

# 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 a parallel I/O approach which has been applied in other  
28 computer science fields but not for existing environmental models and their processes.  
29 Section 2 provides background information about what has been done regarding parallel I/O  
30 applications. Section 3 describes the current implementation of I/O in CMAQ through IOAPI.  
31 Section 4 depicts our application level data aggregation technique to enhance I/O performance  
32 using pnetCDF and demonstrates I/O enhancement through testing with a smaller scale

1 model. This technique was applied to CMAQ and preliminary results are presented in Section  
2 5 while Section 6 summarizes the main findings and presents discussion of future work.

## 3 **2 Previous Work in Parallel I/O**

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

13 To optimize the I/O performance, software is designed to access non-contiguous patterns by  
14 implementation of collective I/O. Data is rearranged and aggregated in memory prior to  
15 writing to files, which reduces the number of disk accesses and the seek-time overhead due to  
16 large amounts of non-contiguous write requests. Improved I/O efficiency is observed through  
17 split writing and hierarchical striping of data (Yu et al., 2007). The benefits of utilizing the  
18 improved parallel I/O techniques on applications in various research areas have been  
19 recognized (Li et al., 2003; Kordenbrock and Oldfield, 2006; Huang et al., 2013). The  
20 approach to parallelize the I/O by using the network Common Data Form (netCDF), a set of  
21 software libraries and self-describing, machine-independent data formats, has been discussed  
22 regarding the performance of different I/O libraries (Li et al., 2003), including serial netCDF,  
23 parallel-netCDF (pnetCDF) and Hierarchical Data Format (Cheng, 2000).

24 File data striping on parallel file systems also influences I/O performance. Data is distributed  
25 using a fixed block size in a round-robin manner among available I/O servers and disks based  
26 on a simple striping data distribution function. Optimal striping setup on parallel file systems  
27 can significantly reduce the I/O time (Nisar et al., 2012) while inappropriate settings may  
28 incur striping overhead for both metadata and file read/write operations (Yu et al., 2007).  
29 Research work has shown degradation of parallel I/O efficiency when large numbers of  
30 processors are applied to scientific applications such as CMAQ (Kordenbrock and Oldfield,  
31 2006). To overcome these short comings, we re-engineered the current CMAQ I/O module to

1 better utilize more processors on high performance computational machines as well as  
2 quantifying the optimal data striping setup on Lustre file systems.

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

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

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

16 When CMAQ was parallelized in late 1998, a “pseudo” parallel I/O library, PARIO, was  
17 created to enable CMAQ to run on a distributed system. PARIO was built on top of the IOAPI  
18 library to handle regular data operations (read and write) from each MPI-process. Each  
19 individual processor can read its sub-domain portion of data straightly from the input files.  
20 However, for output, PARIO requires all processors send its portion of data to the designated  
21 processor, i.e. processor 0, which will stitch all data together and write en masse to the output  
22 file (Figure 1a). Clearly, there are a few shortcomings of this strategy: (1) as the number of  
23 processors increases, the network will be flooded with more MPI messages and require longer  
24 synchronization time to accomplish an output task, (2) if the domain size remains the same  
25 but the number of processors increases, the output data size in each processor decreases which  
26 will lower the I/O efficiency, and (3) it requires extra memory for processor 0 to hold the  
27 entire dataset before writing to the file.

28

### 29 **4 An Application Level Data Aggregation Approach**

30 Besides the shortcomings mentioned in Section 3, IOAPI has another major drawback, which  
31 is not taking advantage of existing technology advancements such as parallel file systems and

1 parallel I/O framework, for example, pnetCDF. Kordenbrock and Oldfield (2006) have  
2 shown an enhancement of model I/O performance with the adoption of pnetCDF. Our new  
3 approach not only utilizes advanced parallel I/O technology, it also addresses all the  
4 shortcomings directly. This new approach performs I/O through pnetCDF on a parallel file  
5 system basis, thus eliminating the first and third shortcomings discussed above.

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

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

#### 18 **4.1 High-Performance Computational Systems (HPCs)**

19 The experiments were performed on two HPCs to examine the CMAQ I/O performance with  
20 various methods. (1) Edison: a Cray XC30 system with 236 Tflop/sec, 5,200 compute nodes  
21 with 64 GB memory per node, 333 TB of aggregate memory, Cray Aries high-speed  
22 interconnect and 6.4 PB of scratch storage space. Each node has dual-socket twelve-core Intel  
23 Ivy Bridge processors at 2.4GHz. The software packages we used were: cray-mpich/7.0.4,  
24 cray-netcdf/4.3.0, parallel-netcdf/1.3.1, lustre: 2.5.0. (2) Kraken: a Cray XT5 system with the  
25 peak performance of 1.17 Petaflops/sec, 112,896 compute cores, 147 TB of compute memory,  
26 3.3 PB of raw parallel file system disk storage space, and 9,408 compute nodes. Each node  
27 has two 2.6 GHz six-core AMD Opteron processors (Istanbul), 16 GB of memory, and is  
28 connected by Cray SeaStar2+ router. The software packages we used were: Cray MPT 5.3.5,  
29 netcdf 3.6.3, pnetcdf 1.2.0, lustre 2.5.0.

30 The file system on both HPCs is managed by Lustre, a massively parallel-distributed file  
31 system that has the ability to distribute the segments of a single file across multiple object

1 storage targets (OSTs). A striping technique is applied when a file with a linear sequence of  
2 bytes is separated into small chunks. Through this technique, the bandwidth of accessing the  
3 file and the available disk space for storing the file both increase as read and write operations  
4 can access multiple OSTs concurrently. The default value of stripe count is 1 OST of stripe  
5 count and 1 MB of stripe size on both Kraken and Edison.

## 6 **4.2 Experimental Design**

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

```
13     DO I = 1, 3  
14         Read in data  
15         Perform numerical calculation  
16         Output result  
17     END DO
```

18 Three domain sizes were chosen to represent the typical 12-km resolution settings in the  
19 CMAQ community: a small domain that covers the entire State of California and vicinity area  
20 (CA), 89 x 134 x 35 x 146 (column by row by layer by species), a medium-sized domain that  
21 covers the Eastern US (EUS) 279 x 299 x 35 x 146, and a large domain that covers the entire  
22 continental US (CONUS), 459 x 299 x 35 x 146 (Figure 2). Various combinations of stripe  
23 counts (2, 5, 11, 20, and 40) and stripe sizes (1MB, 2MB, 4MB, 8MB and 16MB) are tested  
24 on different processor configurations (4x4, 4x8, 8x8, 8x16, and 16x16 on CA and EUS  
25 domain and 4x8, 8x8, 8x16, 16x16, and 16x32 on CONUS domain). Regular domain  
26 decomposition is applied on the spatial domain (column by row) as in the CMAQ model.  
27 Each experiment was carried out multiple times and the averaged values were reported. Four  
28 different I/O techniques were setup: the serial I/O scheme used in the current CMAQ model  
29 which uses regular Network Common Data Form (rnetCDF), I/O implementation with  
30 straight parallel netCDF (pnetCDF) (Fig. 1b), our new technique with data row-wise

1 aggregation among MPI processes plus I/O through using pnetCDF (pnetCDFcr) (Fig. 1c),  
2 and our new technique with data column-wise aggregation among MPI processes plus I/O  
3 through using pnetCDF (pnetCDFcc) (Fig. 1d). Figure 1 illustrates the essence of these  
4 methods. Timing includes the actual I/O time (disk writing time) plus any additional time  
5 such as data gathering as in the method shown in Figure 1a or additional MPI communication  
6 as needed in data aggregation techniques.

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

## 14 **5 Results**

### 15 **5.1 Small scale model results**

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

$$18 \quad rel.performance(\%) = \frac{(T_{m1} - T_{m2})}{T_{m1}} 100\% \quad (1)$$

19 where  $T_{m1}$  and  $T_{m2}$  denote the averaged maximum I/O time for method 1 and method 2,  
20 respectively. Since, all the runs were done on non-dedicated system environments, the same  
21 setup was run multiple times for consistency purposes and outliers were filtered. To avoid  
22 visual blocking when displaying negative values below the xy-plane, absolute values are  
23 plotted and solid bars represent positive values and checkered bars represent negative values.  
24 In each figure, a larger positive solid bar is the desirable outcome.  
25

26 The CA case, which represents a relatively small domain, shows a general trend that  
27 performance degrades as the stripe count increases and/or the stripe size increases. For this  
28 case, pnetCDF performance can be worse than the serial approach using regular netCDF (Fig.  
29 4). With the data aggregation technique, aggregation along the row dimension is better.  
30 Overall, data aggregation along row dimension, pnetCDFcr outperforms pnetCDF. Setting the

1 stripe count to 5 and stripe size to 1MB seems to be the “optimal” settings on both machines  
2 and among all processor configurations. Furthermore, as the number of processors increase,  
3 the relative performance of pnetCDF drops (from ~50% to ~20% on Kraken (Fig. 3) and from  
4 ~40% to about negative 25% on Edison (Fig. 4). Conversely, with the optimal settings,  
5 relative performance of pnetCDFcr increases as the number of processors increases (increases  
6 from ~20% to 75% except 4x4 case on Kraken (Fig. 3) and increases from ~20% to 80% on  
7 Edison (Fig. 4).

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

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

## 24 **5.2 Stripe Size and Stripe Counts Effect with PnetCDF**

25 Stripe size and stripe count are two of the key factors that affect I/O performance as shown in  
26 Figures 3 - 8. The CONUS domain is chosen with various stripe counts (2, 5, and 11) and  
27 stripe sizes (1MB, 2MB, 4MB, 8MB, 16MB) here to summarize these effects (Fig. 9). Among  
28 all stripe counts, the cases using stripe counts of 11 demonstrate the best performance  
29 compared to other stripe counts; for stripe sizes, the 2MB cases were better than the other  
30 stripe sizes. As more processors were applied, larger stripe sizes resulted in decreasing  
31 performance in writing out data while 2MB cases had relatively better results compared to the



1 other four sizes. Shorter writing time was found when fewer processors were requested. The  
2 stripe count effect showed that stripe counts of 11 had the best performance compared to the  
3 other two stripe count cases. The differences became more significant when more processors  
4 were used.

### 5 **5.3 The Impact of Large Number of Processors on I/O Efficiency**

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

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

### 21 **5.4 Application to CMAQ**

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

29

## 1 **6 Conclusions**

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

17 This data aggregation I/O module was also tested for a one-day, 4 km by 4 km EUS domain  
18 using CMAQ compared to the serial I/O mode, which is currently implemented in CMAQ,  
19 and direct parallel netCDF. The results show significant reduction of I/O writing time when  
20 this new data aggregated pnetCDF (pnetCDFcr) technique is used compared with serial I/O  
21 approach and with application of straight pnetCDF. With this finding, the overall run time of  
22 scientific applications which require I/O will be significantly reduced. A more important  
23 implication is that it allows users to use a large number of processors to run applications and  
24 still maintain a reasonable parallel speedup thereby deferring speedup degradation governed  
25 by the Amdahl’s Law. Furthermore, the technique can be transferred to other environmental  
26 models that have large I/O burdens.

27

## 28 **Acknowledgements and Disclaimer**

29 The views expressed here are those of the authors and do not necessarily reflect the views and  
30 policies of the U.S. Environmental Protection Agency (EPA) or any other organization  
31 participating in the AQMEII project. This paper has been subjected to EPA review and  
32 approved for publication. Yang Gao was partly supported by the Office of Science of the U.S.

1 Department of Energy as part of the Regional and Global Climate Modeling Program. The  
2 Pacific Northwest National Laboratory is operated for DOE by Battelle Memorial Institute  
3 (DE-AC05-76RL01830). The Kraken is a supercomputing facility through National Science  
4 Foundation TeraGrid resources provided by National Institute for Computational Sciences  
5 (NICS) under grant numbers TG-ATM110009 and UT-TENN0006.

6

## 1 **References**

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

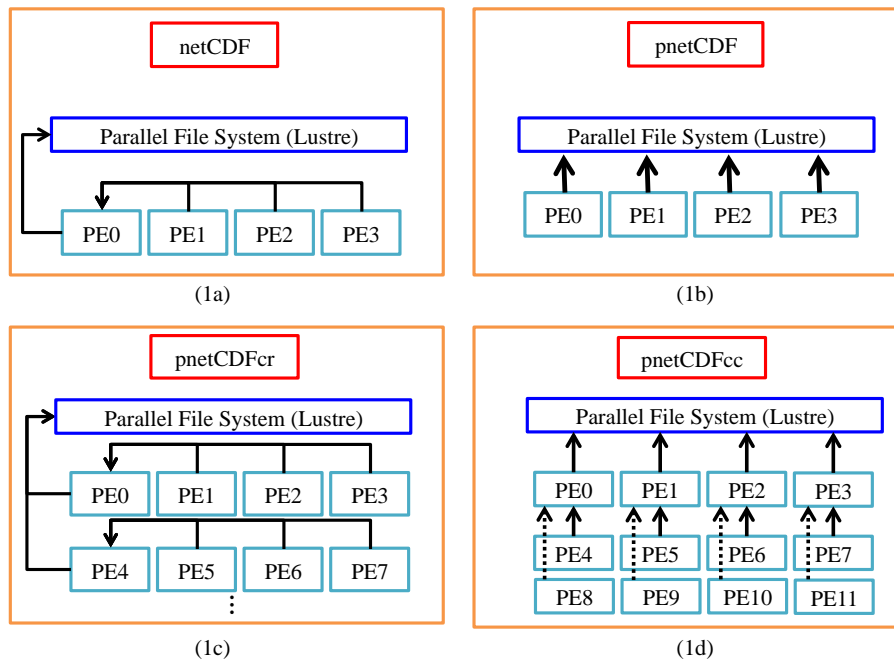
6 Otte, T. L., Pouliot, G., Pleim, J. E., Young, J. O., Schere, K. L., Wong, D. C., Lee, P. C. S.,  
7 Tsidulko, M., McQueen, J. T., Davidson, P., Mathur, R., Chuang, H.-Y., DiMego, G., and  
8 Seaman, N. L.: Linking the Eta Model with the Community Multiscale Air Quality (CMAQ)  
9 modeling system to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

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

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

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

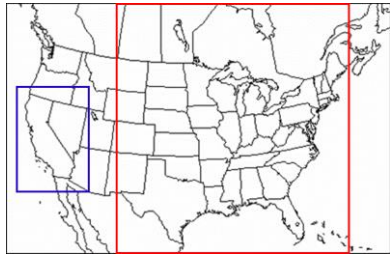
20



1  
2  
3  
4  
5  
6  
7

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

1



2

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

5

6

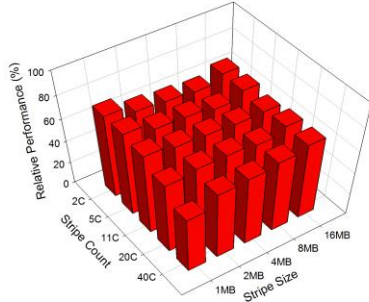
7

8

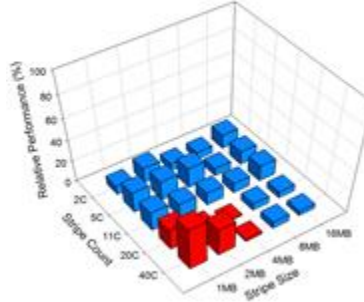
9

10

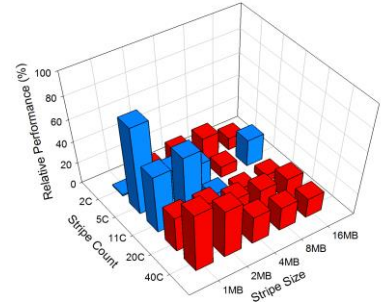
4x4 - CA - metCDF-pnetCDF



4x4 - CA - pnetCDF-pnetCDFcc

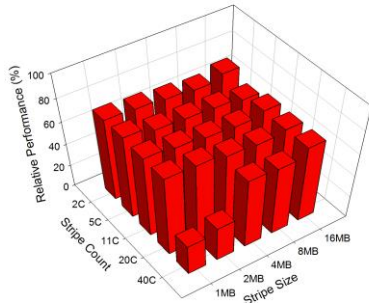


4x4 - CA - pnetCDF-pnetCDFcr

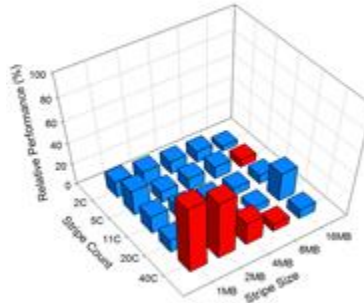


1

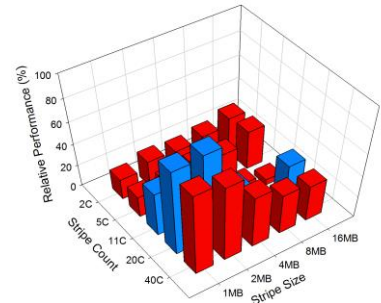
4x8 - CA - metCDF-pnetCDF



4x8 - CA - pnetCDF-pnetCDFcc

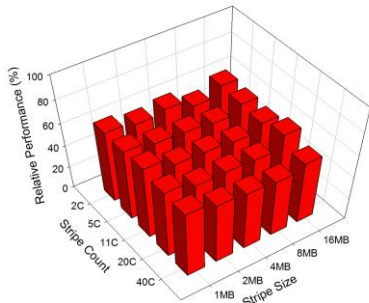


4x8 - CA - pnetCDF-pnetCDFcr

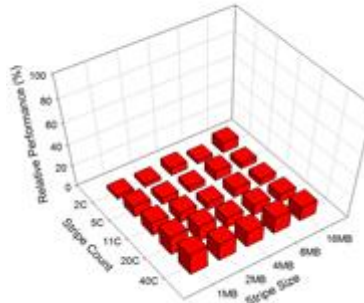


2

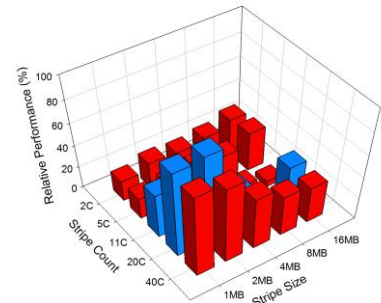
8x8 - CA - metCDF-pnetCDF



8x8 - CA - pnetCDF-pnetCDFcc

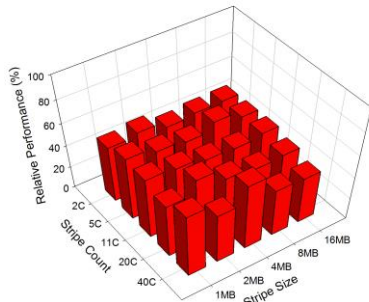


4x8 - CA - pnetCDF-pnetCDFcr

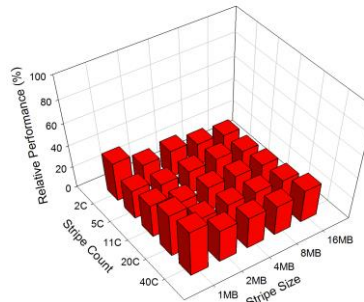


3

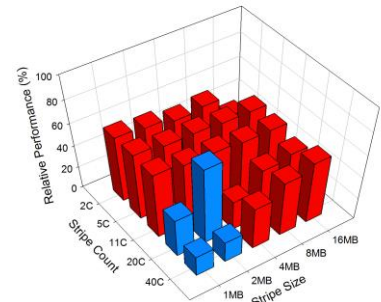
8x16 - CA - metCDF-pnetCDF



8x16 - CA - pnetCDF-pnetCDFcc

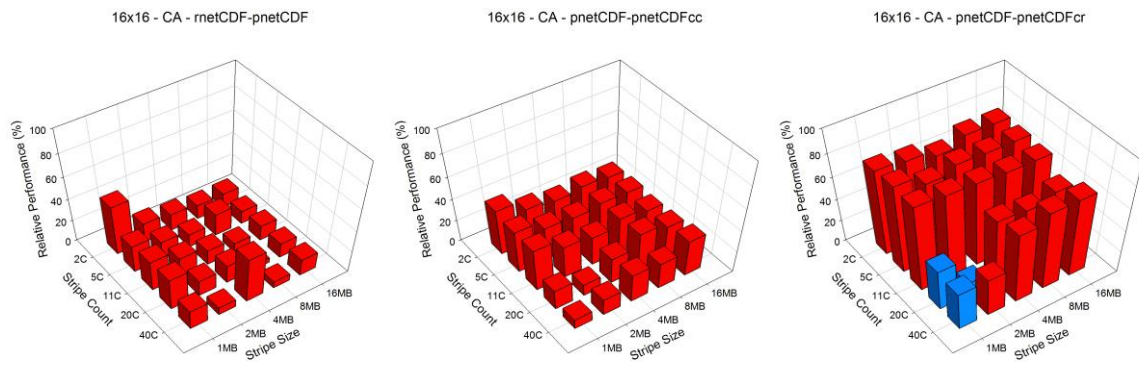


8x16 - CA - pnetCDF-pnetCDFcr



4



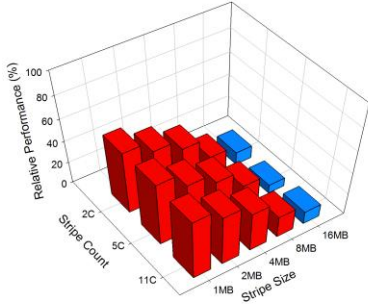


1

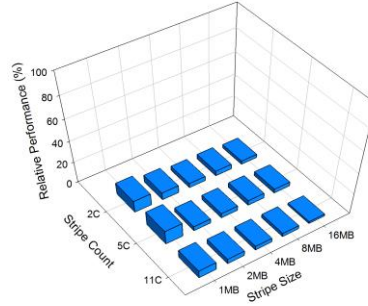
2 Figure 3. Relative I/O performance of pnetCDF to rnetCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CA domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration from  
 4 Kraken. Red colour denotes positive value in relative performance while blue colour denotes  
 5 negative value in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

6

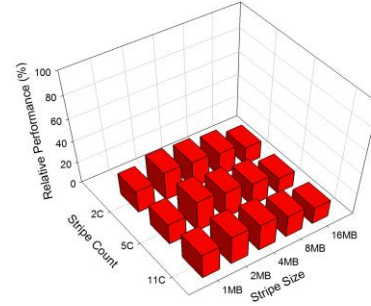
4x4 - CA - metCDF-pnetCDF



4x4 - CA - pnetCDF-pnetCDFcc

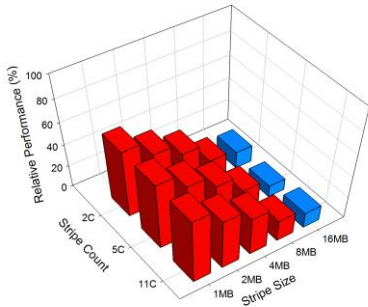


4x4 - CA - pnetCDF-pnetCDFcr

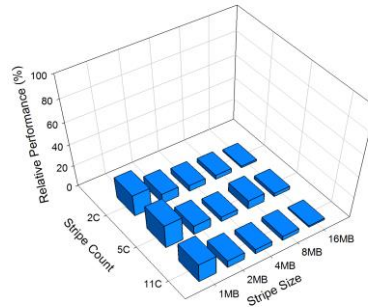


1

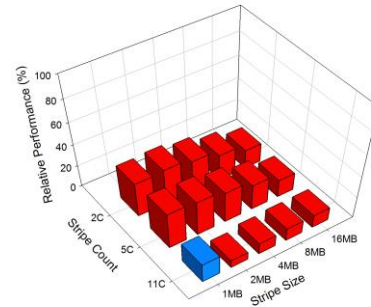
4x8 - CA - metCDF-pnetCDF



4x8 - CA - pnetCDF-pnetCDFcc

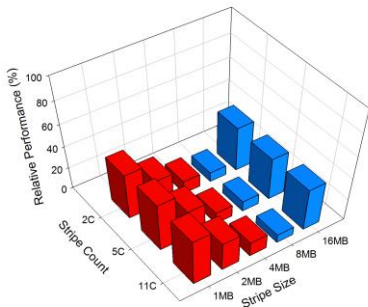


4x8 - CA - pnetCDF-pnetCDFcr

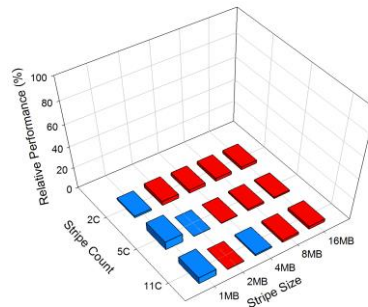


2

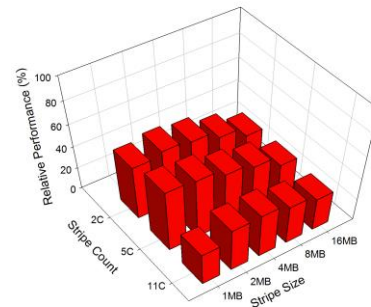
8x8 - CA - metCDF-pnetCDF



8x8 - CA - pnetCDF-pnetCDFcc

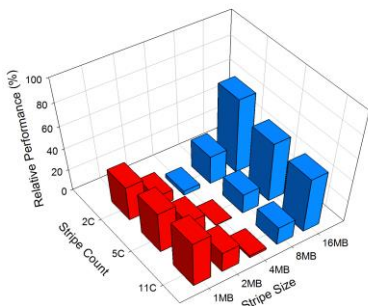


8x8 - CA - pnetCDF-pnetCDFcr

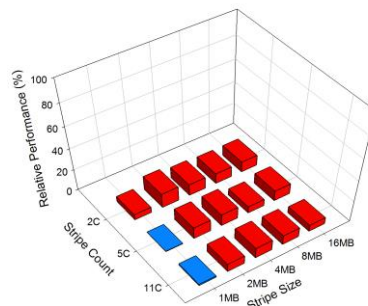


3

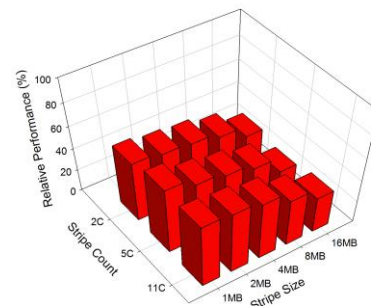
8x16 - CA - metCDF-pnetCDF



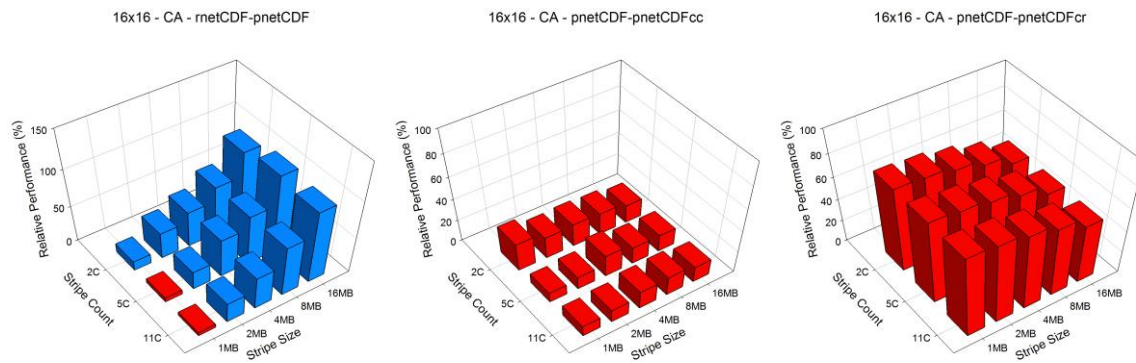
8x16 - CA - pnetCDF-pnetCDFcc



8x16 - CA - pnetCDF-pnetCDFcr



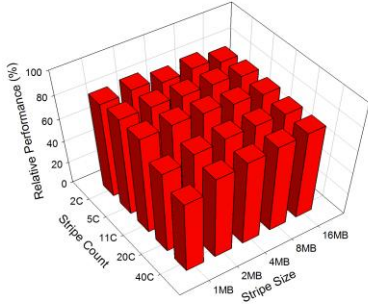
4



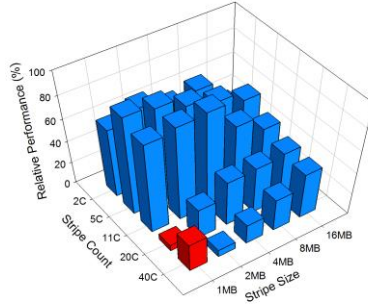
1  
 2 Figure 4. Relative I/O performance of pnetCDF to metCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CA domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration from  
 4 Edison. Red colour denotes positive value in relative performance while blue colour denotes  
 5 negative value in relative performance.

6  
 7  
 8  
 9  
 10  
 11  
 12  
 13  
 14  
 15  
 16  
 17

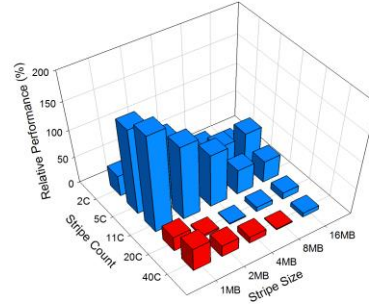
4x4 - EUS - rnetCDF-pnetCDF



4x4 - EUS - pnetCDF-pnetCDFcc

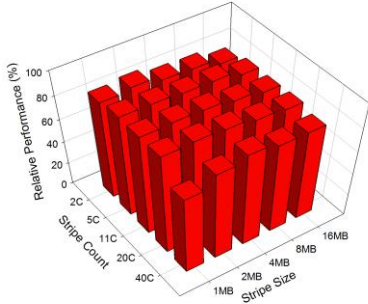


4x4 - EUS - pnetCDF-pnetCDFcr

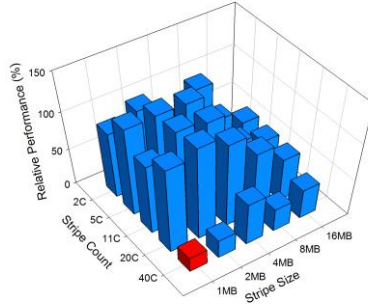


1

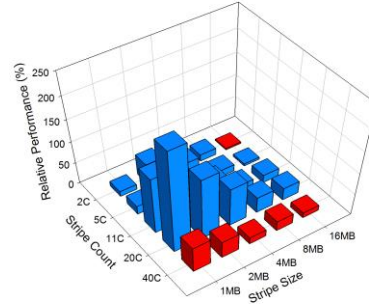
4x8 - EUS - rnetCDF-pnetCDF



4x8 - EUS - pnetCDF-pnetCDFcc

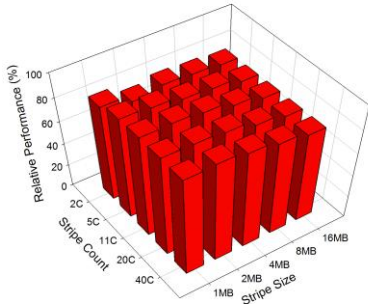


4x8 - EUS - pnetCDF-pnetCDFcr

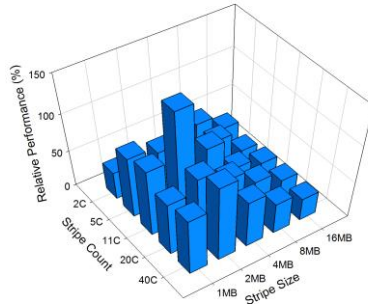


2

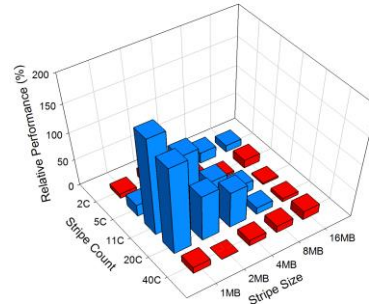
8x8 - EUS - rnetCDF-pnetCDF



8x8 - EUS - pnetCDF-pnetCDFcc

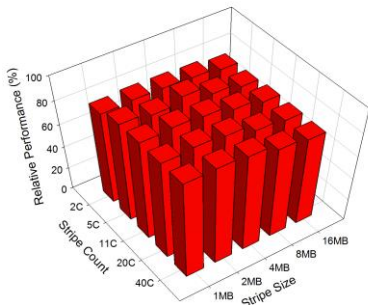


8x8 - EUS - pnetCDF-pnetCDFcr

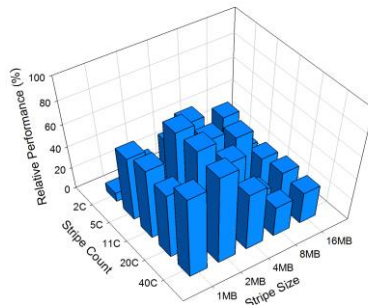


3

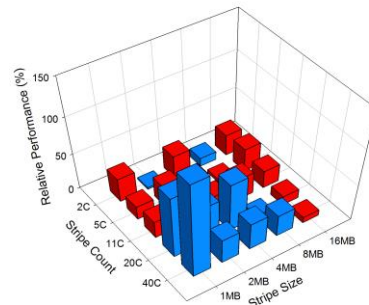
8x16 - EUS - rnetCDF-pnetCDF



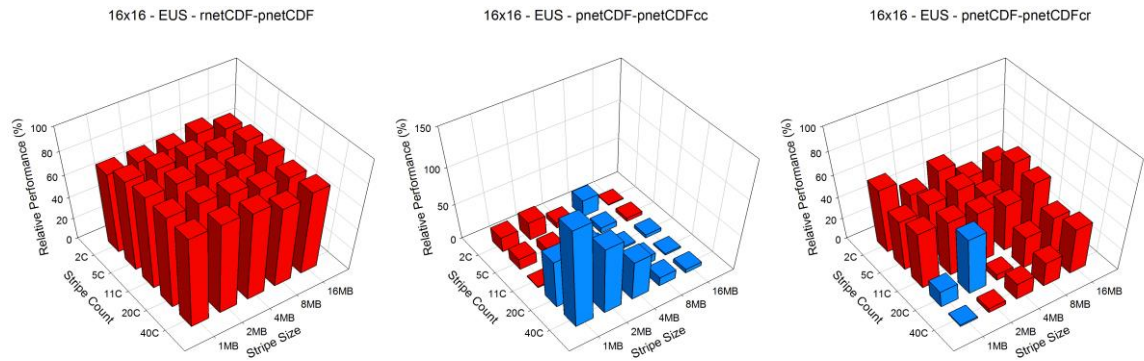
8x16 - EUS - pnetCDF-pnetCDFcc



8x16 - EUS - pnetCDF-pnetCDFcr



4



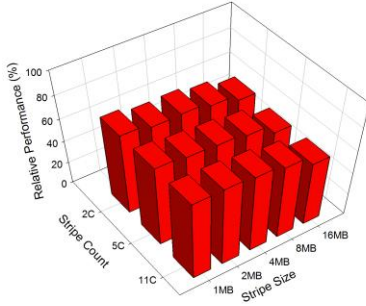
1

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.

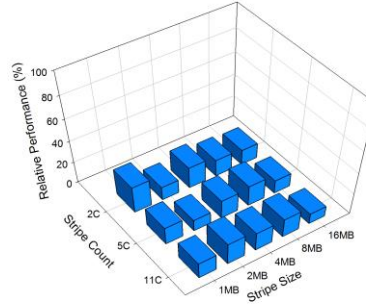
6



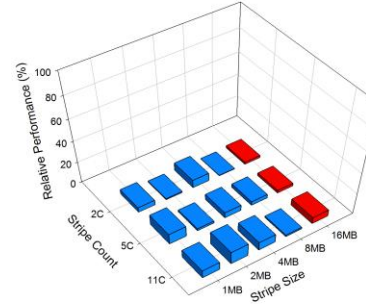
4x4 - EUS - rnetCDF-pnetCDF



4x4 - EUS - pnetCDF-pnetCDFcc

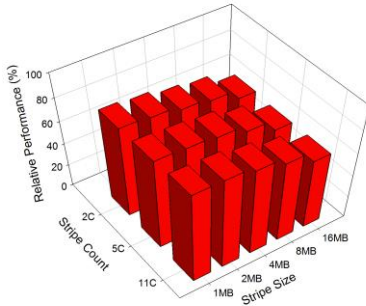


4x4 - EUS - pnetCDF-pnetCDFcr

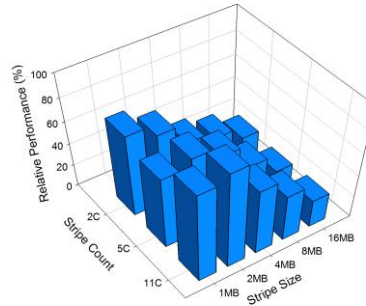


1

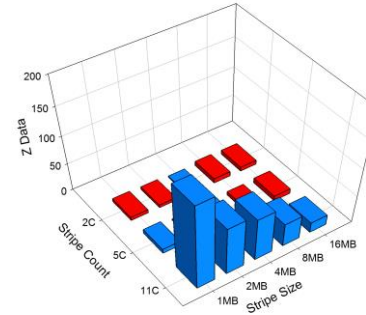
4x8 - EUS - rnetCDF-pnetCDF



4x8 - EUS - pnetCDF-pnetCDFcc

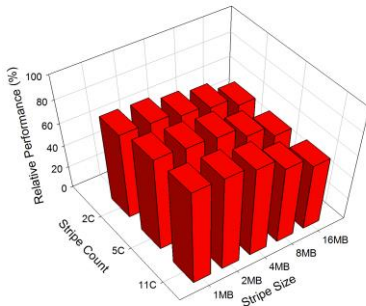


4x8 - EUS - pnetCDF-pnetCDFcr

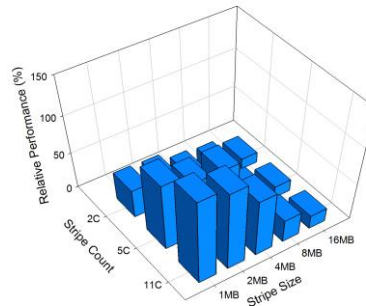


2

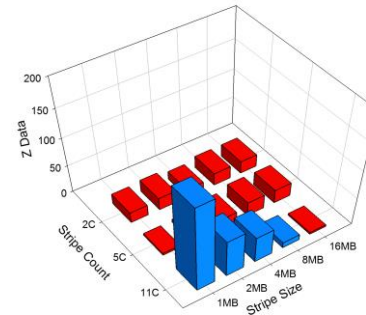
8x8 - EUS - rnetCDF-pnetCDF



8x8 - EUS - pnetCDF-pnetCDFcc

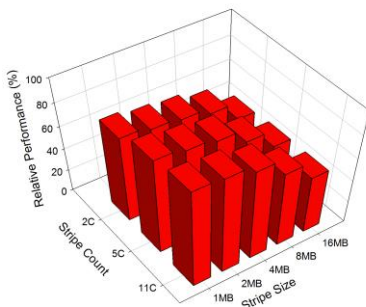


8x8 - EUS - pnetCDF-pnetCDFcr

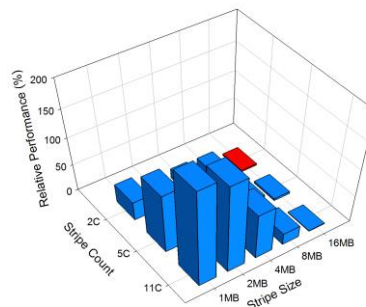


3

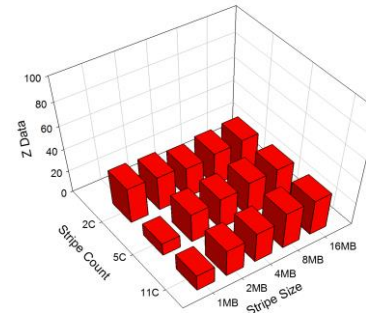
8x16 - EUS - rnetCDF-pnetCDF



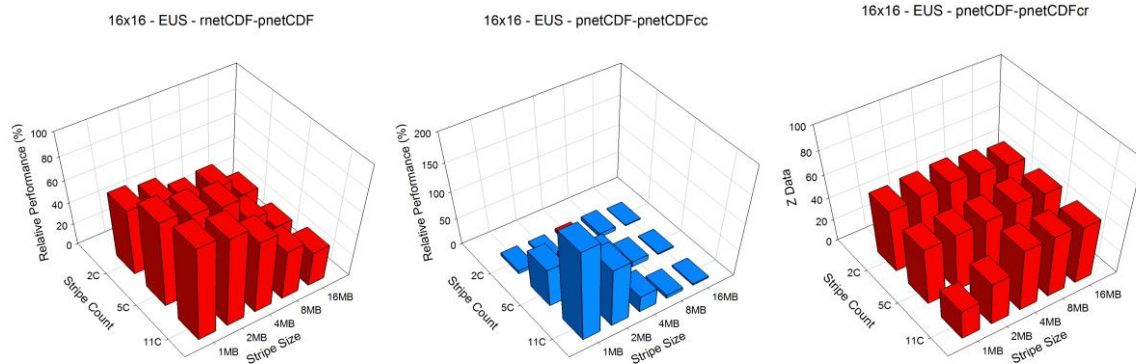
8x16 - EUS - pnetCDF-pnetCDFcc



8x16 - EUS - pnetCDF-pnetCDFcr



4

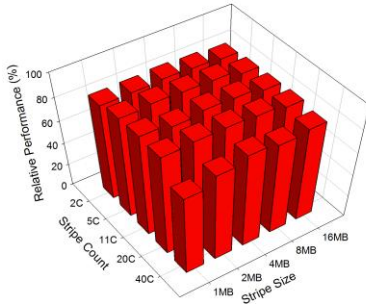


1

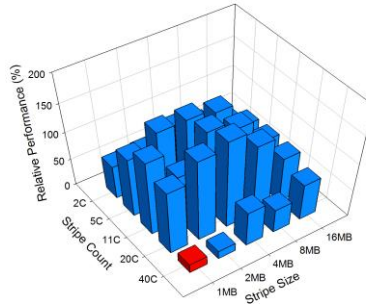
2 Figure 6. Relative I/O performance of pnetCDF to rnetCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on EUS domain with 4x4, 4x8, 8x8, 8x16, and 16x16 processor configuration  
 4 from Edison. Red colour denotes positive value in relative performance while blue colour  
 5 denotes negative value in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

6

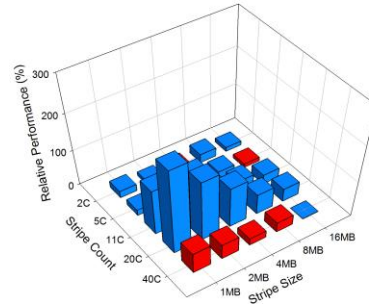
4x8 - CONUS - pnetCDF-pnetCDF



4x8 - CONUS - pnetCDF-pnetCDFcc

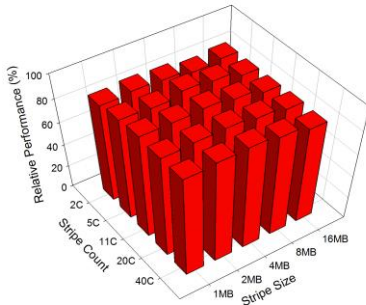


4x8 - CONUS - pnetCDF-pnetCDFcr

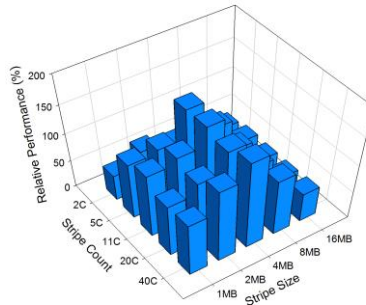


1

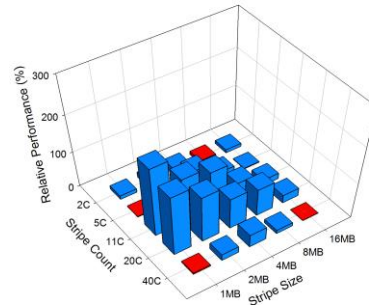
8x8 - CONUS - pnetCDF-pnetCDF



8x8 - CONUS - pnetCDF-pnetCDFcc

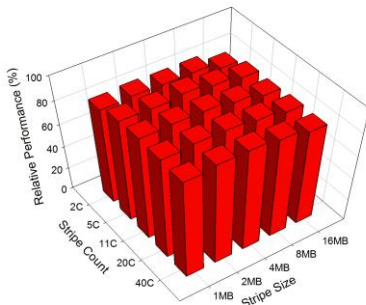


8x8 - CONUS - pnetCDF-pnetCDFcr

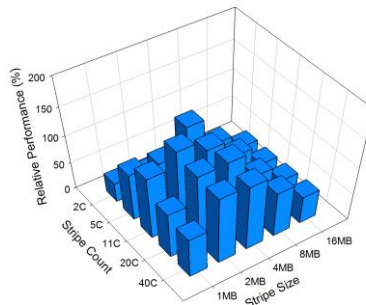


2

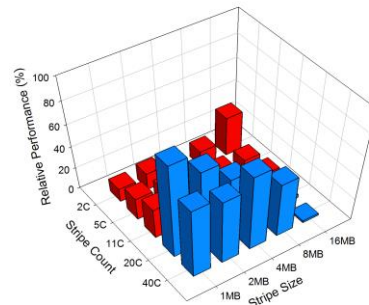
8x16 - CONUS - pnetCDF-pnetCDF



8x16 - CONUS - pnetCDF-pnetCDFcc

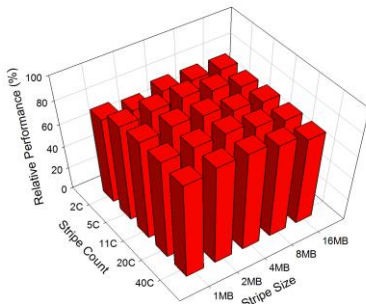


8x16 - CONUS - pnetCDF-pnetCDFcr

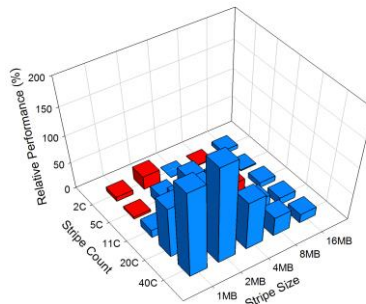


3

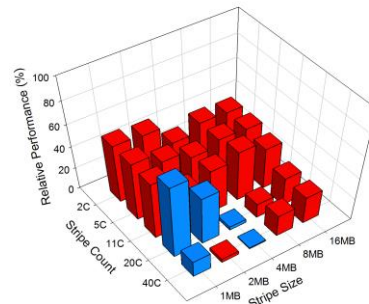
16x16 - CONUS - pnetCDF-pnetCDF



16x16 - CONUS - pnetCDF-pnetCDFcc

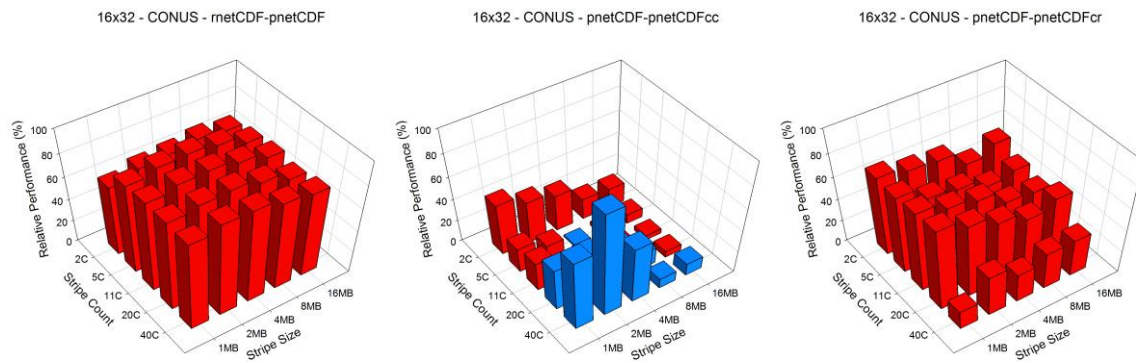


16x16 - CONUS - pnetCDF-pnetCDFcr



4



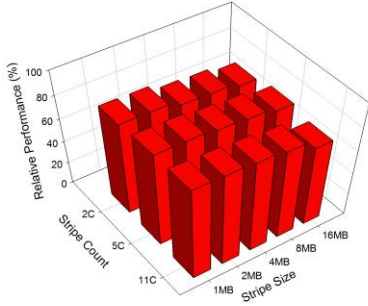


1

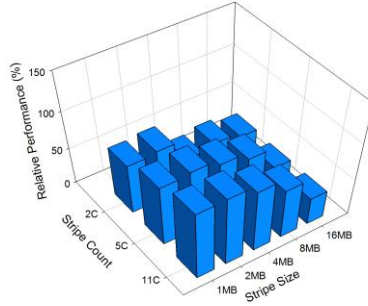
2 Figure 7. Relative I/O performance of pnetCDF to metCDF, and pnetCDFcc and pnetCDFcr  
 3 to pnetCDF on CONUS domain with 4x8, 8x8, 8x16, 16x16, and 16x32 processor  
 4 configuration from Kraken. Red colour denotes positive value in relative performance while  
 5 blue colour denotes negative value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

6

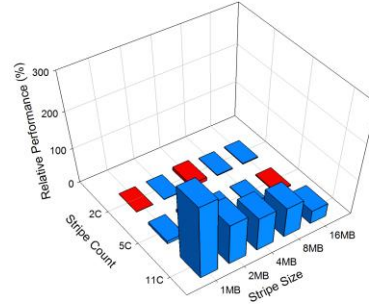
4x8 - CONUS - rnetCDF-pnetCDF



4x8 - CONUS - pnetCDF-pnetCDFcc

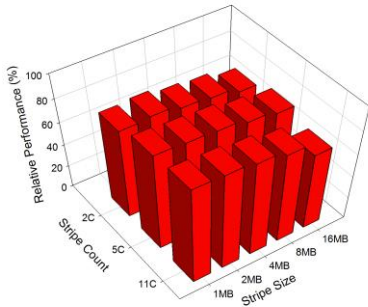


4x8 - CONUS - pnetCDF-pnetCDFcr

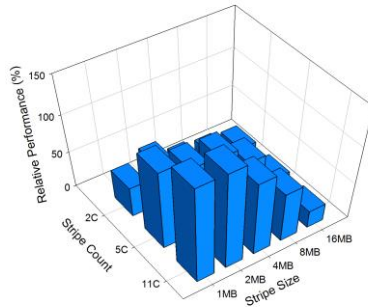


1

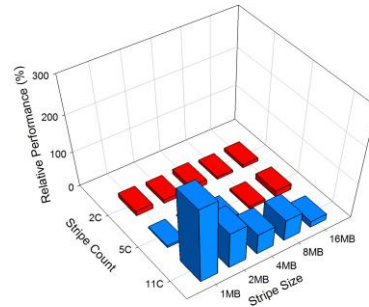
8x8 - CONUS - rnetCDF-pnetCDF



8x8 - CONUS - pnetCDF-pnetCDFcc

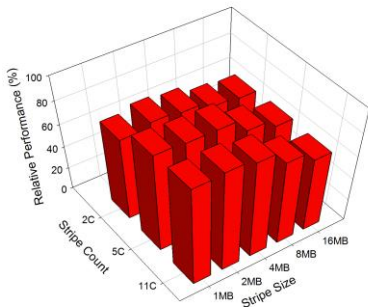


8x8 - CONUS - pnetCDF-pnetCDFcr

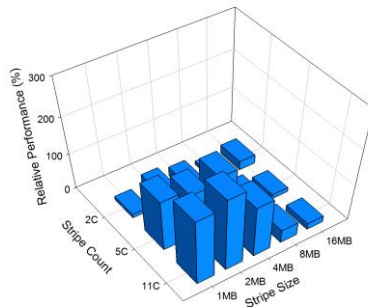


2

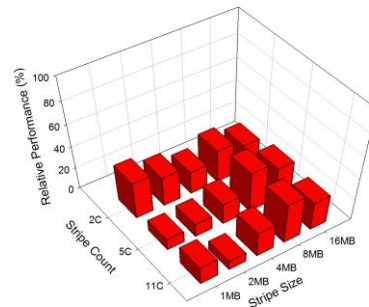
8x16 - CONUS - rnetCDF-pnetCDF



8x16 - CONUS - pnetCDF-pnetCDFcc

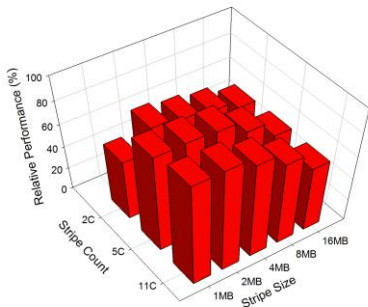


8x16 - CONUS - pnetCDF-pnetCDFcr

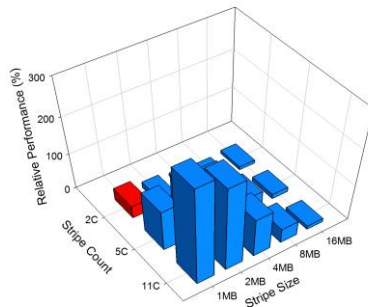


3

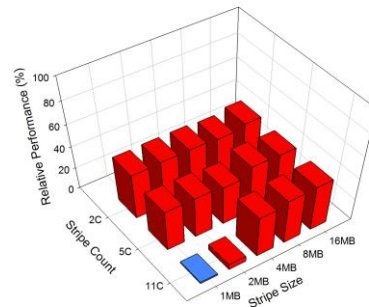
16x16 - CONUS - rnetCDF-pnetCDF



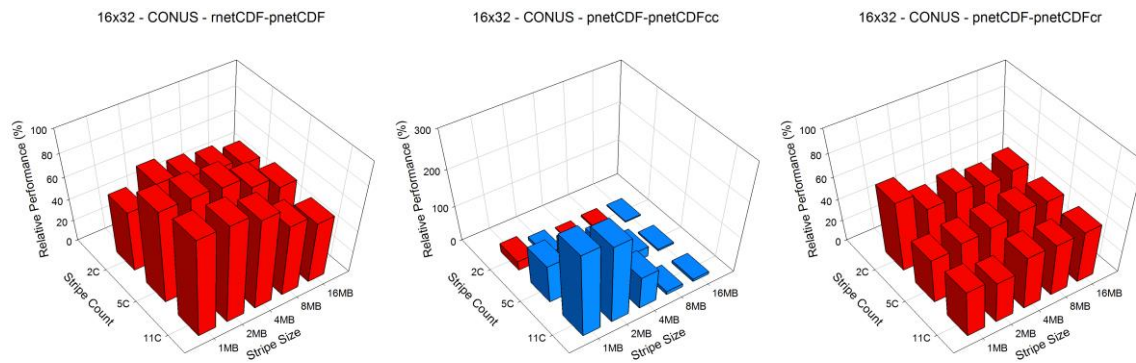
16x16 - CONUS - pnetCDF-pnetCDFcc



16x16 - CONUS - pnetCDF-pnetCDFcr



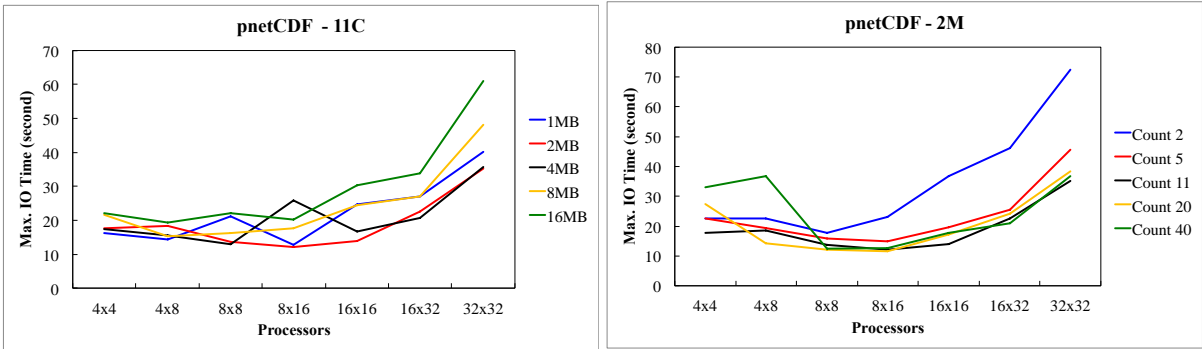
4



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.

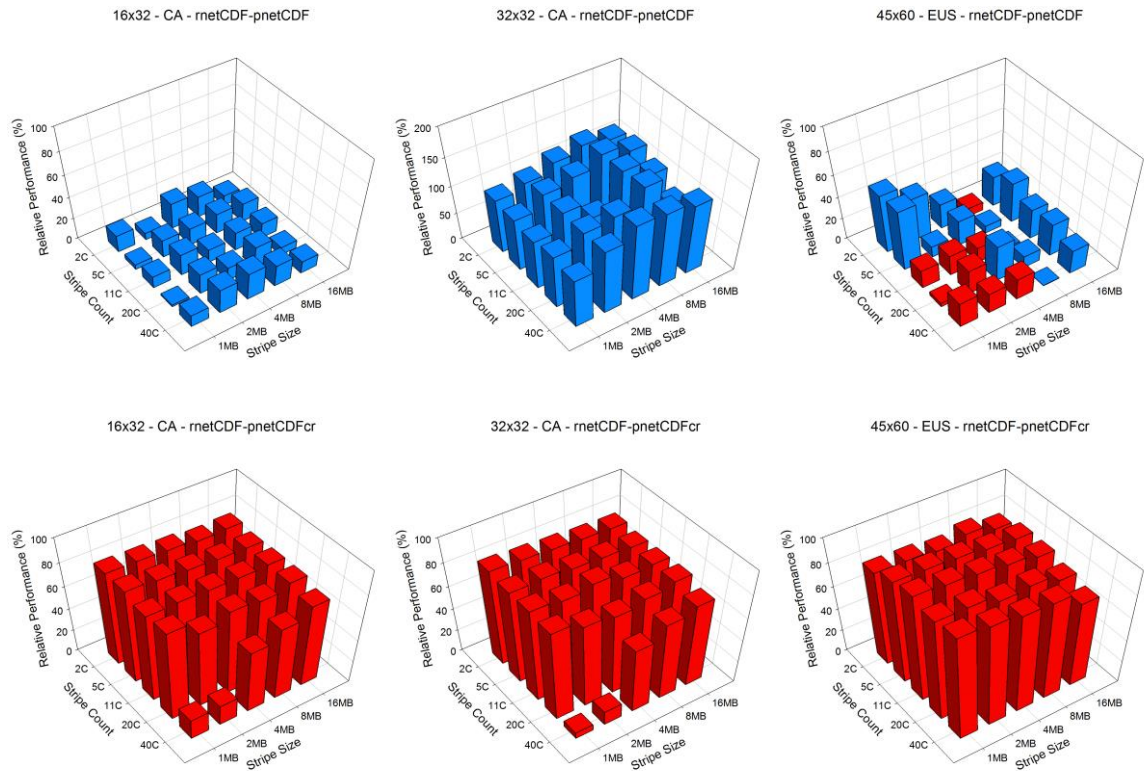
6



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



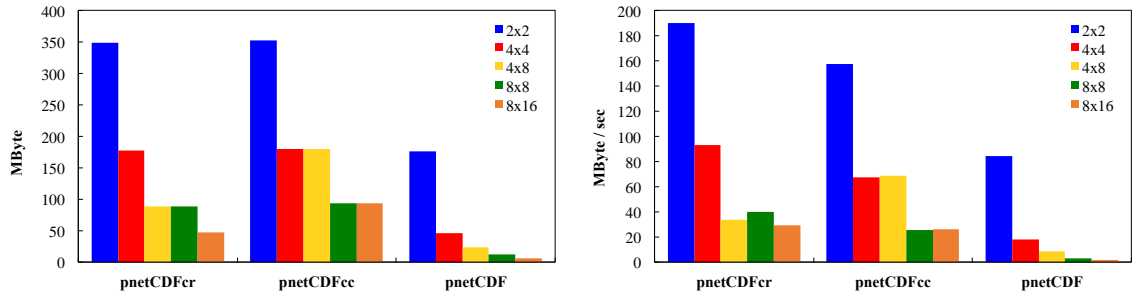
1

2

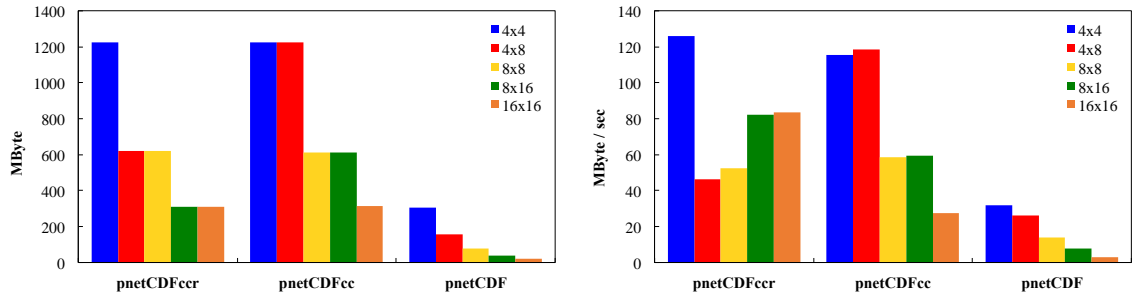
3 Figure 10. Relative performance of pnetCDF and pnetCDFcr with respect to mnetCDF in a  
 4 large number of processors scenario on Kraken. Red colour denotes positive value in relative  
 5 performance while blue colour denotes negative value in relative performance.

6

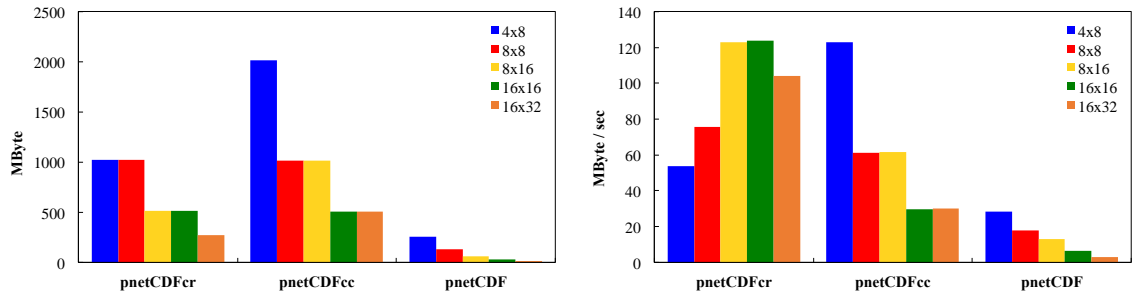
1



2

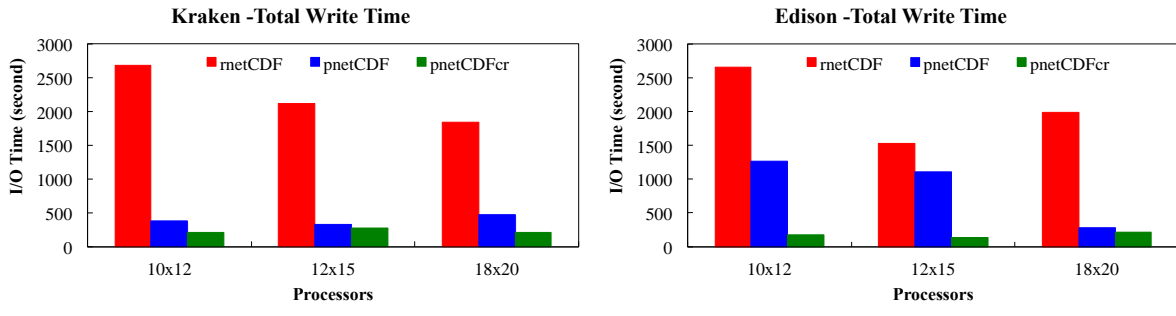


3



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).