

1 **Improving Data Transfer for Model Coupling**

2 **C. Zhang^{2,1}, L. Liu^{1,3}, G. Yang^{2,1,3}, R. Li^{2,1}, and B. Wang^{1,3,4}**

3 [1]{Ministry of Education Key Laboratory for Earth System Modeling, Center for Earth
4 System Science (CESS), Tsinghua University, Beijing, China}

5 [2]{Department of Computer Science and Technology, Tsinghua University, Beijing, China}

6 [3]{Joint Center for Global Change Studies (JCGCS), Beijing, China}

7 [4]{State Key Laboratory of Numerical Modelling for Atmospheric Sciences and Geophysical
8 Fluid Dynamics (LASG), Institute of Atmospheric Physics, Chinese Academy of Sciences,
9 Beijing, China.}

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

11

12 **Abstract**

13 Data transfer means transferring data fields between two component models or rearranging
14 data fields among processes of the same component model. It is a fundamental and most
15 frequently used operation of a coupler. Most versions of state-of-the-art couplers currently use
16 an implementation based on the point-to-point (P2P) communication of the Message Passing
17 Interface (MPI) (refer such an implementation as “P2P implementation” for short). In this
18 paper, we reveal the drawbacks of the P2P implementation, including low communication
19 bandwidth due to small message size, variable and big number of MPI messages, as well as
20 network contention. To overcome these drawbacks, we propose a butterfly implementation
21 for data transfer. Although the butterfly implementation can outperform the P2P
22 implementation in many cases, it degrades the performance in some cases because the total
23 message size transferred by the butterfly implementation is larger than the total message size
24 transferred by the P2P implementation. To further improve data transfer, we design and
25 implement an adaptive data transfer library that combines the advantages of both butterfly
26 implementation and P2P implementation. Performance evaluation shows that the adaptive
27 data transfer library significantly improves the performance of data transfer in most cases, and
28 does not decrease the performance in any cases. Now, the adaptive data transfer library is

1 open to the public and has been imported into a coupler version C-Coupler1 for performance
2 improvement of data transfer. We believe that other coupler versions can also benefit from it.

3

4 **1 Introduction**

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

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

26 A coupler is an important component in a coupled system. It links component models together
27 to construct a coupled model, and controls the integration of the whole coupled model
28 (Valcke, 2012). A number of couplers now are available, e.g., the Model Coupling Toolkit
29 (MCT; Jacob et al., 2005), the Ocean Atmosphere Sea Ice Soil coupling software (OASIS)
30 coupler (Redler et al., 2010; Valcke, 2013), the Earth System Modelling Framework (ESMF;
31 Hill et al., 2004), the CPL6 coupler (Craig et al., 2005), the CPL7 coupler (Craig et al., 2012),
32 the Flexible Modelling System (FMS) coupler (Balaji et al., 2006), the Bespoke Framework

1 Generator (BFG; Ford et al., 2006; Armstrong et al., 2009) and the community coupler
2 version 1 (C-Coupler1; Liu et al., 2014).

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

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

22 In this study, we first propose a butterfly implementation of data transfer. Since the P2P
23 implementation and the butterfly implementation can outperform each other in different cases
24 (Section 5), we next develop an adaptive data transfer library that includes both
25 implementations and can adaptively use the better one for data transfer. Performance
26 evaluation demonstrates that such a library significantly improves the performance of data
27 transfer in most cases and does not degrade the performance in any case. This library has been
28 imported into C-Coupler1 with slight code modification. We believe that other coupler
29 versions can also benefit from it.

30 The remainder of this paper is organized as follows. We briefly introduce the implementation
31 of data transfer in existing couplers in Section 2. Details of the butterfly implementation and
32 the adaptive data transfer library are presented in Sections 3 and 4, respectively. The

1 performances of data transfer implementations are evaluated in Section 5. Conclusions are
2 given in Section 6.

3 **2 Data transfer implementations in existing couplers**

4 **2.1 P2P implementation**

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

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

17 **2.2 Performance bottlenecks of the P2P implementation**

18 Although the P2P implementation can achieve good performance when rearranging data
19 fields for a parallel interpolation in a component model, it is not efficient enough when
20 transferring data between component models (Craig et al., 2012; Valcke, 2013; Valcke et al.,
21 2013; Liu et al., 2014). To reveal why the P2P implementation is not efficient enough, we
22 first derive a benchmark from a real coupled model GAMIL2-CLM3, which includes
23 GAMIL2 (Li et al., 2013) that is an atmosphere model and CLM3 (Oleson et al., 2004;
24 Dickinson et al., 2006) that is a land surface model. GAMIL2 and CLM3 share the same
25 horizontal grid of 7,680 (128×60) grid points, but have different parallel decompositions:
26 GAMIL2 uses a regular 2-D parallel decomposition, while CLM3 uses an irregular 2-D
27 parallel decomposition where the grid points are assigned to the processes in a round-robin
28 fashion.

29 In this benchmark, there is only the data transfer with P2P implementation between two data
30 models with the same horizontal grid of GAMIL2-CLM3. The parallel decomposition of the

1 source data model is derived from CLM3, and the parallel decomposition of target data model
2 is derived from GAMIL2. A high-performance computer named Tansuo100 at Tsinghua
3 University, China is used for the performance tests. It has 700 computing nodes, each of
4 which contains two six-core Intel Xeon X5670 CPUs and 32 GB main memory. All
5 computing nodes are connected by a high-speed InfiniBand network with peak
6 communication bandwidth of 5 GB/s.

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

14 Figure 1 demonstrates the poor performance of the P2P implementation. It is well known that
15 the communication performance heavily depends on message size. As shown in Fig. 2, the
16 P2P communication bandwidth achieved generally increases with message size. So when the
17 message size is small (for example, smaller than 4 KB), the communication bandwidth
18 achieved is very low. The message size in the P2P implementation decreases with increment
19 of process number of models (Fig. 3), indicating that the communication bandwidth becomes
20 lower with the increment of process number. The performance of data transfer also heavily
21 depends on the MPI message number. As shown in Fig. 4, the message number in the P2P
22 implementation increases with increment of process number. Here, we may conclude that the
23 decrease of message size and the increase of message number are primary reasons for the poor
24 performance of the P2P implementation when increasing the process number. However, the
25 ideal performance shown in Fig. 5 is much better than the actual performance. The ratio
26 between the ideal performance and the actual performance significantly increases with the
27 increment of processor number. The significant gap between the ideal performance and the
28 actual performance is due to network contention. For example, when multiple P2P
29 communications share the same source process or target process (Fig. 6), they must wait in an
30 order.

1 **3 Butterfly implementation for better performance of data transfer**

2 The drawbacks of the P2P implementation can be concluded as low communication
3 bandwidth due to small message size, variable and big number of MPI messages, as well as
4 network contention. To overcome these drawbacks, a prospective solution is to organize the
5 communication for data transfer using a better structure, so that we investigate the butterfly
6 structure (Fig. 7), which has already been used in the field of computer (Chong et al., 1994;
7 Foster, 1995; Heckbert et al., 1995; Hemmert et al., 2005; Kim et al., 2007; Jan et al., 2013;
8 Petagon et al, 2016). For example, in hardware aspect, the traditional butterfly structure and
9 its transformation have been used to design networks (Chong et al., 1994; Kim et al., 2007);
10 in software aspect, the butterfly structure has been used to improve the parallel algorithms
11 with all-to-all communications (Foster, 1995), e.g., Fast Fourier Transform (FFT; Heckbert et
12 al., 1995; Hemmert et al., 2005), matrix transposition (Petagon et al, 2016) and sorting (Jan et
13 al., 2013).

14 Unfortunately, the improved all-to-all communication with the butterfly structure cannot be
15 used to improve data transfer, because it requires that one process must communicate with
16 every other process, that the communication load among processes is balanced and that the
17 number of processes must be a power of 2, while the data transfer for model coupling has
18 different characteristics, i.e., one process needs to communicate with a part of other processes
19 (Fig. 6), the communication load among processes is always unbalanced (Fig. 3) and the
20 process number cannot be restricted to a power of 2. Therefore, to benefit from the butterfly
21 structure, we should design a new implementation of data transfer, which is called the
22 butterfly implementation hereafter.

23 The butterfly implementation uses a butterfly structure to transfer data from the sender with
24 the source parallel decomposition to the receiver with the target parallel decomposition. We
25 call the communication following the butterfly structure “the butterfly kernel”. As the process
26 number of the butterfly kernel must be a power of 2, while the process number of the sender
27 or the receiver need not be a power of 2, the butterfly implementation (see Fig. 8) has a
28 process mapping from the sender onto the butterfly kernel and a process mapping from the
29 butterfly kernel onto the receiver, and the butterfly kernel has its own source parallel
30 decomposition and target parallel decomposition, which are determined by the process
31 mappings. Next, we will present the butterfly kernel and the process mappings, respectively.

1 3.1 Butterfly kernel

2 The first question for the butterfly kernel is how to decide its process number. Any process of
3 the sender or the receiver can be used as a process of the butterfly kernel. Given that the total
4 number of unique processes of the sender and receiver is N_T , the process number of the
5 butterfly kernel (N_B) can be any power of 2, which is no larger than N_T . We propose to select
6 the maximum number in order for maximum utilization of resources. We prefer to pick out
7 unique processes first from the sender, and then from the receiver if the sender does not have
8 enough processes.

9 The butterfly kernel is responsible for rearranging the distribution of data among the
10 processes from the source parallel decomposition to the target parallel decomposition. Given
11 the process number $N=2^n$, there are n stages in the butterfly kernel. In a stage, all processes
12 are divided into a number of pairs and the two processes of a pair uses MPI P2P
13 communication to exchange data. Given a process P in the butterfly kernel, after each stage,
14 the number of the processes that may have the data of P on the target parallel decomposition
15 will become a half. Figure 7 is an example for further illustration, where D_j^i means the data is
16 originally in process P_i according to the source parallel decomposition and is finally in
17 process P_j according to the target parallel decomposition. Before the first stage, all processes
18 ($P_0\sim P_7$) may have the data of P_0 on the target parallel decomposition. After the first stage,
19 only four processes (P_0, P_2, P_4 and P_6) may have; and after the second stage, only two
20 processes (P_0 and P_4) may have.

21 To reveal the advantages and disadvantages of the two implementations, we measure the
22 characteristics of the two implementations based on the benchmark introduced in Section 2.2.
23 The results show the total message size transferred by the butterfly implantation is larger than
24 that by the P2P implementation (Fig. 9), which is the major disadvantage of the butterfly
25 implementation. Meanwhile, comparing with the P2P implementation, the butterfly
26 implementation has the following advantages:

- 27 1) bigger message size for better communication bandwidth (Fig. 10);
- 28 2) balanced number of MPI messages among processes (Fig. 11);
- 29 3) ordered communications among processes and fewer communications operated
30 concurrently (Fig. 11), which can dramatically reduce network contention.

1 3.2 Process mapping

2 In this subsection, we will introduce the process mappings from the sender to the butterfly
3 kernel and from the butterfly kernel to the receiver. To minimize the overhead of process
4 mapping from the butterfly kernel to the receiver, we map one or multiple processes of the
5 butterfly kernel onto a process of the receiver if the butterfly kernel has more processes than
6 the receiver; otherwise, we map a process of the butterfly kernel onto one or multiple
7 processes of the receiver. In other words, there is no multiple-to-multiple process mapping
8 between the butterfly kernel and the receiver. Similarly, there is no multiple-to-multiple
9 process mapping between the sender and the butterfly kernel.

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

23 The iterative pairing also requires the number of processes to be a power of 2. Given that the
24 process number of the sender (or receiver) is N_c and the process number of the butterfly
25 kernel is N_B , we first pad empty processes (whose data size is zero) before the iterative
26 pairing to make the process number of the sender (or receiver) be a power of 2 (denoted N_P),
27 which is no smaller than N_B . Therefore, the reordered N_P processes after the iterative pairing
28 can be divided into N_B groups, each of which contains N_P/N_B processes with consecutive
29 reordered indexes and maps onto a unique process of the butterfly kernel.

30 Figure 12 shows an example of the process mapping, where the sender has five processes (S_0 -
31 S_4 in Fig. 12a), the receiver has 10 processes (R_0 - R_9 in Fig. 12b), and the butterfly kernel uses
32 eight processes (B_0 - B_7 in Fig. 12c). At the first, empty processes are padded to the sender (S_5 -

1 S_7 in Fig. 12a) and the receiver (R_{10} - R_{15} in Fig. 12b). Next, the iterative pairing is conducted
 2 for the sender and the receiver, respectively. The iterative pairing has three stages for the
 3 sender. At the first stage, the eight processes of the sender are divided into four groups
 4 $\{S_1, S_7\}$, $\{S_0, S_6\}$, $\{S_2, S_5\}$ and $\{S_4, S_3\}$ (Fig. 12a), according to the data size corresponding to
 5 each process. These four process groups are divided into two bigger groups ($\{\{S_4, S_3\}, \{S_2, S_5\}\}$
 6 and $\{\{S_1, S_7\}, \{S_0, S_6\}\}$) at the second stage (Fig. 12a). Finally, one process group
 7 $\{\{\{S_4, S_3\}, \{S_2, S_5\}\}, \{\{S_1, S_7\}, \{S_0, S_6\}\}\}$ is obtained at the third stage (Fig. 12a), and the eight
 8 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 which is
 9 mapped onto one process of the butterfly kernel (Fig. 12c). Similarly, the iterative pairing has
 10 four stages for the receiver, and the 16 processes of the receiver are reordered as $R_9, R_{15}, R_7,$
 11 $R_{12}, R_4, R_8, R_3, R_{10}, R_1, R_{14}, R_5, R_{13}, R_0, R_6, R_2$ and R_{11} finally, each two of which are mapped
 12 onto one process of the butterfly kernel (Fig. 12c).

13 **4 Adaptive data transfer library**

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

25 As introduced in Section 3.1, the butterfly implementation is divided into multiple stages.
 26 Actually, the data transfer in one stage can be viewed as a P2P implementation with only one
 27 MPI message per process. Inspired by this fact, we try to design an adaptive approach that can
 28 combine the butterfly and P2P implementations, where some stages in the butterfly
 29 implementation are skipped with the P2P implementations of more MPI messages per process.
 30 If all stages of the butterfly implementation are skipped, the adaptive data transfer library will
 31 switch to the P2P implementation. Figure 13 shows an example of the adaptive data transfer

1 library with eight processes, where Stage 2 of the butterfly implementation is skipped with the
2 P2P implementation of three MPI messages per process.

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

26 **5 Performance evaluation**

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

1 Next, we will evaluate the overhead of initialization, the performance in data transfer and the
2 performance in data rearrangement for parallel interpolation.

3 **5.1 Overhead of initialization**

4 We first evaluate the initialization overhead of data transfer implementations. As shown in
5 Fig. 15, the initialization overhead of each implementation increases with the increment of
6 core number. The initialization overhead of the butterfly implementation is a little higher than
7 that of the P2P implementation, while the initialization overhead of the adaptive data transfer
8 library is 2-3 folds higher than that of the P2P implementation, because the adaptive data
9 transfer library uses extra time on the performance profiling (please refer to Section 4).
10 Considering that one data transfer instance should only be initialized at the beginning and
11 executed many times in a coupled model, we can conclude that the initialization overhead of
12 the adaptive data transfer library is reasonable, especially when the simulation is executed for
13 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 factors can impact the performance of a data
17 transfer implementation including MPI message number, the size of data to be transferred
18 (also known as the number of fields in this evaluation) and the number of cores used. In this
19 subsection, we evaluate the performance of data transfer affected by each of these factors. We
20 first build two toy models that use the same logically rectangular grid (of 192×96 grid points).
21 Coupling fields are transferred between the two toy models. In each test, the two toy models
22 use the same number of cores, and each process has the same MPI message number. Next, we
23 evaluate the performance of data transfer through varying one factor and fixing the other
24 factors.

25 In the first experiment, we fix the number of cores to be 192 and the coupling field number to
26 be 10, and vary MPI message number per process. Figure 16 shows the execution time of one
27 data transfer with different implementations when varying the MPI message number per
28 process from 1 to 96. The P2P implementation can outperform the butterfly implementation
29 when the MPI message number is small (say, smaller than 12 in Fig. 16), while the butterfly
30 implementation can outperform the P2P implementation when the MPI message number is

1 big (say, bigger than 12 in Fig. 16). The adaptive data transfer library has the best
2 performance. Moreover, it improves the performance based on the butterfly implementation
3 when the MPI message number is big, because some butterfly stages in the adaptive data
4 transfer library have been skipped with the P2P implementation. When the MPI message
5 number per process is 96, the adaptive data transfer library can achieve a 13.9-fold
6 performance speedup compared to the P2P implementation.

7 In the second experiment, we fix the number of cores and MPI message number per process,
8 and vary the coupling field number transferred. Figure 17 shows the execution time of one
9 data transfer with different implementations in this experiment. The results show that the
10 execution time of each implementation increases with the increment of data size. When MPI
11 message number per process is small (Figs. 17a and 17b), the performance of the butterfly
12 implementation is poorer than that of the P2P implementation, especially when the number of
13 2-D coupling fields gets bigger. The adaptive data transfer library achieves similar
14 performance as the P2P implementation, because it switches to the P2P implementation.
15 When the MPI message number per process is big (Figs. 17c and 17d), both the butterfly
16 implementation and adaptive data transfer library significantly outperform the P2P
17 implementation, and the adaptive data transfer library achieves better performance than the
18 butterfly implementation.

19 In the third experiment, we fix MPI message number per process to be 24 and the coupling
20 field number transferred to be 10, and vary the number of cores. Figure 18 shows the
21 execution time of one data transfer with different implementations when varying the number
22 of cores. The results show that both the butterfly implementation and adaptive data transfer
23 library achieve better parallel scalability than the P2P implementation. The execution time of
24 the P2P implementation slightly increases with the increment of the number of cores used.
25 However, the execution times of the butterfly implementation and adaptive data transfer
26 library slightly decrease with the increment of the number of the cores used. The butterfly
27 implementation outperforms the P2P implementation, while the adaptive data transfer library
28 achieves better performance than the butterfly implementation.

29 The resolution of models becomes higher and higher these days. How about the performance
30 of the data transfer implementations when model resolution becomes higher? Higher model
31 resolution means that a model will use more processor cores for accelerating simulation,
32 while the average number of grid points per processor core can remain constant. Considering

1 that the numbers of grid points are always balanced among the processes of a model, we make
2 each process (which runs on a unique processor core) of the toy models have 96 grid points in
3 this evaluation, while enabling processes to have different message numbers and different
4 message sizes. As shown in Fig. 19, although the execution times of all data transfer
5 implementations increase with the increment of processor core number (from 64 to 1024),
6 both the butterfly implementation and the adaptive data transfer library significantly
7 outperform the P2P implementation, and the adaptive data transfer library achieves the best
8 performance. These results indicate that our proposed implementations can significantly
9 improve the performance of data transfer for higher model resolution.

10 **5.3 Performance of data transfer between realistic models**

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

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

25 For GAMIL2-CLM3, we use the data transfer from CLM3 to GAMIL2 where 14 2-D
26 coupling fields on the GAMIL2 horizontal grid (whose grid size is $128 \times 60 = 7680$) are
27 transferred. Figure 21 shows the execution time of one data transfer of each implementation
28 when increasing the process number of both GAMIL2 and CLM3 from 6 to 192. The results
29 in Fig. 21 confirm that the adaptive data transfer library can constantly show the best
30 performance. Compared to the P2P implementation, the adaptive data transfer library
31 achieves an 11.68-fold performance speedup when the process number is 96, but achieves a

1 much lower speedup (only 3.48-fold) when the process number is 192. This is because the
2 average MPI message number per process reduces from 32 to 18 when the number of process
3 increases from 96 to 192.

4 **5.4 Performance of data rearrangement for interpolation**

5 Besides data transfer between different component models, there is another kind of data
6 transfer in model coupling that rearranges data inside a model for parallel interpolation of
7 fields between different grids. Here, we use the data rearrangement for the parallel
8 interpolation from the atmosphere grid (whose grid size is $144 \times 96 = 13824$) to the ocean grid
9 (whose grid size is $320 \times 384 = 122880$) in the coupled model CESM for further evaluation. As
10 mentioned above, the P2P implementation is sufficient for data rearrangement. However, the
11 butterfly implementation is much poorer than the P2P implementation (Fig. 22). This is
12 because the MPI message number is very small (for example, average MPI message number
13 per process is only 6.49 when each model uses 96 cores) for data rearrangement. On the other
14 hand, the adaptive data transfer library achieves almost the same performance as the P2P
15 implementation, because it switches to the P2P implementation. Therefore, the adaptive data
16 transfer library can always show the best performance.

17 **5.5 Performance improvement for a coupled model**

18 With the performance improvement of data transfer, we expect that the adaptive data transfer
19 library will improve the performance of coupled models. For this evaluation, we first import
20 the adaptive data transfer library into C-Coupler1 and then use the coupled model GAMIL2-
21 CLM3 that uses C-Coupler1 for coupling to measure performance results. As shown in Fig.
22 23, the adaptive data transfer library achieves higher performance improvement (when the
23 P2P implementation is used as the baseline) for GAMIL2-CLM3 when using more processor
24 cores. When each component model uses 128 processor cores, the adaptive data transfer
25 library achieves ~7% performance improvement. This performance improvement would not
26 be low because the model coupling only takes a very small proportion of execution time in the
27 simple coupled model GAMIL2-CLM3 and the parallel scalability of the two coupled models
28 GAMIL2 and CLM3 is not good.

29

1 **6 Conclusions**

2 Data transfer is the fundamental and most frequently used operation in a coupler. This paper
3 demonstrated that the current P2P implementation of data transfer in most state-of-the-art
4 couplers is inefficient for transferring data between two component models. To improve the
5 parallel performance of data transfer, we proposed a butterfly implementation. However,
6 compared to the P2P implementation, the butterfly implementation has both advantages and
7 disadvantages. The evaluation results showed that the butterfly implementation did not always
8 outperform the P2P implementation. To achieve better parallel performance of data transfer,
9 we built an adaptive data transfer library, which combines the advantages of both butterfly
10 implementation and P2P implementation. The evaluation results demonstrated that, the
11 adaptive data transfer library can significantly improve the performance of data transfer so as
12 to improve a coupled model.

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

17 **Code availability**

18 The source code of the adaptive data transfer library is available at [https://github.com/zhang-](https://github.com/zhang-cheng09/Data_transfer_lib)
19 [cheng09/Data_transfer_lib](https://github.com/zhang-cheng09/Data_transfer_lib).

20 **Acknowledgements**

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

24

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
10 expanders, *Parallel Computer Routing and Communication*, Springer Berlin Heidelberg,
11 1994:200-214.
- 12 Craig, A. P., Vertenstein, M., and Jacob, R.: A new flexible coupler for Earth system
13 modelling developed for CCSM4 and CESM1, *Int. J. High Perform. C.*, 26, 31-42,
14 doi:10.1177/1094342011428141, 2012.
- 15 Craig, A. P., Jacob, R., Kauffman, B., Bettge, T., Larson, J., Ong, E., Ding, C., and He, Y.:
16 CPL6: the New Extensible, High Performance Parallel Coupler for the Community Climate
17 System Model, *Int. J. High Perform. C.*, 19, 309–327, 2005.
- 18 Dennis, J. M.: Inverse space-filling curve partitioning of a global ocean model, In *IEEE*
19 *International Parallel & Distributed Processing Symposium*, Long Beach, CA, 2007.
- 20 Dennis, J. M. and Tufo, H. M.: Scaling climate simulation applications on the IBM Blue
21 Gene/L system, *IBM J. Res. Dev.*, 52, 117-126, DOI:10.1147/rd.521.0117, 2008.
- 22 Dennis, J. M., Edwards, J., Evans, K. J., Guba, O., Lauritzen, P. H., Mirin, A. A., St-Cyr, A.,
23 Taylor, M. A., and Worley, P. H.: CAM-SE: a scalable spectral element dynamical core for
24 the Community Atmosphere Model, *Int. J. High Perform. C.*, 26, 74-89,
25 doi:10.1177/1094342011428142, 2012.
- 26 Dickinson, R. E., Oleson, K. W., Bonan, G., Hoffman, F., Thornton, P., Vertenstein, M.,
27 Yang, Z.-L., and Zeng X.: The Community Land surface model and its climate statistics as a
28 component of the Community Climate System Model, *Journal of Climate*, 19(11), 2302–2324,
29 2006.

1 Ford, R. W., Riley, G. D., Bane, M. K., Armstrong, C. W., and Freeman, T. L.: GCF: a
2 general coupling framework, *Concurrency and Computation: Practice and Experience*, 18(2),
3 163–181, 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
9 FFT 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.: $M \times N$ 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,
7 X., Pu, Y., G. and Yang.: Evaluation of Grid-point Atmospheric Model of IAP LASG version
8 2 (GAMIL2), Advances in Atmospheric Sciences, 30, 855–867, doi:10.1007/s00376-013-
9 2157-5, 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.,
17 Williamson, D. L., Rasch, P. J., Vavrus, S. J., Taylor, M. A., Collins, W. D., Zhang, M., and
18 Lin, S.: Description of the NCAR Community Atmosphere Model (CAM 4.0), National
19 Center for Atmospheric Research Ncar Koha Openpat, TN-485+STR, 222p., 2010.

20 Neale, R. B., Chen, C. C., Gettelman, A., Lauritzen, P. H., Park, S., Williamson, D. L.,
21 Conley, A. J., Garcia, R., Kinnison, D., Lamarque, J. F., Marsh, D., Mills, M., Smith, A. K.,
22 Tilmes, S., Vitt, F., Morrison, H., Cameron-Smith, P., Collins, W. D., Iacono, M. J., Easter, R.
23 C., Ghan, S. J., Liu, X., Rasch, P. J., and Taylor, M. A.: Description of the NCAR
24 Community Atmosphere Model (CAM 5.0), National Center for Atmospheric Research Ncar
25 Koha Openpat, 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 Ncar Koha Openpat, 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
5 Earth System Modelling, *Geoscientific Model Development*, 3(1), 87–104, doi:10.5194/gmd-
6 3-87-2010, 2010.

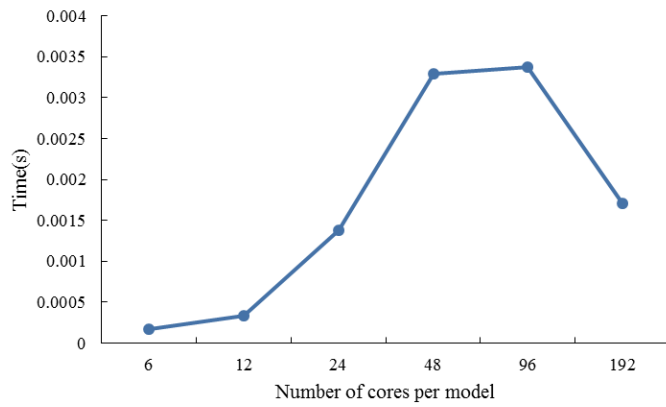
7 Smith, R., Jones, P., Briegleb, B., Bryan, F., Danabasoglu, G., Dennis, J., Dukowicz, J., Eden,
8 C., Fox-Kemper, B., Gent, P., Hecht, M., Jayne, S., Jochum, M., Large, W., Lindsay, K.,
9 Maltrud, M., Norton, N., Peacock, S., Vertenstein, M., and Yeager, S.: The Parallel Ocean
10 Program (POP) reference 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
17 System Modelling, *Geoscientific Model Development*, 5(6), 1589–1596, doi:10.5194/gmd-5-
18 1589-2012, 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 (last access: 15
24 October 2015), 2013.

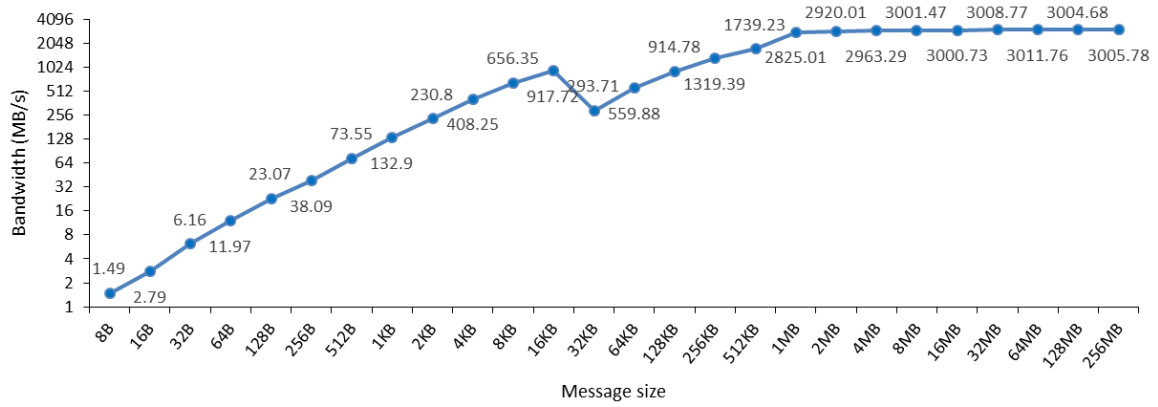
25



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 node. The
5 horizontal grid of the 14 2-D fields contains 7680 (128×60) grid points.

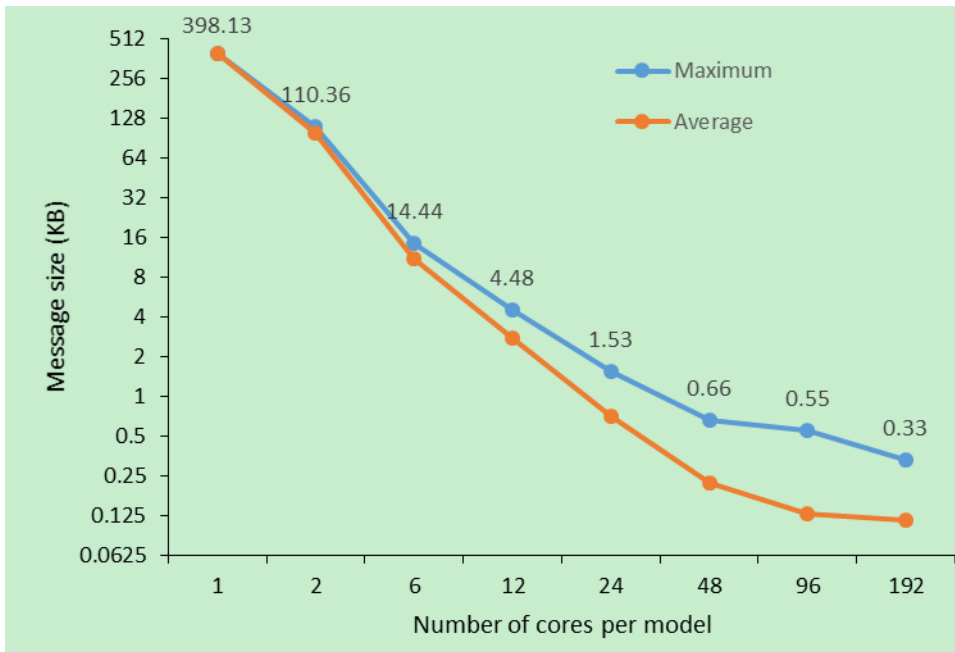
6



1

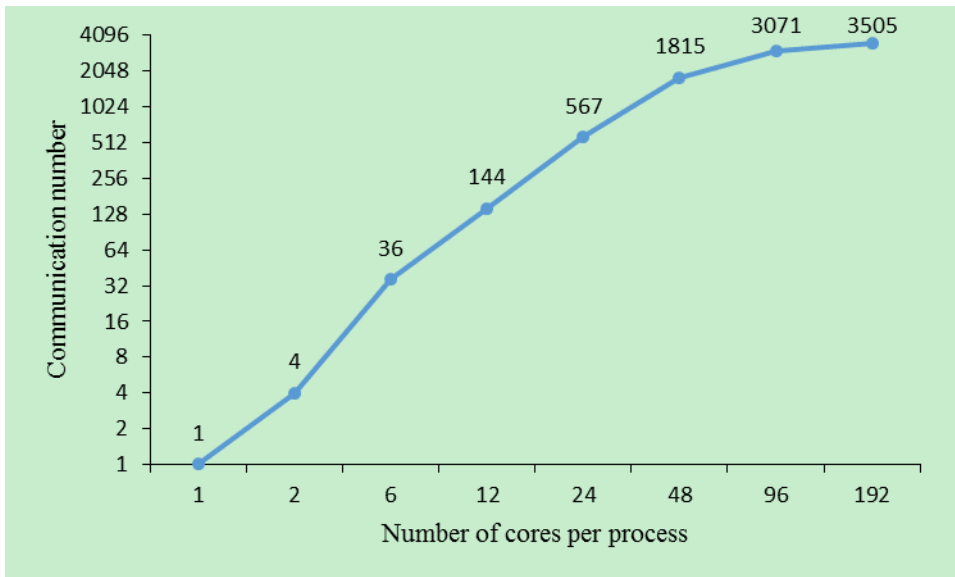
2 Figure 2. Variation of bandwidth (y-axis) of an MPI P2P communication with respect to the
 3 increment of message size (x-axis). The results are generated from our benchmark. In the
 4 benchmark, one process sends messages with different sizes to the other process. The two
 5 processes of the P2P communication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

6



1
2
3
4
5

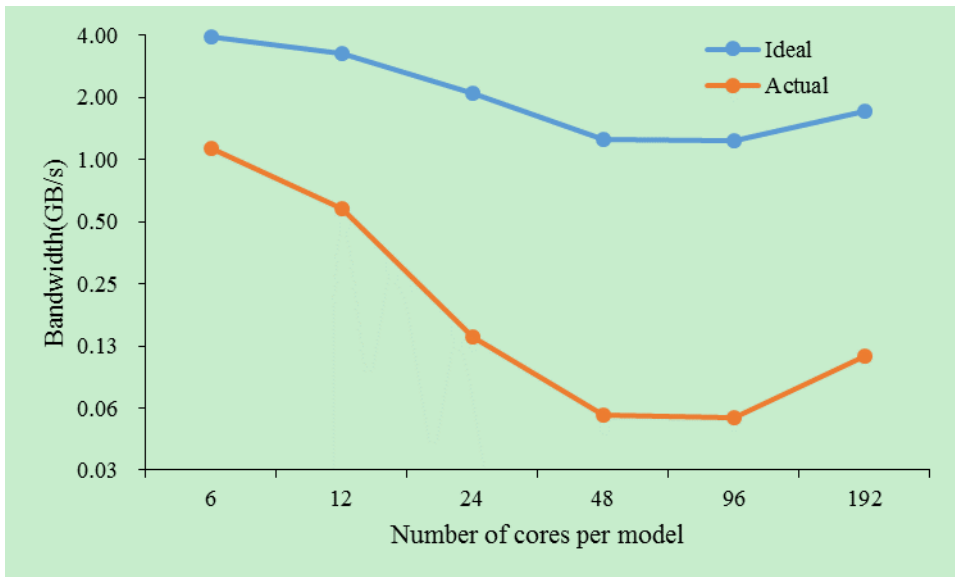
Figure 3. Variation of message size of the P2P implementation (y-axis) in GAMIL2-CLM3 with respect to the increment of core number (x-axis). The experimental setup is similar to that shown in Fig. 1.



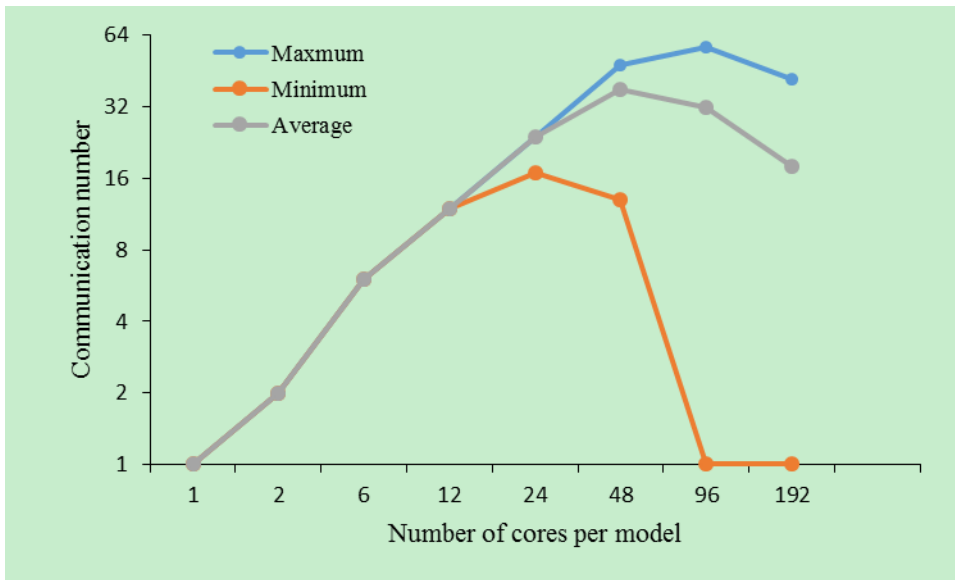
1

2 Figure 4. Variation of total MPI message number (y-axis) of the P2P implementation in
3 GAMIL2-CLM3 with respect to the increment of core number (x-axis). The experimental
4 setup is similar to that shown in Fig. 1.

5



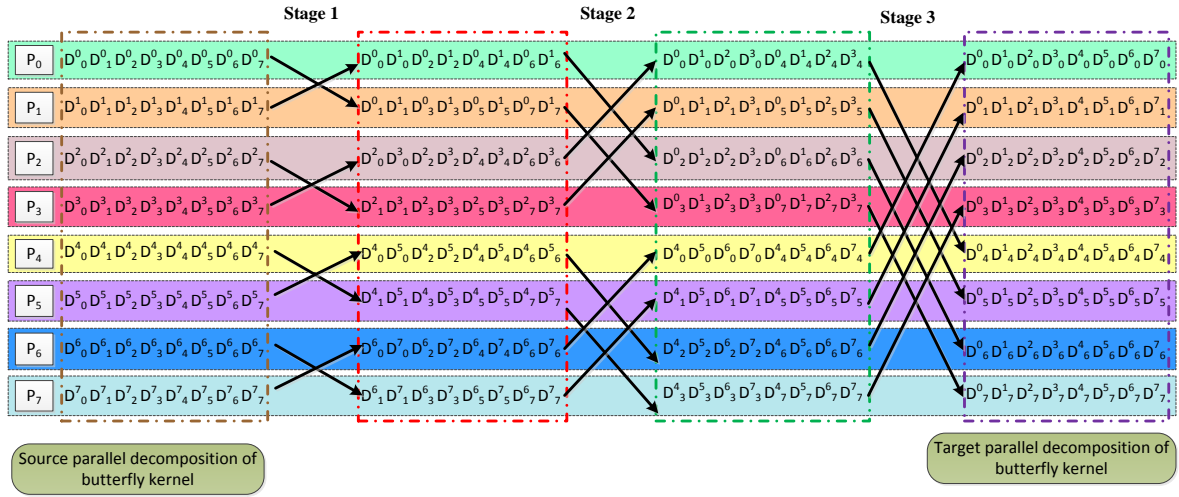
1
 2 Figure 5. Ideal and actual bandwidths of the P2P implementation (y-axis) in GAMIL2-CLM3
 3 when gradually increasing the number of cores used by each component model (x-axis). The
 4 experimental setup is similar to that shown in Fig. 1. The ideal bandwidth is calculated from
 5 the message size and the MPI bandwidth measured in Fig. 2; and the actual bandwidth is
 6 calculated from Fig. 1.
 7



1

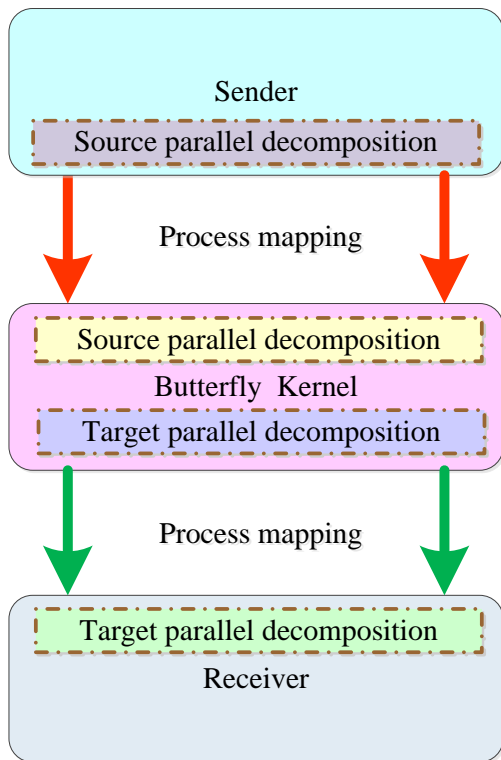
2 Figure 6. Variation of message number of one process (y-axis) using the P2P implementation
 3 in GAMIL2-CLM3 with respect to the increment of core number (x-axis). The experimental
 4 setup is similar to that shown in Fig. 1.

5



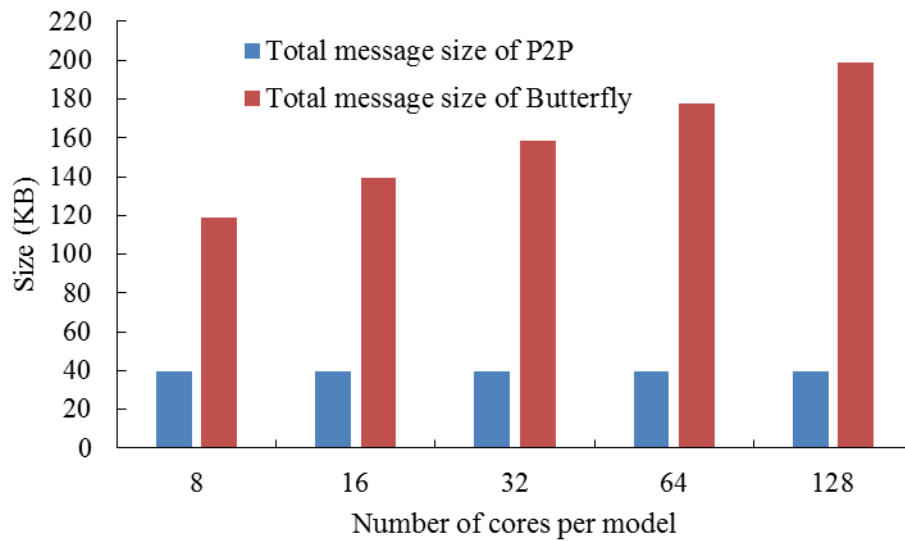
1
 2 Figure 7. An example of the butterfly kernel with eight processes. Each colored row stands
 3 for one process (P_0 - P_7). There are multiple stages (each column of arrows represents a stage
 4 (Stage 1 to Stage 3)) in the butterfly kernel. Each arrow stands for an MPI P2P
 5 communication from one process to another. D_j^i means the data is originally in process P_i
 6 according to the source parallel decomposition and is finally in process P_j according to the
 7 target parallel decomposition.

8

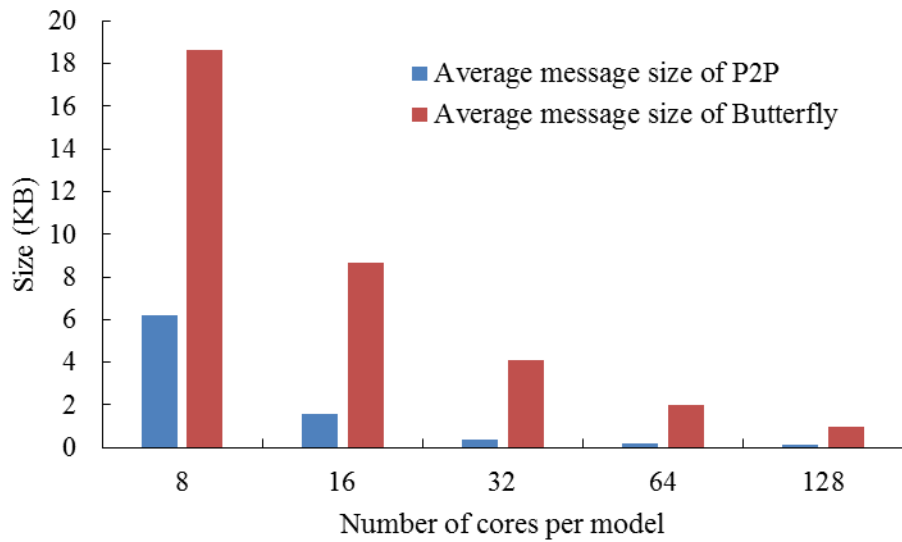


1
2
3
4
5
6

Figure 8. The butterfly implementation, which is composed of three parts: the butterfly kernel; process mapping from the sender to the butterfly kernel; and process mapping from the butterfly kernel to the receiver.



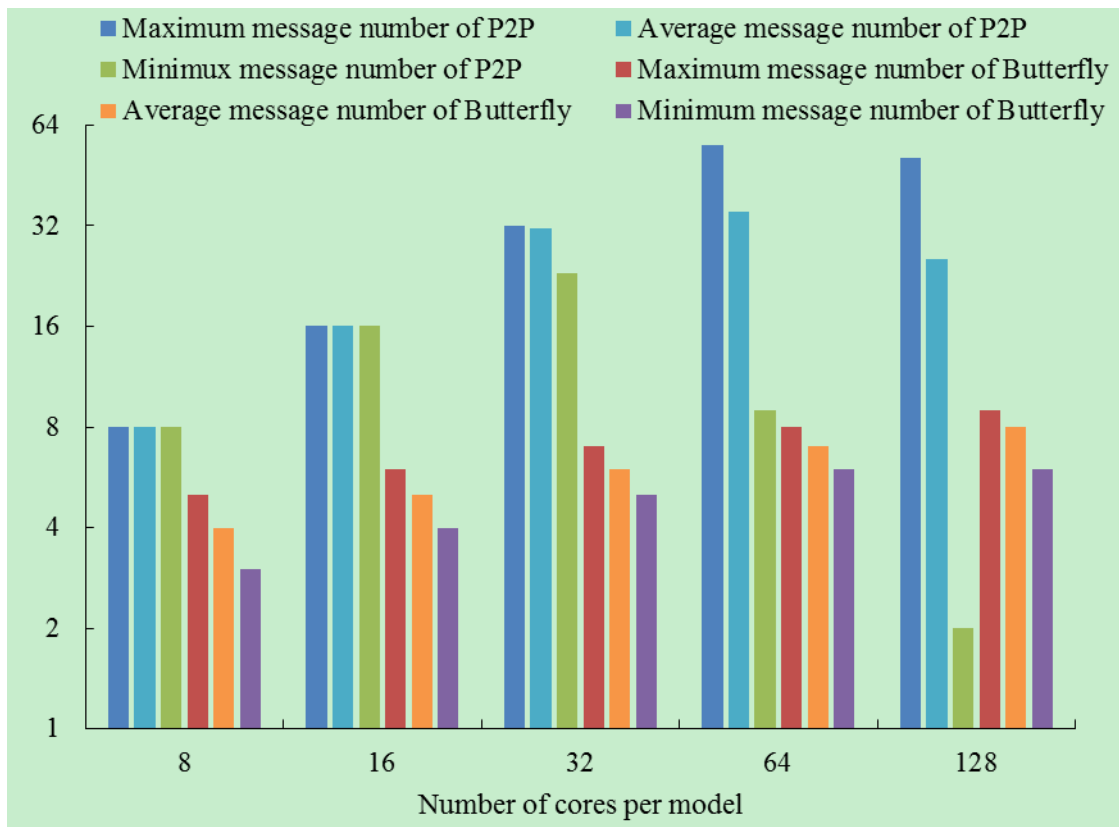
1
 2 Figure 9. Total message size transferred by P2P implementation and butterfly implementation
 3 (y-axis) in GAMIL2-CLM3, when varying the number of cores used by each model (x-axis).
 4 The experimental setup is similar to that shown in Fig. 1.
 5



1

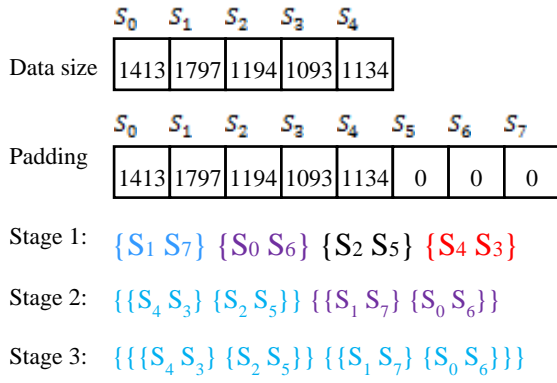
2 Figure 10. Average message size transferred by P2P implementation and butterfly
 3 implementation (y-axis) in GAMIL2-CLM3, when varying the number of cores used by each
 4 model (x-axis). The experimental setup is similar to that shown in Fig. 1.

5

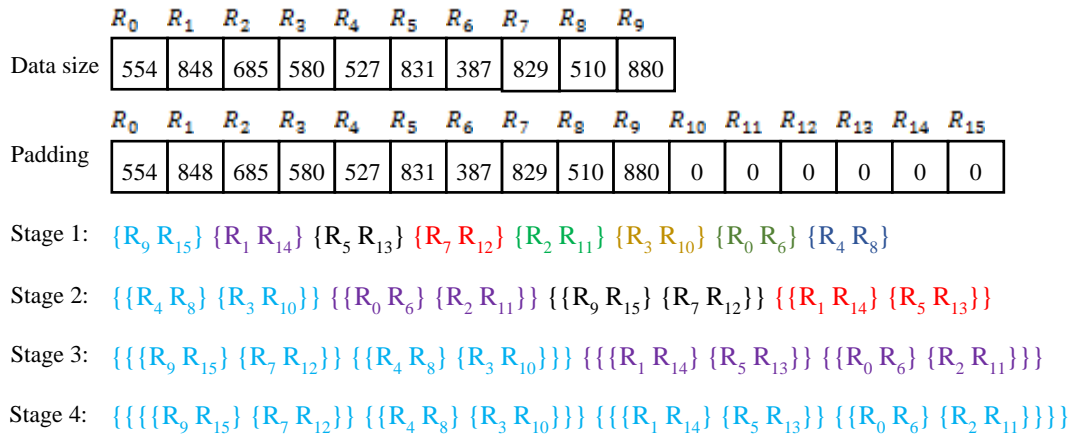


1
2
3
4
5
6

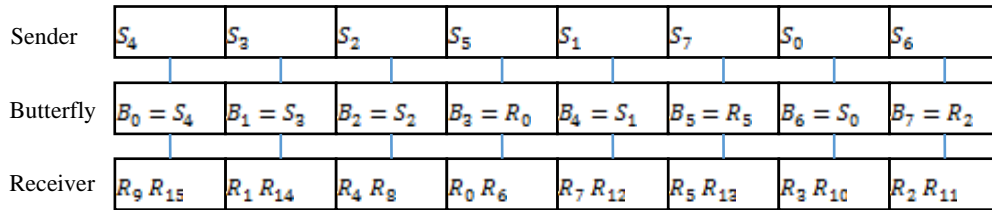
Figure 11. Maximum message number, average message number and minimum message number of processes in P2P implementation and butterfly implementation (y-axis), when varying the number of cores used by each model (x-axis) in GAMIL2-CLM3. The experimental setup is similar to that shown in Fig. 1.



(a)



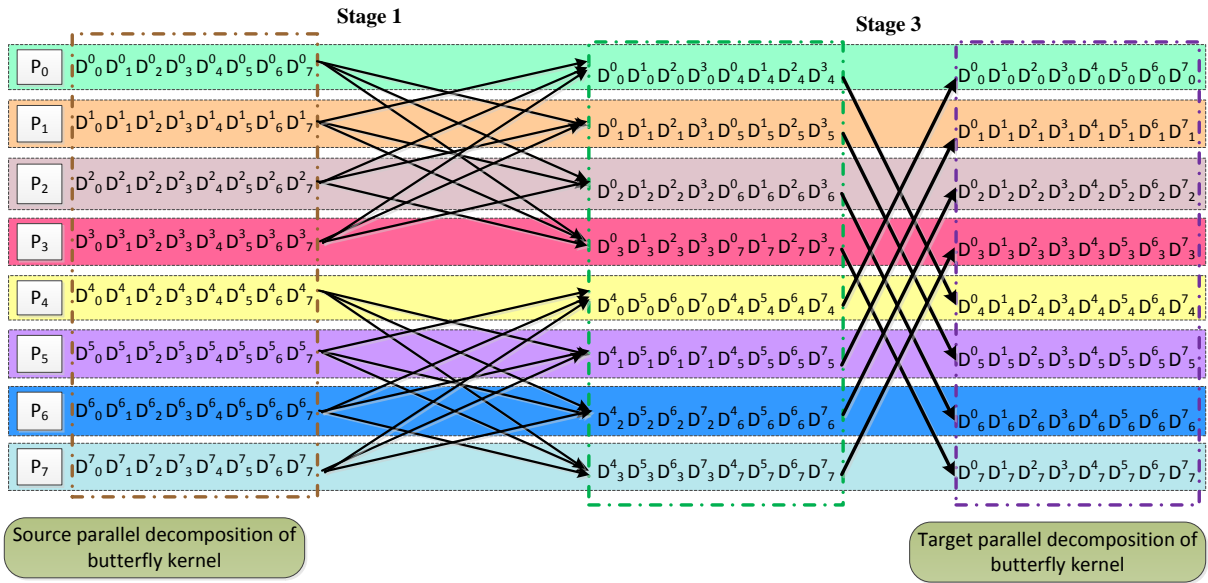
(b)



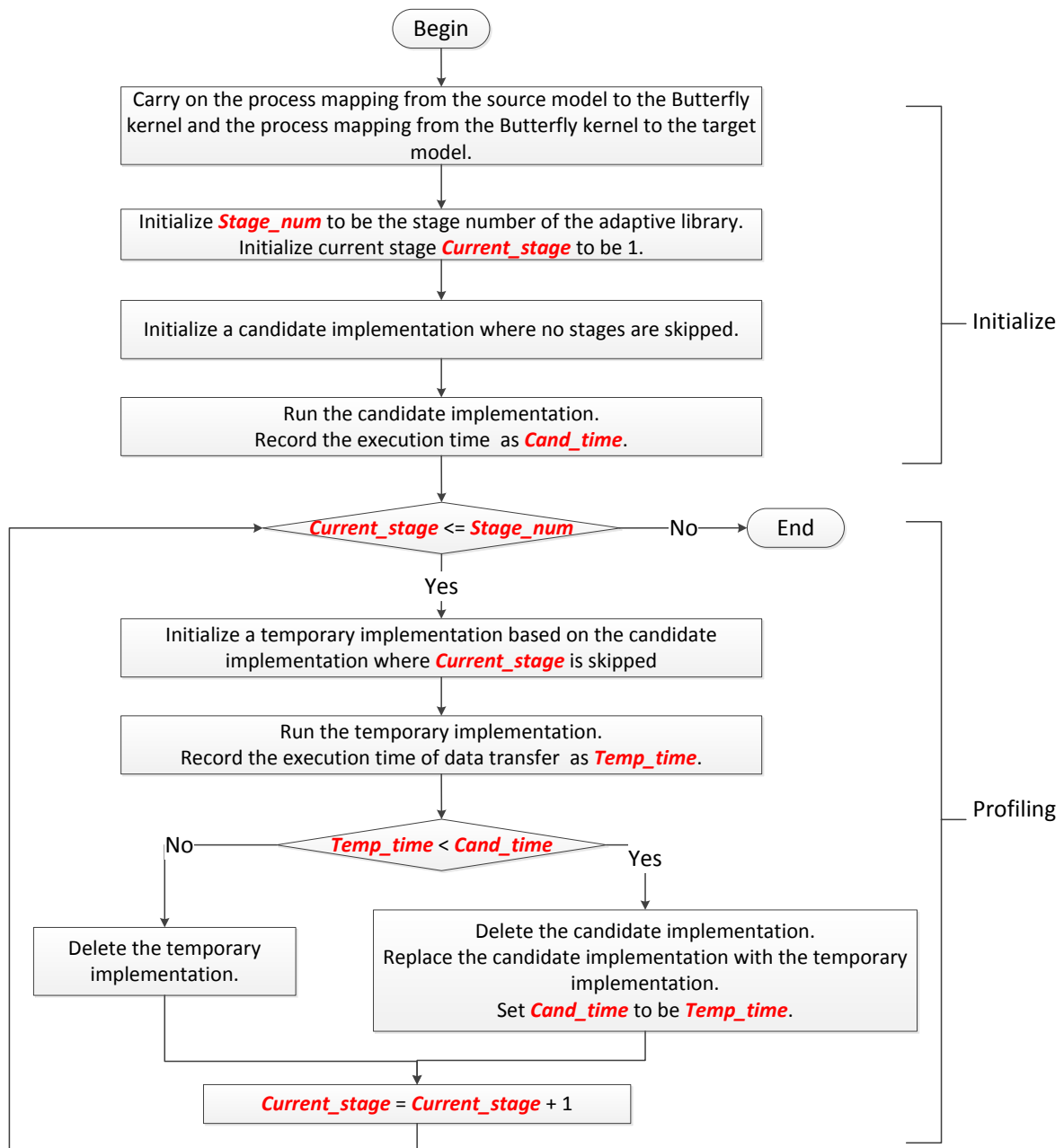
(c)

- 1
- 2 Figure 12. An example of process mappings, given that the sender has five processes (S_0 - S_4),
- 3 the receiver has 10 processes (R_0 - R_9) (there is no common process between the sender and
- 4 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

- 1 kernel, while every two processes of the receiver are mapped onto one process of the butterfly
- 2 kernel.



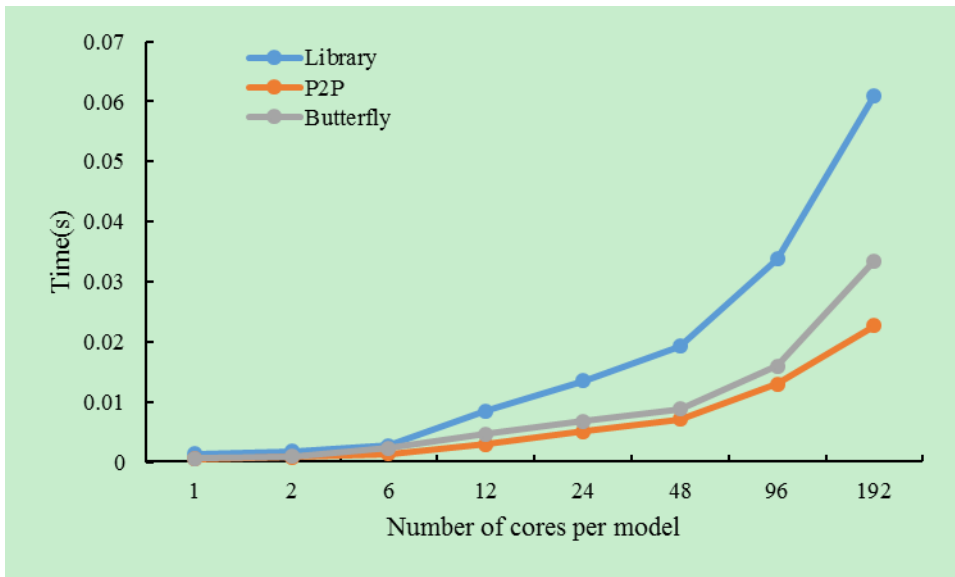
- 3
- 4 Figure 13. An example of the adaptive data transfer library with eight processes, where Stage
- 5 2 of the butterfly implementation is skipped with the P2P implementation of three MPI
- 6 messages per process.
- 7



1

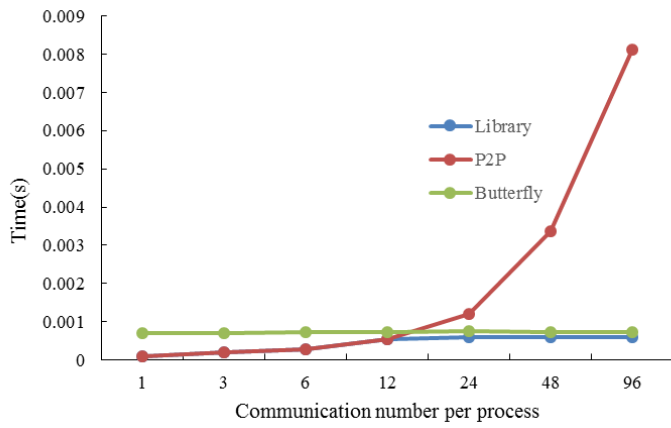
2 Figure 14. A flowchart for determining an appropriate implementation of the adaptive data
 3 transfer library.

4



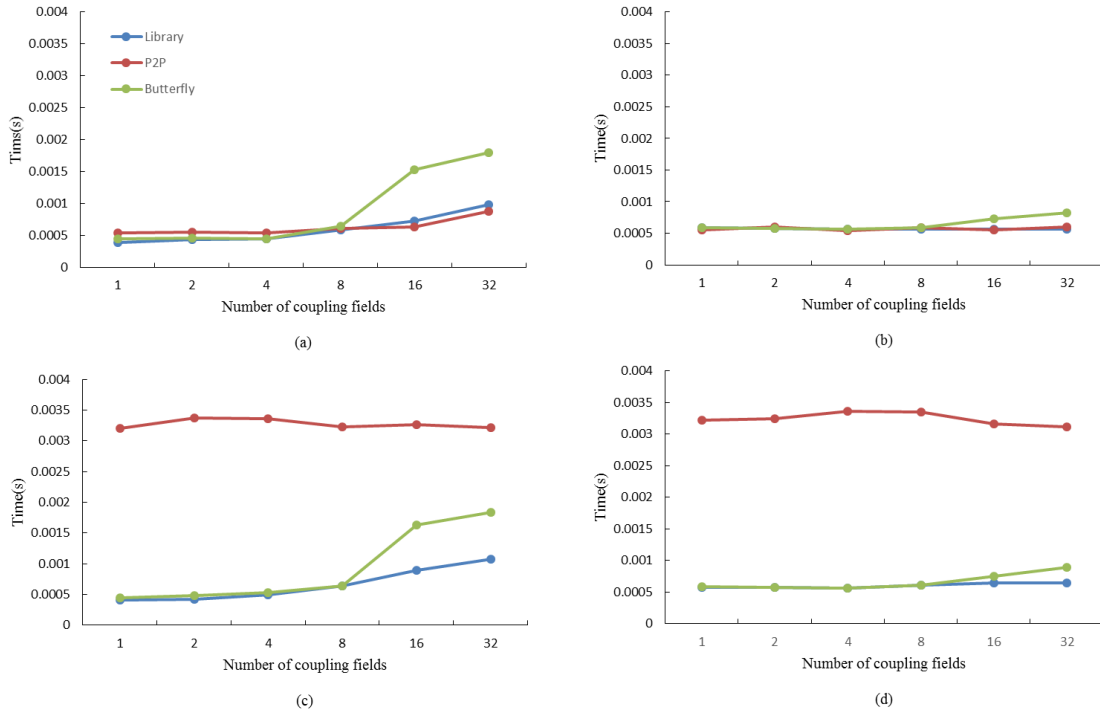
1
2
3
4
5
6
7
8

Figure 15. Initialization time (y-axis) of one data transfer between two toy models using a rectangular grid (of 192×96 grid points) when varying the number of cores used by each toy model (x-axis). There are 10 2-D coupling fields transferred from the source toy model to the target toy model. 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.



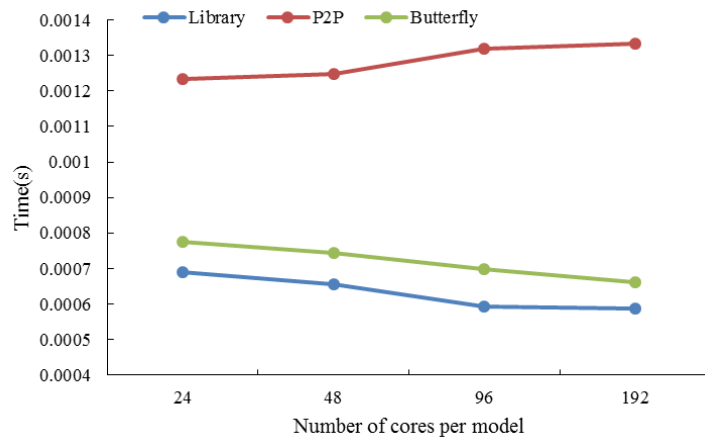
1
 2 Figure 16. Average execution time (y-axis) of one data transfer between two toy models with
 3 the same rectangular grid (of 192×96 grid points) when varying the MPI message number per
 4 process (x-axis). Each toy model is run with 192 cores. There are 10 2-D coupling fields
 5 transferred from the source toy model to the target toy model.

6

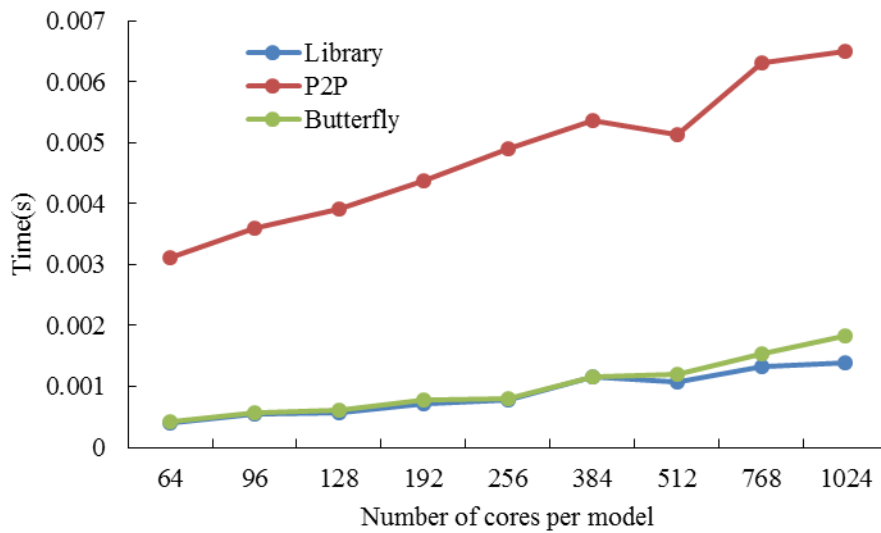


1
2 Figure 17. Average execution time (y-axis) of one data transfer between two toy models with
3 the same rectangular grid (of 192×96 grid points) when varying the number of coupling fields
4 transferred (x-axis). There are four simulation tests for the evaluation. In simulation (a), each
5 toy model is run with 48 cores, and the MPI message number per process is 12. In simulation
6 (b), each toy model is run with 192 cores, and the MPI message number per process is 12. In
7 simulation (c), each toy model is run with 48 cores, and the MPI message number per process
8 is 48. In simulation (d), each toy model is run with 192 cores (or processes), and the MPI
9 message number per process is 48.

10

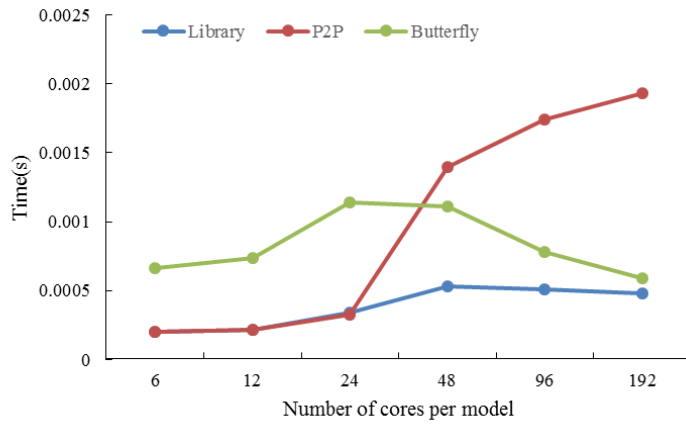


1
 2 Figure 18. Average execution time (y-axis) of one data transfer between two toy models with
 3 the same rectangular grid (of 192×96 grid points) when varying the number of cores used by
 4 each toy model (x-axis). There are 10 2-D coupling fields transferred from the source toy
 5 model to the target toy model. In each test, the MPI message number per process is set to 24.
 6

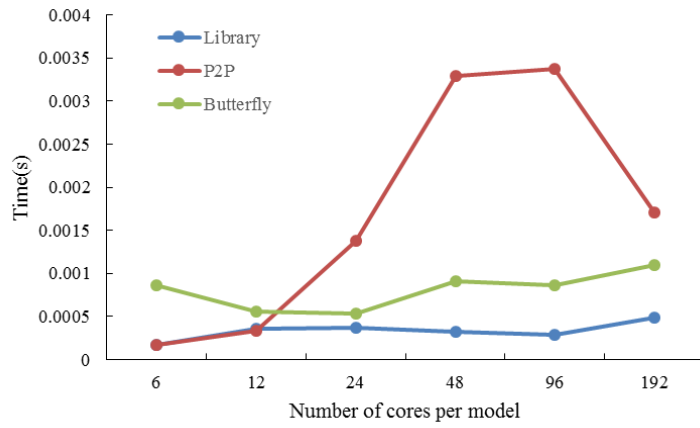


1
2
3
4
5
6

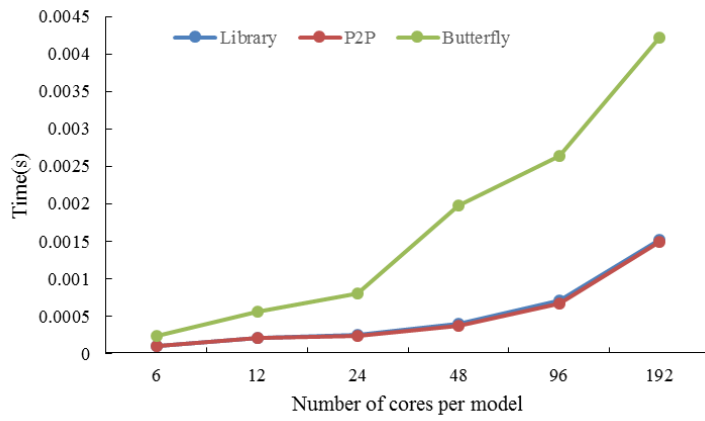
Figure 19. Average execution time (y-axis) of one data transfer between two toy models. In this evaluation, each process (running on a unique processor core) of the toy models have 96 grid points, while different processes have different message numbers and different message sizes. The number of coupling fields transferred is set to 20.



1
 2 Figure 20. 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 results are
 6 from the adaptive data transfer library which switches to the P2P implementation.
 7

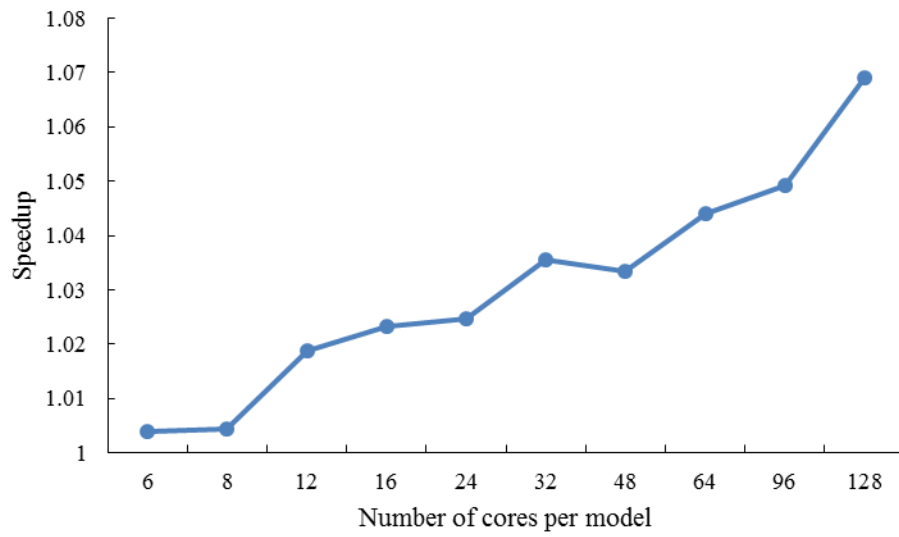


1
 2 Figure 21. Average execution time (y-axis) of one data transfer between the atmosphere
 3 model GAMIL2 and the land surface model CLM3 in GAMIL2-CLM3 when varying the
 4 number of cores used by each model (x-axis): 14 coupling fields on the GAMIL2 horizontal
 5 grid (the grid size is $128 \times 60 = 7680$) are transferred from the land surface model CLM3 to the
 6 atmosphere model GAMIL2.
 7



1
 2 Figure 22. Average execution time (y-axis) of one data rearrangement for the parallel
 3 interpolation from the atmosphere grid (the grid size is $144 \times 96 = 13824$) to the ocean grid (the
 4 grid size is $320 \times 384 = 122880$) in CESM when varying the number of cores used by each
 5 model (x-axis).

6



1

2 Figure 23. Performance improvement of the coupled model GAMIL2-CLM3 achieved by the
3 adaptive data transfer library, with the performance of GAMIL2-CLM3 using the P2P
4 implementation as the baseline.