

We thank Dr. Sophie Valecke very much for the comments and for the help of improving the English grammar and syntax. We'd like to reply the comments one by one as follows.

#### Major comments

1. p.5 L14: You conclude «Figure 1 demonstrates the poor performance of the P2P implementation» but you do not address the referee 2 comments that :
  - 1- the results can be affected by the decomposition strategy and
  - 2- the decrease in time at the end of the graph in Fig 1 (i.e. for more than 96 cores) should be remarked upon as one may wonder if it will continue to go down.Please modify your text to address these two comments.

Response: The corresponding context has been modified. Please refer to P5 L10-L11 for the first comment of the referee 2, and P5 L26-L30 for the second comment.

2. I do not understand the sentences on Fig. 15 «If the number of cores per toy model is less than 24, the MPI message number per process is set to be the number of cores. Otherwise, the MPI message number per process is set to 24. ». Similarly, I do not understand how you can choose the number of MPI messages per process and make it vary (see last paragraph on p.11 and Fig.16, or 2nd paragraph on p.12 and Fig.17, and many other places in the text). Can you explain ?

Response: We replaced “the number of MPI messages per process” with “the communication depth per sender process in the P2P implementation”. The communication depth of one process is defined as the number of the communications that a process is associated with (please refer to P5 L17-L22). For how to make the communication depth per sender process in the P2P implementation vary, please refer to Fig. 14 and Algorithm 1.

3. Section 5.2: I appreciate that you added a test with up to 1024 cores to answer referee1's related comment. But the results illustrated in Fig. 16, 17 and 18 with a toy model still go only up to 192 cores. As referee1, I would have expected that you perform these analysis with higher number of cores; this should be quite straightforward as they are based on toy models and it is easy to increase the size of the grid of a toy model and the number of cores used. As referee1, I do not understand why you did not produce those tests at higher number of cores. If there are sound arguments not to do so, please explain.

Response: In the revised version, up to 1024 cores are used in Fig. 15, 16 and 17.

4. Last paragraph of section 5.2 and Fig 19: The results on Fig 19 show that the adaptive library is only marginally better for 512, 768, 1024 cores. Therefore, I think the sentence «These results indicate that our proposed implementations can significantly improve the performance of data transfer for higher model resolution.» is misleading and should be rephrased.

Response: In our mind, we always refer the P2P implementation as the default baseline for evaluating the adaptive data transfer library, because the paper proposes the butterfly implementation. This misleading sentence has been rephrased in the revised version, please refer to P13 L14-L17.

5. Section 5.4: you explain that the P2P implementation is sufficient in the interpolation case because the number of MPI messages is small but you should clarify that this is probably linked to the fact that the two grids have close parallel decompositions as stressed by referee2.

Response: The fact that the P2P implementation can significantly outperform the butterfly implementation because the parallel decompositions are similar has been stated in the revised version (please refer to P14 L20-L24).

6. Section 5.4: again the conclusion of the paragraph is somewhat misleading. In this case, there is no real benefit brought by the adaptive library and it should be expressed clearly. Please rewrite the last sentence of section 5.4 accordingly.

Response: The conclusion has been corrected. Please refer to P14 L25-L28.

7. Section 5.5 : this section describes the comparison between the P2P and the adaptive library but not with the butterfly implementation so one cannot really conclude on the real improvement brought by the adaptive library. The same comparison should be done with the butterfly implementation as a baseline. Also it is please clarify if Figure 23 shows the gain in data transfer time only or the gain with respect to the whole model time. Finally, I do not understand the sentence «This performance improvement would not be low because the model coupling only takes a very small proportion of execution time in the simple coupled model GAMIL2-CLM3 and the parallel scalability of the two coupled models GAMIL2 and CLM3 is not good. “ Please clarify and rephrase.

Response: The performance speedup corresponding to the butterfly implementation has been added. Please refer to Fig. 22 and P15 L8-L12.

8. Section 6; as underlined by referee1's comment #14, the conclusions are still weak if not misleading. Please rework on the conclusions according to the comments made above.

Response: The conclusions as well as the abstract have been improved. Please refer to P1 L26-L28 and P15 L24-L26.

#### Other comments

1. Grammar and syntax: please consider the attached version of your current manuscript with some propositions to improve the english grammar and syntax.

Response: Thanks a lot for the help of improving the English grammar and syntax. Your modifications have been merged into the current revised version.

2. Title of the paper : Please give a more precise title and mention «adaptive library » in it, something like : «A new data transfer adaptive library improving model coupling. »

Response: The title has been modified according to your suggestion.

3. In your reply, you mention that the «version number has been added into the software» but I do not see any specific number or name for your library in the text and it is indeed missing, also in the title.

Response: Thanks a lot for this suggestion. We have added the version number of the adaptive data transfer library (please refer to P16 L1).

4. p.2 : when you refer to OASIS, please use also Valcke et al 2015, which describes the latest OASIS3-MCT version (see the modified document).

Response: We have referred to the latest OASIS3-MCT version. Please refer to P2 L30 and P20 L25-L26.

5. p. 3, L16-17: Your statement “can only scale to about 100 processor cores when using OASIS3 (Valcke, 2013) and to about 1000 processor cores when using OASIS3-MCT (Valcke et al., 2013);” is not correct. The reference does not show that the data transfer does not scale for more than 1000 cores; instead it shows that the data transfer does scale for up to 1000 cores. Please remove the sentence or rephrase it.

Response: This sentence has been removed.

6. p.4, L18-20 : It seems awkward to conclude here that «P2P implementation can achieve good performance when rearranging data fields for a parallel interpolation in a component model» as this is indeed shown in section 5.4 and Fig. 22; this should be removed at this point in the text.

Response: This sentence has been removed.

7. p.4, L21: The references Valcke, 2013; Valcke et al., 2013 do not show that P2P is “not efficient enough when transferring data between component models” as stated. Please remove those references there.

Response: Those references have been removed.

8. Fig 4 : I think the legend of the x axis should be changed from «number of cores per process » for «number of cores per model »; can you confirm and make the change ?

Response: The original Fig. 4 about the total number of messages has been replaced by the original Fig. 6 (now Fig. 4 in this revised version) about the average communication depth, because the average communication depth is more relative to the performance of the P2P implementation.

9. p.6, 2nd paragraph : I do not understand why you refer to Fig 6 and Fig 3. Please explain or remove the reference.

Response: Those references have been removed, please refer to P6 L21-L24.

10. Fig 12 : I do not understand the meaning of the last sentence «Each process of the sender is mapped onto a process of the butterfly kernel, while every two processes of the receiver are mapped onto one process of the butterfly kernel». I think it does not bring any additional information. In particular, if the receiver has 10 processes and if only 3 processes of the receiver are used for the butterfly kernel, how can this be 'every two processes of the receiver' . Could you please explain or remove the sentence?

Response: This sentence has been removed (please refer to Fig. 11).

11. 2nd paragraph of 5.3 and Fig. 20: The last sentence “When each component uses 192 cores, the adaptive data transfer library is 4.01 times faster than the P2P implementation” is right but Fig. 20 also shows that butterfly and adaptive seem to converge when increasing the number of cores per model. This should also be described in the text.

Response: The corresponding context has been improved according to this suggestion (please refer to P13 L31- P4 L2).

12. Figures 15 to 22, it would be better to change «Library »for «Adaptive library »or «Adaptive »

Response: These figures have been improved accordingly (please refer to Fig. 14 to 21).

13. Figure 20 and your response to Referee2's comment #7: I do not understand the sentence «The P2P results are from the adaptive data transfer library which switches to the P2P implementation. » and this is not mentioned in the text. Please clearly explain what it means.

Response: The corresponding context has been improved (please refer to Fig. 19 and P9 L30- P10 L1)

# Improving Data Transfer for Model Coupling

## A new adaptive data transfer library for model coupling

C. Zhang<sup>2,1</sup>, L. Liu<sup>1,3</sup>, G. Yang<sup>2,1,3</sup>, R. Li<sup>2,1</sup>, and B. Wang<sup>1,3,4</sup>

[1]{Ministry of Education Key Laboratory for Earth System Modeling, Center for Earth System Science (CESS), Tsinghua University, Beijing, China }

[2]{Department of Computer Science and Technology, Tsinghua University, Beijing, China }

[3]{Joint Center for Global Change Studies (JCGCS), Beijing, China }

[4]{State Key Laboratory of Numerical Modelling for Atmospheric Sciences and Geophysical Fluid Dynamics (LASG), Institute of Atmospheric Physics, Chinese Academy of Sciences, Beijing, China. }

Correspondence to: L. Liu (liuli-cess@tsinghua.edu.cn), G. Yang (ygw@tsinghua.edu.cn)

### Abstract

Data transfer means transferring data fields ~~between two component models or rearranging data fields among processes of the same component model from a sender to a receiver~~. It is a fundamental and most frequently used operation of a coupler. Most versions of state-of-the-art couplers currently use an implementation based on the point-to-point (P2P) communication of the Message Passing Interface (MPI) (refer such an implementation as “P2P implementation” for short). In this paper, we ~~reveal~~revealed the drawbacks of the P2P implementation when the parallel decompositions of the sender and the receiver are different, including low communication bandwidth due to small message size, variable and big ~~number of MPI messages~~communication depth, as well as network contention. To overcome these drawbacks, we ~~propose~~proposed a butterfly implementation for data transfer. Although the butterfly implementation can outperform the P2P implementation in many cases, it degrades the performance ~~in some cases because when~~ the total message size transferred by sender and the butterfly implementation is larger than receiver have similar parallel decompositions or the total message size transferred by the P2P implementation number of processor cores used for running models is small. To ~~further improve~~make data transfer always keep the optimal performance, we ~~design~~designed and ~~implement~~implemented an adaptive data transfer library

1 that combines the advantages of both butterfly implementation and P2P implementation.  
2 ~~Performance evaluation shows that~~As the adaptive data transfer library significantly  
3 ~~improves~~can adaptively use the ~~performance of~~better implementation for data transfer, it  
4 outperforms the P2P implementation in ~~most~~many cases, ~~and~~ while does not decrease the  
5 performance in any cases. Now, the adaptive data transfer library is open to the public and has  
6 been imported into a coupler version C-Coupler1 for performance improvement of data transfer.  
7 We believe that other ~~coupler versions~~couplers can also benefit from it.

## 9 **1 Introduction**

10 Climate System Models (CSMs) and Earth System Models (ESMs) are fundamental tools for  
11 simulating, predicting and projecting climate. A CSM or an ESM generally integrates several  
12 component models, such as an atmosphere model, a land surface model, an ocean model and a  
13 sea-ice model, into a coupled system to simulate the behaviours of ~~of~~ the climate system,  
14 including the interactions between components of the climate system. More and more coupled  
15 models have sprung up in the world. For example, the number of coupled model configurations  
16 in the Coupled Model Intercomparison Project (CMIP) has increased from less than 30 (used  
17 for CMIP3) to more than 50 (used for CMIP5).

18 High-performance computing is an essential technical support for model development,  
19 especially for higher and higher resolutions of models. Modern high-performance computers  
20 integrate an increasing number of processor cores for higher and higher computation  
21 performance. Therefore, efficient parallelization, which enables a model to utilize more  
22 processor cores for acceleration, becomes a technical focus in model development; and a  
23 number of component models with efficient parallelization have sprung up. For example, the  
24 Community Ice Code (CICE; Hunke et al., 2008, 2013) at 0.1 °horizontal resolution can scale  
25 to 30,000 processor cores on the IBM Blue Gene/L (Dennis et al., 2008); the Parallel Ocean  
26 Program (POP; Kerbyson, 2005; Smith et al., 2010) at 0.1 °horizontal resolution can also scale  
27 to 30,000 processor cores on the IBM Blue Gene/L and 10,000 processor cores on a Cray XT3  
28 (Dennis, 2007); the Community Atmosphere Model (CAM; Morrison et al., 2008; Neale et al.,  
29 2010, 2012) with a spectral element dynamical core (CAM-SE) at 0.25 °horizontal resolution  
30 can scale to 86,000 processor cores on a Cray XT5 (Dennis et al., 2012).

31 A coupler is an important component in a coupled system. It links component models together  
32 to construct a coupled model, and controls the integration of the whole coupled model (Valcke

1 [et al.](#), 2012). A number of couplers now are available, e.g., the Model Coupling Toolkit (MCT;  
2 Jacob et al., 2005), the Ocean Atmosphere Sea Ice Soil coupling software (OASIS) coupler  
3 (Redler et al., 2010; Valcke, 2013; [Valcke et al., 2015](#)), the Earth System Modelling Framework  
4 (ESMF; Hill et al., 2004), the CPL6 coupler (Craig et al., 2005), the CPL7 coupler (Craig et al.,  
5 2012), the Flexible Modelling System (FMS) coupler (Balaji et al., 2006), the Bespoke  
6 Framework Generator (BFG; Ford et al., 2006; Armstrong et al., 2009) and the community  
7 coupler version 1 (C-Coupler1; Liu et al., 2014).

8 A coupler generally has much smaller overhead than the component models in ~~acurrent~~ coupled  
9 ~~systems~~systems. However, it is potentially a time-consuming component ~~of ain future~~ coupled  
10 ~~model in future~~models. This is because more and more component models (such as land-ice  
11 model, chemistry model and biogeochemical model) will be coupled into a coupled model, and  
12 the coupling frequency between component models will be higher and higher. Data transfer is  
13 a fundamental and ~~most~~ frequently used operation in a coupler. It is responsible for transferring  
14 data fields between the processes of two component models and for rearranging data fields  
15 among processes of the same component model for parallel data interpolation.

16 A coupler may become a bottleneck for efficient parallelization of future coupled models. The  
17 most obvious reason is that the current implementation of data transfer in a state-of-the-art  
18 coupler ~~is not efficient enough for transferring data fields between component models. For~~  
19 ~~example, the data transfer from a component with a logically rectangular grid (of 1021×1442~~  
20 ~~grid points) to a component with a Gaussian Reduced T799 grid (with 843,000 grid points) can~~  
21 ~~only scale to about 100 processor cores when using OASIS3 (Valeke, 2013) and to about 1000~~  
22 ~~processor cores when using OASIS3 MCT (Valeke et al., 2013); the data transfer may be not~~  
23 ~~efficient enough. For example, due to the low efficiency of data transfer, the coupling~~ from a  
24 component model with a horizontal grid (of 576×384 grid points) to another component model  
25 with another horizontal grid (of 3600×2400 grid points) can only scale to about 500 processor  
26 cores when using the CPL7 coupler (Craig et al., 2012). Therefore, it is highly desirable to  
27 improve the ~~parallelization~~parallel data transfer of couplers.

28 In this study, we first propose a butterfly implementation of data transfer. Since the P2P  
29 implementation and the butterfly implementation can outperform each other in different cases  
30 (~~Section~~Sect. 5), we next develop an adaptive data transfer library that includes both  
31 implementations and can adaptively use the better one for data transfer. Performance evaluation  
32 demonstrates that such a library significantly ~~improves the performance of data~~

1 ~~transfer~~outperforms the P2P implementations in most cases and does not degrade the  
2 performance in any case. This library has been imported into C-Coupler1 with slight code  
3 modification. We believe that other ~~coupler versions~~couplers can also benefit from it.

4 The ~~reminder~~remainder of this paper is organized as follows. We briefly introduce the  
5 implementation of data transfer in existing couplers in Section 2. Details of the butterfly  
6 implementation and the adaptive data transfer library are presented in Sections 3 and 4,  
7 respectively. The performances of data transfer implementations are evaluated in Section 5.  
8 Conclusions are given in Section 6.

## 9 **2 Data transfer implementations in existing couplers**

### 10 **2.1 P2P implementation**

11 Almost all state-of-the-art couplers use a similar implementation for data transfer. To achieve  
12 parallel data transfer, MCT first generates a communication router (known as the data mapping  
13 between processes) according to the parallel decompositions (the distribution of grid points  
14 among the processes) of ~~two component models~~the sender and the receiver, and then uses the  
15 point-to-point (P2P) communication of the Message Passing Interface (MPI) to transfer the data.  
16 A data field will be transferred from a process of the ~~source component model~~sender to a  
17 process of the ~~target component model~~receiver, only when the two processes have common  
18 grid points. In the following context, we call this “P2P implementation” for short.

19 Since MCT has already been imported into OASIS3-MCT, the CPL6 coupler and the CPL7  
20 coupler, these couplers also use the P2P implementation for data transfer. Although the other  
21 couplers such as ESMF, OASIS4, the FMS coupler and C-Coupler1 do not directly import MCT,  
22 they also use the P2P implementation for data transfer.

### 23 **2.2 Performance bottlenecks of the P2P implementation**

24 ~~Although the P2P implementation can achieve good performance when rearranging data fields~~  
25 ~~for a parallel interpolation in a component model, it is not efficient enough when transferring~~  
26 ~~data between component models (Craig et al., 2012; Valeke, 2013; Valeke et al., 2013; Liu et~~  
27 ~~al., 2014). To reveal why the P2P implementation is not efficient enough, we first~~To motivate  
28 this work, we first investigate the performance characteristics of the P2P implementation, and  
29 therefore derive a benchmark from a real coupled model GAMIL2-CLM3, which includes  
30 GAMIL2 (Li et al., 2013) that is an atmosphere model and CLM3 (Oleson et al., 2004;



1 Dickinson et al., 2006) that is a land surface model. GAMIL2 and CLM3 share the same  
2 horizontal grid of 7,680 (128×60) grid points, but have different parallel decompositions:  
3 GAMIL2 uses a regular 2-D parallel decomposition, while CLM3 uses an irregular 2-D parallel  
4 decomposition where the grid points are assigned to the processes in a round-robin fashion.

5 In this benchmark, there is only the data transfer with the P2P implementation between two  
6 data models the sender and the receiver with the same horizontal grid of GAMIL2-CLM3. The  
7 parallel decomposition of the source data model sender is derived from CLM3, and the parallel  
8 decomposition of target data model the receiver is derived from GAMIL2. A high-performance  
9 computer named Tansuo100 at Tsinghua University, China is used for the performance tests. It  
10 has 700 computing nodes, each of which contains two six-core Intel Xeon X5670 CPUs and 32  
11 GB main memory. All computing nodes are connected by a high-speed InfiniBand network  
12 with peak communication bandwidth of 5 GB/s.

13 To evaluate the parallel performance of the P2P implementation, 14 2-D coupling fields are  
14 transferred between the two data models sender and the receiver. In each test, the two data  
15 models sender and the receiver use the same number of processes. Since there are 12  
16 CPU processor cores on each computing node, the number of processes is set to be an integral  
17 multiple of 12. When the process number of processes is less than 12, the two data models sender  
18 and the receiver are located on two different computing nodes. The two data models sender and  
19 the receiver do not share the same computing node, so the communication of the P2P  
20 implementation must go through the InfiniBand network.

21 Figure 1 demonstrates that poor performance parallel scalability of the P2P implementation  
22 can be obtained when the parallel decompositions of the sender and receiver are different. It is  
23 well known that the communication performance heavily depends on message size. As shown  
24 in Fig. 2, the P2P communication bandwidth achieved generally increases with message size.  
25 So when the message size is small (for example, smaller than 4 KB), the communication  
26 bandwidth achieved is very low. The message size in the P2P implementation decreases with  
27 the increment of process number of processes of models (Fig. 3), indicating that the  
28 communication bandwidth becomes lower with the increment of process number of processes.  
29 The performance of data transfer also heavily depends on the MPI message number another  
30 term of communication depth, which is defined as the number of the communications that a  
31 process is associated with. The communication depth is determined by the parallel  
32 decompositions of the sender and the receiver. In the P2P implementation, if one process of the

1 sender/receiver has common grid points with  $N$  processes of the receiver/sender, the  
2 communication depth of this process is  $N$ . As shown in Fig. 4, the ~~message number variation~~  
3 of average communication depth in the P2P implementation increases with increment of  
4 process is consistent with the variation of the execution time of the P2P implementation in Fig.  
5 1: both the average communication depth and the execution time of the P2P implementation  
6 increase with the number of cores from 6 to 48, and go down with the number. Here, we may  
7 conclude that the decrease of message size and the increase of message of cores from 96 to 192.  
8 Lower execution time of the P2P implementation will be obtained if more cores are used (the  
9 maximum number are primary of cores in both Fig. 1 and Fig. 4 is limited to 192 because  
10 GAMIL2-CLM3 will not be further accelerated when using more cores) since the average  
11 communication depth will further go down. To further reveal possible reasons for the poor  
12 performance of the P2P implementation when increasing the process number. However, the  
13 parallel scalability, we evaluate the ideal performance ~~shown~~ and actual performance in Fig. 5.  
14 The ideal performance is much better than the actual performance. ~~The, and the~~ ratio between  
15 the ideal performance and the actual performance significantly increases with the increment of  
16 ~~processor~~the number of processes. The significant gap between the ideal performance and the  
17 actual performance is due to network contention. For example, when multiple P2P  
18 communications share the same ~~source~~sender process or ~~target~~receiver process (Fig. 6),<sub>2</sub> they  
19 must wait in ~~an~~ order.

### 20 **3 Butterfly implementation for better performance of data transfer**

21 The drawbacks of the P2P implementation when the sender and the receiver use different  
22 parallel decompositions can be ~~concluded~~identified as low communication bandwidth due to  
23 small message size, variable and big ~~number of MPI messages~~communication depth, as well as  
24 network contention. To overcome these drawbacks, a prospective solution is to organize the  
25 ~~communication for data transfer~~ of data using a better ~~structure, so that we investigate~~algorithm,  
26 e.g., the butterfly ~~structure~~algorithm (Fig. 76), which has already been ~~used~~studied in ~~the field~~  
27 of ~~computer~~computing sciences (Chong et al., 1994; Foster, 1995; Heckbert et al., 1995;  
28 Hemmert et al., 2005; Kim et al., 2007; Jan et al., 2013; Petagon et al, 2016). ~~For example, in~~In  
29 hardware aspect, the traditional butterfly ~~structure~~algorithm and its transformation have been  
30 used to design networks (Chong et al., 1994; Kim et al., 2007); in software aspect, the butterfly  
31 ~~structure~~algorithm has been used to improve the parallel algorithms with all-to-all

1 communications (Foster, 1995), e.g., Fast Fourier Transform (FFT; Heckbert et al., 1995;  
2 Hemmert et al., 2005), matrix transposition (Petagon et al, 2016) and sorting (Jan et al., 2013).

3 Unfortunately, the ~~improved all-to-all communication with the~~classical butterfly  
4 ~~structure~~algorithm cannot be used as is to improve data transfer, because it requires that one  
5 process ~~must communicate~~communicates with every other process, that the communication  
6 load among processes is balanced and that the number of processes must be a power of 2, ~~while~~  
7 ~~the.~~ In practice, data transfer for model coupling has different characteristics, i.e., one process  
8 needs to communicate with a part of other processes (~~Fig. 6~~), the communication load among  
9 processes is always unbalanced (~~Fig. 3~~) and the ~~process~~ number of processes cannot be  
10 restricted to a power of 2. Therefore, ~~to benefit from the butterfly structure,~~ we should  
11 design~~propose here~~ a new implementation of data transfer, ~~which is called the~~ involving an  
12 additional butterfly ~~implementation hereafter~~.

13 ~~The butterfly implementation uses a butterfly structure~~kernel to transfer data from the sender  
14 with the source parallel decomposition to the receiver with the target parallel decomposition.  
15 ~~We call the communication following the butterfly structure “the butterfly kernel”.~~ As the  
16 ~~process~~ number of processes of the butterfly kernel must be a power of 2, while the ~~process~~  
17 number of processes of the sender or the receiver ~~need are~~ not be a power of 2, the butterfly  
18 ~~implementation (see Fig. 8) has a process mapping from the sender onto the butterfly kernel~~  
19 ~~and a process mapping from the butterfly kernel onto the receiver, and~~necessarily, the butterfly  
20 kernel has its own source parallel decomposition and target parallel decomposition, ~~which are~~  
21 ~~determined by the~~and process mappings: are needed from the sender onto the butterfly kernel  
22 and from the butterfly kernel onto the receiver (see Fig. 7). Next, we ~~will~~ present the butterfly  
23 kernel and the process mappings, respectively.

### 24 **3.1 Butterfly kernel**

25 The first question for the butterfly kernel is how to decide its ~~process~~ number of processes.  
26 Any process of the sender or ~~the~~ receiver can be used as a process ~~offor~~ the butterfly kernel.  
27 Given that the total number of unique processes of the sender and receiver is  $N_T$ , the ~~process~~  
28 number of processes of the butterfly kernel ( $N_B$ ) can be any power of 2, which is no larger than  
29  $N_T$ . We propose to select the maximum number in order for maximum utilization of resources.  
30 We prefer to pick out unique processes first from the sender, and then from the receiver if the  
31 sender does not have enough processes.

1 The butterfly kernel is responsible for rearranging the distribution of data among the processes  
2 from the source parallel decomposition to the target parallel decomposition. Given the ~~process~~  
3 number of processes  $N=2^n$ , there are  $n$  stages in the butterfly kernel. In a stage, all processes  
4 are divided into a number of pairs and the two processes of a pair uses MPI P2P communication  
5 to exchange data. ~~Given a process P in the butterfly kernel, after~~After each stage, the number  
6 of ~~the butterfly kernel~~ processes that may have the data ~~of P that will finally belong to any one~~  
7 process on the target parallel decomposition will become a half. Figure ~~76~~ is an example for  
8 further illustration, where  $D_i^j$  means the data is originally in process  $P_i$  according to the source  
9 parallel decomposition and is finally in process  $P_j$  according to the target parallel decomposition.  
10 Before the first stage, all processes ( $P_0 \sim P_7$ ) may have the data of  $P_0$  on the target parallel  
11 decomposition. After the first stage, only four processes ( $P_0, P_2, P_4$  and  $P_6$ ) may have that data;  
12 and after the second stage, only two processes ( $P_0$  and  $P_4$ ) may have it.

13 To reveal the advantages and disadvantages of the two implementations, we measure the  
14 characteristics of the two implementations based on the benchmark introduced in Section 2.2.  
15 The results show that the total message size transferred by the butterfly implantation is larger  
16 than that by the P2P implementation (Fig. ~~98~~), which is the major disadvantage of the butterfly  
17 implementation. Meanwhile, comparing with the P2P implementation, the butterfly  
18 implementation has can have the following advantages:

- 19 1)- bigger message size for better communication bandwidth (Fig. ~~109~~);
- 20 2) -balanced number of MPI messages and smaller communication depth among processes (Fig.  
21 ~~110~~);
- 22 3) ordered communications among processes and fewer communications operated concurrently  
23 (Fig. ~~110~~), which can dramatically reduce network contention.

### 24 **3.2 Process mapping**

25 In this subsection, we will introduce the process mappings from the sender to the butterfly  
26 kernel and from the butterfly kernel to the receiver. To minimize the overhead of process  
27 mapping from the butterfly kernel to the receiver, we map one or multiple processes of the  
28 butterfly kernel onto a process of the receiver if the butterfly kernel has more processes than  
29 the receiver; otherwise, we map a process of the butterfly kernel onto one or multiple processes  
30 of the receiver. In other words, there is no multiple-to-multiple process mapping between the

1 butterfly kernel and the receiver. Similarly, there is no multiple-to-multiple process mapping  
2 between the sender and the butterfly kernel.

3 Processes of the sender or the receiver may be unbalanced in terms of the size of the data  
4 transferred, which may result in unbalanced communications among processes of the butterfly  
5 kernel. As mentioned in Section 3.1, at each stage of the butterfly kernel, all processes are  
6 divided into a number of pairs, each of which is involved in P2P communications. To improve  
7 the balance of communications among the processes in the butterfly kernel, one solution is to  
8 try to make the process pairs at each stage more balanced in terms of data size of P2P  
9 communications, so we propose to reorder the processes of the sender or the receiver according  
10 to data size. At the first stage, each time we pick out the process with the largest data size and  
11 the process with the smallest data size from the remaining processes that have not been paired,  
12 to generate a process group. For the next stage, the outputs of two process groups from the  
13 previous stage are paired into a bigger process groups in a similar way. After finishing the  
14 iterative pairing throughout all stages, all processes of the sender or the receiver are reordered.

15 The iterative pairing also requires the number of processes to be a power of 2. Given that the  
16 ~~process~~-number of processes of the sender (or receiver) is  $N_C$  and the ~~process~~-number of  
17 processes of the butterfly kernel is  $N_B$ , we first pad empty processes (whose data size is zero)  
18 before the iterative pairing to make the ~~process~~-number of processes of the sender (or receiver)  
19 be a power of 2 (denoted  $N_P$ ), which is no smaller than  $N_B$ . Therefore, the reordered  $N_P$   
20 processes after the iterative pairing can be divided into  $N_B$  groups, each of which contains  $N_P/N_B$   
21 processes with consecutive reordered indexes and maps onto a unique process of the butterfly  
22 kernel.

23 Figure 4211 shows an example of the process mapping, where the sender has five processes  
24 ( $S_0$ - $S_4$  in Fig. 42a11a), the receiver has 10 processes ( $R_0$ - $R_9$  in Fig. 42b11b), and the butterfly  
25 kernel uses eight processes ( $B_0$ - $B_7$  in Fig. 42e11c). At the first, empty processes are padded to  
26 the sender ( $S_5$ - $S_7$  in Fig. 42a11a) and the receiver ( $R_{10}$ - $R_{15}$  in Fig. 42b11b). Next, the iterative  
27 pairing is conducted for the sender and the receiver, respectively. The iterative pairing has three  
28 stages for the sender. At the first stage, the eight processes of the sender are divided into four  
29 groups  $\{S_1, S_7\}$ ,  $\{S_0, S_6\}$ ,  $\{S_2, S_5\}$  and  $\{S_4, S_3\}$  (Fig. 42a11a), according to the data size  
30 corresponding to each process. These four process groups are divided into two bigger groups  
31 ( $\{\{S_4, S_3\}, \{S_2, S_5\}\}$  and  $\{\{S_1, S_7\}, \{S_0, S_6\}\}$  at the second stage (Fig. 42a11a). Finally, one process  
32 group  $\{\{\{S_4, S_3\}, \{S_2, S_5\}\}, \{\{S_1, S_7\}, \{S_0, S_6\}\}\}$  is obtained at the third stage (Fig. 42a11a), and

1 the eight processes of the sender are reordered as  $S_4, S_3, S_2, S_5, S_1, S_7, S_0$  and  $S_6$ , each one ~~of~~  
2 ~~which is being~~ mapped onto one process of the butterfly kernel (Fig. 42e11c). Similarly, the  
3 iterative pairing has four stages for the receiver, and the 16 processes of the receiver are  
4 reordered as  $R_9, R_{15}, R_7, R_{12}, R_4, R_8, R_3, R_{10}, R_1, R_{14}, R_5, R_{13}, R_0, R_6, R_2$  and  $R_{11}$  finally, ~~each~~  
5 ~~two with pairs~~ of ~~which are these being~~ mapped onto one process of the butterfly kernel (Fig.  
6 42e11c).

#### 7 4 Adaptive data transfer library

8 Now, we have two kinds of implementations (the P2P implementation and the butterfly  
9 implementation) for data transfer. Although the butterfly implementation can effectively  
10 improve the performance of data transfer in many cases (examples are given in Section 5), it  
11 has some drawbacks: 1) it generally has a larger total message size ~~of communications~~ than the  
12 P2P implementation; 2) its ~~stage-number~~ ~~of stages~~ is  $\log_2 N$  (where  $N$  is the number of processes  
13 for the butterfly kernel) (Foster, 1995), which may be bigger than the average ~~number of MPI~~  
14 ~~messages per process~~ ~~communication depth~~ in the P2P implementation in some cases (for  
15 example, ~~the data rearrangement for when the sender and the receiver use the similar~~ parallel  
16 ~~interpolation decompositions~~). Therefore, it is possible that the P2P implementation  
17 outperforms the butterfly implementation in some cases ~~(examples are given in Section 5)~~. To  
18 achieve optimal performance for data transfer, we propose an adaptive data transfer library that  
19 can take the advantages of the two implementations in all cases.

20 As introduced in Section 3.1, the butterfly implementation is divided into multiple stages.  
21 Actually, the data transfer in one stage can be viewed as a P2P implementation with only one  
22 MPI message per process. Inspired by this fact, we try to design an adaptive approach that can  
23 combine the butterfly and P2P implementations, where some stages in the butterfly  
24 implementation are skipped ~~with the and replaced by~~ P2P ~~implementations~~ ~~communication~~ of  
25 more MPI messages per process. ~~If~~ ~~When~~ all stages of the butterfly implementation are skipped,  
26 the adaptive data transfer library ~~will switch~~ ~~completely~~ ~~switches~~ to the original P2P  
27 implementation. That is to say, the adaptive data transfer can adaptively choose the optimal  
28 implementation from the P2P implementation and the butterfly implementation. Figure 4312  
29 shows an example of the adaptive data transfer library with eight processes, where Stage 2 of  
30 the butterfly implementation is skipped ~~with the and replaced by~~ P2P  
31 ~~implementation~~ ~~communication~~ of three MPI messages per process.

1 The most significant challenge ~~to~~of such an adaptive approach is ~~how~~ to determine which  
2 stage(s) of the butterfly implementation should be skipped. The first attempt was to design a  
3 cost model that can accurately predict the performance of data transfer in various  
4 implementations. We eventually gave up this ~~because~~approach as it was almost impossible to  
5 accurately predict the performance of the communications on a high-performance computer,  
6 especially when a lot of users share the computer to run various applications. Performance  
7 profiling which means directly measuring the performance of data transfer is more practical to  
8 determine an appropriate implementation, because the simulation of Earthearth system  
9 modelling always takes a long time to run. Figure 1413 shows our flowchart of how the adaptive  
10 data transfer library determines an appropriate implementation. It consists of an initialization  
11 segment and a profiling segment. The initialization segment generates the process mappings  
12 and a candidate implementation that is a butterfly implementation with no skipped stages. The  
13 profiling segment iterates through each stage of the butterfly implementation to determine  
14 whether the current stage should be skipped or kept. In an iteration, the profiling segment first  
15 generates a temporary implementation based on the candidate implementation where the current  
16 stage is skipped, and then runs the temporary implementation to get the time the data transfer  
17 takes. When the temporary implementation is more efficient than the candidate implementation,  
18 the current stage is skipped and the temporary implementation replaces the candidate  
19 implementation. When the profiling segment finishes, the appropriate implementation is set to  
20 be the candidate implementation. To reduce the overhead introduced by the adaptive data  
21 transfer library, the profiling segment truly transfers the data for model coupling. In other words,  
22 before obtaining an appropriateoptimal implementation, the data is transferred by the profiling  
23 segment.

## 24 **5 Performance evaluation**

25 In this section, we empirically evaluate the adaptive data transfer library, through comparing it  
26 to the P2P implementation and the butterfly implementation ~~and the P2P implementation~~. Both  
27 toy models and realistic models (GAMIL2-CLM3 and CESM) are used for the performance  
28 evaluation. GAMIL2-CLM3 has been introduced in Section 2.2. CESM (Hurrell et al., 2013) is  
29 a state-of-the-art ESM developed by the National Center for Atmospheric Research (NCAR).  
30 All the experiments are run on the high performance computer Tansuo100.

1 Next, we will evaluate the overhead of initialization, the performance ~~of transferring~~ data  
2 ~~transfer~~fields between two different component models and the performance ~~of rearranging~~  
3 data ~~rearrangement~~fields intra a component model for parallel interpolation.

## 4 5.1 Overhead of initialization

5 We first evaluate the initialization overhead of data transfer implementations. As shown in Fig.  
6 ~~1514~~, the initialization overhead of each implementation increases with the increment of core  
7 number. The initialization overhead of the butterfly implementation is a little higher than that  
8 of the P2P implementation, while the initialization overhead of the adaptive data transfer library  
9 is 2-3 folds higher than that of the P2P implementation, because the adaptive data transfer  
10 library uses extra time on the performance profiling (~~please refer to see~~ Section 4). Considering  
11 that one data transfer instance should only be initialized at the beginning and executed many  
12 times in a coupled model, we can conclude that the initialization overhead of the adaptive data  
13 transfer library is reasonable, especially when the simulation is executed for a very long time.

## 14 5.2 Performance of data transfer between toy models

15 ~~As mentioned in Section 3, the butterfly implementation has different characterizations~~  
16 ~~compared to the P2P implementation. Many~~The factors ~~that~~ can impact the performance of a  
17 data transfer implementation ~~including MPI message number~~generally include the  
18 ~~communication depth~~, the size of ~~the~~ data to be transferred (also ~~known~~referred to as the number  
19 of fields in this evaluation) and the number of cores used. In this subsection, we evaluate the  
20 ~~impact of each factor on the~~ performance of data transfer ~~affected by each of these factors. for~~  
21 ~~different implementations.~~ We first build two toy models that ~~both~~ use the same logically  
22 rectangular grid ~~(of 192×96480 grid points).~~. Coupling fields are transferred between the two  
23 toy models. ~~In each~~For any test, the two toy models use the same number of cores, ~~and each~~  
24 ~~process has the same MPI message number.~~ Next, we evaluate the performance of data transfer  
25 through varying one factor ~~and~~while fixing the other ~~two~~ factors.

26 In the first experiment, we fix the number of cores to be ~~192~~1024 and the ~~number of~~ coupling  
27 ~~field number~~fields to be 10, ~~and~~while only vary ~~the communication depth in the P2P~~  
28 ~~implementation. In each test, all processes of the sender have the same communication depth.~~  
29 ~~As the communication depth is determined by the parallel decompositions of the sender and the~~  
30 ~~receiver, we design an algorithm (Algorithm 1) that can generate the parallel decompositions~~



1 of the two toy models according to the average communication depth of the sender in the P2P  
2 implementation. MPI message number per process. Figure 4615 shows the execution time of  
3 one data transfer with different implementations when ~~varying~~increasing the MPI message  
4 number communication depth per sender process in the P2P implementation from 1 to ~~96~~90.  
5 The P2P implementation can outperform the butterfly implementation when the MPI message  
6 number communication depth is small (say, smaller than 12 in Fig. 4615), while the butterfly  
7 implementation can outperform the P2P implementation when the MPI message  
8 number communication depth is big (say, bigger than 12 in Fig. 4615). The adaptive data  
9 transfer library ~~has the best performance. Moreover can~~ adaptively choose the optimal  
10 implementation from the P2P implementation and the butterfly implementation, and moreover,  
11 it improves the performance based on the butterfly implementation when the MPI message  
12 number communication depth is big, because some butterfly stages ~~in the adaptive data transfer~~  
13 ~~library have been~~ of the butterfly implementation are skipped ~~with the P2P implementation.~~  
14 When the MPI message number per process is 96 communication depth is 90, the adaptive data  
15 transfer library can achieve a ~~13.9~~19.2-fold performance speedup compared to the P2P  
16 implementation.

17 In the second experiment, we fix the number of cores and ~~MPI message number the~~  
18 communication depth per sender process in the P2P implementation, and vary the number of  
19 coupling ~~field number fields~~ transferred. Figure 4716 shows the execution time of one data  
20 transfer with different implementations in this experiment. The results show that the execution  
21 time of each implementation increases with the increment of data size. When MPI message  
22 number the communication depth per sender process in the P2P implementation is small (Figs.  
23 47a16a and 47b16b), the performance of the butterfly implementation is poorer than that of the  
24 P2P implementation, especially when the number of 2-D coupling fields gets bigger. ~~The~~  
25 ~~adaptive data transfer library achieves similar performance as~~ When the communication depth  
26 per sender process in the P2P implementation, ~~because it switches to the P2P implementation.~~  
27 ~~When the MPI message number per process is~~ big (Figs. 47e16c and 47d), ~~both 16d~~, the  
28 butterfly implementation ~~and adaptive data transfer library~~ significantly  
29 ~~outperform~~outperforms the P2P implementation, however, the advantage of the butterfly  
30 implementation decreases with the increment of the number of coupling fields. The results also  
31 demonstrate that the adaptive data transfer library can adaptively choose the optimal  
32 implementation from the P2P implementation and the butterfly implementation, and can further

1 ~~improve~~ the ~~adaptive data transfer library achieves better~~ performance ~~than based on~~ the  
2 butterfly implementation.

3 In the third experiment, we fix ~~MPI message number~~ the communication depth per sender  
4 process in the P2P implementation to be 24 and the number of coupling field number  
5 transferred to be 10, and vary the number of cores used. Figure ~~18~~17 shows the execution time  
6 of one data transfer with different implementations when varying the number of cores. The  
7 ~~results show that both~~ P2P implementation outperforms the butterfly implementation ~~and~~  
8 ~~adaptive data transfer library achieve better parallel scalability, when small number of cores are~~  
9 ~~used (say, smaller than the P2P implementation. The execution time of the P2P implementation~~  
10 ~~slightly increases with the increment of the number of cores used. However, the execution times~~  
11 ~~of the butterfly implementation and adaptive data transfer library slightly decrease with the~~  
12 ~~increment of the number of the cores used. The~~ 256 in Fig. 17); while the butterfly  
13 implementation outperforms the P2P implementation, ~~while~~ when large number of cores are  
14 used (say, larger than 256 in Fig.17). Similar to above two experiments, the adaptive data  
15 transfer library ~~achieves better performance than~~ can adaptively choose the optimal  
16 implementation from the P2P implementation and the butterfly implementation.

17 The ~~resolution~~ resolutions of models ~~becomes~~ become higher and higher these days. How about  
18 the performance of the data transfer implementations when model resolution  
19 ~~becomes~~ resolutions become higher? Higher model ~~resolution means~~ resolutions mean that a  
20 model will use more processor cores for accelerating a simulation, while the average number  
21 of grid points per processor core can remain constant. Considering that the numbers of grid  
22 points are always balanced among the processes of a model, we make each process (which runs  
23 on a unique processor core) of the toy models evenly have around 96 grid points in this  
24 evaluation, while enabling processes to have different ~~message numbers~~ communication depth  
25 and different message sizes: (the average communication depth of the sender in P2P  
26 implementation is 34). As shown in Fig. ~~19~~18, although the execution times of all data transfer  
27 implementations increase ~~with~~ when increasing the ~~increment~~ number of processor ~~core~~  
28 ~~number~~ cores (from 64 to 1024), ~~both~~ the butterfly implementation ~~and~~ significantly  
29 outperforms the P2P implementation. So the adaptive data transfer library significantly  
30 outperform the P2P implementation, and the adaptive data transfer library achieves the best  
31 performance. These results indicate that our proposed implementations can significantly  
32 improve adaptively chooses the ~~performance of data transfer for higher~~ butterfly implementation,

1 and further slightly outperforms the butterfly implementation when each model resolution uses  
2 more than 512 cores because some butterfly stages are skipped.

### 3 **5.3 Performance of data transfer between realistic models**

4 In this subsection, we evaluate the performance using two realistic models: GAMIL2-CLM3  
5 (horizontal resolution of  $2.8^\circ \times 2.8^\circ$ ) and CESM (resolution of  $1.9 \times 2.5_{\text{gx1v6}}$ ).

6 For CESM, we use the data transfer between the coupler CPL7 (Craig et al., 2012) and the land  
7 surface model CLM4 (Oleson et al., 2004), where 32 2-D coupling fields on the CLM4  
8 horizontal grid (the grid size is  $144 \times 96 = 13824$ ) are transferred. Figure 2019 shows the  
9 performance of one data transfer of different implementations when increasing the ~~process~~  
10 number of processes of both CPL7 and CLM4 from 6 to 192. When the ~~process~~-number of  
11 processes is small (say, smaller than 24 in Fig. 2019), the butterfly implementation is much  
12 poorer than the P2P implementation, ~~and. In this case,~~ the adaptive data transfer library achieves  
13 similar performance as chooses the P2P implementation because it switches to as the P2P optimal  
14 implementation. However, when the ~~process~~-number of processes gets bigger (say, larger than  
15 24 in Fig. 2019), the ~~adaptive data transfer library dramatically~~ butterfly implementation  
16 outperforms the P2P implementation with more speedup and also. In this case, the adaptive  
17 data transfer library based on the butterfly implementation skips some stages, so it  
18 outperforms the butterfly implementation. Figure 19 also shows that the  
19 butterfly implementaion and the adaptive transfer library seem to converge when increasing the  
20 number of cores per model. When each ~~component~~model uses 192 cores, the adaptive data  
21 transfer library is 4.01 times faster than the P2P implementation.

22 For GAMIL2-CLM3, we use the data transfer from CLM3 to GAMIL2 where 14 2-D coupling  
23 fields on the GAMIL2 horizontal grid (whose grid size is  $128 \times 60 = 7680$ ) are transferred. Figure  
24 2120 shows the execution time of one data transfer of each implementation when increasing  
25 the ~~process~~-number of processes of both GAMIL2 and CLM3 from 6 to 192. The results in Fig.  
26 2120 confirm that the adaptive data transfer library can constantly show adaptively choose the  
27 best performance optimal implementation from the P2P implmentation and the butterfly  
28 implementation. Compared to the P2P implementation, the adaptive data transfer library  
29 achieves an 11.68-fold performance speedup when the ~~process~~-number of processes is 96, but  
30 achieves a much lower speedup (only 3.48-fold) when the ~~process~~-number of processes is 192.

1 This is because the average MPI message number communication depth per process in the P2P  
2 implementation reduces from 32 to 18 when the number of process increases from 96 to 192.

### 3 **5.4 Performance of data rearrangement for interpolation**

4 Besides data transfer between different component models, there is another kind of data transfer  
5 in model coupling that rearranges data inside a model for parallel interpolation of fields between  
6 different grids. Here, we use the data rearrangement for the parallel interpolation from the  
7 atmosphere grid (whose grid size is  $144 \times 96 = 13824$ ) to the ocean grid (whose grid size is  
8  $320 \times 384 = 122880$ ) in the coupled model CESM for further evaluation. As mentioned  
9 aboveshown on Fig. 21, the P2P implemtation is sufficient for data rearrangement.  
10 However, implementation significantly outperforms the butterfly implementation is much  
11 poorer than the P2P implementation (Fig. 22). This is because the MPI message number is very  
12 small (for corresponding parallel decompositions for data rearrangement are always similar  
13 while similar parallel decompositions generally introduce small communication depth. For  
14 example, average MPI message number per process communication depth in the P2P  
15 implementation corresponding to Fig. 21 is only 6.49 when eachthe model uses 96 cores) for  
16 data rearrangement. On the other hand. In this case, the adaptive P2P implementation is chosen  
17 as the optimal implementation of the data transfer library, so the data transfer library achieves  
18 almost the same performance as the P2P implementation, because it switcheslibrary does not  
19 provide real benefit compared to the P2P implementation. Therefore, the adaptive data transfer  
20 library can always show the best performance.

### 21 **5.5 Performance impevementimprovement for a coupled model**

22 With the performance improvement of data transfer, we expect that the adaptive data transfer  
23 library will improve the performance of coupled models. For this evaluation, we first  
24 importimported the adaptive data transfer library into C-Coupler1 and then useused the coupled  
25 model GAMIL2-CLM3 that uses C-Coupler1 for coupling to measure performance results. As  
26 shown in Fig. 2322, the adaptive data transfer library achieves higher performance  
27 improvement (when the P2P implementation is used as the baseline) for GAMIL2-CLM3 when  
28 using more processor cores. When each component model uses 128 processor cores, the  
29 butterfly implementation achieves ~4.6% performance improvement, and the adaptive data  
30 transfer library achieves ~7% performance improvement. This6.9% performance improvement  
31 would not be low because the model coupling only takes a very small proportion of execution

1 ~~time in. So the simple coupled model GAMIL2-CLM3 and the parallel scalability data transfer~~  
2 ~~library can improve the performance of data transfer, and then improve the performance~~ of the  
3 ~~two whole coupled models GAMIL2 and CLM3 is not good.~~

## 5 **6 Conclusions**

6 Data transfer is ~~the~~ a fundamental and most frequently used operation in a coupler. This paper  
7 ~~demonstrated showed~~ that the ~~current~~ P2P implementation ~~of data transfer currently used~~ in  
8 most state-of-the-art couplers for data transfer is inefficient ~~for transferring data between two~~  
9 ~~component models. To improve when~~ the parallel decompositions of the sender and the receiver  
10 are different, and further revealed the corresponding performance ~~of data transfer~~ bottlenecks.  
11 To overcome these bottlenecks, we proposed a butterfly implementation. ~~However, compared~~  
12 ~~to the P2P implementation, the butterfly implementation has both advantages and disadvantages.~~  
13 ~~The evaluation results showed that the the~~ butterfly implementation ~~did not always that can~~  
14 outperform the P2P implementation. ~~To achieve better parallel implementation in many cases,~~  
15 however, degrades the performance in some cases, for example, when a small number of cores  
16 are used to run models or the parallel decompositions of data transfer, we built the sender and  
17 receiver are similar. We therefore further designed and implemented an adaptive data transfer  
18 library, which combines the advantages of both butterfly implementation and P2P  
19 implementation. The evaluation results demonstrated that can not only adaptively choose an  
20 optimal implementation from the P2P implementation and the butterfly implemtation, but also  
21 further improve the performance based on the butterfly implementation through skipping some  
22 butterfly stages. Compared to the P2P implementation, the adaptive data transfer library can  
23 significantly improve the performance of data transfer ~~so as to improve a coupled model when~~  
24 the parallel decompositions of the sender and the receiver are different.

25 The initialization overhead for the adaptive data transfer library could become expensive when  
26 using a large number of processor cores. In the future version, the adaptive data transfer will  
27 allow users to record the results of performance profiling offline to save the time used for  
28 performance profiling in next runs of the same coupled model.

## 29 **Code availability**

30 The source code of the adaptive data transfer library version 1.0 is available at  
31 [https://github.com/zhang-cheng09/Data\\_transfer\\_lib](https://github.com/zhang-cheng09/Data_transfer_lib).

1 **Acknowledgements**

2 This work is supported in part by the Natural Science Foundation of China (no. 41275098), the  
3 National Grand Fundamental Research 973 Program of China (no. 2013CB956603) and the  
4 Tsinghua University Initiative Scientific Research Program (no. 20131089356).

5

## 1 **References**

- 2 Armstrong, C. W., Ford, R. W., and Riley, G. D.: Coupling integrated Earth System Model  
3 components with BFG2, *Concurrency and Computation: Practice and Experience*,  
4 2009;21;767–791, doi:10.1002/cpe.1348, 2009.
- 5 Balaji, V., Anderson, J., Held, I., Winton, M., Durachta, J., Malyshev, S., and Stouffer, R. J.:  
6 The Exchange Grid: a mechanism for data exchange between Earth system components on  
7 independent grids, In *Parallel Computational Fluid Dynamics 2005 Theory and Applications*,  
8 2006, 179-186, doi: 10.1016/B978-044452206-1/50021-5, 2006.
- 9 Chong, F. T., and Brewer, E. A.: Packaging and multiplexing of hierarchical scalable expanders,  
10 *Parallel Computer Routing and Communication*, Springer Berlin Heidelberg, 1994:200-214.
- 11 Craig, A. P., Vertenstein, M., and Jacob, R.: A new flexible coupler for Earth system modelling  
12 developed for CCSM4 and CESM1, *Int. J. High Perform. C.*, 26, 31-42,  
13 doi:10.1177/1094342011428141, 2012.
- 14 Craig, A. P., Jacob, R., Kauffman, B., Bettge, T., Larson, J., Ong, E., Ding, C., and He, Y.:  
15 CPL6: the New Extensible, High Performance Parallel Coupler for the Community Climate  
16 System Model, *Int. J. High Perform. C.*, 19, 309–327, 2005.
- 17 Dennis, J. M.: Inverse space-filling curve partitioning of a global ocean model, In *IEEE*  
18 *International Parallel & Distributed Processing Symposium*, Long Beach, CA, 2007.
- 19 Dennis, J. M. and Tufo, H. M.: Scaling climate simulation applications on the IBM Blue Gene/L  
20 system, *IBM J. Res. Dev.*, 52, 117-126, DOI:10.1147/rd.521.0117, 2008.
- 21 Dennis, J. M., Edwards, J., Evans, K. J., Guba, O., Lauritzen, P. H., Mirin, A. A., St-Cyr, A.,  
22 Taylor, M. A., and Worley, P. H.: CAM-SE: a scalable spectral element dynamical core for the  
23 Community Atmosphere Model, *Int. J. High Perform. C.*, 26, 74-89,  
24 doi:10.1177/1094342011428142, 2012.
- 25 Dickinson, R. E., Oleson, K. W., Bonan, G., Hoffman, F., Thornton, P., Vertenstein, M., Yang,  
26 Z.-L., and Zeng X.: The Community Land surface model and its climate statistics as a  
27 component of the Community Climate System Model, *Journal of Climate*, 19(11), 2302–2324,  
28 2006.

1 Ford, R. W., Riley, G. D., Bane, M. K., Armstrong, C. W., and Freeman, T. L.: GCF: a general  
2 coupling framework, *Concurrency and Computation: Practice and Experience*, 18(2), 163–181,  
3 2006.

4 Foster I.: *Designing and building parallel programs: concepts and tools for parallel software*  
5 *engineering*, Addison-Wesley, 1995.

6 Heckbert P.: *Fourier Transforms and the Fast Fourier Transform (FFT) Algorithm*, *Computer*  
7 *Graphics*, 2: 15-463, 1995.

8 Hemmert, K. S., and K. D. Underwood.: *An analysis of the double-precision floating-point FFT*  
9 *on FPGAs*. *Field-Programmable Custom Computing Machines*, 2005. FCCM 2005. 13th  
10 *Annual IEEE Symposium on IEEE*, 2005:171-180.

11 Hill, C., DeLuca, C., Balaji, V., Suarez, M., and da Silva, A.: *The Architecture of the Earth*  
12 *System Modelling Framework*, *Computing in Science & Engineering*, 6(1), 18–28, 2004.

13 Hurrell, J. W., Holland, M. M., Gent, P. R., Ghan, S., Kay, J. E., Kushner, P. J., Lamarque, J.-  
14 F., Large, W. G., Lawrence, D., Lindsay, K., Lipscomb, W. H., Long, M. C., Mahowald, N.,  
15 Marsh, D. R., Neale, R. B., Rasch, P., Vavrus, S., Vertenstein, M., Bader, D., Collins, W. D.,  
16 Hack, J. J., Kiehl, J., and Marshall, S.: *The Community Earth System Model: a framework for*  
17 *collaborative research*, *Bulletin of the American Meteorological Society*, 94(9), 1339–1360,  
18 2013.

19 Hunke, E. C. and Lipscomb W. H.: *CICE: the Los Alamos Sea Ice Model Documentation and*  
20 *Software User’s Manual 4.0*, Technical Report LA-CC-06-012, Los Alamos National  
21 Laboratory, T-3 Fluid Dynamics Group, 2008.

22 Hunke, E. C., Lipscomb, W. H., Turner, A. K., Jeffery, N., and Elliott, S.: *CICE: the Los*  
23 *Alamos Sea Ice Model Documentation and Software User’s Manual Version 5.0*, LA-CC-06-  
24 012, Los Alamos National Laboratory, Los Alamos NM, 87545, 115, 2013.

25 Jacob, R., Larson, J., and Ong, E.: *M × N Communication and Parallel Interpolation in*  
26 *Community Climate System Model version 3 using the Model Coupling Toolkit*, *International*  
27 *Journal of High Performance Computing Applications*, 19(3), 293–307, 2005.

28 Jan, B., Montrucchio, B., Ragusa, C., Khan, F. G., and Khan, O.: *Parallel butterfly sorting*  
29 *algorithm on gpu*, Acta Press, 2013.



1 Kerbyson, D. J., and Jones, P. W.: A performance model of the parallel ocean program,  
2 International Journal of High Performance Computing Applications, 19(3), 261-276,  
3 doi:10.1177/1094342005056114, 2005.

4 Kim J., Dally W. J., and Abts D.: Flattened butterfly: A cost-efficient topology for high-radix  
5 networks, ISCA, 2007, 35(2):126-137.

6 Li, L. J., Wang, B., Dong, L., Liu, L., Shen, S., Hu, N., Sun, W., Wang, Y., Huang, W., Shi, X.,  
7 Pu, Y., G. and Yang.: Evaluation of Grid-point Atmospheric Model of IAP LASG version 2  
8 (GAMIL2), Advances in Atmospheric Sciences, 30, 855–867, doi:10.1007/s00376-013-2157-  
9 5, 2013.

10 Liu, L., Yang, G., Wang, B., Zhang, C., Li, R., Zhang, Z., Ji, Y., and Wang, L.: C-Coupler1: a  
11 Chinese community coupler for Earth system modeling, Geoscientific Model Development,  
12 7(5), 2281-2302, doi:10.5194/gmd-7-2281-2014, 2014.

13 Morrison, H., and A. Gettelman: A new two-moment bulk stratiform cloud microphysics  
14 scheme in the Community Atmosphere Model, version 3 (CAM3). Part I: Description and  
15 numerical tests, Journal of Climate, 21(15), 3642–3659, doi:10.1175/2008JCLI2105.1, 2008.

16 Neale, R. B., Richter, J. H., Conley, A. J., Park, S., Lauritzen, P. H., Gettelman, A., Williamson,  
17 D. L., Rasch, P. J., Vavrus, S. J., Taylor, M. A., Collins, W. D., Zhang, M., and Lin, S.:  
18 Description of the NCAR Community Atmosphere Model (CAM 4.0), National Center for  
19 Atmospheric Research Near Koha Opencat, TN-485+STR, 222p., 2010.

20 Neale, R. B., Chen, C. C., Gettelman, A., Lauritzen, P. H., Park, S., Williamson, D. L., Conley,  
21 A. J., Garcia, R., Kinnison, D., Lamarque, J. F., Marsh, D., Mills, M., Smith, A. K., Tilmes, S.,  
22 Vitt, F., Morrison, H., Cameron-Smith, P., Collins, W. D., Iacono, M. J., Easter, R. C., Ghan,  
23 S. J., Liu, X., Rasch, P. J., and Taylor, M. A.: Description of the NCAR Community  
24 Atmosphere Model (CAM 5.0), National Center for Atmospheric Research Near Koha  
25 Opencat, TN-486+STR, 289p., 2012

26 Oleson, K. W., Dai, Y., Bonan, G., Bosilovich, M., Dickinson, R., Dirmeyer, P., Hoffman, F.,  
27 Houser, P., Levis, S., Niu, G. Y., Thornton, P., Vertenstein, M., Yang, Z. L., and Zeng, X.:  
28 Technical Description of the Community Land Surface Model (CLM), National Center for  
29 Atmospheric Research Near Koha Opencat, TN-461+STR, 186p., 2004.

1 Petagon, R., and Werapun, J.: Embedding the optimal all-to-all personalized exchange on  
2 multistage interconnection networks + + mathContainer Loading Mathjax, *Journal of Parallel  
3 & Distributed Computing* 88(2016):16-30.

4 Redler, R., Valcke, S., and Ritzdorf, H.: OASIS4—a coupling software for next generation Earth  
5 System Modelling, *Geoscientific Model Development*, 3(1), 87–104, doi:10.5194/gmd-3-87-  
6 2010, 2010.

7 Smith, R., Jones, P., Briegleb, B., Bryan, F., Danabasoglu, G., Dennis, J., Dukowicz, J., Eden,  
8 C., Fox-Kemper, B., Gent, P., Hecht, M., Jayne, S., Jochum, M., Large, W., Lindsay, K.,  
9 Maltrud, M., Norton, N., Peacock, S., Vertenstein, M., and Yeager, S.: The Parallel Ocean  
10 Program (POP) reference manual ocean component of the Community Climate System Model  
11 (CCSM) and Community Earth System Model (CESM), Los Alamos National Laboratory,  
12 LAUR-10-01853, available at  
13 <http://www.cesm.ucar.edu/models/cesm1.1/pop2/doc/sci/POPRefManual.pdf> (last access: 15  
14 October 2015), 141 p., 2010.

15 Valcke, S., Balaji, V., Craig, A., DeLuca, C., Dunlap, R., Ford, R. W., Jacob, R., Larson, J.,  
16 O’Kuinghttons, R., Riley, G. D., and Vertenstein, M.: Coupling technologies for Earth System  
17 Modelling, *Geoscientific Model Development*, 5(6), 1589–1596, doi:10.5194/gmd-5-1589-  
18 2012, 2012.

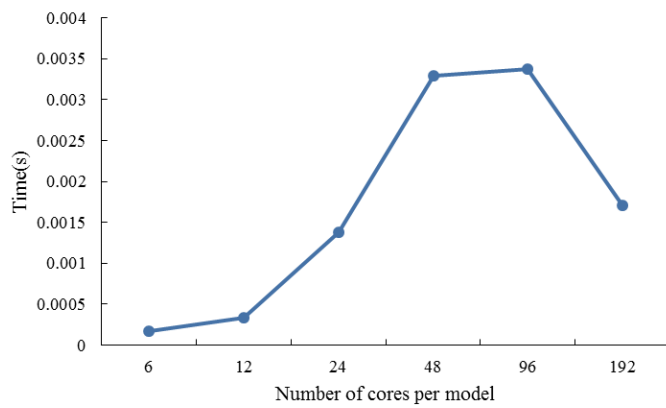
19 Valcke, S.: The OASIS3 coupler: a European climate modelling community software,  
20 *Geoscientific Model Development*, 6(2), 373–388, doi:10.5194/gmd-6-373-2013, 2013.

21 Valcke, S., Craig, T., and Coquart, L.: The OASIS3-MCT parallel coupler, in: The Second  
22 Workshop on Coupling Technologies for Earth System Models (CW2013), available at:  
23 [https://wiki.cc.gatech.edu/CW2013/images/a/a0/OASIS\\_MCT\\_abstract.pdf](https://wiki.cc.gatech.edu/CW2013/images/a/a0/OASIS_MCT_abstract.pdf) (last access: 15  
24 October 2015), 2013.

25 [Valcke, S., Craig, T. and Coquart, L.: OASIS3-MCT User Guide, OASIS3-MCT 3.0,](#)  
26 [Technical Report TR/CMGC/15/38, Cerfacs, France, 2015](#)  
27

Algorithm 1. Generating the parallel decompositions of the sender and the receiver according to an average communication depth of the sender in the P2P implementation.

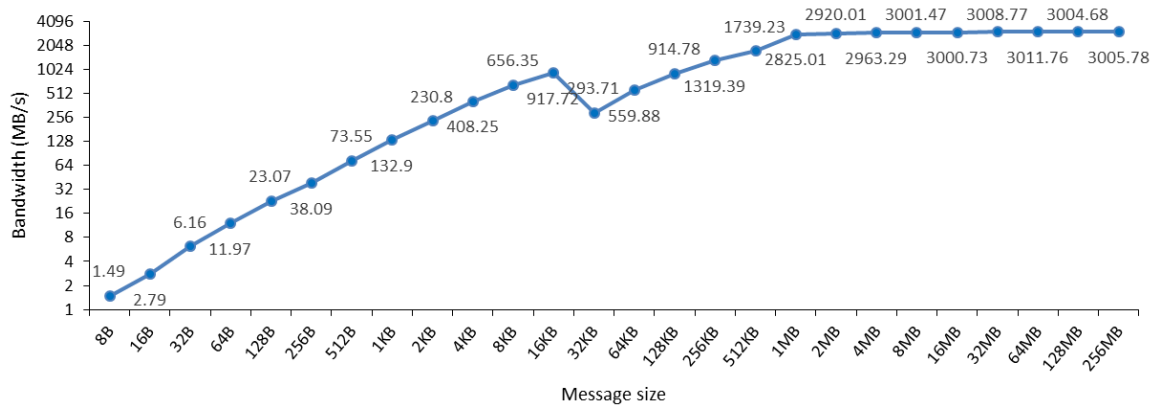
<u>Input</u>	<u>Number of processes of the sender: <math>M</math></u> <u>Number of processes of the receiver: <math>N</math></u> <u>Number of points in the grid: <math>Grid\_pnts</math></u> <u>Average communication depth per process of the sender in the P2P implementation: <math>Avg\_send\_depth, Avg\_send\_depth \leq N</math></u> <u>The flag that specifies whether the communication depths among processes are the same: <math>Is\_balanced</math></u>
<u>Output</u>	<u>Parallel decomposition of the sender</u> <u>Parallel decomposition of the receiver</u>
<u>1</u>	<u>Determine the parallel decomposition of the sender</u> <u>Considering that the numbers of grid points are always balanced among the processes of a model, assign around <math>Grid\_pnts/M</math> grid points to each process of the sender.</u>
<u>2</u>	<u>Determine the communication depth of each process of the sender</u>
<u>2.1</u>	<u>If the flag <math>Is\_balanced</math> is set to true, set the communication depth of each process of the sender to be <math>Avg\_send\_depth</math>;</u>
<u>2.2</u>	<u>Otherwise, randomly determine the communication depth of each process of the sender</u>
<u>2.2.1</u>	<u>Initialize the communication depth of each process of the sender to be 1</u>
<u>2.2.2</u>	<u>Randomly select a process of the sender whose communication depth does not exceed <math>N</math> and <math>Grid\_pnts/M</math>, and then increase its communication depth by 1, until the average communication depth of all processes of the sender reaches <math>Avg\_send\_depth</math>.</u>
<u>3</u>	<u>Determine the grid points of each communication</u> <u>For each process of the sender, assign the corresponding grid points to all communications of this process (a grid point belongs to only one communication)</u>
<u>3.1</u>	<u>If the flag <math>Is\_balanced</math> is set to true, assign the grid points to all communications evenly.</u>
<u>3.2</u>	<u>Otherwise, assign the grid points to each communication randomly</u>
<u>3.2.1</u>	<u>Assign one grid point to each communication</u>
<u>3.2.2</u>	<u>For each of remaining grid points, randomly select a communication for it</u>
<u>4</u>	<u>Determine the parallel decomposition of the receiver through assigning the grid points in each communication to a process of the receiver</u> <u>For each process of the sender, assign the grid points in each communication of it to a distinct receiver process: to make the numbers of grid points balance among the processes of the receiver in the final parallel decomposition, a communication with bigger number of grid points will be assigned to a receiver process with smaller total number of grid points that have been assigned to it.</u>



1

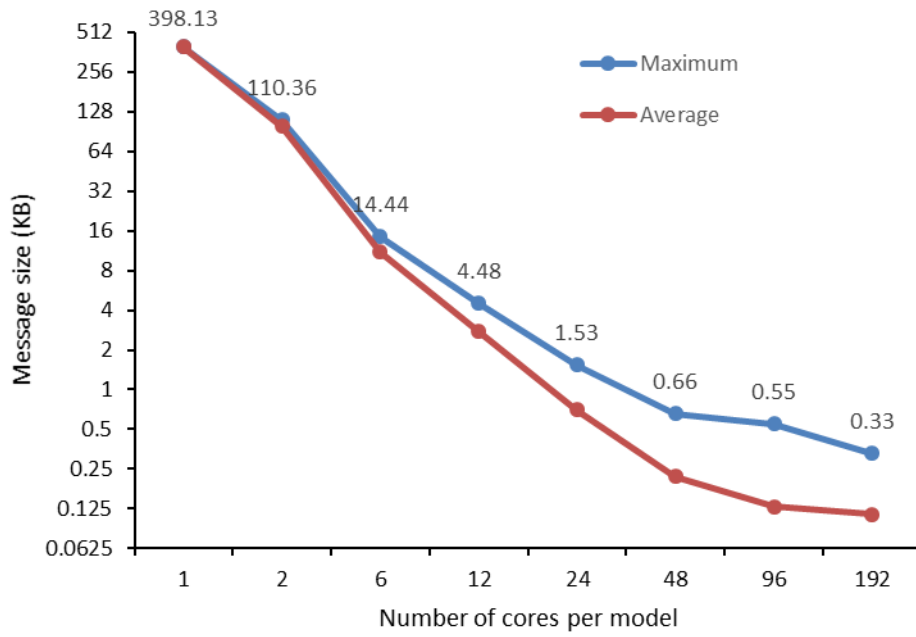
2 Figure 1. Average execution time of the P2P implementation when transferring 14 2-D fields  
 3 from CLM3 to GAMIL2. In each test, the atmosphere model GAMIL2 and the land surface  
 4 model CLM3 use the same number of cores; they do not share the same computing ~~nodes~~nodes.  
 5 The horizontal grid of the 14 2-D fields contains 7680 (128×60) grid points.

6



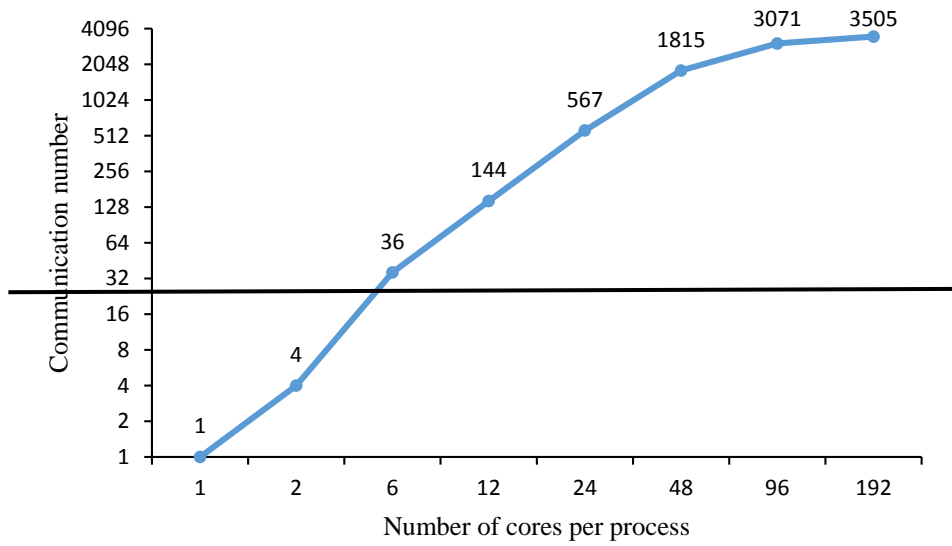
1  
2  
3  
4  
5  
6

Figure 2. Variation of bandwidth (y-axis) of an MPI P2P communication with respect to the ~~increment of~~ message size (x-axis). The results are generated from our benchmark. In the benchmark, one process sends messages with different sizes to the other process. The two processes of the P2P communication run on two different computing nodes of Tansuo100.

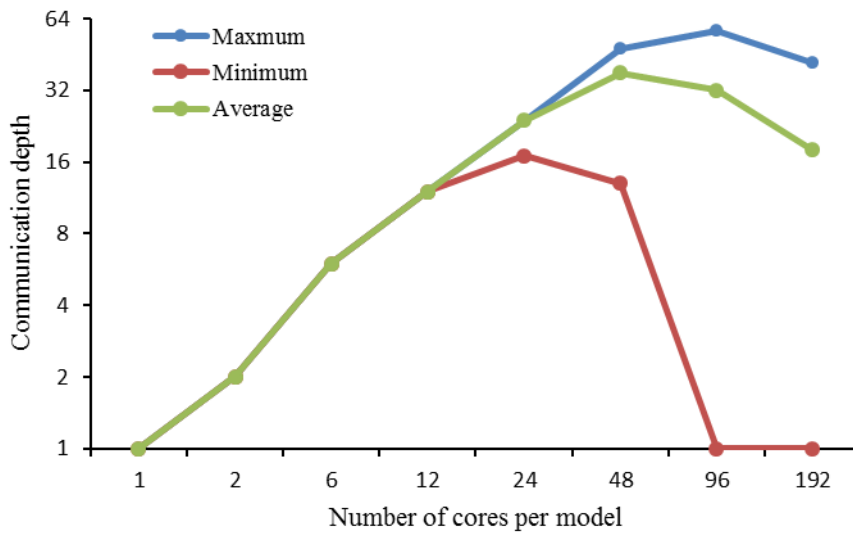


1  
2  
3  
4  
5

Figure 3. Variation of message size of the P2P implementation (y-axis) in GAMIL2-CLM3 with respect to the ~~increment of core~~ number of cores per model (x-axis). The experimental setup is similar to that shown in Fig. 1.



1

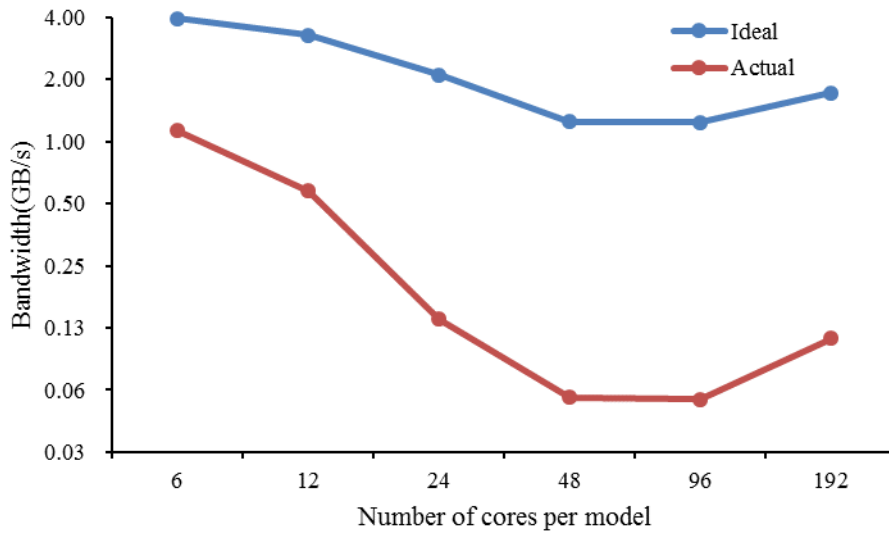


2



3 Figure 4. Variation of ~~total MPI message number~~ the communication depth of one process (y-  
 4 axis) ~~of using~~ the P2P implementation in GAMIL2-CLM3 with respect to the increment number  
 5 of ~~core number~~ cores per model (x-axis). The experimental setup is similar to that shown in Fig.

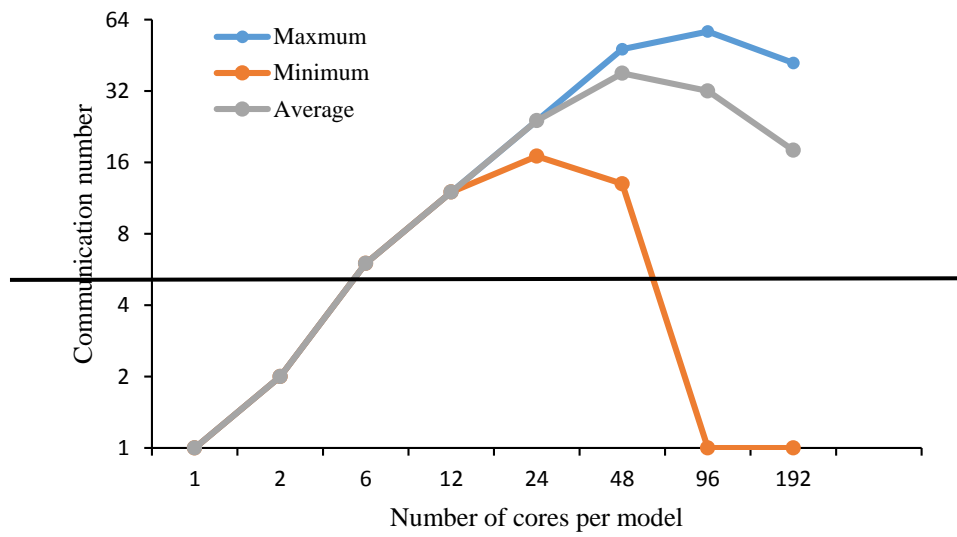
6 1.  
 7



1  
 2 Figure 5. Ideal and actual bandwidths of the P2P implementation (y-axis) in GAMIL2-CLM3  
 3 when gradually increasing the number of cores used by each component model (x-axis). The  
 4 experimental setup is similar to that shown in Fig. 1. The ideal bandwidth is calculated from  
 5 the message size and the MPI bandwidth measured in Fig. 2; and the actual bandwidth is  
 6 calculated from Fig. 1.

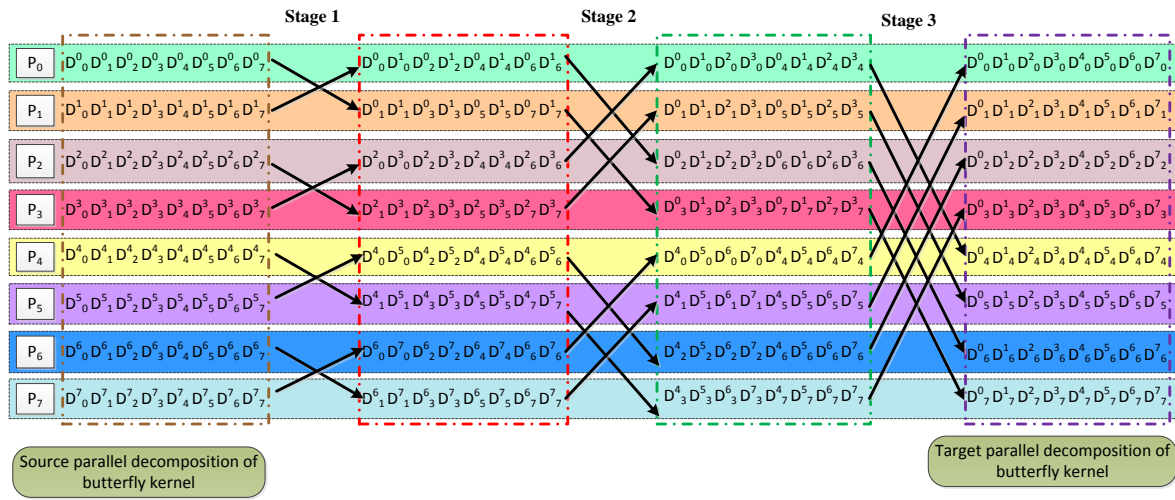
7





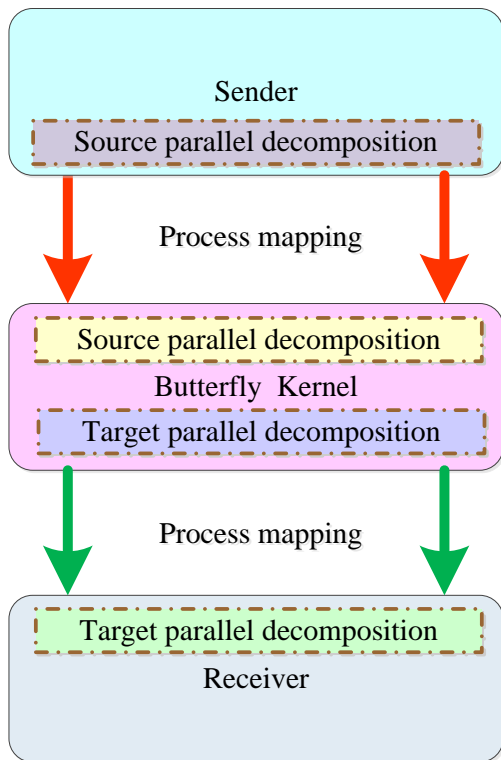
1  
2  
3  
4  
5  
6

Figure 6. Variation of message number of one process (y axis) using the P2P implementation in GAMIL2-CLM3 with respect to the increment of core number (x axis). The experimental setup is similar to that shown in Fig. 1.



1  
2  
3  
4  
5  
6  
7  
8

Figure 76. An example of the butterfly kernel with eight processes. Each colored row stands for one process ( $P_0$ - $P_7$ ). There are multiple stages (each column of arrows represents a stage (Stage 1 to Stage 3)) in the butterfly kernel. Each arrow stands for an MPI P2P communication from one process to another.  $D^i_j$  means the data is originally in process  $P_i$  according to the source parallel decomposition and is finally in process  $P_j$  according to the target parallel decomposition.

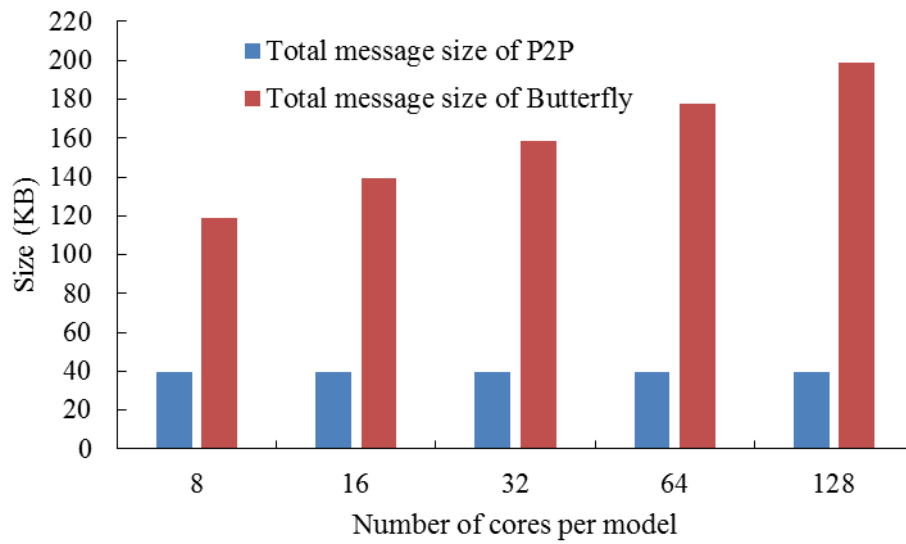


1

2 Figure 87. The butterfly implementation, which is composed of three parts: the butterfly kernel;  
 3 process mapping from the sender to the butterfly kernel; and process mapping from the butterfly  
 4 kernel to the receiver.

5

6

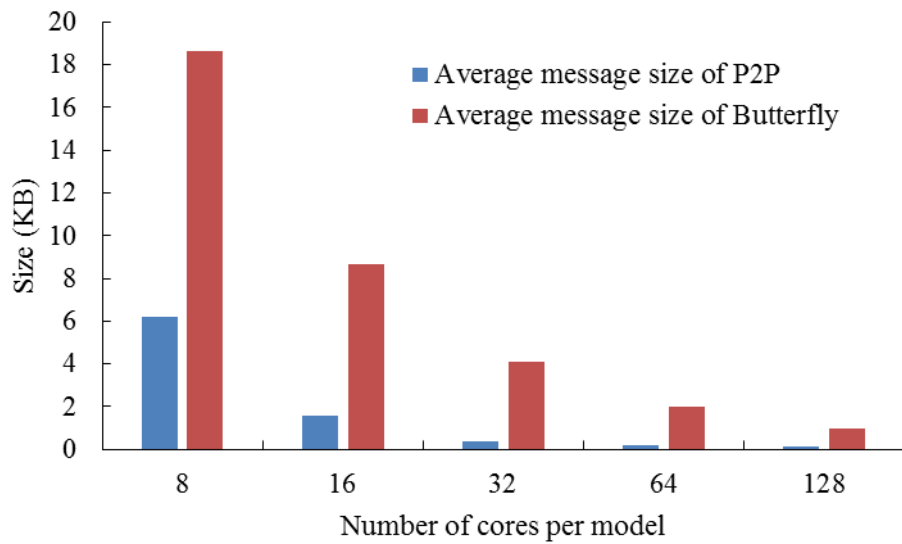


1

2 Figure 98. Total message size transferred by P2P implementation and butterfly implementation  
 3 (y-axis) in GAMIL2-CLM3, when varying the number of cores used by each model (x-axis).

4 The experimental setup is similar to that shown in Fig. 1.

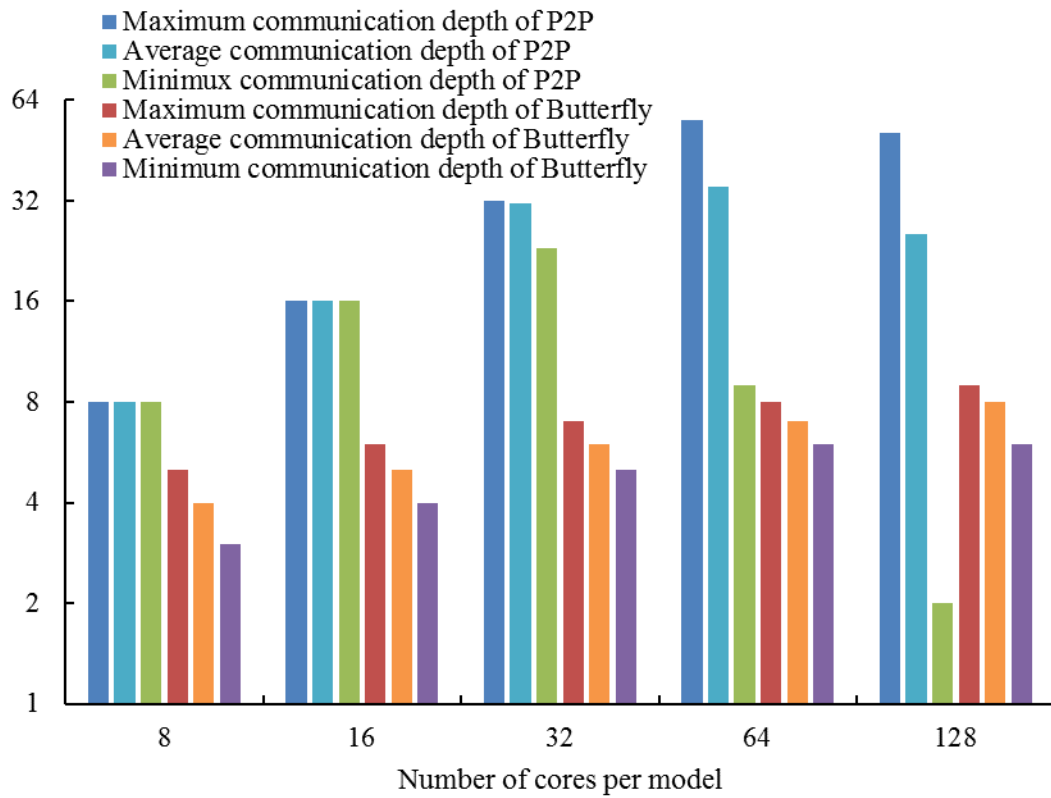
5



1

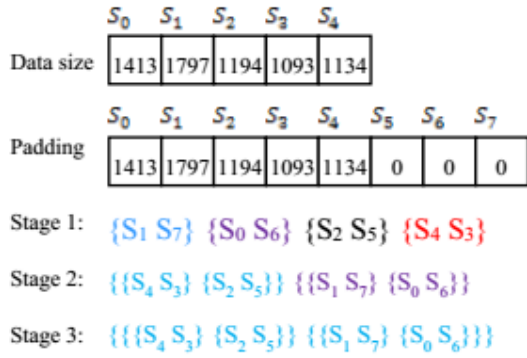
2 Figure 109. Average message size transferred by P2P implementation and butterfly  
 3 implementation (y-axis) in GAMIL2-CLM3, when varying the number of cores used by each  
 4 model (x-axis). The experimental setup is similar to that shown in Fig. 1.

5

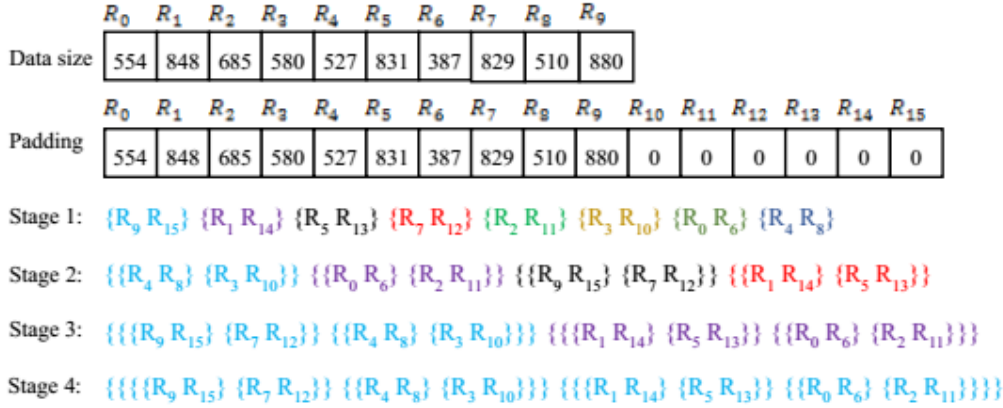


1  
2  
3  
4  
5  
6  
7

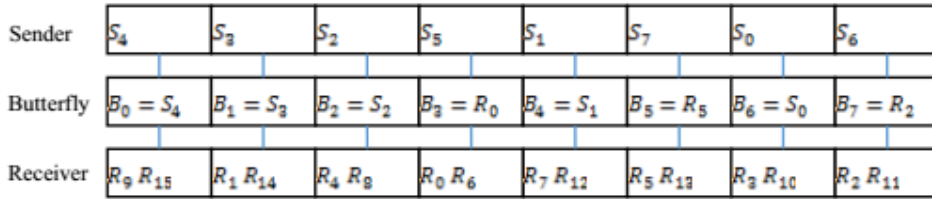
Figure 4.10. Maximum ~~message number~~communication depth, average ~~message number~~communication depth and minimum ~~message number of processes~~communication depth in P2P implementation and butterfly implementation (y-axis), when varying the number of cores used by each model (x-axis) in GAMIL2-CLM3. The experimental setup is similar to that shown in Fig. 1.



(a)



(b)



(c)

1

2 Figure 12.11. An example of process mappings, given that the sender has five processes ( $S_0$ -3  $S_4$ ), the receiver has 10 processes ( $R_0$ - $R_9$ ) (there is no common process between the sender and4 receiver), and the butterfly kernel contains eight processes ( $B_0$ - $B_7$ ). Panels (a) and (b) show

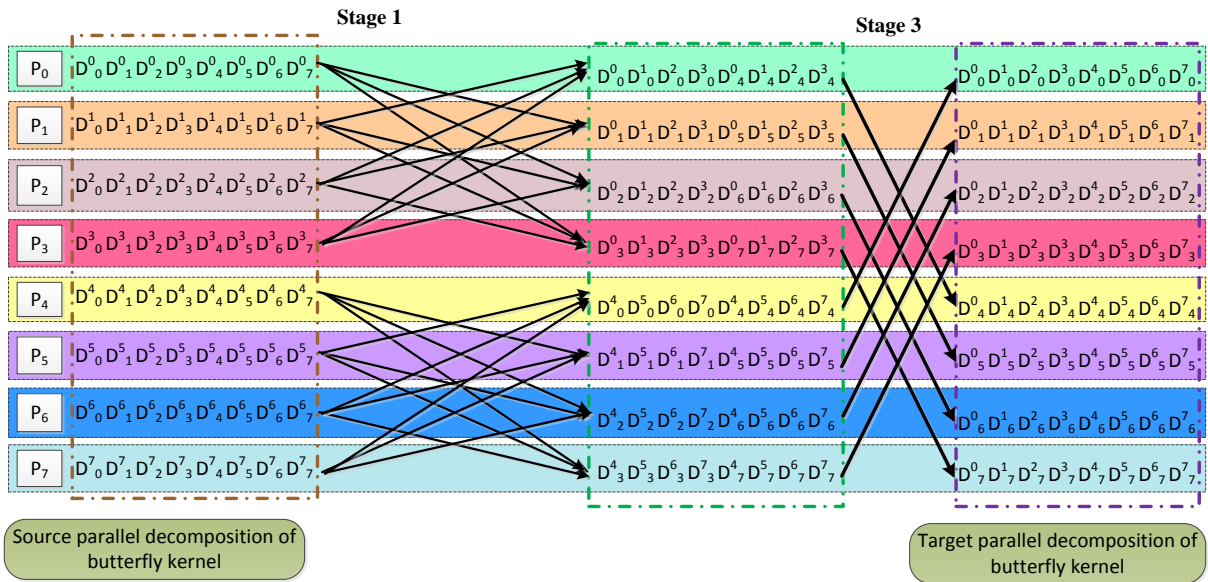
5 how to iteratively pair processes of the sender and receiver, respectively. There are

6 multiple stages in the iterative pairing of processes of the sender and receiver. In each stage,

7 the processes in the same color are grouped into one process pair. Panel (c) shows how to map

8 the reordered processes of the sender and receiver onto the processes of the butterfly kernel.

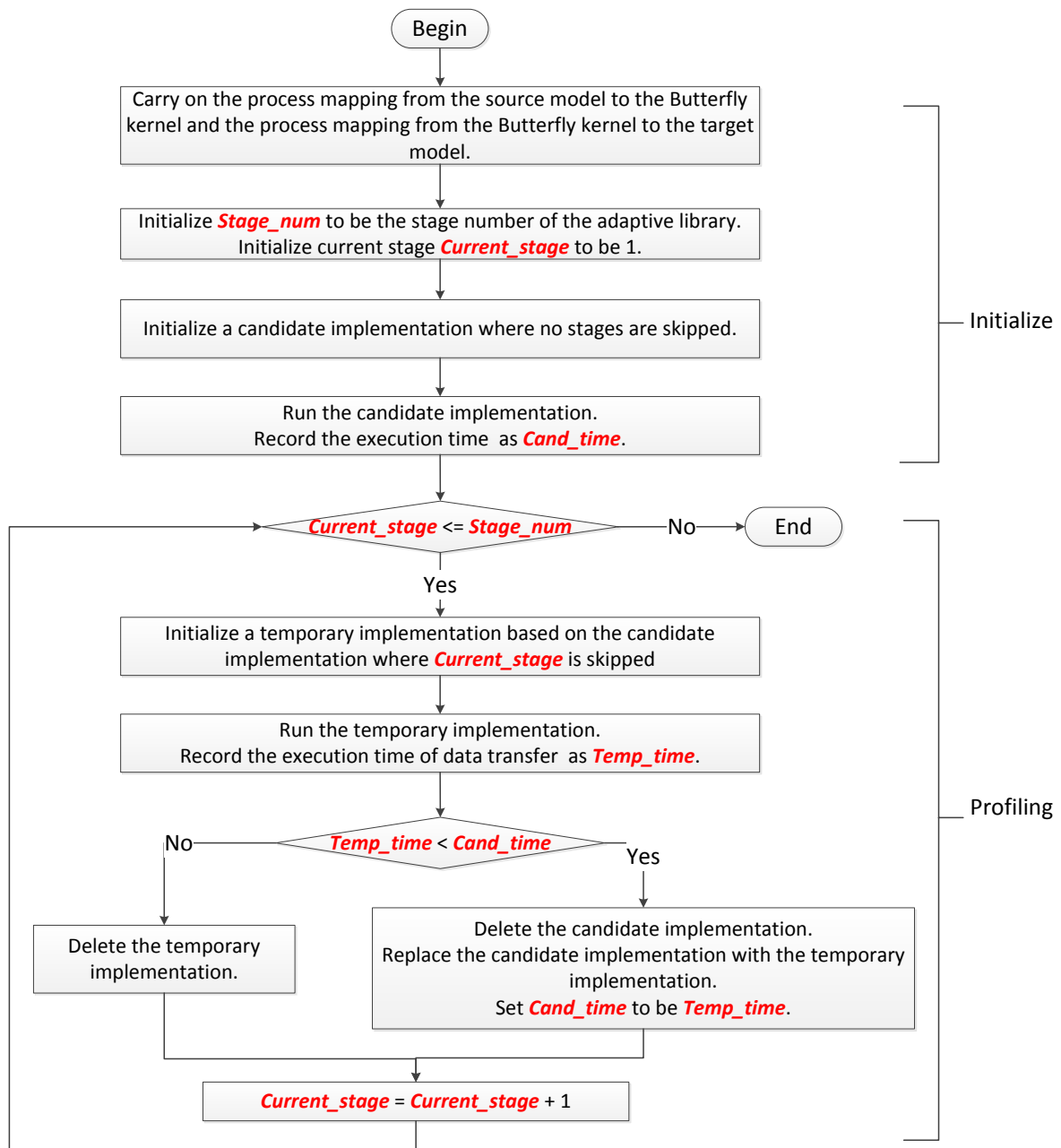
9 ~~All five processes of the sender and three processes of the receiver are used as the processes~~10 ~~of the butterfly kernel. Each process of the sender is mapped onto a process of the butterfly~~11 ~~kernel, while every two processes of the receiver are mapped onto one process of the butterfly~~12 ~~kernel.~~



1  
2  
3  
4  
5

Figure 13.12. An example of the adaptive data transfer library with eight processes, where Stage 2 of the butterfly implementation is skipped ~~with the~~ and replaced by P2P implementation communication of three MPI messages per process.

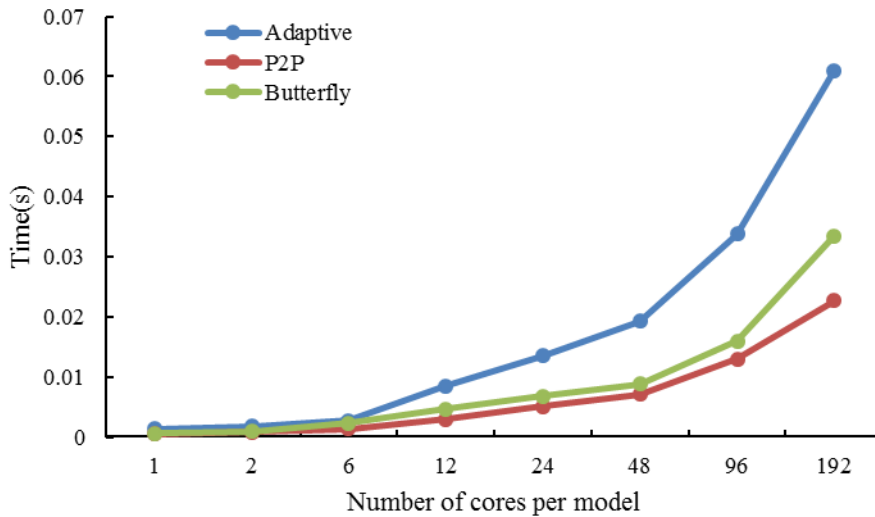




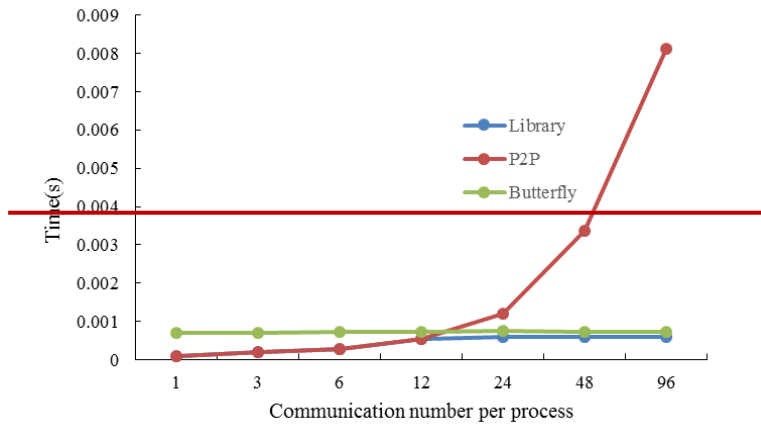
1

2 Figure 1413. A flowchart for determining an appropriate implementation of the adaptive data  
 3 transfer library.

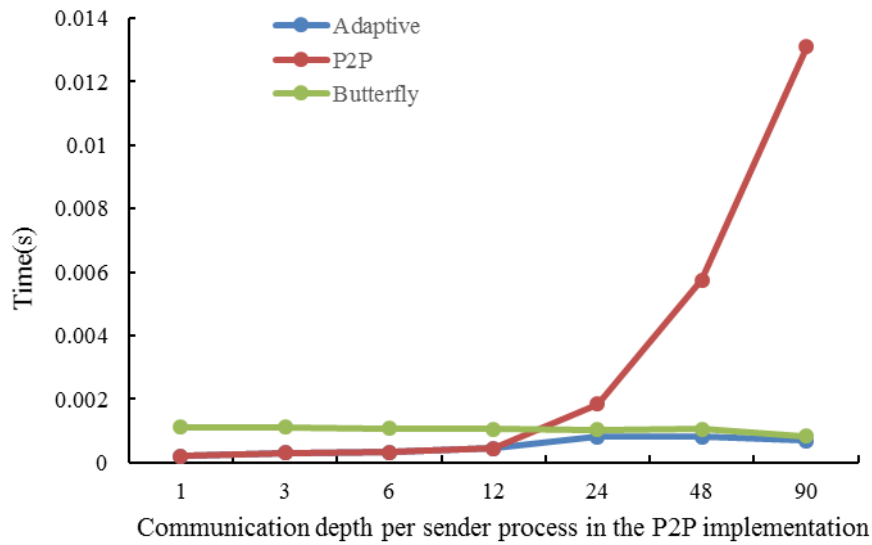
4



1  
2 Figure 4514. Initialization time (y-axis) of one data transfer between two toy models using a  
3 rectangular grid (of  $192 \times 96$  grid points) when varying the number of cores used by each toy  
4 model (x-axis). There are 10 2-D coupling fields transferred from the source toy model to the  
5 target toy model. In each test, all processes of the sender in the P2P implementation have the  
6 same communication depth. If the number of cores per toy-model used is less than 24, the MPI  
7 message number communication depth per sender process in the P2P implementation is set equal  
8 to be the number of cores. Otherwise, the MPI message number per model; otherwise, the  
9 communication depth per sender process in the P2P implementation is set to 24. The parallel  
10 decompositions of the sender and the receiver for a given setting of communication depth are  
11 generated by Algorithm 1.  
12



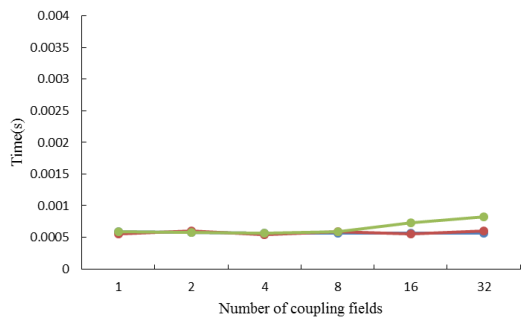
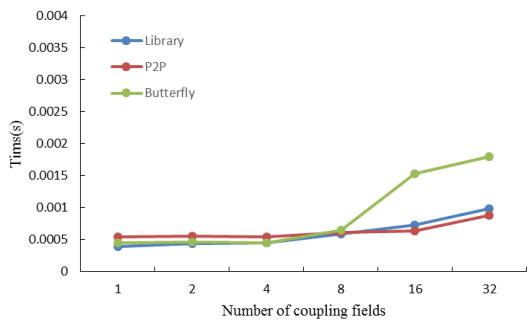
1



2

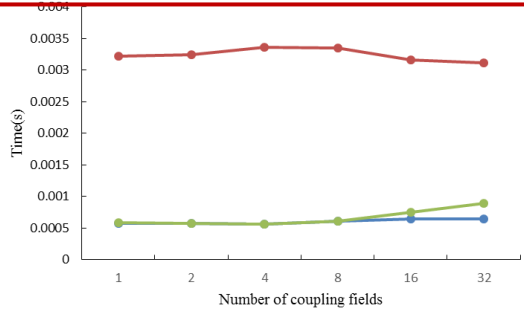
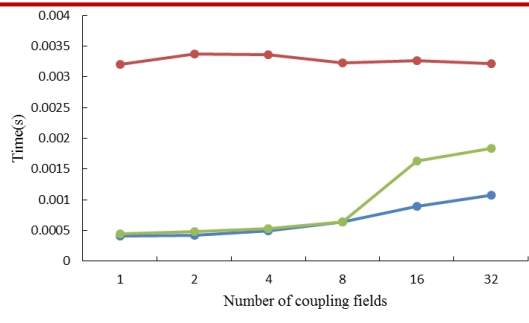
3 Figure 1615. Average execution time (y-axis) of one data transfer between two toy models with  
 4 the same rectangular grid (of  $192 \times 96480$  grid points) when varying the MPI message  
 5 number/communication depth per sender process in the P2P implementation (x-axis). Each toy  
 6 model is run with  $1921024$  cores. There are 10 2-D coupling fields transferred from the source  
 7 toy model to the target toy model.

8



(a)

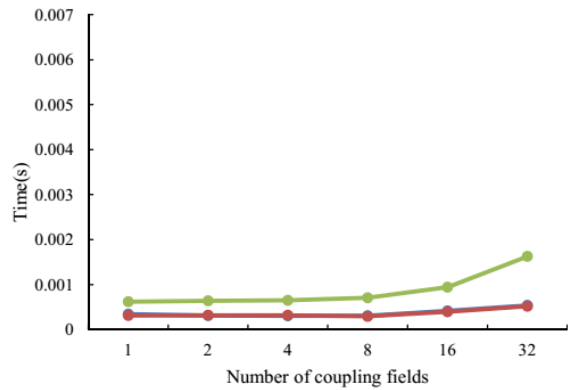
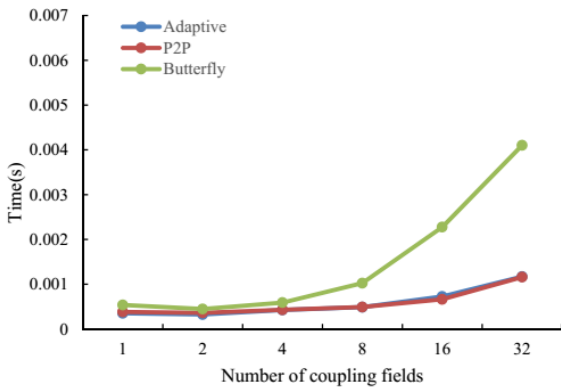
(b)



(c)

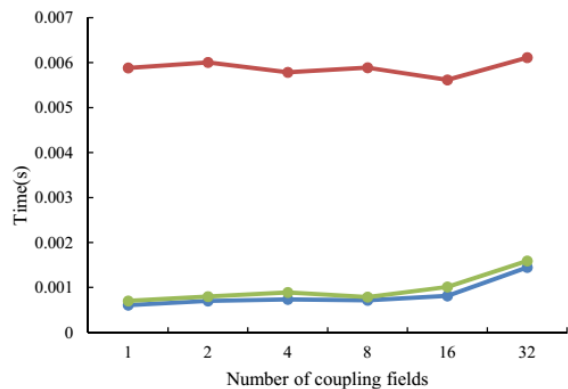
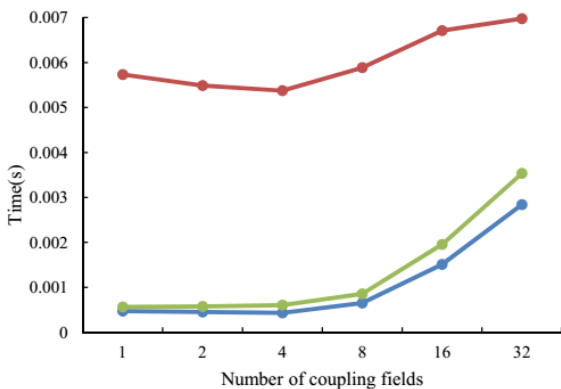
(d)

1



(a)

(b)

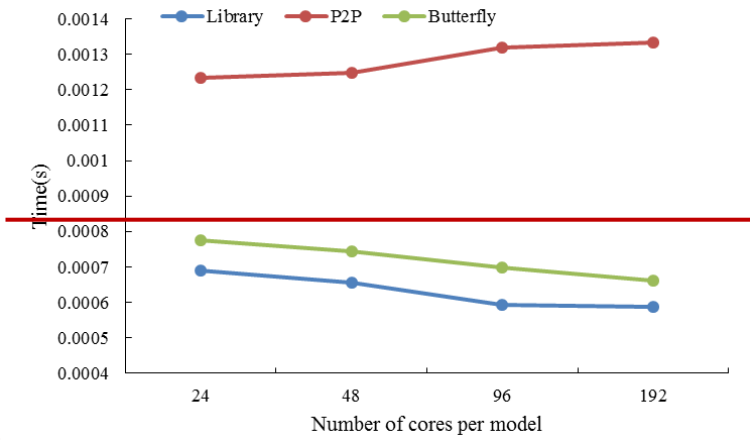


(c)

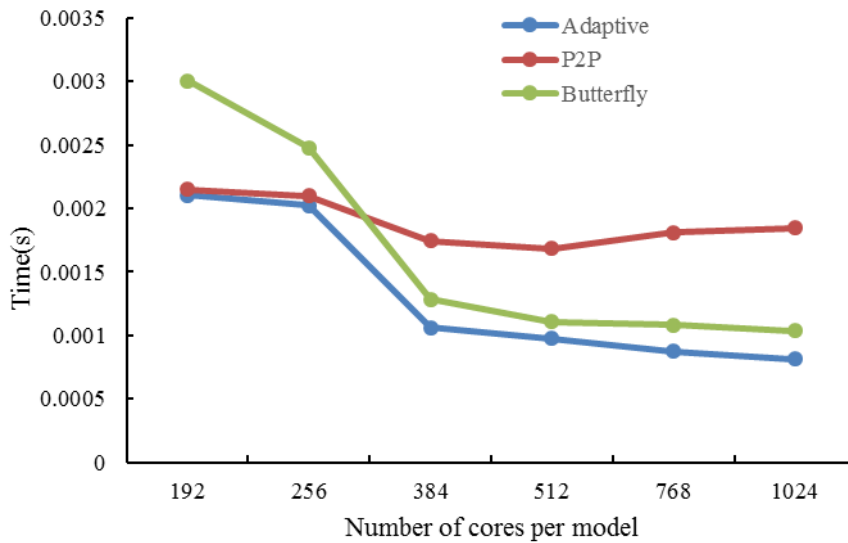
(d)

2

1 Figure ~~4716~~. Average execution time (y-axis) of one data transfer between two toy models with  
2 the same rectangular grid (of  $192 \times 96480$  grid points) when varying the number of coupling  
3 fields transferred (x-axis). There are four simulation tests for the evaluation. In simulation (a),  
4 each toy model is run with ~~48256~~ cores, and the ~~MPI message number communication depth~~  
5 per ~~sender~~ process ~~in the P2P implementation~~ is 12. In simulation (b), each toy model is run  
6 with ~~1921024~~ cores, and the ~~MPI message number communication depth~~ per ~~sender~~ process is  
7 ~~in the P2P implementation~~ 12. In simulation (c), each toy model is run with ~~48256~~ cores, and  
8 the ~~MPI message number communication depth~~ per ~~sender~~ process ~~in the P2P implementation~~  
9 is 48. In simulation (d), each toy model is run with ~~1921024~~ cores (or processes), and the ~~MPI~~  
10 ~~message number communication depth~~ per ~~sender~~ process ~~in the P2P implementation~~ is 48.  
11



1

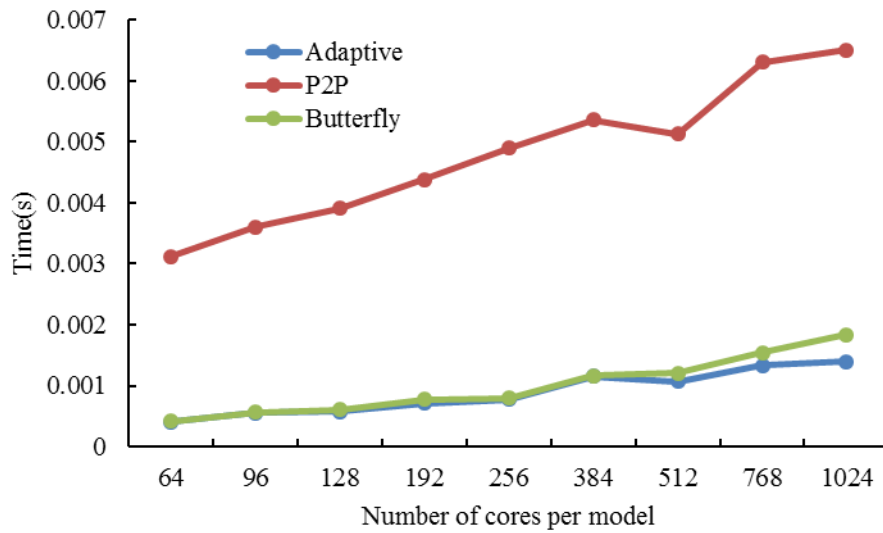


2

3 Figure 1817. Average execution time (y-axis) of one data transfer between two toy models with  
 4 the same rectangular grid (of  $192 \times 96480$  grid points) when varying the number of cores used  
 5 by each toy model (x-axis). There are 10 2-D coupling fields transferred from the source toy  
 6 model to the target toy model. In each test, the MPI message number communication depth per  
 7 sender process in the P2P implementation is ~~set to~~ 24.

8

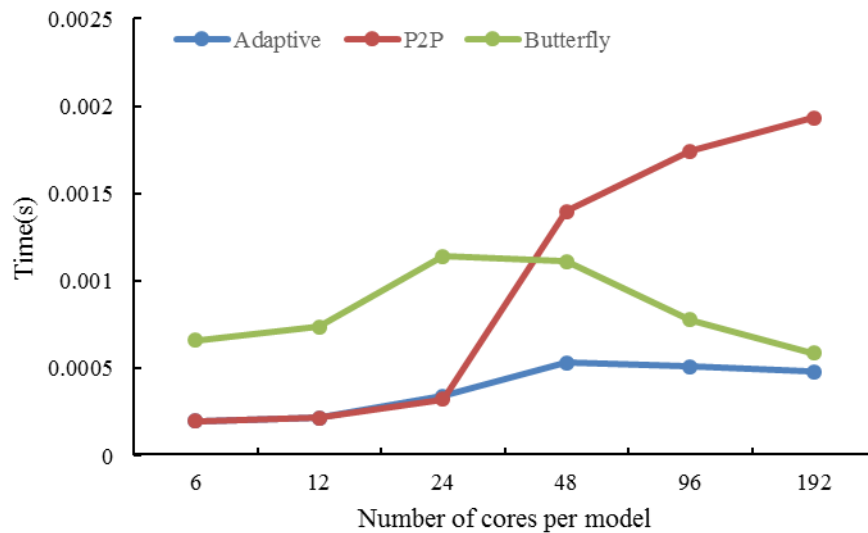
9



1

2 Figure 1918. Average execution time (y-axis) of one data transfer between two toy models. In  
 3 this evaluation, each process (running on a unique processor core) of the toy models have 96  
 4 grid points, while different processes have different message numberscommunication depth and  
 5 different message sizes in the P2P implementation. The number of coupling fields transferred  
 6 is set to 20.

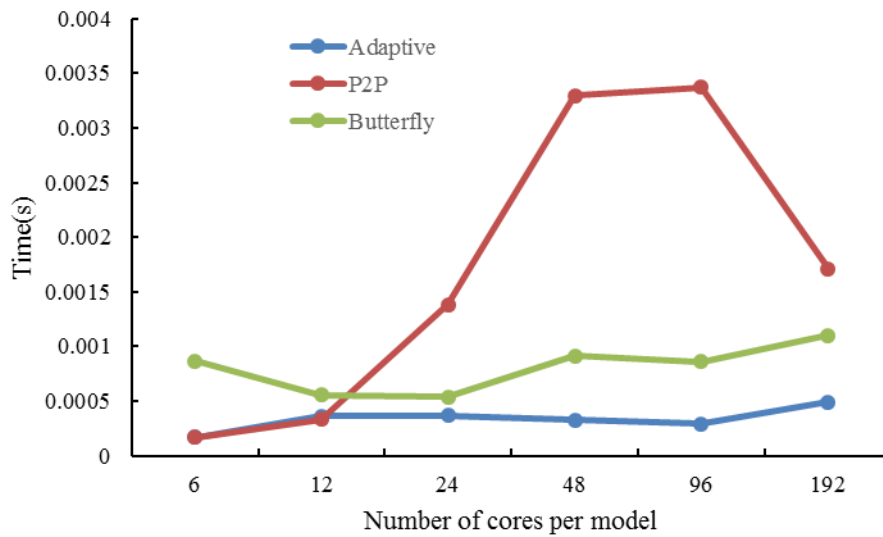
7



1  
 2 Figure 2019. Average execution time (y-axis) of one data transfer between the land surface  
 3 model CLM4 and the coupler CPL7 in CESM when varying the number of cores used by each  
 4 model (x-axis): 32 coupling fields on the CLM horizontal grid (the grid size is  $144 \times 96 = 13824$ )  
 5 are transferred from the land surface model CLM4 to the coupler CPL7. The P2P performance  
 6 results of the P2P implementation are ~~from~~ obtained through running the adaptive data transfer  
 7 library ~~which~~ when it completely switches to the original P2P implementation.

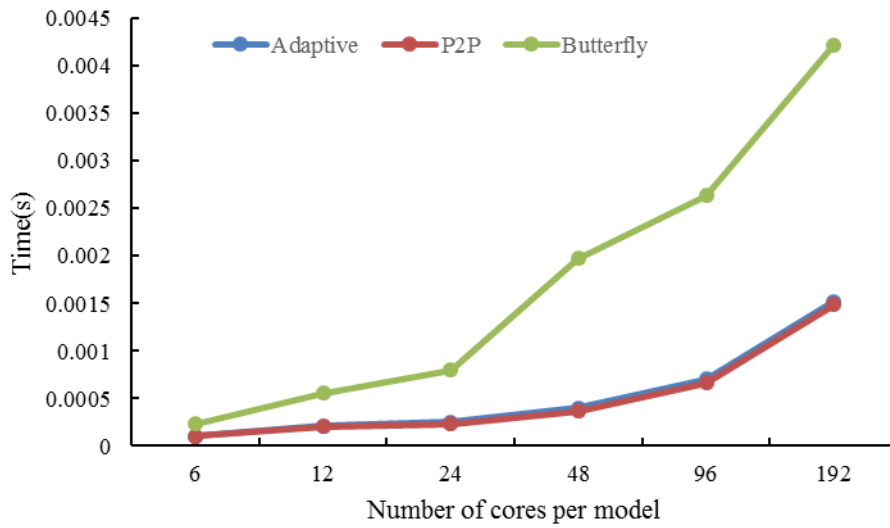
8





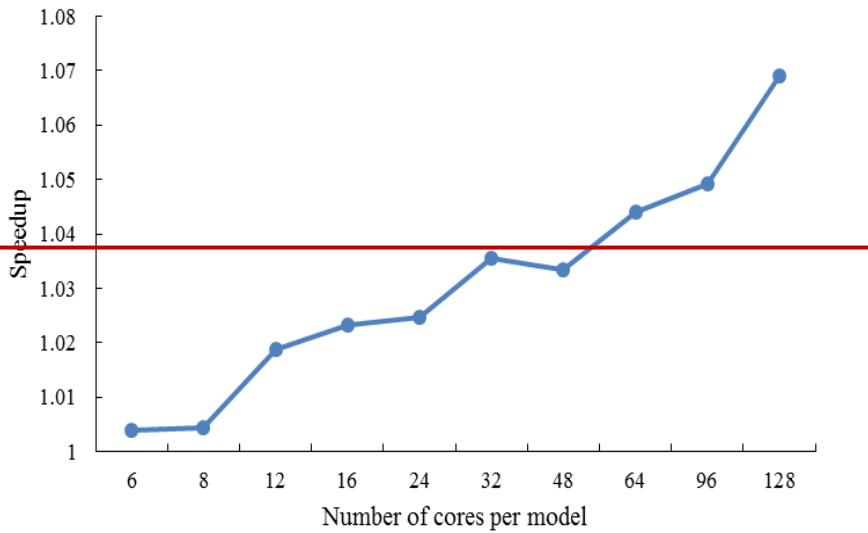
1  
 2 Figure 2420. Average execution time (y-axis) of one data transfer between the atmosphere  
 3 model GAMIL2 and the land surface model CLM3 in GAMIL2-CLM3 when varying the  
 4 number of cores used by each model (x-axis): 14 coupling fields on the GAMIL2 horizontal  
 5 grid (the grid size is  $128 \times 60 = 7680$ ) are transferred from the land surface model CLM3 to the  
 6 atmosphere model GAMIL2.

7

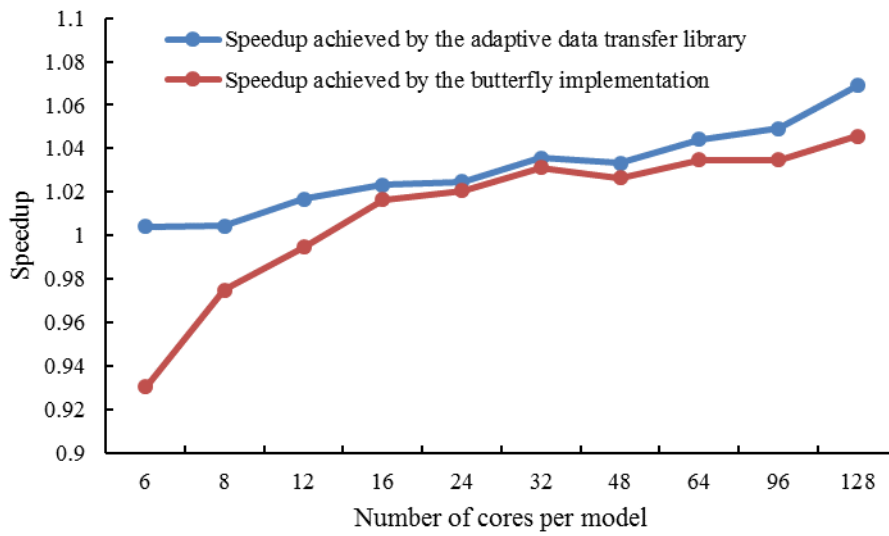


1  
2  
3  
4  
5  
6  
7

Figure 2221. Average execution time (y-axis) of one data rearrangement for the parallel interpolation from the atmosphere grid (the grid size is  $144 \times 96 = 13824$ ) to the ocean grid (the grid size is  $320 \times 384 = 122880$ ) in CESM when varying the number of cores used by each model (x-axis).



1



2

3 Figure 2322. Performance improvement of the coupled model GAMIL2-CLM3 achieved by  
 4 the butterfly implementation and the adaptive data transfer library, with the performance whole  
 5 model time of GAMIL2-CLM3 using the P2P implementation as the baseline.