

# 1 An Approach to Enhance pnetCDF Performance in 2 Environmental Modeling Applications

3 D. C. Wong<sup>1</sup>, C. E. Yang<sup>2,‡</sup>, J. S. Fu<sup>2</sup>, K. Wong<sup>2</sup>, Y. Gao<sup>2,\*</sup>

4 [1]{U.S. Environmental Protection Agency, Research Triangle Park, NC, USA}

5 [2]{University of Tennessee, Knoxville, TN, USA}

6 [\*]{now at: Pacific Northwest National Laboratory, Richland, WA, USA}

7 Correspondence to: D. C. Wong ([wong.david-c@epa.gov](mailto:wong.david-c@epa.gov))

8

## 9 Abstract

10 Data intensive simulations are often limited by their I/O performance, and “novel” techniques  
11 need to be developed in order to overcome this limitation. The software package, pnetCDF  
12 which works with parallel file systems, was developed to address this issue by providing  
13 parallel I/O capability. This study examines the performance of an application level data  
14 aggregation approach which performs data aggregation along either row or column dimension  
15 of MPI processes on a spatially decomposed domain, and then applies the pnetCDF parallel  
16 I/O paradigm. The test was done with three different domain sizes which represent small,  
17 moderately large, and large data domains, using a small-scale Community Multi-scale Air  
18 Quality model (CMAQ) mock-up code. The examination includes comparing I/O  
19 performance with traditional serial I/O technique, straight application of pnetCDF, and the  
20 data aggregation along row and column dimension before applying pnetCDF. After the  
21 comparison, “optimal” I/O configurations of this application level data aggregation approach  
22 were quantified. Data aggregation along the row dimension (pnetCDF<sub>cr</sub>) works better than  
23 along the column dimension (pnetCDF<sub>cc</sub>) although it may perform slightly worse than the  
24 straight pnetCDF method with a small number of processors. When the number of processors  
25 becomes larger, pnetCDF<sub>cr</sub> out performs pnetCDF significantly. If the number of processors  
26 keeps increasing, pnetCDF reaches a point where the performance is even worse than the

---

‡ Co-first Author

1 serial I/O technique. This new technique has also been tested for a real application where it  
2 performs two times better than the straight pnetCDF paradigm.

### 3 **1 Introduction**

4 The Community Multiscale Air Quality (CMAQ) model (Byun and Schere, 2006) is a  
5 regional air quality model which is widely used in air quality research and regulatory  
6 applications (e.g. Fu et al., 2012). This model was developed in the 1990s by the U.S.  
7 Environmental Protection Agency (US EPA) and it has continued to evolve. Recently, CMAQ  
8 was combined with WRF to form a WRF-CMAQ two-way coupled model (Wong et al.,  
9 2012) with direct aerosol effects on radiation. CMAQ has been and continues to be  
10 extensively used to provide guidance in rulemaking such as CSAPR (Cross-State Air  
11 Pollution Rule, <http://www.epa.gov/airtransport/CSAPR/>), used by state and local agencies  
12 for air quality management analyses such as SIP (State Implementation Plan), and also used  
13 by academia and industry for studying relevant atmospheric processes and model  
14 applications. CMAQ has also been adapted into the real-time US National Air Quality  
15 Forecasting system (AQF) (Otte et al., 2005) operationally at National Weather Service since  
16 2003 and was recently deployed for forecasting air quality for the 2010 Shanghai World  
17 Expo.

18 CMAQ uses IOAPI (Input/Output Applications Programming Interface  
19 <http://www.cmascenter.org/ioapi>) to handle I/O since the initial model inception. In recent  
20 years, I/O has been shown as one of the bottlenecks in the CMAQ model system. I/O  
21 consumes about 24% - 35% of the entire simulation time. For many applications model run  
22 time is critically important such air quality forecasting which requires meeting operational  
23 time constraints, studies of relationships between climate change and air quality that involve  
24 decadal scale simulations (e.g. Gao et al., 2013), or multiscale model simulations involving  
25 multiple coarse and fine grid resolutions. For such runtime critical applications improving  
26 efficiency of I/O becomes an even more important element that needs to be addressed. To  
27 increase the I/O efficiency, we turn to data aggregation technique (Palmera et al, 2011) which  
28 has been used to achieve high I/O bandwidth as well as a parallel I/O approach which has  
29 been applied in other computer science fields but not for existing environmental models and  
30 their processes. Section 2 provides background information about what has been done  
31 regarding parallel I/O applications. Section 3 describes the current implementation of I/O in  
32 CMAQ through IOAPI. Section 4 depicts our application level data aggregation technique to

1 enhance I/O performance using pnetCDF (Li et al., 2003,  
2 <http://trac.mcs.anl.gov/projects/parallel-netcdf/wiki/Download>) and demonstrates I/O  
3 enhancement through testing with a smaller scale model. This technique was applied to  
4 CMAQ and preliminary results are presented in Section 5 while Section 6 summarizes the  
5 main findings and presents discussion of future work.

## 6 **2 Previous Work in Parallel I/O**

7 The independent I/O and collective I/O are the two most common I/O strategies in parallel  
8 applications. However, a shortcoming of the independent I/O approach is the servicing of the  
9 I/O requests of each process individually (Chen et al., 2010). The collective I/O provides a  
10 better solution of managing non-contiguous portions of a file with multiple processes  
11 interleaved (Thakur et al., 1999). Several collective I/O techniques are hence developed to  
12 improve the parallel I/O performance at various levels by enabling the compute nodes to  
13 cooperate with efficient parallel access to the storage system. Examples include, two-phase  
14 I/O (del Rosario et al., 1993), data sieving (Thakur et al., 1999), and the collective buffering  
15 (Nitzberg and Lo, 1997).

16 To optimize the I/O performance, software is designed to access non-contiguous patterns by  
17 implementation of collective I/O. Data is rearranged and aggregated in memory prior to  
18 writing to files, which reduces the number of disk accesses and the seek-time overhead due to  
19 large amounts of non-contiguous write requests. Improved I/O efficiency is observed through  
20 split writing and hierarchical striping of data (Yu et al., 2007). The benefits of utilizing the  
21 improved parallel I/O techniques on applications in various research areas have been  
22 recognized (Li et al., 2003; Kordenbrock and Oldfield, 2006; Huang et al., 2013). The  
23 approach to parallelize the I/O by using the network Common Data Form (netCDF,  
24 <http://www.unidata.ucar.edu/software/netcdf/>), a set of software libraries and self-describing,  
25 machine-independent data formats, has been discussed regarding the performance of different  
26 I/O libraries (Li et al., 2003), including serial netCDF, parallel-netCDF (pnetCDF) and  
27 Hierarchical Data Format (Cheng, 2000). An auto-tuning framework (behzad and Luu et al.,  
28 2013) has been developed to attempt to provide the best I/O setting with respect to the entire  
29 I/O stack automatically. It used genetic algorithm to search the “optimal” solution from the  
30 parameter-space. Our approach gears toward application level rather than the I/O stack level  
31 but we also dealing with two parameters, stripe count and stripe size in the parallel file system  
32 level.

1 File data striping on parallel file systems also influences I/O performance. Data is distributed  
2 using a fixed block size in a round-robin manner among available I/O servers and disks based  
3 on a simple striping data distribution function. Optimal striping setup on parallel file systems  
4 can significantly reduce the I/O time (Nisar et al., 2012) while inappropriate settings may  
5 incur striping overhead for both metadata and file read/write operations (Yu et al., 2007).  
6 Research work has shown degradation of parallel I/O efficiency when large numbers of  
7 processors are applied to scientific applications such as CMAQ (Kordenbrock and Oldfield,  
8 2006). To overcome these shortcomings, we re-engineered the current CMAQ I/O module to  
9 better utilize more processors on high performance computational machines as well as  
10 quantifying the optimal data striping setup on Lustre file systems.

### 11 **3 I/O in CMAQ**

12 The Community Multiscale Air Quality (CMAQ) modeling system, an active open-source  
13 development project of the U.S. Environmental Protection Agency, is an air quality model for  
14 regulatory and policy analysis. The interactions of atmospheric chemistry and physics are  
15 studied through this three-dimensional Eulerian atmospheric chemistry and transport  
16 modeling system. The primary goal for CMAQ is to simulate ozone, particulate matter, toxic  
17 airborne pollutants, visibility, and acidic and nutrient pollutant species within the troposphere  
18 and across spatial scales ranging from local to hemispheric.

19 IOAPI, a third-party software, was created concurrently with the initial development of the  
20 CMAQ model. It provides a simple interface to handle read and write data in netCDF format  
21 in CMAQ. It originally operated in serial mode and was later expanded to run on SMP  
22 machines using OpenMP. It has never been implemented with capability to run on a  
23 distributed system.

24 When CMAQ was parallelized in late 1998, a “pseudo” parallel I/O library, PARIO, was  
25 created to enable CMAQ to run on a distributed system. PARIO was built on top of the IOAPI  
26 library to handle regular data operations (read and write) from each MPI-process. Each  
27 individual processor can read its sub-domain portion of data straightly from the input files.  
28 However, for output, PARIO requires all processors send its portion of data to the designated  
29 processor, i.e. processor 0, which will stitch all data together and write en masse to the output  
30 file (Figure 1a). Clearly, there are a few shortcomings of this strategy: (1) as the number of  
31 processors increases, the network will be flooded with more MPI messages and require longer  
32 synchronization time to accomplish an output task, (2) if the domain size remains the same

1 but the number of processors increases, the output data size in each processor decreases which  
2 will lower the I/O efficiency, and (3) it requires extra memory for processor 0 to hold the  
3 entire dataset before writing to the file.

4

## 5 **4 An Application Level Data Aggregation Approach**

6 Besides the shortcomings mentioned in Section 3, IOAPI has another major drawback, which  
7 is not taking advantage of existing technology advancements such as parallel file systems and  
8 parallel I/O framework, for example, pnetCDF. Kordenbrock and Oldfield (2006) have  
9 shown an enhancement of model I/O performance with the adoption of pnetCDF. Our new  
10 approach not only utilizes advanced parallel I/O technology, it also addresses all the  
11 shortcomings directly. This new approach performs I/O through pnetCDF using collective  
12 parallel netCDF API on a parallel file system basis, thus eliminating the first and third  
13 shortcomings discussed above.

14 Spatial domain decomposition is widely used in parallelizing scientific models such as  
15 CMAQ. The key characteristic of this new technique is data aggregation which can be  
16 considered as mitigation for the second shortcoming described above. Generally speaking,  
17 data can be either aggregated along the row dimension or column dimension of a rectangular  
18 horizontal grid to enhance the I/O efficiency. During aggregation, a small number of localized  
19 MPI communication processes were introduced which does not diminish the overall benefit of  
20 the technique.

21 In order to determine the performance of this new technique, a small-scale code was devised.  
22 This smaller code, which is designed to mimic the CMAQ model, contains a time step loop  
23 with artificial workload. Data is output at the end of each time step.. This small scaled code  
24 was tested with three time steps and was run on two different machines. The following two  
25 subsections provide brief information about the machines as well as how the test was setup.

### 26 **4.1 High-Performance Computational Systems (HPCs)**

27 The experiments were performed on two HPCs to examine the CMAQ I/O performance with  
28 various methods. (1) Edison: a Cray XC30 system with 236 Tflop/sec, 5,200 compute nodes  
29 with 64 GB memory per node, 333 TB of aggregate memory, Cray Aries high-speed  
30 interconnect and 6.4 PB of scratch storage space. Each node has dual-socket twelve-core Intel

1 Ivy Bridge processors at 2.4GHz. The software packages we used were: cray-mpich/7.0.4,  
2 cray-netcdf/4.3.0, parallel-netcdf/1.3.1, lustre: 2.5.0. (2) Kraken: a Cray XT5 system with the  
3 peak performance of 1.17 Petaflops/sec, 112,896 compute cores, 147 TB of compute memory,  
4 3.3 PB of raw parallel file system disk storage space, and 9,408 compute nodes. Each node  
5 has two 2.6 GHz six-core AMD Opteron processors (Istanbul), 16 GB of memory, and is  
6 connected by Cray SeaStar2+ router. The software packages we used were: Cray MPT 5.3.5,  
7 netcdf 3.6.3, pnetcdf 1.2.0, lustre 2.5.0.

8 The file system on both HPCs is managed by Lustre, a massively parallel-distributed file  
9 system that has the ability to distribute the segments of a single file across multiple object  
10 storage targets (OSTs). A striping technique is applied when a file with a linear sequence of  
11 bytes is separated into small chunks. Through this technique, the bandwidth of accessing the  
12 file and the available disk space for storing the file both increase as read and write operations  
13 can access multiple OSTs concurrently. The default value of stripe count is 1 OST of stripe  
14 count and 1 MB of stripe size on both Kraken and Edison.

## 15 **4.2 Experimental Design**

16 To examine the I/O performance of each module, a small scaled model (pseudo code I/O  
17 module) written in Fortran90 which includes the basic functions, reading data, writing data  
18 and performing arithmetic in between read and write operations, was tested to imitate the  
19 complex CMAQ model with the emphasis on the I/O behavior. The code loops for three times  
20 to represent three time steps as in regular CMAQ simulations. The pseudo code of this small  
21 scaled model (is available on request) looks like this:

```
22     DO I = 1, 3  
23         Read in data  
24         Perform numerical calculation  
25         Output result  
26     END DO
```

27 Three domain sizes were chosen to represent the typical 12-km resolution settings in the  
28 CMAQ community: a small domain that covers the entire State of California and vicinity area  
29 (CA), 89 x 134 x 35 x 146 (column by row by layer by species), a medium-sized domain that  
30 covers the Eastern US (EUS) 279 x 299 x 35 x 146, and a large domain that covers the entire

1 continental US (CONUS), 459 x 299 x 35 x 146 (Figure 2). Various combinations of stripe  
2 counts (2, 5, 11, 20, and 40) and stripe sizes (1MB, 2MB, 4MB, 8MB and 16MB) are tested  
3 on different processor configurations (4x4, 4x8, 8x8, 8x16, and 16x16 on CA and EUS  
4 domain and 4x8, 8x8, 8x16, 16x16, and 16x32 on CONUS domain). Regular domain  
5 decomposition is applied on the spatial domain (column by row) as in the CMAQ model.  
6 Each experiment was carried out multiple times and the averaged values were reported. Four  
7 different I/O techniques were setup: the serial I/O scheme used in the current CMAQ model  
8 which uses regular Network Common Data Form (rnetCDF), I/O implementation with  
9 straight parallel netCDF (pnetCDF) (Fig. 1b), our new technique with data row-wise  
10 aggregation among MPI processes plus I/O through using pnetCDF (pnetCDFcr) (Fig. 1c),  
11 and our new technique with data column-wise aggregation among MPI processes plus I/O  
12 through using pnetCDF (pnetCDFcc) (Fig. 1d). Figure 1 illustrates the essence of these  
13 methods. Timing includes the actual I/O time (disk writing time) plus any additional time  
14 such as data gathering as in the method shown in Figure 1a or additional MPI communication  
15 as needed in data aggregation techniques.

16 The results provided by the small scaled model serve as a basis to determine the optimal  
17 striping information (count and size) for further experiments with the pre-released CMAQ  
18 version 5.0.2. One-day simulations of CMAQ on a 4-km resolution EUS domain (423 x 594 x  
19 14 x146) were run to evaluate the differences among rnetCDF, pnetCDF, and the data  
20 aggregation schemes using pnetCDF with respect to I/O performance. These tests were  
21 conducted on Kraken and Edison with three different processor configurations: 10x12, 12x15,  
22 and 18x20.

## 23 5 Results

### 24 5.1 Small scale model results

25 In Figures 3 - 8 and 10, a relative performance calculation shown in formula (1) is plotted  
26 against the stripe counts and sizes:

$$rel.performance(\%) = \frac{(T_{m1} - T_{m2})}{T_{m1}} 100\% \quad (1)$$

28

29 where  $T_{m1}$  and  $T_{m2}$  denote the averaged maximum I/O time for method 1 and method 2,  
30 respectively. Since, all the runs were done on non-dedicated system environments, the same

1 setup was run multiple times for consistency purposes and outliers were filtered. To avoid  
2 visual blocking when displaying negative values below the xy-plane, absolute values are  
3 plotted and solid bars represent positive values and checkered bars represent negative values.  
4 In each figure, a larger positive solid bar is the desirable outcome.

5 The CA case, which represents a relatively small domain, shows a general trend that  
6 performance degrades as the stripe count increases and/or the stripe size increases. For this  
7 case, pnetCDF performance can be worse than the serial approach using regular netCDF (Fig.  
8 4). With the data aggregation technique, aggregation along the row dimension is better.  
9 Overall, data aggregation along row dimension, pnetCDFcr outperforms pnetCDF. Setting the  
10 stripe count to 5 and stripe size to 1MB seems to be the “optimal” settings on both machines  
11 and among all processor configurations. Furthermore, as the number of processors increase,  
12 the relative performance of pnetCDF drops (from ~50% to ~20% on Kraken (Fig. 3) and from  
13 ~40% to about negative 25% on Edison (Fig. 4). Conversely, with the optimal settings,  
14 relative performance of pnetCDFcr increases as the number of processors increases (increases  
15 from ~20% to 75% except 4x4 case on Kraken (Fig. 3) and increases from ~20% to 80% on  
16 Edison (Fig. 4).

17 The EUS case, which represents a moderately large domain, shows similar result as in the CA  
18 case. Relative performance of aggregation along row dimension is much better than along  
19 column dimension (Fig. 5 - 6). With small number of processors scenarios, 4x4 and 4x8,  
20 pnetCDF performs better than pnetCDFcr. At 8x8, pnetCDFcr performs better than pnetCDF  
21 slightly, ~10%. As the number of processors grows, the enhancement becomes more  
22 significant. Overall, the optimal setting on Kraken is stripe count equals to 11 and stripe size  
23 equals to 2MB and on Edison is stripe count equals to 5 and stripe size equals to 2MB. Again,  
24 with pnetCDF, the relative performance drops as the number of processors increases  
25 (decreases from ~90% to ~75% on Kraken and ~50% to ~40% on Edison). However, the  
26 pnetCDFcr shows the opposite behavior: as the number of processors increases, the relative  
27 performance increases significantly.

28 The CONUS case represents a relatively large domain, showing similar results as the CA and  
29 EUS cases (Fig. 7 - 8). When the number of processors increases, the relative performance of  
30 pnetCDF decreases (from ~80 down to ~60% on Kraken and from ~75% down to ~50% on  
31 Edison). However, the relative performance of the pnetCDFcr scheme increases dramatically.  
32 Overall the “optimal” settings are stripe count equals to 11 and stripe size equals to 2MB.

## 1 **5.2 Stripe Size and Stripe Counts Effect with PnetCDF**

2 Stripe size and stripe count are two of the key factors that affect I/O performance as shown in  
3 Figures 3 - 8. The CONUS domain is chosen with various stripe counts (2, 5, and 11) and  
4 stripe sizes (1MB, 2MB, 4MB, 8MB, 16MB) here to summarize these effects (Fig. 9). Among  
5 all stripe counts, the cases using stripe counts of 11 demonstrate the best performance  
6 compared to other stripe counts; for stripe sizes, the 2MB cases were better than the other  
7 stripe sizes. As more processors were applied, larger stripe sizes resulted in decreasing  
8 performance in writing out data while 2MB cases had relatively better results compared to the  
9 other four sizes. Shorter writing time was found when fewer processors were requested. The  
10 stripe count effect showed that stripe counts of 11 had the best performance compared to the  
11 other two stripe count cases. The differences became more significant when more processors  
12 were used.

## 13 **5.3 The Impact of Large Number of Processors on I/O Efficiency**

14 Section 5.1 has shown pnetCDF performance decreases as the number of processors  
15 increases. When the number of processors continues to increase, the performance of pnetCDF  
16 reaches a point that's worse than the serial I/O scheme (Fig. 10). In contrast, pnetCDFcr  
17 scheme continues to improve significantly as the number of processors increases.

18 The I/O efficiency is defined as the rate of data being output. In parallel applications with a  
19 spatial domain decomposition strategy, the domain size in each processor becomes smaller as  
20 the number of processors increase (Fig. 11 left panel). It is known that the I/O rate is higher if  
21 a large chunk of data is being output. Figure 11 (right panel) reflects this assertion which was  
22 tested on Kraken. When the data is aggregated, no matter whether it is along row or column  
23 dimension, it will increase the data size in the processor which is responsible for the I/O. This  
24 is clearly shown in Figure 11 left panel. With data aggregation (pnetCDFcc or pnetCDFcr),  
25 the data size decreases slower than the pnetCDF approach as the number of processors  
26 increases. This translates into a higher I/O rate in aggregated schemes than the pnetCDF  
27 approach with respect to the same number of processors. pnetCDFcc is worse than pnetCDFcr  
28 due to the internal data alignment in the netCDF internal format (row major).

## 1 **5.4 Application to CMAQ**

2 Based on this small-scale code experiment, the setting of 11 stripe count and 2MB stripe size  
3 is selected to employ in a real CMAQ application: one-day simulation on a 4 km resolution  
4 EUS domain (423x594x14x142). Figure 12 shows the overall writing time recorded on  
5 Kraken and Edison with respect to three different ways to perform I/O: the current way using  
6 regular netCDF (rnetCDF), using straight pnetCDF with 11 stripe count and 2MB stripe size  
7 (pnetCDF), and our new approach (pnetCDFcr) with 11 stripe counts and 2MB stripe size.  
8 Clearly pnetCDFcr shortens the writing time significantly.

9

## 10 **6 Conclusions**

11 We performed a series of experiments with four different I/O modules to examine their I/O  
12 efficiencies in CMAQ. First, a small scaled code was tested on three different domains: CA,  
13 EUS and CONUS which represents small, moderately large, and large data sizes of CMAQ  
14 outputs. The I/O modules include serial mode which is currently used in CMAQ, direct  
15 application of parallel netCDF (pnetCDF), and a new technique based on data aggregation  
16 which can be along row or column dimension (pnetCDFcr and pnetCDFcc) before applying  
17 the parallel netCDF technique. The experiment results show: (1) pnetCDFcr performs better  
18 than pnetCDFcc; (2) pnetCDF performance deteriorates as the number of processors increases  
19 and becomes worse than serial mode when certain large numbers of processors are used; (3)  
20 even though pnetCDFcr does not perform as well as pnetCDF in the small number of  
21 processors scenarios, it does out-perform pnetCDF once the number of processors becomes  
22 larger. In addition, an overall “optimal” setting has been shown based on the experiments: 5  
23 stripe count and 1MB stripe size for small domain, 11 stripe count and 2MB stripe size or 5  
24 stripe count and 2MB stripe size for the moderately large domain, and 11 stripe count and  
25 2MB stripe size for the large domain.

26 This data aggregation I/O module was also tested for a one-day, 4 km by 4 km EUS domain  
27 using CMAQ compared to the serial I/O mode, which is currently implemented in CMAQ,  
28 and direct parallel netCDF. The results show significant reduction of I/O writing time when  
29 this new data aggregated pnetCDF (pnetCDFcr) technique is used compared with serial I/O  
30 approach and with application of straight pnetCDF. With this finding, the overall run time of  
31 scientific applications which require I/O will be significantly reduced. A more important  
32 implication is that it allows users to use a large number of processors to run applications and

1 still maintain a reasonable parallel speedup thereby deferring speedup degradation governed  
2 by the Amdahl's Law. Furthermore, the technique can be transferred to other environmental  
3 models that have large I/O burdens.

4

## 5 **Acknowledgements and Disclaimer**

6 The views expressed here are those of the authors and do not necessarily reflect the views and  
7 policies of the U.S. Environmental Protection Agency (EPA) or any other organization  
8 participating in the AQMEII project. This paper has been subjected to EPA review and  
9 approved for publication. Yang Gao was partly supported by the Office of Science of the U.S.  
10 Department of Energy as part of the Regional and Global Climate Modeling Program. The  
11 Pacific Northwest National Laboratory is operated for DOE by Battelle Memorial Institute  
12 (DE-AC05-76RL01830). The Kraken is a supercomputing facility through National Science  
13 Foundation TeraGrid resources provided by National Institute for Computational Sciences  
14 (NICS) under grant numbers TG-ATM110009 and UT-TENN0006.

15

## 1 **References**

- 2 Behzad, B., Luu, H. V. T., Huchette, J., Byna, S., Prabhat, Aydt, R. A., Koziol, Q., and Snir,  
3 M., “Taming parallel I/O complexity with auto-tuning”, SC 2013, page 68. ACM, 2013.
- 4 Byun, D. W. and Schere, K. L.: Review of the governing equations, computational  
5 algorithms, and other components of the Models-3 Community Multiscale Air Quality  
6 (CMAQ) Modeling System. *Appl. Mech. Rev.*, 59, 51-77, 2006.
- 7 Chen, Y., Sun, X.-H., Thakur, R., Song, H., and Jin, H.: Improving Parallel I/O Performance  
8 with Data Layout Awareness. *Cluster '10: Proceedings of the IEEE International Conference  
9 on Cluster Computing 2010*, Washington, DC, USA: IEEE Computer Society, 2010.
- 10 Cheng, A. and Folk, M.: HDF5: High performance science data solution for the new  
11 millennium. *Proceedings of SC2000: High Performance Networking and Computing*, Dallas,  
12 TX, ACM Press and IEEE Computer Society Press, 2000.
- 13 del Rosario, J., Brodawekar, R., and Choudhary, A.: Improved Parallel I/O via a Two-Phase  
14 Run-time Access Strategy. *Workshop on I/O in Parallel Computer Systems at IPPS '93*, Apr.  
15 1993, pp. 56-70, 1993.
- 16 Fu, J. S., Dong, X., Gao, Y., Wong, D., and Lam Y. F.: Sensitivity and linearity analysis of  
17 ozone in East Asia: The effects of domestic emission and intercontinental transport, *J. Air  
18 Waste Manag. Assoc.*, 62(9), 1102-1114, 2012.
- 19 Gao, Y., Fu, J. S., Drake, J. B., Lamarque, J.-F., and Liu, Y.: The impact of emission and  
20 climate change on ozone in the United States under representative concentration pathways  
21 (RCPs), *Atmos. Chem. Phys.*, 13, 9607-9621, doi:10.5194/acp-13-9607-2013, 2013.
- 22 Huang, X. M., Wang, W. C., Fu, H. H., Yang, G. W., Wang, B., and Zhang, C.: A fast  
23 input/output library for high resolution climate models. *Geosci. Model Dev.*, 7, 93–103, 2014,  
24 doi:10.5194/gmd-7-93-2014, 2014.
- 25 Kordenbrock, T. H. and Oldfield, R. A.: Parallel I/O Advancements in Air Quality Modeling  
26 Systems. Poster on 5th annual CMAS conference, 2006.
- 27 Li, J., Liao, W.-K., Choudhary, A., Ross, R., Thakur, R., Gropp, W., Latham, R., Siegel, A.,  
28 Gallagher, B., and Zingale, M.: Parallel netCDF: A High-Performance Scientific I/O  
29 Interface, *Proceedings of the 2003 ACM/IEEE conference on Supercomputing*, p.39,  
30 November 15-21, 2003. DOI: 10.1145 /1048935.1050189.

1 Nisar, A., Liao, W.-K., and Choudhary, A.: Delegation-Based I/O Mechanism for High  
2 Performance Computing Systems. *IEEE Transactions on Parallel and Distributed Systems*,  
3 vol. 23, no. 2, pp. 271-279, Feb. 2012, doi:10.1109/TPDS.2011.166, 2012.

4 Nitzberg, B. and Lo, V. M.: Collective Buffering: Improving Parallel I/O Performance.  
5 *HPDC*, pp. 148, 1997.

6 Palmera, B., Koontza, A., Schuchardta, K., Heikesb, R., and Randallb, D., "Efficient data IO  
7 for a Parallel Global Cloud Resolving Model", *Environmental Modelling & Software*,  
8 Volume 26, Issue 12, December 2011, Pages 1725–1735.

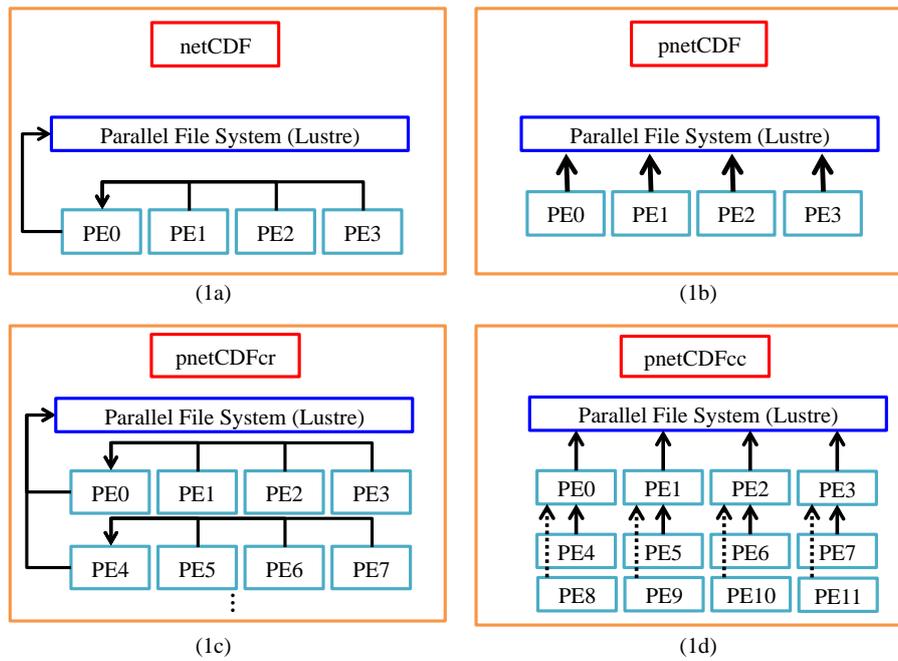
9 Otte, T. L., Pouliot, G., Pleim, J. E., Young, J. O., Schere, K. L., Wong, D. C., Lee, P. C. S.,  
10 Tsidulko, M., McQueen, J. T., Davidson, P., Mathur, R., Chuang, H.-Y., DiMego, G., and  
11 Seaman, N. L.: Linking the Eta Model with the Community Multiscale Air Quality (CMAQ)  
12 modeling system to build a national air quality forecasting system. *Weather Forecast*, 20, 367-  
13 384, 2005.

14 Thakur, R., Gropp, W., and Lusk, E.: Data sieving and collective I/O in ROMIO. *Proceedings*  
15 *of the Seventh Symposium on the Frontiers of Massively Parallel Computation*, IEEE  
16 *Computer Society Press*, 182-189, Feb. 1999.

17 Wong, D. C., Pleim, J., Mathur, R., Binkowski, F., Otte, T., Gilliam, R., Pouliot, G., Xiu, A.,  
18 Young, J. O., and Kang, D.: WRF-CMAQ two-way coupled system with aerosol feedback:  
19 software development and preliminary results. *Geosci. Model Dev.*, 5, 299-312, 2012.

20 Yu, W., Vetter, J., Canon, R. S., and Jiang, S.: Exploiting Lustre File Joining for Effective  
21 Collective IO. *The Seventh IEEE International Symposium on Cluster Computing and the*  
22 *Grid*, 2007.

23



1  
 2 Figure 1. Conceptual diagrams of four I/O modules: serial I/O used in current CMAQ with  
 3 netCDF data format (1a), straight pnetCDF implementation (1b), with data aggregation along  
 4 row dimension and then use pnetCDF (1c), and with data aggregation along column  
 5 dimension and then use pnetCDF (1d). PE denotes a processor. Arrows show the direction of  
 6 data movement.

7

1



2

3 Figure 2: Regional representation of the CA (blue box), EUS (red box) and CONUS (entire)  
4 domains.

5

6

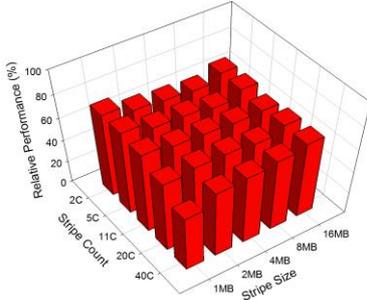
7

8

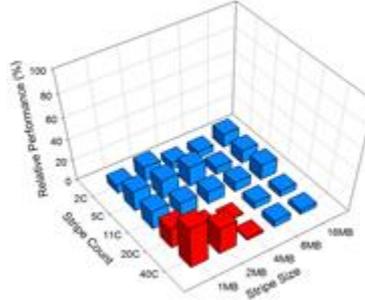
9

10

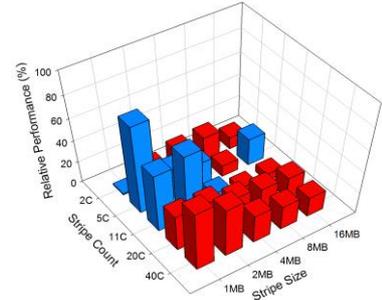
4x4 - CA - metCDF-pnetCDF



4x4 - CA - pnetCDF-pnetCDFcc

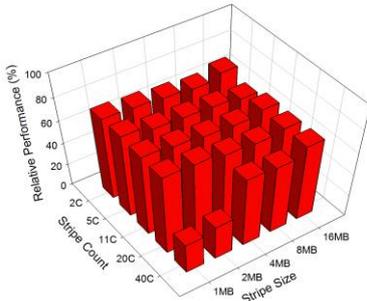


4x4 - CA - pnetCDF-pnetCDFcr

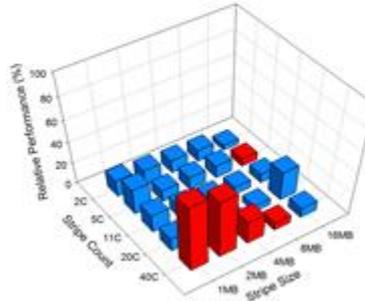


1

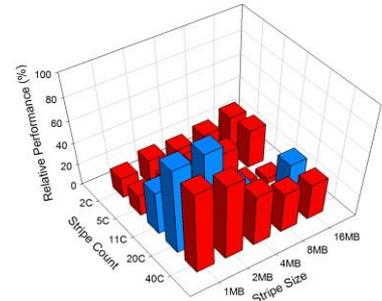
4x8 - CA - metCDF-pnetCDF



4x8 - CA - pnetCDF-pnetCDFcc

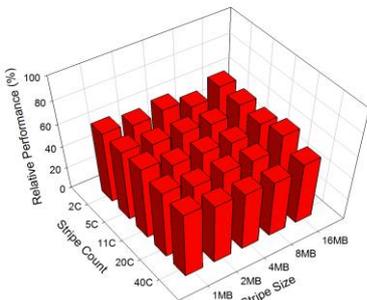


4x8 - CA - pnetCDF-pnetCDFcr

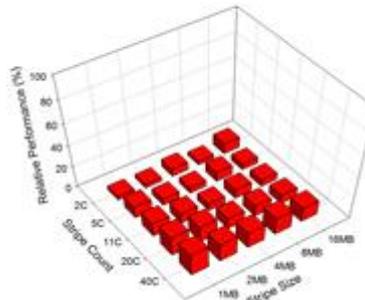


2

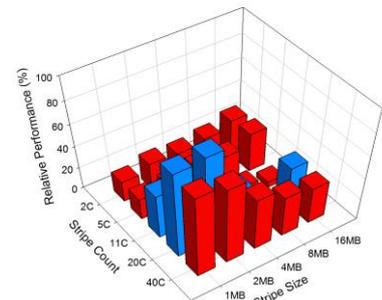
8x8 - CA - metCDF-pnetCDF



8x8 - CA - pnetCDF-pnetCDFcc

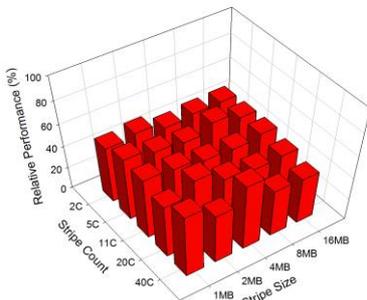


4x8 - CA - pnetCDF-pnetCDFcr

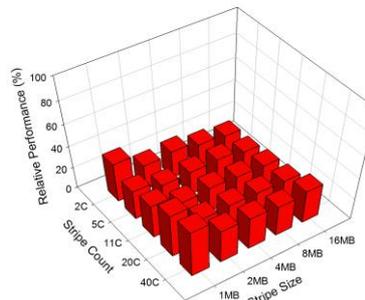


3

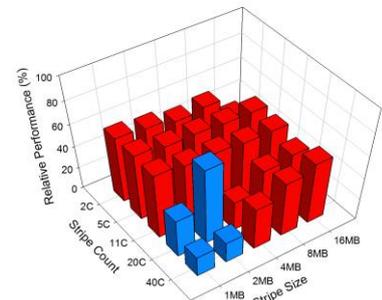
8x16 - CA - metCDF-pnetCDF



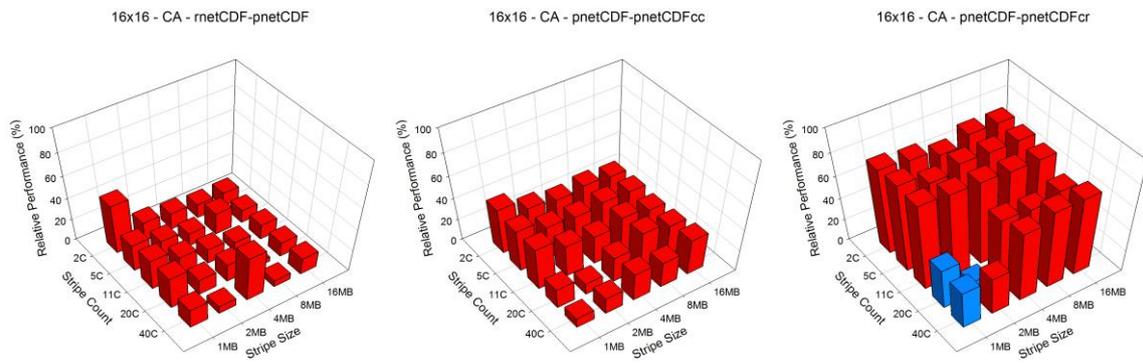
8x16 - CA - pnetCDF-pnetCDFcc



8x16 - CA - pnetCDF-pnetCDFcr



4

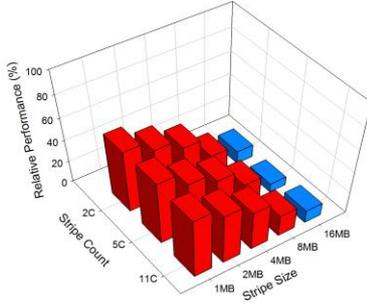


1

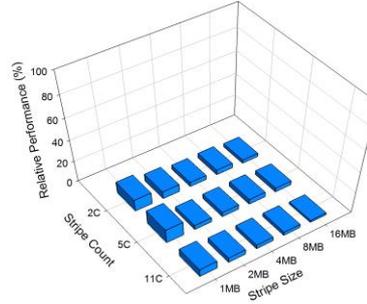
2 Figure 3. Relative I/O performance of pnetCDF to rnetCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CA domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration from  
 4 Kraken. Red colour denotes positive value in relative performance while blue colour denotes  
 5 negative value in relative performance.

6

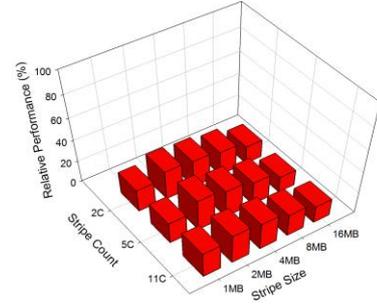
4x4 - CA - metCDF-pnetCDF



4x4 - CA - pnetCDF-pnetCDFcc

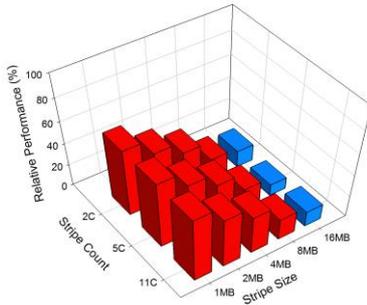


4x4 - CA - pnetCDF-pnetCDFcr

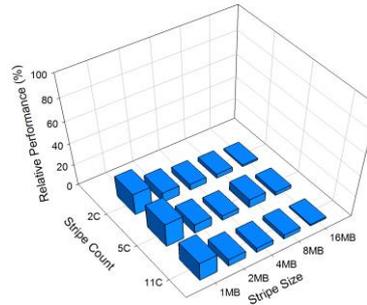


1

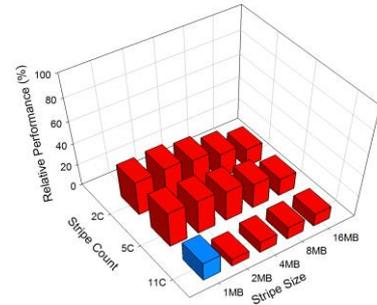
4x8 - CA - metCDF-pnetCDF



4x8 - CA - pnetCDF-pnetCDFcc

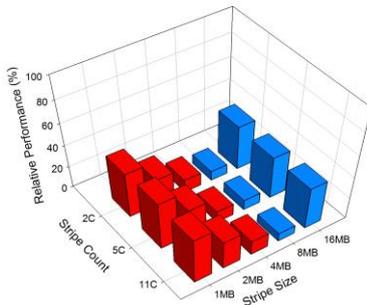


4x8 - CA - pnetCDF-pnetCDFcr

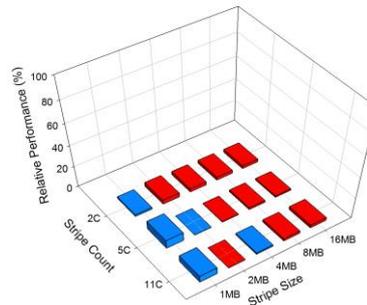


2

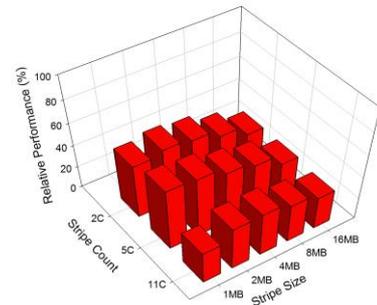
8x8 - CA - metCDF-pnetCDF



8x8 - CA - pnetCDF-pnetCDFcc

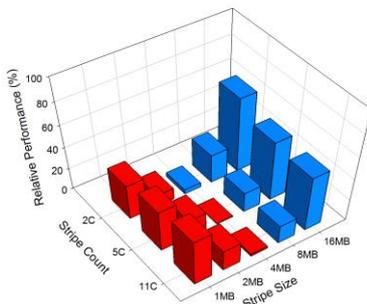


8x8 - CA - pnetCDF-pnetCDFcr

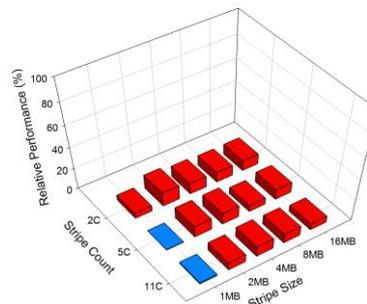


3

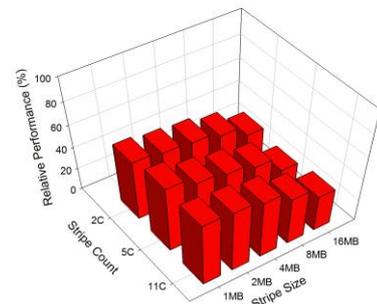
8x16 - CA - metCDF-pnetCDF



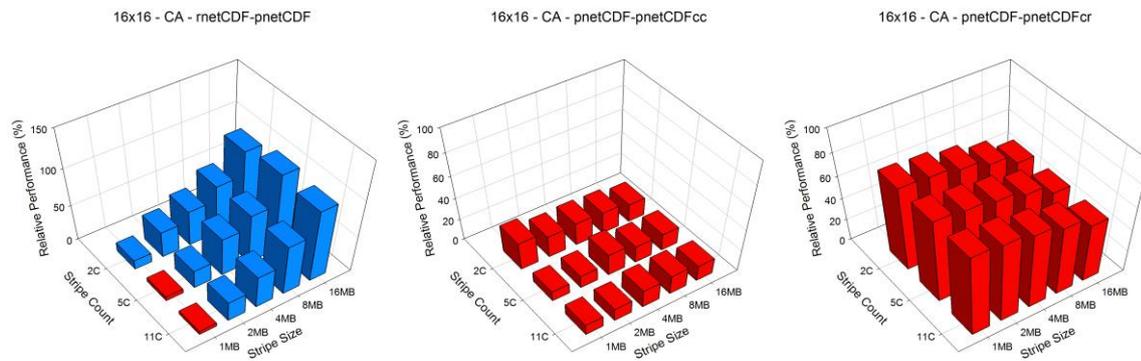
8x16 - CA - pnetCDF-pnetCDFcc



8x16 - CA - pnetCDF-pnetCDFcr



4



1

2 Figure 4. Relative I/O performance of pnetCDF to metCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CA domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration from  
 4 Edison. Red colour denotes positive value in relative performance while blue colour denotes  
 5 negative value in relative performance.

6

7

8

9

10

11

12

13

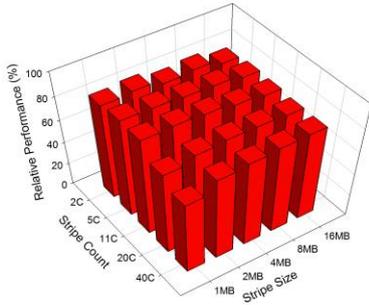
14

15

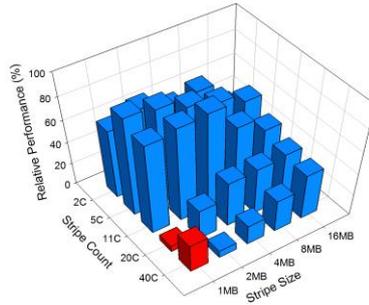
16

17

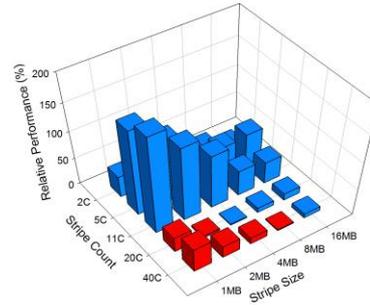
4x4 - EUS - rnetCDF-pnetCDF



4x4 - EUS - pnetCDF-pnetCDFcc

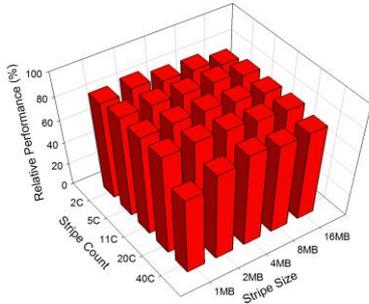


4x4 - EUS - pnetCDF-pnetCDFcr

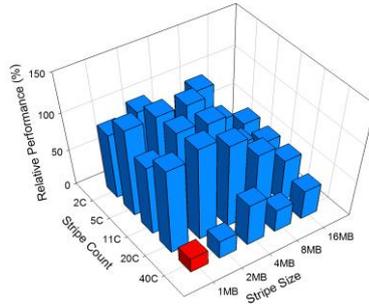


1

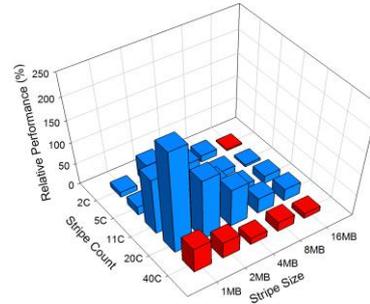
4x8 - EUS - rnetCDF-pnetCDF



4x8 - EUS - pnetCDF-pnetCDFcc

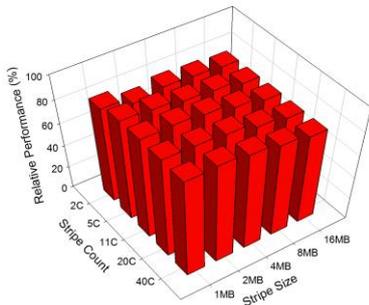


4x8 - EUS - pnetCDF-pnetCDFcr

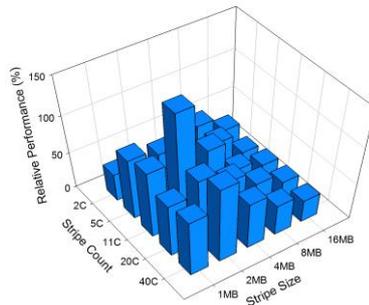


2

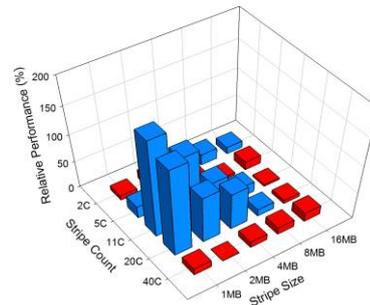
8x8 - EUS - rnetCDF-pnetCDF



8x8 - EUS - pnetCDF-pnetCDFcc

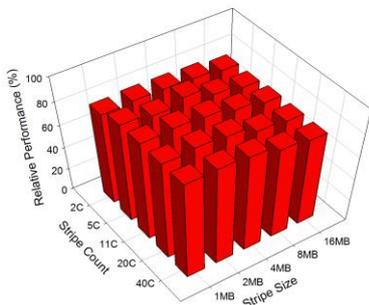


8x8 - EUS - pnetCDF-pnetCDFcr

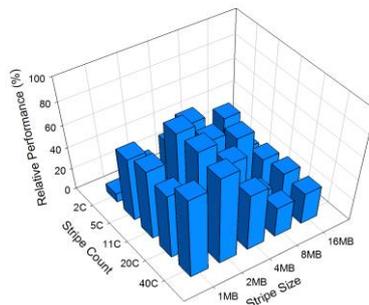


3

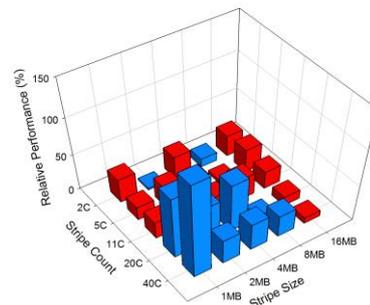
8x16 - EUS - rnetCDF-pnetCDF



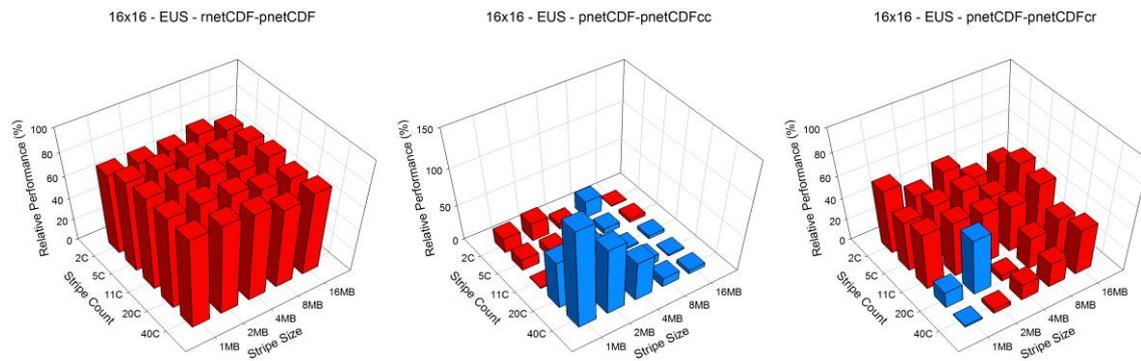
8x16 - EUS - pnetCDF-pnetCDFcc



8x16 - EUS - pnetCDF-pnetCDFcr



4

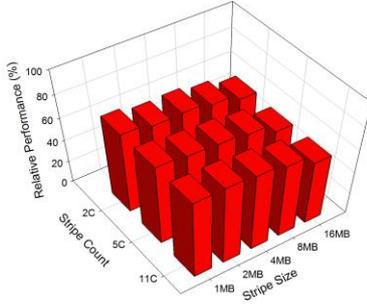


1

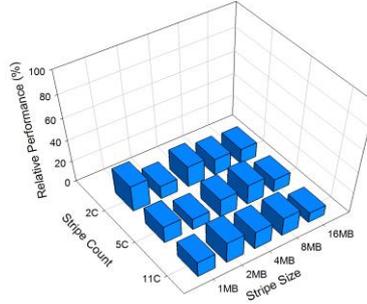
2 Figure 5. Relative I/O performance of pnetCDF to rnetCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on EUS domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration  
 4 from Kraken. Red colour denotes positive value in relative performance while blue colour  
 5 denotes negative value in relative performance.

6

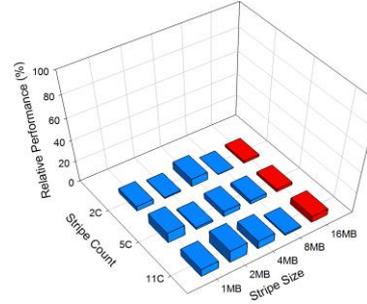
4x4 - EUS - rnetCDF-pnetCDF



4x4 - EUS - pnetCDF-pnetCDFcc

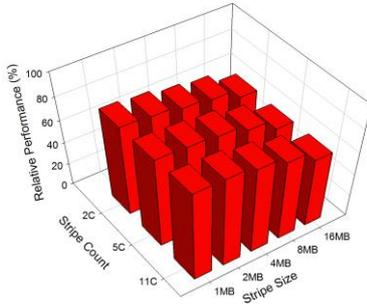


4x4 - EUS - pnetCDF-pnetCDFcr

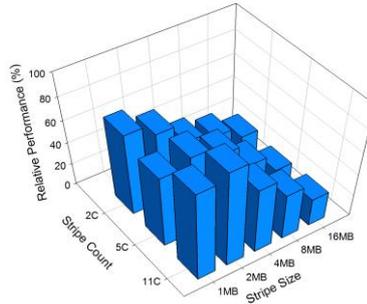


1

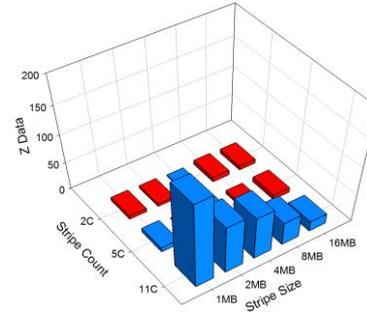
4x8 - EUS - rnetCDF-pnetCDF



4x8 - EUS - pnetCDF-pnetCDFcc

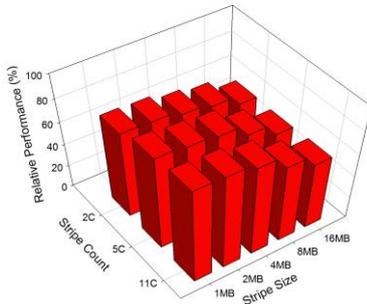


4x8 - EUS - pnetCDF-pnetCDFcr

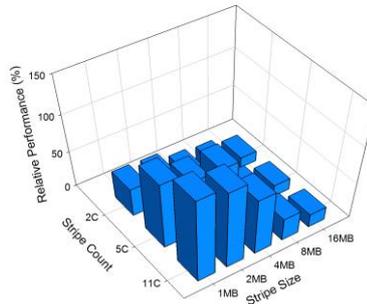


2

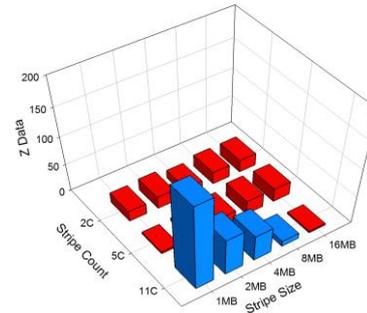
8x8 - EUS - rnetCDF-pnetCDF



8x8 - EUS - pnetCDF-pnetCDFcc

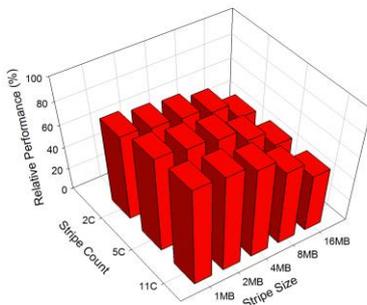


8x8 - EUS - pnetCDF-pnetCDFcr

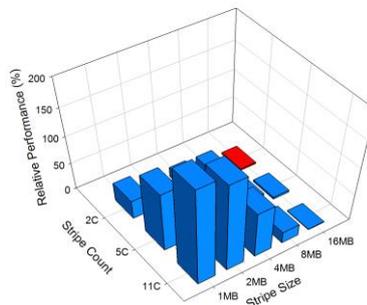


3

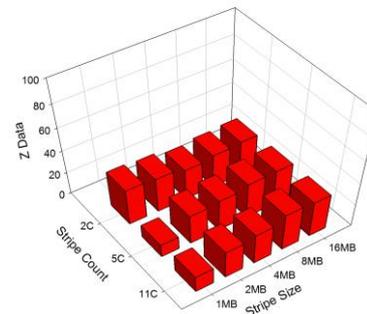
8x16 - EUS - rnetCDF-pnetCDF



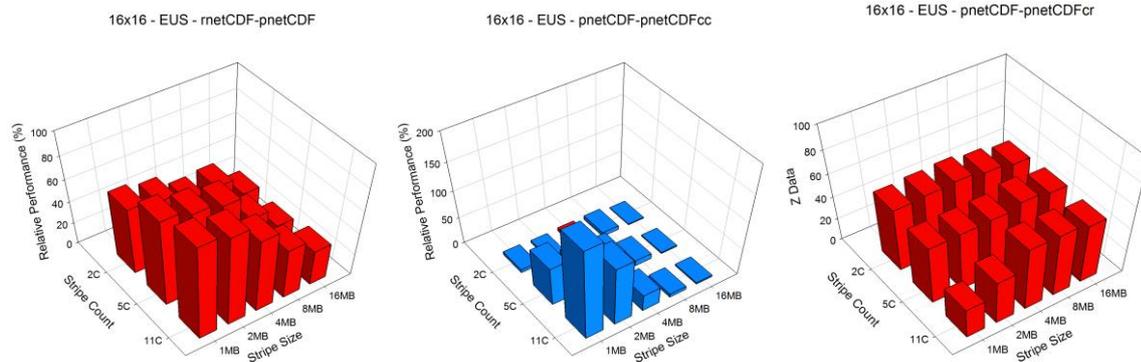
8x16 - EUS - pnetCDF-pnetCDFcc



8x16 - EUS - pnetCDF-pnetCDFcr



4

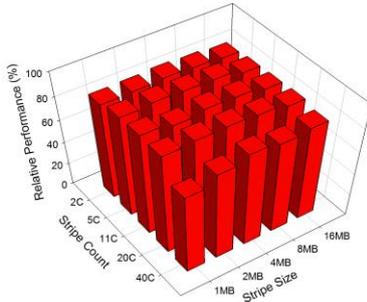


1

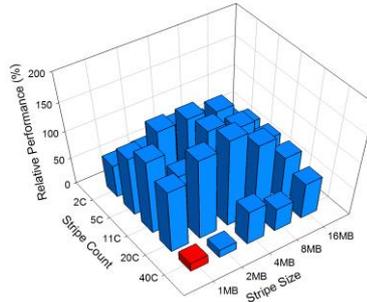
2 Figure 6. Relative I/O performance of pnetCDF to rnetCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on EUS domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration  
 4 from Edison. Red colour denotes positive value in relative performance while blue colour  
 5 denotes negative value in relative performance.

6

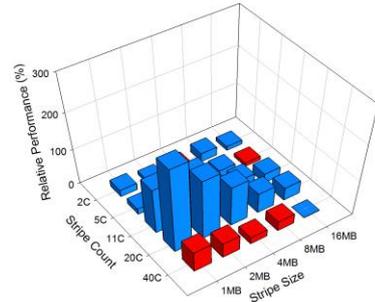
4x8 - CONUS - pnetCDF-pnetCDF



4x8 - CONUS - pnetCDF-pnetCDFcc

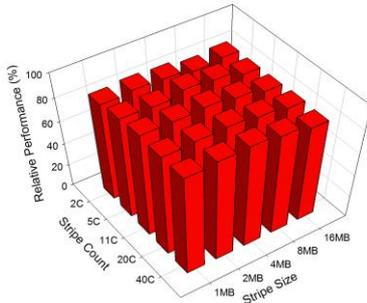


4x8 - CONUS - pnetCDF-pnetCDFcr

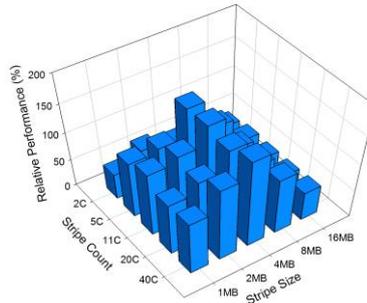


1

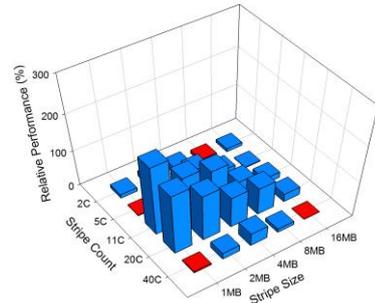
8x8 - CONUS - pnetCDF-pnetCDF



8x8 - CONUS - pnetCDF-pnetCDFcc

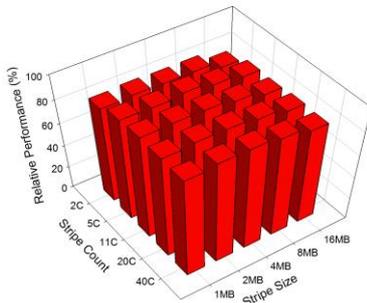


8x8 - CONUS - pnetCDF-pnetCDFcr

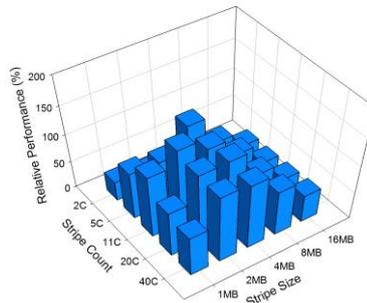


2

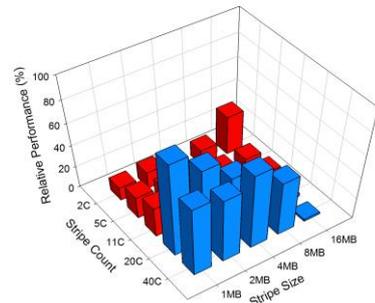
8x16 - CONUS - pnetCDF-pnetCDF



8x16 - CONUS - pnetCDF-pnetCDFcc

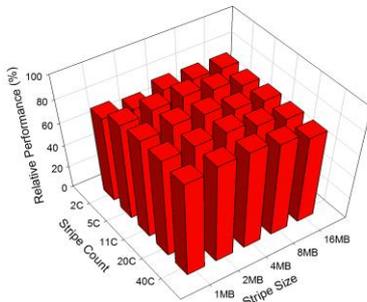


8x16 - CONUS - pnetCDF-pnetCDFcr

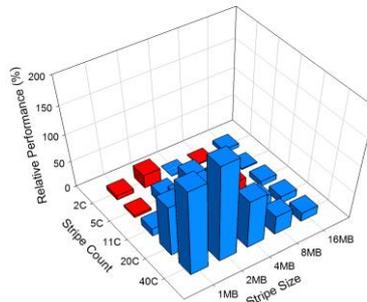


3

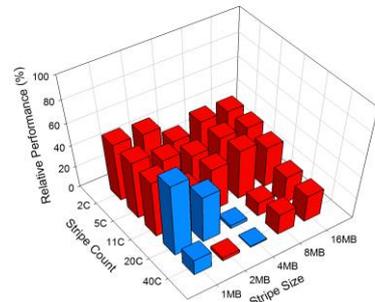
16x16 - CONUS - pnetCDF-pnetCDF



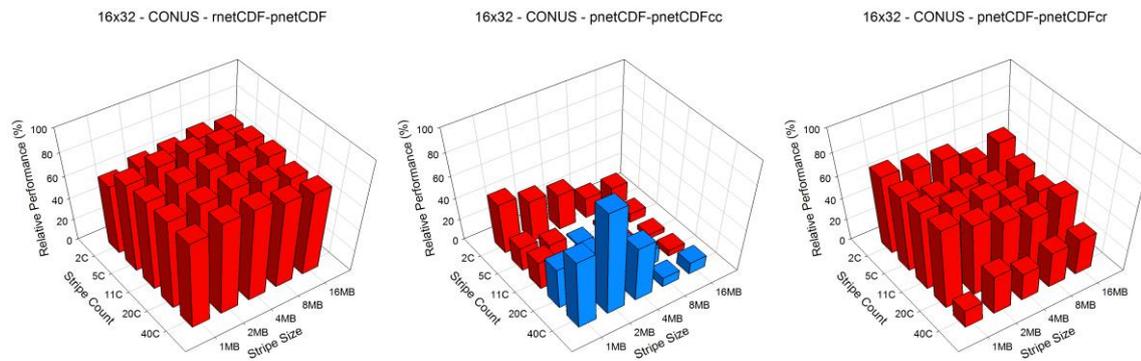
16x16 - CONUS - pnetCDF-pnetCDFcc



16x16 - CONUS - pnetCDF-pnetCDFcr



4

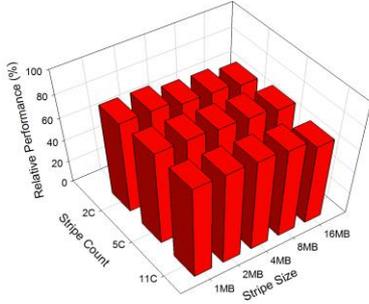


1

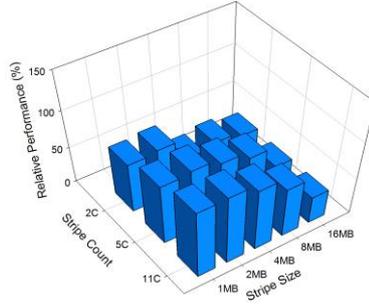
2 Figure 7. Relative I/O performance of pnetCDF to metCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CONUS domain with 4x8, 8x8, 8x16, 16x16, and 16x32 processor  
 4 configuration from Kraken. Red colour denotes positive value in relative performance while  
 5 blue colour denotes negative value in relative performance.

6

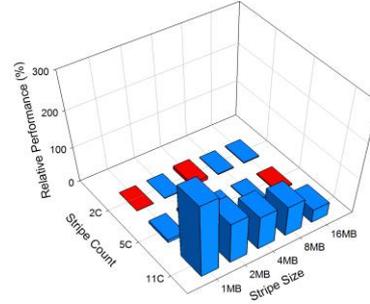
4x8 - CONUS - rnetCDF-pnetCDF



4x8 - CONUS - pnetCDF-pnetCDFcc

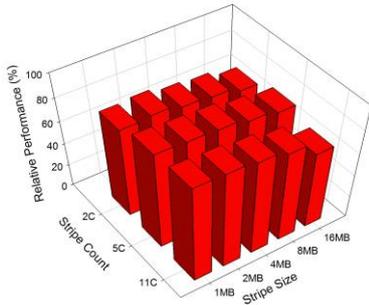


4x8 - CONUS - pnetCDF-pnetCDFcr

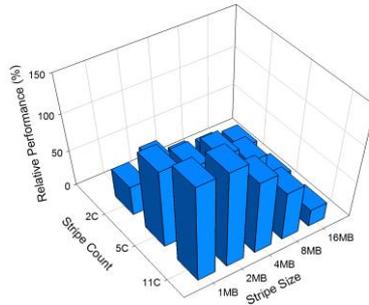


1

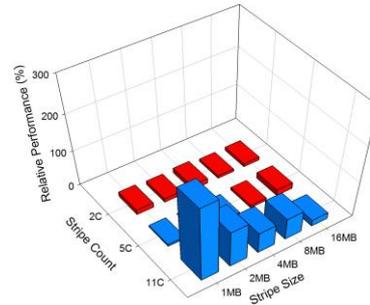
8x8 - CONUS - rnetCDF-pnetCDF



8x8 - CONUS - pnetCDF-pnetCDFcc

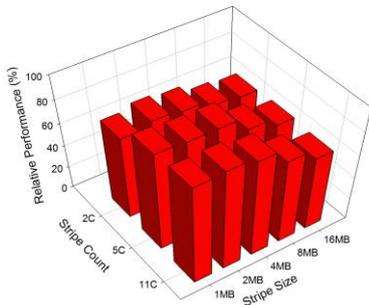


8x8 - CONUS - pnetCDF-pnetCDFcr

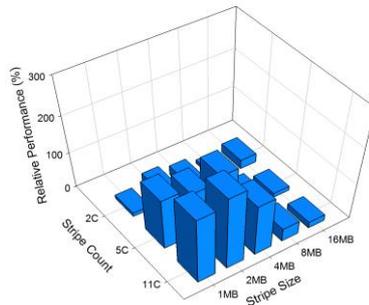


2

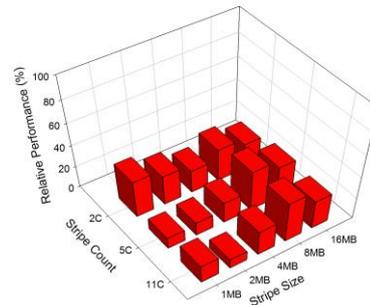
8x16 - CONUS - rnetCDF-pnetCDF



8x16 - CONUS - pnetCDF-pnetCDFcc

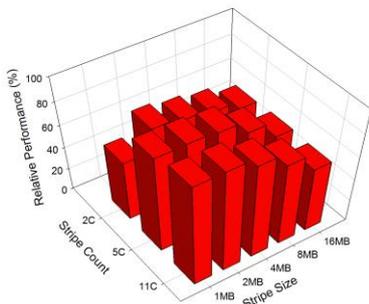


8x16 - CONUS - pnetCDF-pnetCDFcr

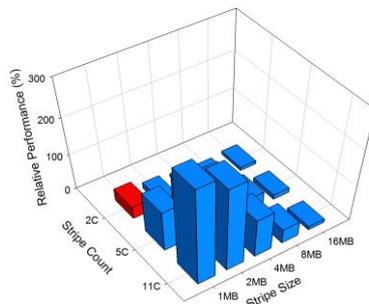


3

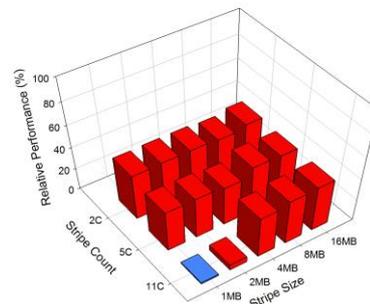
16x16 - CONUS - rnetCDF-pnetCDF



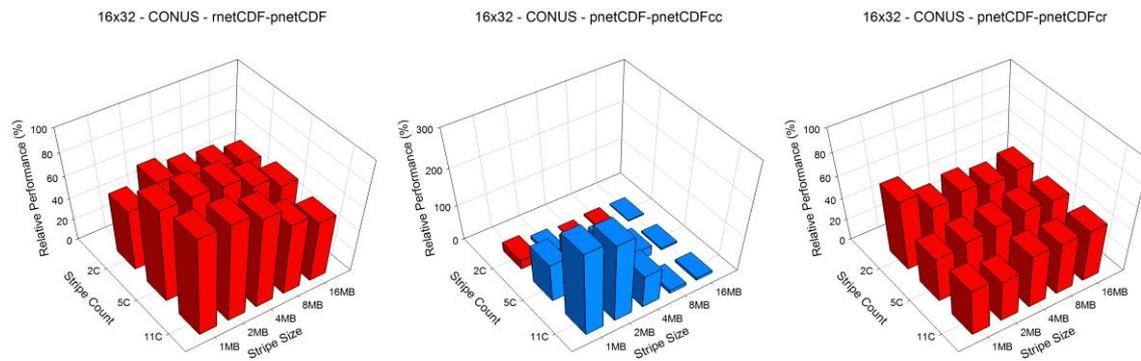
16x16 - CONUS - pnetCDF-pnetCDFcc



16x16 - CONUS - pnetCDF-pnetCDFcr



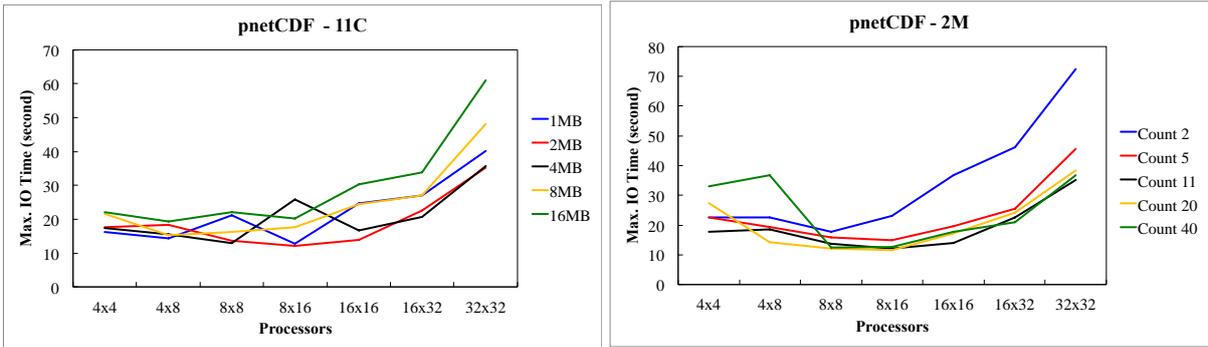
4



1

2 Figure 8. Relative I/O performance of pnetCDF to metCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CONUS domain with 4x8, 8x8, 8x16, 16x16, and 16x32 processor  
 4 configuration from Edison. Red colour denotes positive value in relative performance while  
 5 blue colour denotes negative value in relative performance.

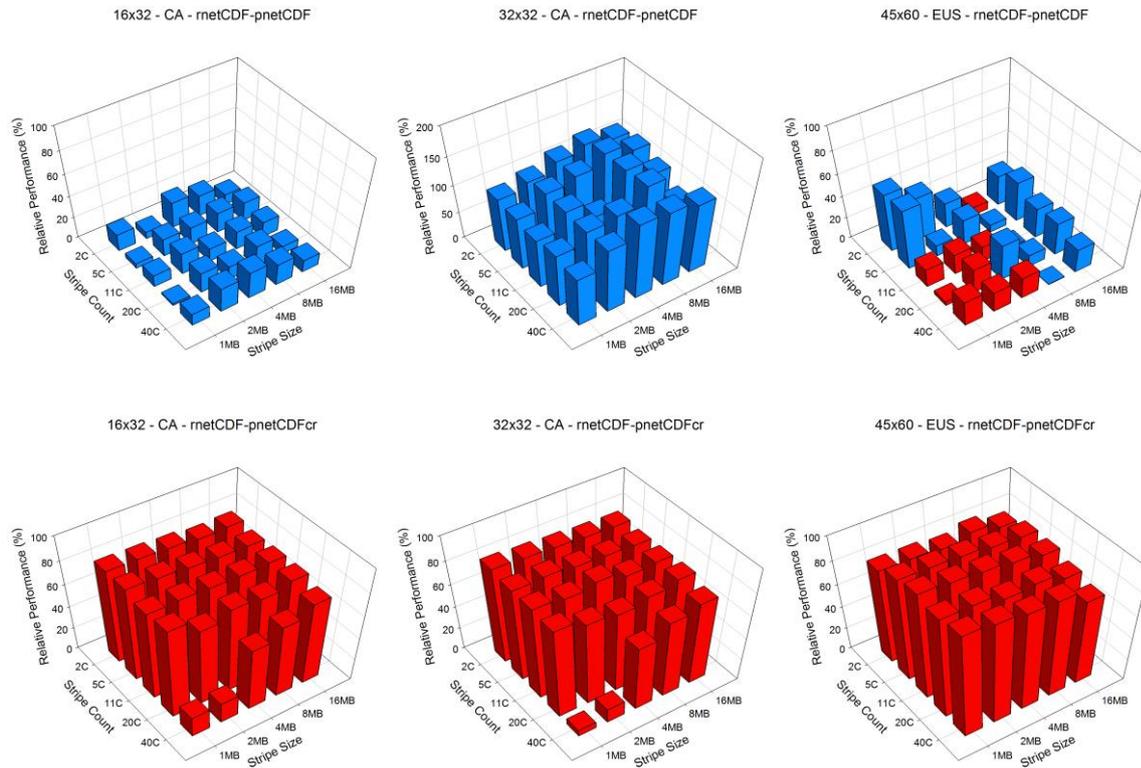
6



1

2 Figure 9. The impact of stripe count and size on parallel netCDF I/O performance on CONUS  
 3 domain. Left: various stripe sizes with fixed 11 stripe counts. Right: various stripe counts with  
 4 fixed 2MB stripe size.

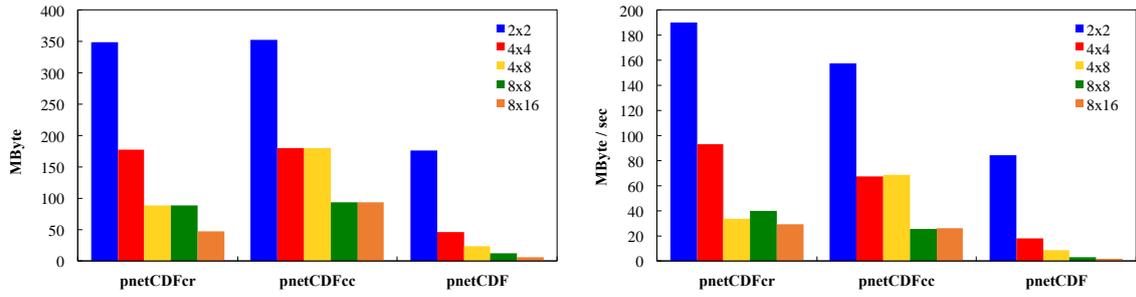
5



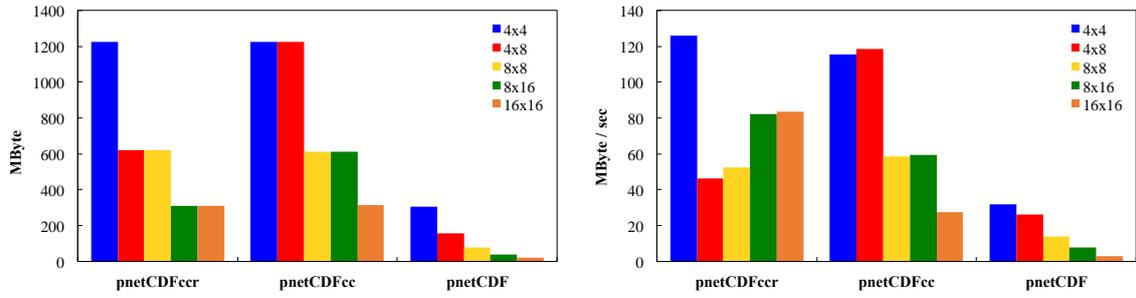
3 Figure 10. Relative performance of pnetCDF and pnetCDFcr with respect to rnetCDF in a  
 4 large number of processors scenario on Kraken. Red colour denotes positive value in relative  
 5 performance while blue colour denotes negative value in relative performance.

6

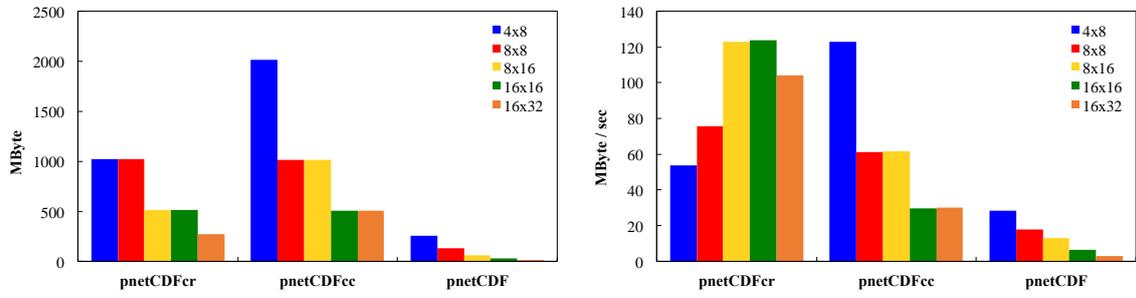
1



2

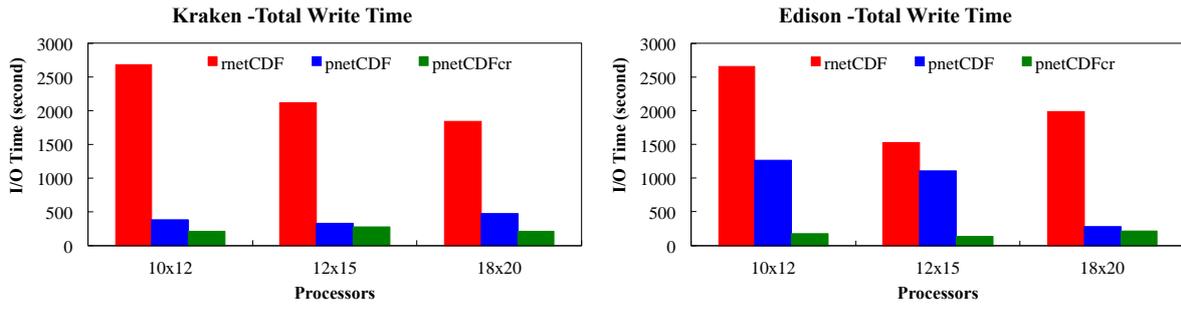


3



4 Figure 11: The maximum data size (left panel) and I/O rate (right panel) among all I/O  
 5 processors in CA (top), EUS (middle) and CONUS domain (bottom), respectively, in the  
 6 pseudo code experiment running on Edison.

7



1  
 2 Figure 12. Total write time in one-day CMAQ simulation by different I/O approaches on a  
 3 4km EUS domain with stripe size 2MB and stripe count 11 (Kraken left and Edison right).