial loops.
 - suggests OPENMP directive that can be inserted to a loop.
 - can attach the performance data collected during execution to identify profile of loops.
 - requires knowledge of OPENMP to resolve conflicts and issues.
 - works only with Cray compiling environment.
 - does not support tasks, barrier, critical or atomic regions.

For more details - refer Cray documentation (not much)

CRAY Reveal

The screenshot displays the CRAY Reveal application interface. The main window shows a source code file with line numbers and code snippets. A navigation pane on the left lists various functions and loops. An 'OpenMP Scoping' dialog box is open, showing a table of variable scopes and their status.

Navigation Pane (Function View):

- 20.07% UM_MAIN
- 7.59% RAD_CTL
- 5.58% EG_CUBIC_LAGRANGE
- 2.59% LIDAR_SIMULATOR
- 0.89% NI_CONV_CTL
- 0.84% ATMOS_PHYSICS2
- 0.63% MONO_ENFORCE
- 0.56% EG_VERT_WEIGHTS_ETA
- 0.49% NI_IMP_CTL
- 0.45% BDY_IMPL3
- 0.30% PARCEL_ASCENT
- 0.28% MICROPHYS_CTL
 - Loop@735
 - Loop@737
 - Loop@739
 - Loop@849**
 - Loop@851
 - Loop@853
 - Loop@1263
 - Loop@1303
 - Loop@1305
 - Loop@1318
 - Loop@1320
 - Loop@1330
 - Loop@1418
 - Loop@1420
 - Loop@1422

Source Code (microphys_ctl.f90):

```

847 ! Dry level T values (only required for di
848
849 DO k = qdims%k_end + 1, tdims%k_end
850
851   DO j = qdims%j_start, qdims%j_end
852
853     DO i = qdims%i_start, qdims%i_end
854       t_work(i,j,k) = t_n(i,j,k)
855     END DO
856   END DO
857 END DO
858
859 END DO
860
861 END DO
862
863 l_aggfr_diag = sf(100,4)
864 l_point_diag = sf(101,4)
865 l_vm_cry_diag = sf(102,4)
866 l_vm_agg_diag = sf(103,4)
867 l_vtbranch_diag = sf(104,4)
868 l_vm_used_diag = sf(105,4)
869 l_nifrw_diag = sf(106,4)
    
```

Reveal OpenMP Scoping Dialog (microphys_ctl.f90: Loop@849):

Name	Type	Scope	Info
QDIMS	Scalar	Unresolved	FAIL: No scoping information available
TDIMS	Scalar	Unresolved	FAIL: No scoping information available
i	Scalar	Private	
j	Scalar	Private	
k	Scalar	Private	
t_n	Array	Shared	
t_work	Array	Shared	

Info - Line 849:

- A loop starting at line 849 was not vectorized because a better candidate was found at line 853.
- The loop is flat.
- The loop is flat.

System Path: /home2/n02/n02/karthee/um/xjlek/umatmos/tmp/xjkel.pl loaded

CRAY Reveal

The screenshot displays the CRAY Reveal application interface. The main window shows a source code file named `microphys_ctl.f90` with a tree view on the left. The tree view lists various loops and their percentages, with `Loop@849` selected. The source code editor shows a loop starting at line 849, with variables `ILS`, `IL`, and `AI` highlighted. An `OpenMP Directive` dialog box is open, showing the following text:

```
! Directive inserted by Cray Reveal. May be incomplete.
!$OMP parallel do default(none) &
!$OMP& private (i,j,k) &
!$OMP& shared (t_n,t_work,QDIMS,TDIMS)
```

The `OpenMP Scoping` dialog box is also open, showing the following table:

Name	Type	Scope	Info
i	Scalar	Private	
j	Scalar	Private	
k	Scalar	Private	
QDIMS	Scalar	Shared	FAIL: No scoping information available
TDIMS	Scalar	Shared	FAIL: No scoping information available
t_n	Array	Shared	
t_work	Array	Shared	

The `OpenMP Scoping` dialog box also includes a `Reduction` dropdown menu set to `None`, a `Find Name:` input field, and buttons for `Insert Directive`, `Show Directive`, and `Close`. The `OpenMP Directive` dialog box has buttons for `Copy Directive` and `Close`.

CRAY Reveal

The screenshot displays the CRAY Reveal application interface. On the left is a 'Navigation' pane with a 'Program View' showing a tree of source files. The main window shows the source code for 'balance_lbc_values_4A.f90' with line numbers 447 to 469. A 'Reveal OpenMP Scoping' dialog box is open, showing a table of loops to be scoped. The dialog has tabs for 'Scope Loops' and 'Scoping Results'. The 'Scoping Results' tab is active, displaying a table with columns 'Scope?', 'Line #', and 'File or Source Line'. The table lists several loops with their line numbers and corresponding code snippets. At the bottom of the dialog, there are fields for 'Time' (0.000) and 'Trips' (0), and buttons for 'Apply Filter', 'Start Scoping', 'Cancel', and 'Close'.

Navigation Pane:

- atm_step_local_mod.f90
- atm_step_phys_init.f90
- atm_step_phys_reset.f90
- atm_step_stash.f90
- atm_step_swapbnds.f90
- atm_step_timestep_init.f90
- atmos_constants_mod.f90
- atmos_max_sizes.f90
- atmos_physics1.f90
- atmos_physics2.f90
- atmos_physics2_alloc.f90
- augment_radiance.f90
- augment_tiled_radiance.f90
- backscatter_spectrum.f90
- balance_lbc_values.f90
- balance_lbc_values_4A.f90 (selected)
- EG_BALANCE_LBC_VALUES
 - Loop@279
 - Loop@280
 - Loop@286
 - Loop@287
 - Loop@294
 - Loop@295
 - Loop@302
 - Loop@303
 - Loop@309
 - Loop@310

Source Code:

```
447
448 ! NB. Broken into two separate loops to make vectorisable
449
450 ILS
451
452 ILS
453
454 ILS
455
456 !-----
457 ! Calculat
458 !-----
459
460
461
462
463
464
465
466
467
468
469
```

Reveal OpenMP Scoping Dialog:

Scope?	Line #	File or Source Line
<input checked="" type="checkbox"/>		/home2/n02/n02/karthee/um/xjlek/umatmos/ppsrc/UM/at
<input checked="" type="checkbox"/>	279	D0 k = tdims_s%k_start, tdims_s%k_end
<input checked="" type="checkbox"/>	280	D0 i = 1, lenrim
<input checked="" type="checkbox"/>	286	D0 k = tdims_s%k_start, tdims_s%k_end
<input checked="" type="checkbox"/>	287	D0 i = 1, lenrim
<input checked="" type="checkbox"/>	294	D0 k = tdims_s%k_start, tdims_s%k_end
<input checked="" type="checkbox"/>	295	D0 i = 1, lenrim
<input checked="" type="checkbox"/>	302	D0 k = tdims_s%k_start, tdims_s%k_end
<input checked="" type="checkbox"/>	303	D0 i = 1, lenrim
<input checked="" type="checkbox"/>	309	D0 k = tdims_s%k_start, tdims_s%k_end
<input checked="" type="checkbox"/>	310	D0 i = 1, lenrim
<input checked="" type="checkbox"/>	324	D0 j = pdims%j_start, pdims%j_end

Time: 0.000 Trips: 0

Buttons: Start Scoping, Cancel, Close

home2/n02/n02/karthee/um/xjlek/umatmos/tmp/xjkel.pl loaded

CRAY Reveal

The screenshot displays the Cray Reveal IDE interface. The main window shows the source code for `balance_lbc_values_4A.f90`. The code includes a loop structure with `DO` statements for `k`, `j`, and `i`. A status bar at the bottom indicates "The loop is flat."

An **OpenMP Directive** dialog box is open, showing the directive inserted by Cray Reveal:

```
! Directive inserted by Cray Reveal. May be incomplete.
!$OMP parallel do default(none) &
!$OMP& unresolved (K,KP,ROW_END_PT,LAST_ROW) &
!$OMP& private (i,j,iloc,v_av,u_av,denom1,metric_term,dpdz,temp) &
!$OMP& shared (lenrim,exner_lbc,model_levels,rho_lbc,theta_lbc,w_lbc, &
!$OMP& lenrimu,u_lbc,lenrimv,v_lbc,r_rho_levels, &
!$OMP& r_theta_levels,row_length,halo_i,halo_j,rimwidth, &
!$OMP& f1_comp_eh,f2_comp_eh,g_theta_eh,intw_rho2w, &
!$OMP& intw_w2rho,lbc_address,lbc_address_im1, &
!$OMP& lbc_address_jm1,lbc_address_u,lbc_address_v, &
!$OMP& l_cartesian,ns_weights,theta_wet_lbc,work1)
```

The **Scoping Results** window is also open, showing a table of variables and their scoping information:

Name	Type	Scope	Info
K	Scalar	Unresolved	FAIL: No scoping information
KP	Scalar	Unresolved	FAIL: No scoping information
LAST_ROW	Scalar	Unresolved	FAIL: No scoping information
ROW_END_PT	Scalar	Unresolved	FAIL: No scoping information
denom1	Scalar	Private	
dpdz	Scalar	Private	
i	Scalar	Private	
iloc	Scalar	Private	
j	Scalar	Private	
metric_term	Scalar	Private	

Below the table, there are checkboxes for `Enable FirstPrivate` and `Enable LastPrivate`, and a `Reduction` dropdown menu set to `None`. A `Find Name:` search box and `Insert Directive` / `Show Directive` buttons are also present.

CRAY Reveal

The screenshot displays the CRAY Reveal application window. The title bar reads "X Reveal". The menu bar includes "File", "Edit", "View", and "Help". The window contains two main panes: a "Navigation" pane on the left and a "Source" pane on the right.

Navigation Pane: Shows a "Function View" of the code. The current file is "xjkel.pl". The navigation tree lists various functions and loops with their respective execution percentages and optimization status icons (green for flat, red for non-optimized):

- 20.07% UM_MAIN
- 7.59% RAD_CTL
- 5.58% EG_CUBIC_LAGRANGE
- 2.59% LIDAR_SIMULATOR
 - Loop@271 (Flat)
 - Loop@292 (Flat)
 - Loop@312 (Flat)
 - Loop@331 (Flat)
 - Loop@338 (Flat)
 - Loop@343 (Not Optimized)
 - Loop@353 (Flat)
 - Loop@372 (Flat)
 - Loop@381 (Flat)
- 0.89% NI_CONV_CTL
- 0.84% ATMOS_PHYSICS2
- 0.63% MONO_ENFORCE
- 0.56% EG_VERT_WEIGHTS_ETA (Not Optimized)
- 0.49% NI_IMP_CTL (Not Optimized)
- 0.45% BDY_IMPL3
- 0.30% PARCEL_ASCENT
- 0.28% MICROPHYS_CTL
- 0.25% EXCF_NL_9C
- 0.24% EG_TRI_LINEAR
- 0.23% BDY_IMPL4
- 0.20% PC2_PRESSURE_FORCING
- 0.16% KMKHZ_9C
- 0.16% UM_SHELL

Source Pane: Displays the source code for "lidar_simulator.f90". The code is annotated with optimization status icons in the left margin:

- Line 301: CILV (Flat)
- Line 302: (Flat)
- Line 303: (Flat)
- Line 304: (Flat)
- Line 305: (Flat)
- Line 306: CIVf (Flat)
- Line 307: If (Collapsed)
- Line 308: If (Inlined)
- Line 309: If (Vectorized)
- Line 310: If (Fused)
- Line 311: (Flat)
- Line 312: ILS (Flat)
- Line 313: I (Flat)
- Line 314: IL (Flat)
- Line 315: (Flat)
- Line 316: (Flat)
- Line 317: (Flat)
- Line 318: CILV (Flat)
- Line 319: (Flat)
- Line 320: (Flat)
- Line 321: (Flat)

A tooltip is visible over the "If" statements (lines 307-310), listing optimization options: "Collapsed", "Inlined", "Vectorized", and "Fused". Below the tooltip, it states: "Compiler generated pseudo code for the inlined function can be viewed by clicking the arrow on this line".

Info Pane: Located at the bottom, it provides details for "Line 271":

- A loop starting at line 271 was not vectorized because a recurrence was found on "zheight" at line 272.
- The loop is flat.
- The loop is flat.

CRAY Reveal

- 2389 serial loops are parallelised using OPENMP directives as suggested by CRAY Reveal
- Fortran array notation expressions , for all and where statements are parallelised using

```
!$OMP PARALLEL WORKSHARE
```

```
!$OMP END PARALLEL WORKSHARE
```

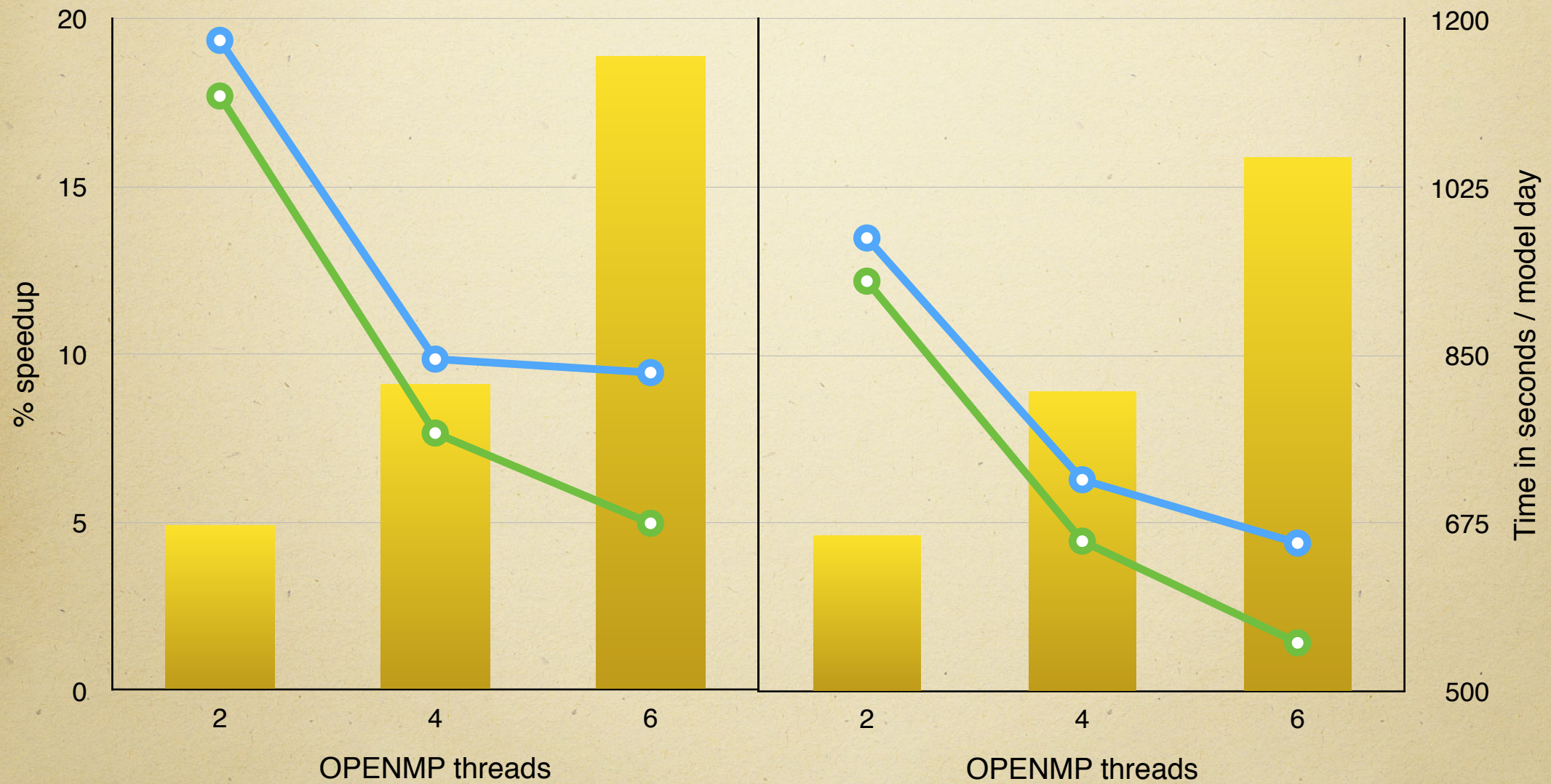
- Bit reproducibility is tested for correctness of OPENMP directives.
- default(none) option is used for all the newly added directives.

Performance speedup

○ UM 8.6 ○ UM 8.6 + Cray Reveal ■ % speedup

N96

N512



Conclusions

- ARCHER (Cray XC30) has good MPI bandwidth and supports hardware threads
- CrayPAT tools used to analyse the performance
- MPI tuning result in 5 to 12% speedup
- UM has poor thread load balance.
- CRAY Reveal tool used to parallelise serial loops
- New OPENMP directives improve performance by 5% up to 20% when using 2 to 6 threads.

THANK YOU

- ARCHER (Cray XC30) has good MPI bandwidth and supports hardware threads
- CrayPAT tools used to analyse the performance
- MPI tuning result in 5 to 12% speedup
- UM has poor thread load balance.
- CRAY Reveal tool used to parallelise serial loops
- New OPENMP directives improve performance by 5% up to 20% when using 2 to 6 threads.