# Part 1: Responses to Moritz Hanke

We thank Dr. Moritz Hanke for the comments and suggestions. We have modified the manuscript accordingly. In the following, we will reply them one by one.

1. The presented algorithm is basically a rendezvous algorithm, which uses a distributed directory. This was first introduced in [1]. I assume that you did not know about this paper and therefore did not reference it. This algorithm works by distributing data among the processes in a globally known decomposition. This is called "regular intermediate distribution" in the manuscript. Using this intermediate decomposition, accessing data without knowing the original decomposition and without gathering all data on a single process is easily possible. The use of distributed directories in a similar context is mentioned in the following references [2], [3], and [4].

**Response**: Thanks a lot for introducing distributed directory and the related works. We now know that this work as well as our future works deeply benefits from the general idea of distributed directory. The manuscript has been modified accordingly. Please refer to P4L95, P8L249, P9L257.

2. The manuscript describes the router network generation based on two predefined decompositions from two component models as being a fundamental functionality of a coupler. However, this routing table can also be a by-product, in case both components have different grids and the coupler generates interpolation weights online.

**Response**: A routing table can also make help in data interpolation. For example, MCT employs an interpolation decomposition and the corresponding data rearrangement that also follows on a routing table.

3. Line 30: "there is almost no evidence of scalable initialization of a coupler" [...] "(Craig et al., 2017; Liu et al., 2018)". Both papers mentioned in this context contain figures with initialisation cost measurements (see figure 2 and figure 8 respectively). Scaling behaviour in these figures is indeed sub-optimal. However, the cause for this scaling behaviour is not explicitly attributed by either paper to the router generation. Your manuscript indicates that in fact this is the case without providing evidence. By recreating figure 8 from the second paper with your new router generation implementation, you could have confirmed (at least for the C-Coupler) this. In Hanke et al., 2016 (cited by the manuscript) figure 3 (b) shows good scalability of the overall coupler initialisation for up to 3072 processes per component.

**Response**: Our statements here are incorrect. They have been modified. Please refer to P2L47 to P2L49. Moreover, experimental results based on MCT are added. Please refer to P2L49 and Fig. 2.

4. Line 56: "almost all existing couplers use the following 4 steps for generating a routing network" This is a very strong claim, which I would not support (see [2], [3], and [4]).

**Response**: This incorrect statement has been modified. Please refer to P2L59~P3L62.

5. Paragraph 3.2 and figure 2: This is the description of a basic parallel sorting algorithm. A shorter paragraph and a reference to a respective paper would have been enough.

**Response**: In the discussions, we have introduced the differences between the algorithm in Section 3.2 and a basic distributed sorting algorithm. In spite of the differences, we tried to make Section 3.2 as shorter as possible.

6. Line 167-168: "it only utilizes point-to-point communications and does not rely on any collective communication". The parallel sorting algorithm has a complexity of $O(log(n))$ and a similar communication pattern as the MPI implementations of various collective communication operations[5]. Therefore, you could argue that you actually did use collective communication, which you implemented using point-to-point communication.

**Response**: This incorrect statement has been modified into "it does not rely on any gather/broadcast communication". Please refer to P6L179.

7. Paragraph 4: This paragraph does not describe how these measurements where generated. Did you do a single run, average of multiple runs or average over multiple executions of the algorithm within a single run? It is possible, that especially the first execution of this algorithm produces for some MPI implementations a much higher run time than the following ones.

**Response**: We used average of multiple runs in our evaluation. Please refer to P7L199 in the revised manuscript.

8. Figure 3: I would have preferred to have absolute runtimes instead of speedups for the evaluation of the individual algorithms and speedups only for the direct comparison between

the two.

**Response**: Figure 3 (current Fig. 4) has been modified accordingly. Please refer to P15.

9. Tables 1 to 6: In my opinion, these tables add no significant value to the understanding of the algorithm.

**Response**: Considering some existing couplers such as MCT, OASIS3-MCT and CPL6/CPL7 can also benefit from DiRong1.0, we propose to reserve these tables that can show more details about the implementation.

10. Table 7: I assume the 1600 cores mean, that you used two toy components with 1600 cores each. I could not find an explicit description of this. The core counts are rather odd. You mentioned that your nodes have 24 processors each. Therefore, I would have assume, that the core counts in your tests are multiples of 24.

**Response**: The first column of Table 7 is the number of cores used in each toy component model. We have corrected it and please refer to P22. We used 1600 cores for each component model, because we can use at most 3200 cores on that computer and try to maximize the core number used by models.

11. Intel MPI library (3.2.2) Why did you use such an old version? In my experiments Intel MPI often performed very poorly. Maybe give the most recent OpenMPI a try.

**Response**: We have corrected for this mistake, the correct MPI version used for the evaluations in the manuscript is Intel MPI library (2018 Update 2). Please refer to P7L198 in the revised manuscript.

12. You propose to use a sorting algorithm with complexity O(log(n)) to generate the distributed directory. In some tests, I have seen that some MPI implementations introduce significant delays with such communication pattern, when being used for the first time in a run.

**Response**: We agree that significant delays can happen, while we would think that is the problem of network. Under an unstable network, any kind of MPI communications may suffer from significant delays.

13. Alternatively, each process could compute for all its local points the destination rank. Using alltoall, you can exchange the number of points that need to be sent, in order to get directly from the original to the intermediate decomposition. Afterwards, alltoallv can redistribute the data in a single communication call. Depending on the MPI implementation alltoall can also have a complexity of O(log(n)). However, very little data is exchanged and it can be highly optimized within the MPI. The communication matrix for the alltoallv is probably very sparse. Hence, it can be implemented by the user using point-to-point communication. In my tests, this approach delivers really good results.

**Response**: It is true that each process can easily compute the destination rank for all its local points, and then can easily prepare the parameters of sendcnts as well as sdisls for using MPI_alltoallv. However, each process cannot compute recvcnts as well as rdispls from its local points without extra communications, while recvcnts is necessary for using MPI_alltoallv. So, extra collective communications will be also required for computing recvcnts.

# Part 2: Responses to Vijay Mahadevan

We thank Dr. Vijay Mahadevan for the comments and suggestions. We have modified the manuscript accordingly. In the following, we will reply them one by one.

1. "First, the authors claim that "existing couplers employ an inefficient and unscalable global implementation for routing network generation that relies on collective communications. That's a main reason why the initialization cost of a coupler increases rapidly when using more processor cores."." "Can you provide some actual timings from the fully coupled climate solver runs to put the actual setup costs in perspective ? The scalability shown in Fig. (3) still indicate about 3.5s of compute time (since speedup for global routing network case is ≈1) for the 16M grid case on 1600 processes. If this accounts for say over 5% of the actual runtime of the solver, or a non-trivial percentage of total time to simulate a year (or days for high-res) of climate interactions, then 20x improvements in the setup cost could be quite impactful. However, such one time costs get amortized with physics setup costs for high-res runs, in addition to long-term temporal integration of the actual coupled simulation. Hence, I think the manuscript lacks a strong motivation, and provides only an incremental update to avoid the collective algorithms in the coupled simulation invoked during the initial setup phase. "

**Response**: This incorrect statement has been modified and please refer to P2L47~P2L49. Moreover, we measured the impact of routing network generation on the total initialization time of MCT in a coupled model (please refer to Figure 2). We did not evaluate the impact of DiRong1.0 on the total time of a model simulation, because this impact can be relative. We discussed about that and introduced more about the motivation in Section 5 (P8L228~P9L258).

2. Secondly, the global ID based partitioning strategy used in the distributed sort with DaRong to determine the communication pattern is not an innovative concept. There have been several algorithmic ideas based on graph partitioning strategies used in the parallel Sparse Matrix-Vector (SpMV) linear algebra context [1]. In a simplified sense, without a constraint on the message volume, data locality or latency of communication (such strategies may require repartitioning and/or task mapping), the globally unique ID space can be used in a round-robin type partitioning scheme. For instance, if the source component data are distributed on M processes, and destination on N processes, an implicit decomposition can be determined a-priori based on the global ID numbering that leads to MxN data redistribution. Such an implicit ID decomposition establishes a direct point-to-point communication pattern after which the CLGMT table can be created on both source and destination processes for further send/receive of DoF data at runtime. There may be a need for multiple rounds of rendezvous communication to establish message size for buffer allocation etc, but such an algorithm can eliminate collectives like broadcast and allreduce operations as necessary for better scalability.

**Response**: The global ID based partitioning strategy can be viewed as a precondition of the MxN data transfer, which is also a precondition of DiRong1.0. We made a slight modification at P4L99 accordingly. The sparse MxN data transfer has been widely used in existing couplers, and this manuscript focuses on how to accelerate the routing network generation for initializing the MxN data transfer. Current implementations of routing network generation rely on gather/scatter communications, while this manuscript might correspond to "There may be a need for multiple rounds of rendezvous communication to establish message size for buffer allocation". There are a lot of algorithms or optimizations based on the butterfly communication network, which has been stated at P5L138 ~ P5L139. It is not an innovative concept indeed. That's why we just give a simple example (Fig. 3) in the manuscript.

3. Line 29-30: "Although most existing couplers achieve scalable data transfer and data interpolation, i.e., the data transfer and data interpolation generally can be faster when using more processor cores, there is almost no evidence of scalable initialization of a coupler." Total cost of a coupled solver includes both the setup/initialization cost and the runtime remap operator application and data-transfer at every time step per coupled component pair/field. Hence, cost of initialization often gets amortized in a climate simulation run. As mentioned previously, please cite or provide data to substantiate such strong claims, preferably with real results using MCT and C-Coupler2 runs.

**Response**: This incorrect statement has been corrected. Please refer to P2L47~P2L49. We have deleted this description and described the scalability of the current coupler more precisely in Section 1 (Line49). Moreover, we measured the impact of routing network generation on the total initialization time of MCT in a coupled model (please refer to Figure 2) and give the related discussions in Section 5 (P8L228~P9L258).

4. Paragraph starting at line 79 is confusing. Please rephrase the sentence better to clearly describe the particular step and its time complexity that leads to the inefficient and non-scalable implementation of routing network generation.

**Response**: The numbers of steps have been added. Please refer to P3L84 and P3L85.

5. Line 104: "Specifically, we employ a regular intermediate distribution that evenly distributes the CLGMT entries among processes based on the ascending order of the global cell index. Such an intermediate distribution is not only simple, but also enables to easily achieve the rearrangement to the intermediate distribution via a sorting procedure similar to distributed sort. " As noted previously, the GSMap and Router infrastructure in MCT already has such options

to redistribute data based on Global ID numbering. This is inherently what has been described here as the intermediate distribution of the CLGMT. The primary difference seems to be that GSMap is a O(P) data structure that grows with core counts, and is accumulated on all process through a gather on root and a subsequent broadcast. The authors of the current paper are trying to avoid this one-time collective operation, which could be an over optimization considering the total runtime of the climate solver.

**Response**: Yes, we are trying to avoid this "one-time collective operation", because it becomes more important in C-Coupler development, which has been discussed in Section 5 (P8L228~P9L258).

6.  The paragraph starting at Line 171 can be rewritten in the context of the distributed sort workflow shown in Fig. (2).

**Response**: This paragraph is about the whole implementation of DiRong1.0 (corresponding to the major steps introduced in Section 3.1), while the original Fig. 2 (current Figure 3) is only about the first major step (Section 3.2) and the fourth major step (Section 3.5).

7.  It will be helpful to explicitly mention the time and memory complexity for each stage in a table format, for both the global and the DaRong algorithm so that the reader can immediately get a sense of the actual improvement.

**Response**: Thanks a lot for this suggestion. We do not provide a table in the revised manuscript currently, because the paragraph starting from P3L84 states the complexity of each step in the global implantation. We will add a table if it is still required.

8.  The weak scalability results shown in Table. (7) are not uniform since the grid sizes are doubled, but the core counts do not, going from 250 to 450 and 900 to 1600. Please rerun these calculations with P=[200,400,800,1600] instead

**Response**: Table 7 has been updated with new results accordingly. Please refer to P22.

9.  Section (3.2) title: "Rearranging CLGMT entries intra a component model". Please rephrase. Do you mean to say "between component models" ? Same comment applies to Section (3.5) title as well.

**Response**: We confirm that Section 3.2 and 3.5 are for "intra a component model" but not for

"between component models".

10. Line 175, "SPT" should be "SRT" ?

**Response**: Thanks a lot. The current P7L187 has been corrected.

Part 3: a marked-up manuscript version

# DiRong1.0: a distributed implementation for improving routing network generation in model coupling ~~A new distributed algorithm for routing network generation in model coupling and its evaluation based on C-Coupler2~~

Hao Yu[1], Li Liu[1], Chao Sun[1], Ruizhe Li[1], Xinzhu Yu[1], Cheng Zhang[1], Zhiyuan Zhang[2], Bin Wang[1,3]

[1] Ministry of Education Key Laboratory for Earth System Modeling, Department of Earth System Science, Tsinghua University, Beijing, China

[2] Hydro-Meteorological Center of Navy China, Beijing China

[3] State Key Laboratory of Numerical Modeling for Atmospheric Sciences and Geophysical Fluid Dynamics (LASG), Institute of Atmospheric Physics, Chinese Academy of Sciences, Beijing, China

*Correspondence to*: Li Liu (liuli-cess@tsinghua.edu.cn)

**Abstract.** It is a fundamental functionality of a coupler for Earth system modeling to efficiently handle data transfer between component models. As an approach of *MxN* communication following a routing network has been used in existing couplers for achieving data transfer, ~~R~~routing network generation is generally a major step for initializing the *MxN* communication~~data transfer functionality~~. ~~Most~~ Some existing couplers such as MCT and C-Coupler employ an inefficient ~~and unscalable~~ global implementation for routing network generation that relies on ~~collective~~gather/broadcast communications. That's a~~n important~~ ~~main~~ reason why the initialization cost of a coupler increases ~~rapidly~~ when using more processor cores. In this paper, we propose a ~~new~~ **D**istributed ~~i~~mplementation~~algorithm~~ for **Ro**uting **n**etwork **g**eneration, version 1.0 (~~DaRong~~DiRong1.0), which does not introduce any gather/broadcast~~collective~~ communication ~~and achieves much lower complexities than the global implementation~~. The empirical evaluations show that ~~DaRong~~DiRong1.0 is ~~of course~~ much more efficient ~~and scalable~~ than the global implementation~~, which has been further demonstrated via empirical evaluations~~. ~~DaRong~~DiRong1.0 has already been implemented in C-Coupler2. We believe that ~~existing and future~~some other couplers can also benefit from it.

## 1 Introduction

A coupled model regarding Earth system Modelling and numerical weather forecasting generally highly depends on existing couplers (Hill et al., 2004; Craig et al., 2005; Larson et al., 2005; Balaji et al., 2006; Redler et al., 2010; Craig et al., 2012; Valcke, 2013; Liu et al., 2014; Hanke et al., 2016; Craig et al., 2017; Liu et al., 2018), each of which can combine different

component models into a whole system and handle data interpolation between different model grids and data transfer between
30  component models (Valcke, 2012).

The functionality of data interpolation generally takes two major steps, i.e., preparing remapping weights that are from an
offline file or from online calculation when initializing the coupler, and conducting parallel interpolation calculation based on
the sparse matrix-vector multiplication with the remapping weights throughout the coupled model integration. ~~In response to~~
35  ~~more and more computation resulting from higher and higher resolutions in model development, the parallel efficiency of a~~
~~coupled model on modern high-performance computers becomes more and more critical. Any module in a coupled model,~~
~~including the coupler, can impact or even may damage the parallel efficiency of the whole coupled model. Although most~~
~~existing couplers achieve scalable data transfer and data interpolation, i.e., the data transfer and data interpolation generally~~
~~can be faster when using more processor cores, there is almost no evidence of scalable initialization of a coupler. Experiences~~
40  ~~from OASIS3-MCT and C-Coupler2 have revealed that the initialization cost of a coupler increases rapidly when using more~~
~~processor cores (Craig et al., 2017; Liu et al., 2018). To achieve scalable initialization of couplers, this paper tries to make a~~
~~first step through focusing on the initialization of data transfer.~~

The functionality of data transfer of couplers is transferring scalar variables or fields on a model grid (called gridded fields
45  hereafter) from one component model to another via MPI (Message Passing Interface). A component model generally has been
parallelized through decomposing the cells of a model grid into distinct subsets each of is assigned to an MPI process for
cooperative concurrent computation, e.g., the sample parallel decompositions in Fig. 1a and 1b. To efficiently transfer gridded
fields in parallel, Jacob et al. (2005) proposed an approach of $MxN$ communication (called $MxN$ approach) following ~~the~~ a
routing network where each pair of processes from the two component models should have a communication connection only
50  when they have common grid cells (for example, Fig. 1c). This $MxN$ approach has already been used in existing couplers for
more than ten years. As the parallel decompositions of component models generally keep constant throughout the whole
integration, a routing network can also keep constant. Thus, the $MxN$ approach can be achieved with two major steps:
generating the routing network when initializing the coupler, and transferring gridded fields based on the routing network
throughout the coupled model integration. ~~In spite of the scalability of the second major step, the first major step in existing~~
55  ~~couplers is unscalable, introducing higher cost when using more processor cores.~~

In response to more and more computation resulting from higher and higher resolutions in model development, the parallel
efficiency of a coupled model on modern high-performance computers becomes more and more critical. Any module in a
coupled model, including the coupler, can impact the parallel efficiency of the whole coupled model. Most existing couplers
60  achieve scalable data transfer and data interpolation throughout the coupled model integration, i.e., the data transfer and data
interpolation generally can be faster when using more processor cores, while experiences from OASIS3-MCT and C-Coupler2
have shown that the initialization cost of a coupler can increase rapidly when using more processor cores (Craig et al., 2017;

2

Liu et al., 2018). A further investigation in Fig. 2 based on MCT shows that the initialization of data transfer, i.e., generating routing networks, is an important source of the initialization cost.

65

To lower the initialization cost of a coupler, this paper tries to make a first step through focusing on the generation of routing networks

~~In this paper, we~~, and proposes a new **D**istributed **i**mplementation~~algorithm~~ for **Ro**uting **n**etwork **g**eneration, version 1.0 (~~DaRong~~DiRong1.0). The evaluation based on C-Coupler2 shows that, it ~~, which~~ is much faster ~~and consumes much less~~

70 ~~memory~~ than the existing approach. The remainder of this paper is organized as follows. We ~~reveal the causes of unscalability~~investigate ~~of~~ the existing implementations in Section 2, present and then evaluate ~~DaRong~~DiRong1.0 in Section 3 and 4 respectively, and conclude and discuss this work in Section 5.

**2 Existing implementations of routing network generation**

75 In ~~existing~~some couplers such as MCT and C-Coupler, the global information of a parallel decomposition ~~generally~~ is originally distributed among all processes of a component model, where a process only records its local parallel decomposition corresponding to the grid cells assigned to it. Thus, ~~almost all existing~~these couplers generally use the following 4 steps for generating a routing network between the parallel decompositions of a source (*src*) and a destination (*dst*) component model.

1) Gathering global parallel decomposition: the *src*/*dst* root process gathers the global information of the *src*/*dst* parallel
80   decomposition from all *src*/*dst* processes.

2) Exchanging global parallel decomposition: the *src*/*dst* root process first exchanges the *src*/*dst* global parallel decomposition with the *dst*/*src* root process, and then broadcasts the *dst*/*src* global parallel decomposition to all *src*/*dst* processes.

3) Detecting common grid cells: each *src*/*dst* process detects its common grid cells with each *dst*/*src* process based on its
85   local parallel decomposition and the *dst*/*src* global parallel decomposition.

4) Generating the routing network: each *src*/*dst* process generates its local routing network according to the information about common grid cells.

Given that each of the *src* and *dst* component models uses $K$ processes and the corresponding grid size is $N$ (the grid has $N$
90 cells), the first and second steps correspond to gather/broadcast~~collective~~ communications with the time complexity of at least $O(N*logK)$ and the memory complexity of $O(N)$. Regarding the third step, the average time complexity corresponding to MCT (as well as CPL6/CPL7 and OASIS3-MCT that employ MCT for data transfer) on each *src*/*dst* process can be~~is~~ $O(N*N/K)$, because the main loop of this step consists of two levels, i.e., the first level corresponds to the local parallel decomposition (the average number of cells in the local parallel decomposition is $N/K$) while the second level corresponds to the *dst*/*src* global

3

95　parallel decomposition. The average time complexity of the third step corresponding to C-Coupler is $O(N)$, as C-Coupler first generates a map corresponding to the global parallel decomposition and next detects common cells based on looking up the map. Although this implementation can lower the time complexity, but introduces inefficient irregular memory accesses. As the last step does not depend on any global parallel decomposition, its average time complexity is $O(N/K)$.

100　Determined by the gather/broadcast~~collective~~ communications (the first and second steps) and the corresponding time complexity of $O(N*logK)$, and the time complexity of $O(N*N/K)$ or $O(N)$ corresponding to common grid cells detection (the third step), such existing implementations of routing network generation are of course inefficient ~~and unscalable~~ under the increment of processor cores. Moreover, in response to the memory complexity of $O(N)$, more memory will be consumed ~~when~~with ~~the~~finer model grids ~~get finer~~.

105

In the following context, the existing implementations relying on gather/broadcast communications~~of routing network generation~~ are called global routing network generation.

**3 Design and implementation**

110

**3.1 Overall design**

The design and implementation of DiRong1.0 significantly benefits the generally idea of distributed directory (Pinar and Hendrickson, 2001) that has already been used in coupler development (Theurich et al, 2008, Hanke et al., 2016), and different kinds of specific distributed directories are defined and used in DiRong1.0.

115

To generate the $MxN$ data transfer, ~~E~~each cell of a grid can be numbered with a unique index from 1 to $N$ (called global cell index), while each grid cell assigned to the same process can also be numbered with a unique local cell index. Thus, the information of a given parallel decomposition can be recorded as a Cell Local-Global Mapping Table (CLGMT), each element of which is a triple of global cell index, process ID, and local cell index. For example, Tables 1 and 2 are the CLGMTs

120　corresponding to the parallel decompositions in Fig. 1a and Fig. 1b respectively.

Generally, the CLGMT entries of a parallel decomposition are distributed among the processes of a component model, which means a process only stores a part of the CLGMT. The key idea of the existing global implementation can be summarized as reconstructing the global CLGMT of the peer parallel decomposition in each process for routing network generation. To be a~~n~~ efficient ~~scalable~~ solution, ~~DaRong~~DiRong1.0 should be fully based on distributed CLGMT without reconstructing any global

4

CLGMT. The reason why the existing global implementations have to depend on global CLGMTs is because the distribution of the CLGMT entries is determined by a model and thus a coupler generally has to view any distribution as random.

Motivated by the above analysis, the key challenge to ~~DaRong~~DiRong1.0 becomes how to efficiently rearrange the original
130 distribution of the CLGMT entries of a given parallel decomposition into a regular intermediate distribution and how to efficiently generate the routing network based on the intermediate distribution. Specifically, we employ a regular intermediate distribution that evenly distributes the CLGMT entries among processes based on the ascending order of the global cell index. Such an intermediate distribution is not only simple, but also enables to easily achieve the rearrangement to the intermediate distribution via a sorting procedure similar to distributed sort. With the above preparations, ~~DaRong~~DiRong1.0 is designed
135 with the following major steps for generating a routing network between the *src* and *dst* component models:

1) The *src*/*dst* component model rearranges the original distribution of the CLGMT entries of the *src*/*dst* parallel decomposition into the regular intermediate distribution.

2) The *src* and *dst* component models exchange the CLGMT entries based on the intermediate distributions.

3) Each *src*/*dst* process generates table entries of the sharing relationship about how each grid cell is shared between the
140 processes of the *src* and *dst* component models, based on the *src* and *dst* CLGMT entries on the intermediate distributions.

4) The *src*/*dst* component model rearranges the intermediate distribution of the entries of the sharing relationship table (SRT) into the original distribution of the CLGMT entries of the *src*/*dst* parallel decomposition.

5) Each *src*/*dst* process generates its local routing network based on the local SRT entries.

145 In the following context of this section, we will detail the implementation of each major step except the last one because it is similar to the last major step in the global implementation.

### 3.2 Rearranging CLGMT entries intra a component model

Such rearrangement is achieved via a divide-and-conquer sorting procedure that is similar to a merge sort with the keyword
150 of global cell index. This procedure first sorts the CLGMT entries locally in each process, and next iteratively conducts distributed sort by a main loop of $logK$ iterations ($K$ is the number of processes of the *src*/*dst* component model). In an iteration, processes are divided into distinct pairs and the two processes in each pair swap the CLGMT entries based on a point-to-point communication. Figure 3~~2~~ shows an example of the distributed sort corresponding to the CLGMT entries in Table 1, and Table 3 shows the distributed CLGMT after rearranging the CLGMT entries in Table 2. ~~As shown in Fig. 3, the distributed sort
155 employed in DiRong1.0 uses a similar butterfly communication pattern as various optimized collective communication operations and an optimized matrix-vector multiplication (Brooks, 1986; Hendrickson et al, 1995; Thakur et al., 2005).

5

**3.3 Exchanging CLGMT entries between component models**

After the rearrangement of the CLGMT in a component model, the CLGMT entries are sorted in an ascending order of the global cell indexes and evenly distributed among processes. The CLGMT entries reserved in each process therefore have a determinate and non-overlapping range of global cell indexes, and such a range can be easily calculated from the grid size, the number of total processes, and process ID. Thus, it is easy to calculate the overlapping relationship of global cell index range between a *src* process and a *dst* process. As it is only necessary to exchange CLGMT entries between a pair of *src* and *dst* processes with overlapping ranges, point-to-point communications only are enough for handling the exchange of the CLGMT entries.

**3.4 Generation of SRT**

After the previous major step, each process reserves two sequences of CLGMT entries corresponding to the *src* and *dst* parallel decompositions respectively. Given that the two sequences contain *n1* and *n2* entries respectively, the time complexity of detecting the sharing relationship is $O(n1+n2)$, because the entries in each sequence have already been ordered in ascending global cell indexes, and a procedure similar to the kernel of merge sort that merges two ordered data sequences can handle such detection.

To record the sharing relationship, a SRT entry is designed as a quintuple of global cell index, *src* process ID, *src* local cell index, *dst* process ID, and *dst* local cell index. Given a quintuple $<q_1,q_2,q_3,q_4,q_5>$, it means that number $q_3$ local cell in number $q_2$ process of the *src* component model is number $q_1$ global cell, and the data on it will be transferred to number $q_5$ local cell in number $q_4$ process of the *dst* component model. Table 4 shows the SRT in the *src* component model, calculated from the rearranged distributed CLGMT entries in Fig. 3~~2~~ and Table 3.

It is possible that multiple *src* CLGMT entries correspond to the same global cell index. Under such a case, any *src* CLGMT entry can be used for generating the corresponding SRT entries, because the *src* component model should guarantee that the data copies on the same grid cell are exactly the same. Given a *dst* CLGMT entry, if there is no *src* CLGMT entry with the same global cell index, no SRT entry will be generated. Given that multiple *dst* CLGMT entries correspond to the same global cell index and there is at least one *src* CLGMT entry with the same global cell index, a SRT entry will be generated for each *dst* CLGMT entry.

6

### 3.5 Rearranging SRT entries intra a component model

After the previous major step, the SRT entries are distributed among processes of a component model according to the intermediate distribution. As a process can use only the SRT entries corresponding to its local cells for the last major step of local routing network generation, the SRT entries should be rearranged among the processes of a component model. We find that such rearrangement can also be achieved via a sorting procedure similar to the distributed sort with the keyword of *src/dst* process ID, or even the sorting procedure implemented for the first major step can be reused. Tables 5 and 6 show the SRT entries distributed in the *src* and *dst* component model respectively, after the rearrangement.

### 3.6 Time complexity and memory complexity

As ~~DaRong~~DiRong1.0 does not reconstruct the global CLGMT, it ~~only utilizes point-to-point communications and~~ does not rely on any gather/broadcast~~collective~~ communication, and its average memory complexity is $O(N/K)$ on each process. As the implementation of its most time-consuming major steps are similar to a merge sort, and the time complexity of merge sort is $O(N*logN)$, the average time complexity of ~~DaRong~~DiRong1.0 on each process is $O(N*(logN)/K)$, and average communication complexity is $O(N*(logK)/K)$.

To facilitate the implementation of the sorting procedure, we force the number of processes regarding the $1^{st} \sim 4^{th}$ major steps (Section 3.1) to be the maximum power of 2 ($2^n$) no larger than the total process number of the *src/dst* component model. For a process whose ID $I$ is not smaller than $2^n$, its CLGMT entries will be merged into the process with the ID of $I$-$2^n$ before the first major step, and the S~~R~~PT entries corresponding to it will be obtained from the process with the ID of $I$-$2^n$ after the fourth major step. This strategy will not change the above time complexity and memory complexity of ~~DaRong~~DiRong1.0, as $2^n$ is larger than a half of the total process number.

### 4 Evaluation

For evaluating ~~DaRong~~DiRong1.0, we implemented it in C-Coupler2, which enables us to compare it with the original global routing network generation in C-Coupler2. We developed a toy coupled model consisting of two toy component models and C-Coupler2 for the evaluation, which enables us to flexibly change the model settings in terms of grid size and number of processor cores (processes). The toy coupled model is run on a supercomputer, where each computing node on the supercomputer includes two Intel Xeon E5-2678 v3 CPUs (Intel(R) Xeon(R) CPU (24 processor cores in total)), and all computing nodes were connected with an InfiniBand network. The codes were compiled by an Intel Fortran and C++ compiler at the optimization level O2, using an Intel MPI library (2018 Update 2~~3.2.2~~). A maximum number of 3200 cores are used for running the toy coupled model. All test results are from the average of multiple runs of the toy coupled model.

7

We made an evaluation under the variation of process numbers (Fig. 4~~3~~; two component models use the same number of processor cores). For the grid size of 500,000 (Fig. 4~~3~~a), the execution time of ~~DaRong~~DiRong1.0 does not really decrease when using more processor cores. This result is reasonable although it does not match the time complexity of ~~DaRong~~DiRong1.0. The communication complexity of ~~DaRong~~DiRong1.0 is O(*N\*(logK)/K*), where *logK* stands for the number of point-to-point communications in each process and *N/K* stands for the average message size in each communication. The average message size corresponding Fig. 4~~3~~a is small (about 160KB under 60 cores while about 6KB under 1600 cores for each toy component model), while the execution time of point-to-point communication cannot keep linear to the message size and may be unstable when the message size is small. Different from ~~DaRong~~DiRong1.0, the execution time of the global implementation increases rapidly with the increment of core number. As a result, ~~DaRong~~DiRong1.0 outperforms the global implementation more significantly when using more cores. When the grid size gets larger (e.g., 4,000,000 in Fig. 4~~3~~b and 16,000,000 in Fig. 4~~3~~c), ~~DaRong~~DiRong1.0 still significantly outperforms the global implementation, while with better scalability.

Considering a model can use more processor cores for acceleration when its resolution gets finer, we further evaluated the weak scalability of ~~DaRong~~DiRong1.0, where we concurrently increased the grid size and core number to achieve similar numbers of grid points per process. As shown in Table 7, the execution time of ~~DaRong~~DiRong1.0 increases slowly while the execution time of the global implementation increases rapidly with the increment of grid size and core number. This demonstrates that ~~DaRong~~DiRong1.0 achieves much better weak scalability than the global implementation.

**5 Conclusion and discussion**

In this paper, we propose a new distributed implementation~~algorithm~~, ~~DaRong~~DiRong1.0, for routing network generation. As it does not introduce any gather/broadcast~~collective~~ communication and achieves much lower complexity in terms of time, memory and communication than the global implementation ~~that is widely used in existing couplers~~, it is of course much more efficient ~~and scalable~~ than the global implementation. The evaluation results further demonstrate this conclusion.

~~DaRong~~DiRong1.0 has already been implemented in C-Coupler2. Its code is publicly available in a C-Coupler2 version and will be further used in future C-Coupler versions. We do believe that some existing couplers such as MCT, OASIS3-MCT and CPL6/CPL7, can also benefit from ~~DaRong~~DiRong1.0, for accelerating ~~routing~~ the routing network generation as well as the coupler initialization ~~of couplers~~.

We did not evaluate the impact of DiRong1.0 on the total time of a model simulation, because this impact can be relative. The overhead of routing network generation as well as coupler initialization will be trivial under a long simulation (e.g., hundreds

8

of model days or even hundreds of model years), but may be significant for a short simulation (e.g., several model days or even several model hours in weather forecasting (Palmer et al., 2008; Hoskins, 2013)). In a data assimilation for weather forecasting, it can be required to start a model run just for only several model hours or even shorter. Regarding an operational model, there is generally a time limitation of producing forecasting results (for example, finishing 5-day forecasting in two hours), and thus developers always have to carefully optimize various software modules especially when the model resolution gets finer. In fact, we have been asked to accelerate the initialization of C-Coupler2 for an operational coupled model used in China, and that's a main reason why we developed DiRong1.0.

Another main reason why we developed DiRong1.0 is that, routing network generation will become more important along with the development of C-Coupler. Recently, a new framework for weakly coupled ensemble data assimilation (EDA) based on C-Coupler2, named DAFCC1 (Sun et al., 2020), was developed. DAFCC1 will be an important part in C-Coupler3, the next version of C-Coupler. Given a coupled EDA system and that users want the atmosphere component to perform EDA, DAFCC1 will automatically generate an ensemble component corresponding to all ensemble members of the atmosphere component for calling the DA algorithm, and will automatically conduct routing network generation for the data transfers between the ensemble component and each ensemble member. Thus, routing network generation will be more frequently used in EDA with DAFCC1. For example, given that there are 50 ensemble members, the routing network generation with the ensemble component will be conducted at least 50 times.

We note that, the current sequential read of a remapping weight file is another bottleneck of C-Coupler2. Similar to Hanke et al. (2016), we will design a specific distributed directory for reading in the remapping weights in parallel, while will enable to efficiently redistribute remapping weights among processes based on DiRong1.0. Currently, C-Coupler2 employs a simple global representation for horizontal grids, which means that each process keeps all points of a horizontal grid. The global representation will become a bottleneck in at least two aspects. First, it will consume too much memory to run a model simulation. For example, given a horizontal grid with 16,000,000 points, the memory for keeping it in each process will be about 1.3GB (given that each point has four vertexes and the data type is double precision), which is a large memory requirement. Second, initialization of the data interpolation functionality requires exchanging model grids between different component models, which introduces global communications (e.g., broadcast) for the global grid representations. To address this bottleneck, we will design and develop a distributed grid representation that can be viewed as a specific distributed directory, will enable to efficiently redistribute horizontal grid points among processes based on DiRong1.0.

**References**

Balaji, V., J. Anderson, I. Held, M. Winton, J. Durachta, S. Malyshev, and R. J. Stouffer, 2006: The Exchange Grid: a mechanism for data exchange between Earth System components on independent grids. In: Proceedings of the 2005 International Conference on Parallel Computational Fluid Dynamics, College Park, MD, USA, Elsevier, 2006.

E. D. Brooks,1986: The Butterfly Barrier. International Journal of Parallel Programming,15(4):295–307

Craig, A. P., M. Vertenstein, and R. Jacob, 2012: A New Flexible Coupler for Earth System Modelling developed for CCSM4 and CESM1. Int. J. High Perform. C, 26-1, 31–42, doi:10.1177/1094342011428141.

Craig, A. P., R. L. Jacob, B. Kauffman, T. Bettge, J. W. Larson, E. T. Ong, C. H. Q. Ding, Y. He, 2005: CPL6: The New Extensible, High Performance Parallel Coupler for the Community Climate System Model. International Journal of High Performance Computing Applications, 19(3): 309-327.

Craig, A., Valcke, S., and Coquart, L.: Development and performance of a new version of the OASIS coupler, OASIS3-MCT_3.0, Geosci. Model Dev., 10, 3297-3308, https://doi.org/10.5194/gmd-10-3297-2017, 2017

Hanke, M., Redler, R., Holfeld, T., and Yastremsky, M.: YAC 1.2.0: new aspects for coupling software in Earth system modelling, Geosci. Model Dev., 9, 2755–2769, https://doi.org/10.5194/gmd-9-2755-2016, 2016
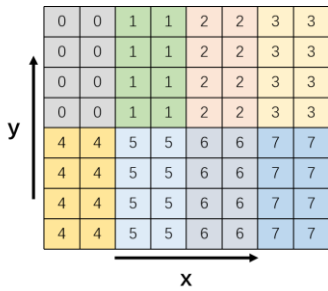
Hendrickson, B., R. Leland, and S. Plimpton. An efficient parallel algorithm for matrix-vector multiplication. International Journal of High Speed Computing 7, no. 01(1995): 73-88.

Hill, C., C. DeLuca, V. Balaji, M. Suarez, and A. da Silva, 2004: Architecture of the Earth System Modelling Framework. Comput. Sci. Eng., 6, 18–28.

Hoskins, B.: The potential for skill across the range of the seamless weather-climate prediction problem: a stimulus for our science, Q. J. Roy. Meteor. Soc., 139, 573-584, 2013.

Jacob, R., J. Larson, and E. Ong, 2005: M x N Communication and Parallel Interpolation in Community Climate System Model Version 3 Using the Model Coupling Toolkit, Int. J. High. Perform. C, 19, 293–307.

Larson, J., R. Jacob, and E. Ong, 2005: The Model Coupling Toolkit: A New Fortran90 Toolkit for Building Multiphysics Parallel Coupled Models. Int. J. High Perform, C, 19, 277–292.

Liu, L., Yang, G., Wang, B., Zhang, C., Li, R., Zhang, Z., Ji, Y., and Wang, L.: C-Coupler1: a Chinese community coupler for Earth system modeling, Geosci. Model Dev., 7, 2281-2302, doi:10.5194/gmd-7-2281-2014, 2014.

Liu, L., Zhang, C., Li, R., Wang, B., and Yang, G.: C-Coupler2: a flexible and user-friendly community coupler for model coupling and nesting, Geosci. Model Dev., 11, 3557-3586, https://doi.org/10.5194/gmd-11-3557-2018, 2018.

Palmer, T. N., Doblas-Reyes, F. J., Weisheimer, A. and Rodwell, M. J.: Toward seamless prediction: Calibration of climate change projections using seasonal forecasts, Bull. Am. Meteorol. Soc., 89, 459-470, 2008.

Pinar, A. and Hendrickson, B.: Communication Support for Adaptive Computation, Proceedings of the Tenth SIAM Conference on Parallel Processing for Scientific Computing, Portsmouth, Virginia, USA, 12–14 March 2001, SIAM, 2001

R. Redler, S. Valcke and H. Ritzdorf. OASIS4 - a coupling software for next generation earth system modeling. Geosci. Model Dev., 2010, 3(1): 87~104.

Sun, C., Liu, L., Li, R., Yu, X., Yu, H., Zhao, B., Wang, G., Liu, J., Qiao, F., and Wang, B.: Developing a common, flexible and efficient framework for weakly coupled ensemble data assimilation based on C-Coupler2.0, Geosci. Model Dev. Discuss., https://doi.org/10.5194/gmd-2020-75, in review, 2020.

Thakur, R., Rabenseifner, R., Gropp, W. (2005). Optimization of Collective Communication Operations in MPICH. The International Journal of High Performance Computing Applications, 19(1), 49–66. https://doi.org/10.1177/1094342005051521

G. Theurich, ESMF Core Team: Performance Optimization of ESMF. https://www.earthsystemcog.org/site_media/projects/esmf/pres_0812_board_gerhard.pdf, 2008

Valcke, S., 2013: The OASIS3 coupler: A European climate modelling community software. Geosci. Model Dev., 6, 373–388, doi:10.5194/gmd-6–373-2013

Valcke, S., Balaji, V., Craig, A., Deluca, C., Dunlap, R., Ford, R. W., Jacob, R., Larson, J., Okuinghttons, R., Riley, G., Vertenstein, M: Coupling technologies for Earth System Modelling, Geosci. Model Dev., 5, 1589-1596, 2012.
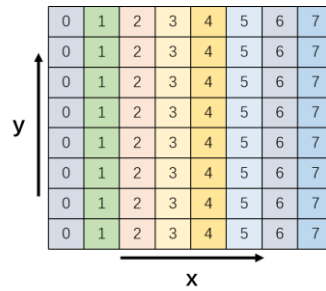
315

320

325

330

335

340

345

11

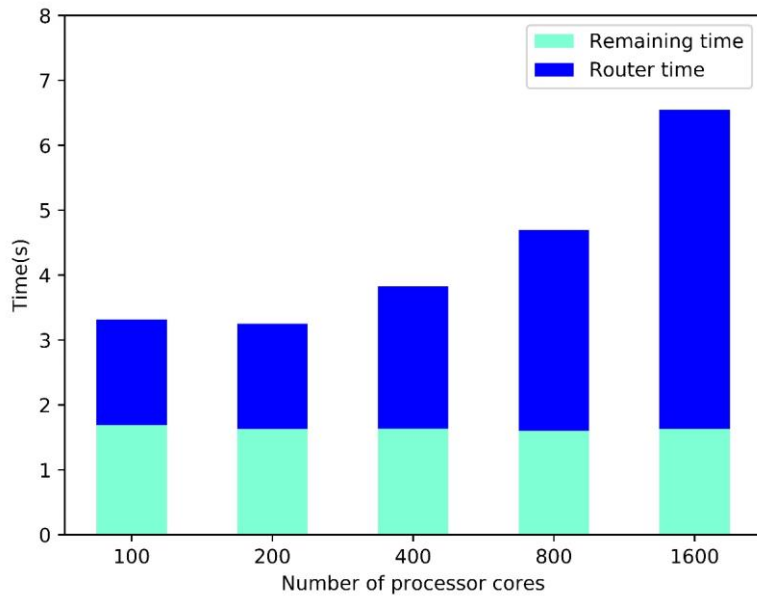(a) A regular 2-D parallel decomposition in both X and Y direction

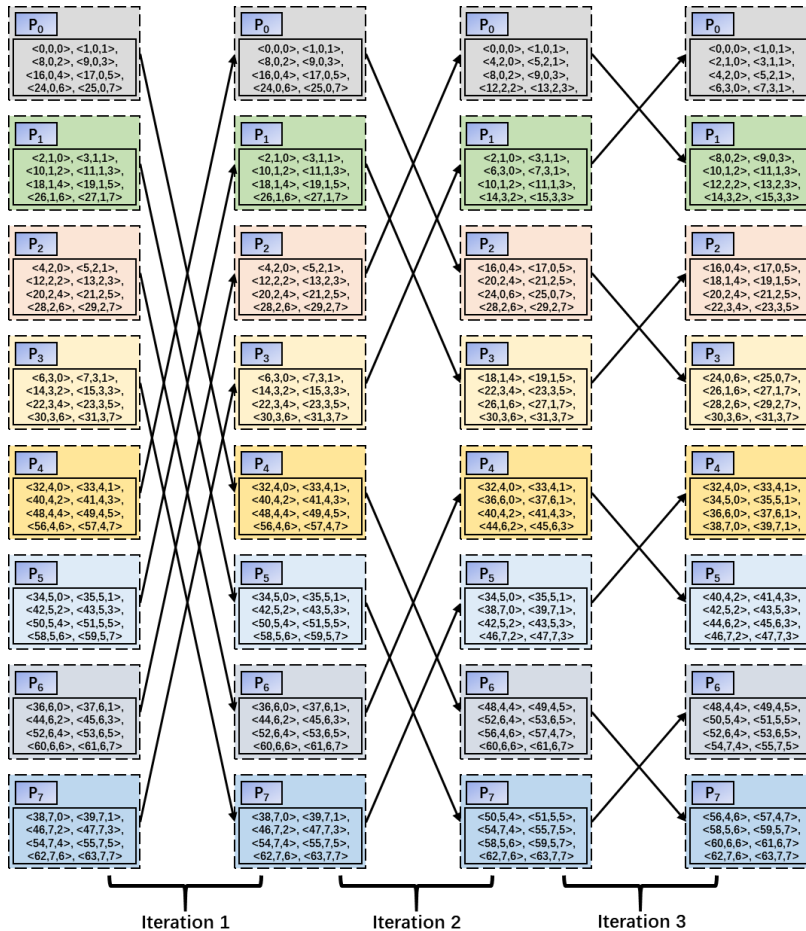(b) A regular 1-D parallel decomposition in only X direction

(c) The routing network from the parallel decomposition in Fig. 1a (**S**ource) to the parallel decomposition in Fig. 1(b) (**D**estination).

**Figure 1. Two sample parallel decompositions of an 8x8 grid under 8 processes (Fig. 1a and 1b) and the routing network between them (Fig. 1c). Each colour corresponds to a process)**
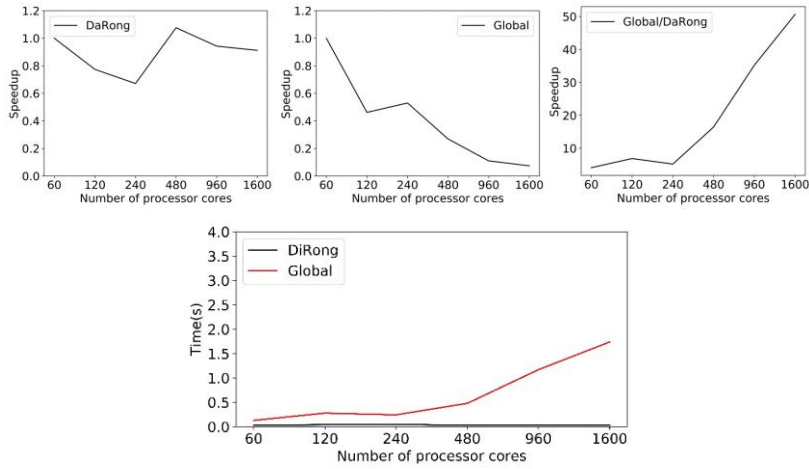
350

**Figure 2.** **The total time of routing network generation (router time) and the remaining time for initializing a two-way MCT coupling between two toy component models. One toy component model uses a longitude-latitude grid with 4 million points and a regular 2-D parallel decomposition, while the other uses a cubed-sphere grid at 0.3 degree and a round-robin parallel decomposition. The time for reading an offline remapping weight file has been taken account in the remaining time, and a regular 1-D parallel decomposition is designed for the data interpolation. The supercomputer as well as the corresponding software stacks described in Section 4 is used for this test.**

355
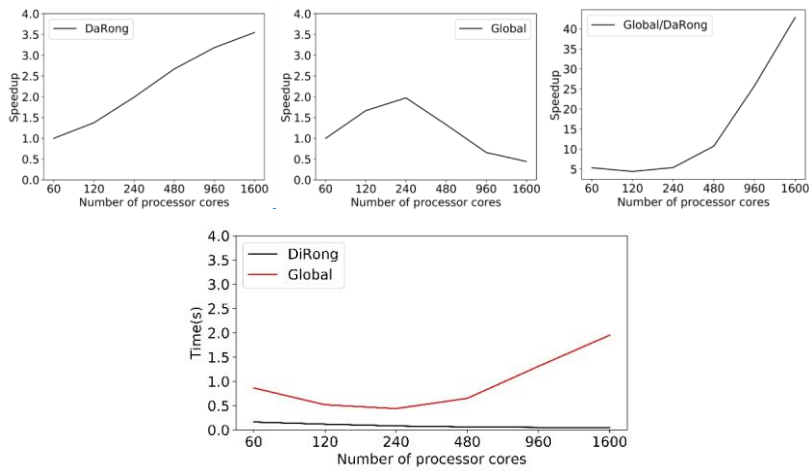
360   **Figure 32. The distributed sort corresponding to the CLGMT entries in Table 1. Each iteration makes the CLGMT entries with larger global cell indexes reserved in the processes with larger IDs. For example, after the first iteration, the CLGMT entries with global cell indexes between 0 and 31 are reserved in P0~P3, while the remaining CLGMT entries are reserved in P4~P7.**

(a) ~~Speedups under the grid size of 500,000.~~ The execution time of ~~DaRong~~DiRong1.0 and the global under ~~60~~ different ~~cores~~ numbers and the grid size of 500,000. ~~is 0.031 s and 0.129 s respectively.~~



(b) ~~Speedups under the grid size of 4,000,000.~~ The execution time of ~~DaRong~~DiRong1.0 and the global under ~~60~~ different core numbers and the grid size of 4,000,000.~~s is 0.161 s and 0.863 s respectively.~~
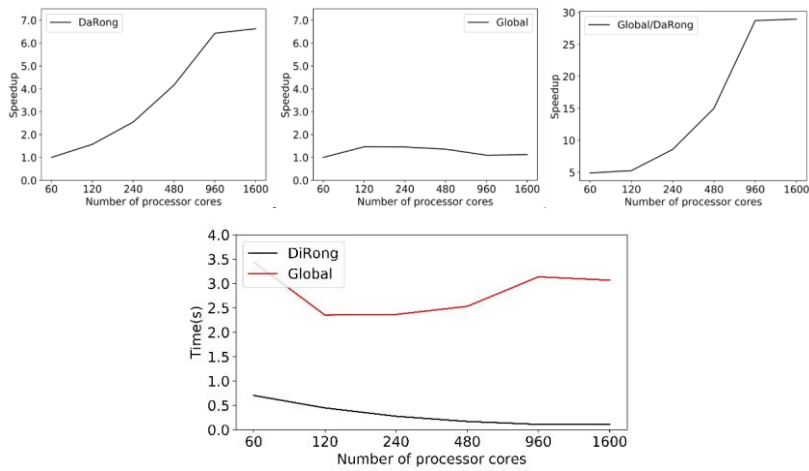
15

375



(c) ~~Speedups under the grid size of 16,000,000.~~ The execution time of ~~DaRong~~DiRong1.0 and the global under ~~60~~ different core numbers and the grid size of 16,000,000. ~~s is 0.702 s and 3.44 s respectively.~~

**Figure 4**~~3~~**. Performance of ~~DaRong~~DiRong1.0 and the comparison with the original global routing network generation**
380 **(Global) under different core numbers and grid sizes. Two toy component models use the same number of processor**
**cores in each test case. The comparison of the two algorithms in these figures shows that the acceleration effect of**
**DiRong1.0 is more obvious when the number of grids and the number of processes is larger, that is, DiRong1.0 has**
**higher parallel efficiency and better scalability.**~~The speedup in the left or middle graph of each sub figure is used for~~
~~evaluating the scalability of DaRong~~DiRong1.0 or the global when increasing processor cores. It is the ratio between~~
385 ~~the execution time under 60 cores and the execution time under another core number. The speedup in the right graph~~
~~of each sub figure is the ratio of the execution time between the global and DaRong~~DiRong1.0, and bigger speedup~~
~~means that is DaRong~~DiRong1.0 is faster.~~

16

**Table 1. The Cell Local-Global Mapping Table (CLGMT) of the parallel decomposition in Fig. 1a**

| Process ID | Cell Local-Global Mapping Table entries |
|---|---|
| 0 | <0,0,0>, <1,0,1>, <8,0,2>, <9,0,3>, <16,0,4>, <17,0,5>, <24,0,6>, <25,0,7> |
| 1 | <2,1,0>, <3,1,1>, <10,1,2>, <11,1,3>, <18,1,4>, <19,1,5>, <26,1,6>, <27,1,7> |
| 2 | <4,2,0>, <5,2,1>, <12,2,2>, <13,2,3>, <20,2,4>, <21,2,5>, <28,2,6>, <29,2,7> |
| 3 | <6,3,0>, <7,3,1>, <14,3,2>, <15,3,3>, <22,3,4>, <23,3,5>, <30,3,6>, <31,3,7> |
| 4 | <32,4,0>, <33,4,1>, <40,4,2>, <41,4,3>, <48,4,4>, <49,4,5>, <56,4,6>, <57,4,7> |
| 5 | <34,5,0>, <35,5,1>, <42,5,2>, <43,5,3>, <50,5,4>, <51,5,5>, <58,5,6>, <59,5,7> |
| 6 | <36,6,0>, <37,6,1>, <44,6,2>, <45,6,3>, <52,6,4>, <53,6,5>, <60,6,6>, <61,6,7> |
| 7 | <38,7,0>, <39,7,1>, <46,7,2>, <47,7,3>, <54,7,4>, <55,7,5>, <62,7,6>, <63,7,7> |

390

**Table 2. The Cell Local-Global Mapping Table (CLGMT) of the parallel decomposition in Fig. 1b**

| Process ID | Cell Local-Global Mapping Table entries |
|---|---|
| 0 | <0,0,0>, <8,0,1>, <16,0,2>, <24,0,3>, <32,0,4>, <40,0,5>, <48,0,6>, <56,0,7> |
| 1 | <1,1,0>, <9,1,1>, <17,1,2>, <25,1,3>, <33,1,4>, <41,1,5>, <49,1,6>, <57,1,7> |
| 2 | <2,2,0>, <10,2,1>, <18,2,2>, <26,2,3>, <34,2,4>, <42,2,5>, <50,2,6>, <58,2,7> |
| 3 | <3,3,0>, <11,3,1>, <19,3,2>, <27,3,3>, <35,3,4>, <43,3,5>, <51,3,6>, <59,3,7> |
| 4 | <4,4,0>, <12,4,1>, <20,4,2>, <28,4,3>, <36,4,4>, <44,4,5>, <52,4,6>, <60,4,7> |
| 5 | <5,5,0>, <13,5,1>, <21,5,2>, <29,5,3>, <37,5,4>, <45,5,5>, <53,5,6>, <61,5,7> |
| 6 | <6,6,0>, <14,6,1>, <22,6,2>, <30,6,3>, <38,6,4>, <46,6,5>, <54,6,6>, <62,6,7> |
| 7 | <7,7,0>, <15,7,1>, <23,7,2>, <31,7,3>, <39,7,4>, <47,7,5>, <55,7,6>, <63,7,7> |

**Table 3. The distributed CLGMT after rearranging the CLGMT entries in Table 2**

| Process ID | CLGMT entries |
|---|---|
| 0 | <0,0,0>, <1,1,0>, <2,2,0>, <3,3,0>, <4,4,0>, <5,5,0>, <6,6,0>, <7,7,0> |
| 1 | <8,0,1>, <9,1,1>, <10,2,1>, <11,3,1>, <12,4,1>, <13,5,1>, <14,6,1>, <15,7,1> |
| 2 | <16,0,2>, <17,1,2>, <18,2,2>, <19,3,2>, <20,4,2>, <21,5,2>, <22,6,2>, <23,7,2> |
| 3 | <24,0,3>, <25,1,3>, <26,2,3>, <27,3,3>, <28,4,3>, <29,5,3>, <30,6,3>, <31,7,3> |
| 4 | <32,0,4>, <33,1,4>, <34,2,4>, <35,3,4>, <36,4,4>, <37,5,4>, <38,6,4>, <39,7,4> |
| 5 | <40,0,5>, <41,1,5>, <42,2,5>, <43,3,5>, <44,4,5>, <45,5,5>, <46,6,5>, <47,7,5> |
| 6 | <48,0,6>, <49,1,6>, <50,2,6>, <51,3,6>, <52,4,6>, <53,5,6>, <54,6,6>, <55,7,6> |
| 7 | <56,0,7>, <57,1,7>, <58,2,7>, <59,3,7>, <60,4,7>, <61,5,7>, <62,6,7>, <63,7,7> |

**Table 4. The Sharing Relationship Table (SRT) calculated from the rearranged distributed CLGMT entries in Fig. 32 and Table 3.**

| Process ID | Sharing Relationship Table entries |
|---|---|
| 0 | <0,0,0,0,0>, <1,0,1,1,0>, <2,1,0,2,0>, <3,1,1,3,0>, <4,2,0,4,0>, <5,2,1,5,0>, <6,3,0,6,0>, <7,3,1,7,0> |
| 1 | <8,0,2,0,1>, <9,0,3,1,1>, <10,1,2,2,1>, <11,1,3,3,1>, <12,2,2,4,1>, <13,2,3,5,1>, <14,3,2,6,1>, <15,3,3,7,1> |
| 2 | <16,0,4,0,2>, <17,0,5,1,2>, <18,1,4,2,2>, <19,1,5,3,2>, <20,2,4,4,2>, <21,2,5,5,2>, <22,3,4,6,2>, <23,3,5,7,2> |
| 3 | <24,0,6,0,3>, <25,0,7,1,3>, <26,1,6,2,3>, <27,1,7,3,3>, <28,2,6,4,3>, <29,2,7,5,3>, <30,3,6,6,3>, <31,3,7,7,3> |
| 4 | <32,4,0,0,4>, <33,4,1,1,4>, <34,5,0,2,4>, <35,5,1,3,4>, <36,6,0,4,4>, <37,6,1,5,4>, <38,7,0,6,4>, <39,7,1,7,4> |
| 5 | <40,4,2,0,5>, <41,4,3,1,5>, <42,5,2,2,5>, <43,5,3,3,5>, <44,6,2,4,5>, <45,6,3,5,5>, <46,7,2,6,5>, <47,7,3,7,5> |
| 6 | <48,4,4,0,6>, <49,4,5,1,6>, <50,5,4,2,6>, <51,5,5,3,6>, <52,6,4,4,6>, <53,6,5,5,6>, <54,7,4,6,6>, <55,7,5,7,6> |
| 7 | <56,4,6,0,7>, <57,4,7,1,7>, <58,5,6,2,7>, <59,5,7,3,7>, <60,6,6,4,7>, <61,6,7,5,7>, <62,7,6,6,7>, <63,7,7,7,7> |

395

**Table 5. The SRT entries distributed in the *src* component model after rearranging the SRT in Table 4**

| Process ID | Sharing Relationship Table entries |
|---|---|
| 0 | <0,0,0,0,0>, <1,0,1,1,0>, <8,0,2,0,1>, <9,0,3,1,1>, <16,0,4,0,2>, <17,0,5,1,2>, <24,0,6,0,3>, <25,0,7,1,3> |
| 1 | <2,1,0,2,0>, <3,1,1,3,0>, <10,1,2,2,1>, <11,1,3,3,1>, <18,1,4,2,2>, <19,1,5,3,2>, <26,1,6,2,3>, <27,1,7,3,3> |
| 2 | <4,2,0,4,0>, <5,2,1,5,0>, <12,2,2,4,1>, <13,2,3,5,1>, <20,2,4,4,2>, <21,2,5,5,2>, <28,2,6,4,3>, <29,2,7,5,3> |
| 3 | <6,3,0,6,0>, <7,3,1,7,0>, <14,3,2,6,1>, <15,3,3,7,1>, <22,3,4,6,2>, <23,3,5,7,2>, <30,3,6,6,3>, <31,3,7,7,3> |
| 4 | <32,4,0,0,4>, <33,4,1,1,4>, <40,4,2,0,5>, <41,4,3,1,5>, <48,4,4,0,6>, <49,4,5,1,6>, <56,4,6,0,7>, <57,4,7,1,7> |
| 5 | <34,5,0,2,4>, <35,5,1,3,4>, <42,5,2,2,5>, <43,5,3,3,5>, <50,5,4,2,6>, <51,5,5,3,6>, <58,5,6,2,7>, <59,5,7,3,7> |
| 6 | <36,6,0,4,4>, <37,6,1,5,4>, <44,6,2,4,5>, <45,6,3,5,5>, <52,6,4,4,6>, <53,6,5,5,6>, <60,6,6,4,7>, <61,6,7,5,7> |
| 7 | <38,7,0,6,4>, <39,7,1,7,4>, <46,7,2,6,5>, <47,7,3,7,5>, <54,7,4,6,6>, <55,7,5,7,6>, <62,7,6,6,7>, <63,7,7,7,7> |

**Table 6. The SRT entries distributed in the *dst* component model after rearranging the SRT in Table 4**

| Process ID | Sharing Relationship Table entries |
|---|---|
| 0 | <0,0,0,0,0>, <8,0,2,0,1>, <16,0,4,0,2>, <24,0,6,0,3>, <32,4,0,0,4>, <40,4,2,0,5>, <48,4,4,0,6>, <56,4,6,0,7> |
| 1 | <1,0,1,1,0>, <9,0,3,1,1>, <17,0,5,1,2>, <25,0,7,1,3>, <33,4,1,1,4>, <41,4,3,1,5>, <49,4,5,1,6>, <57,4,7,1,7> |
| 2 | <2,1,0,2,0>, <10,1,2,2,1>, <18,1,4,2,2>, <26,1,6,2,3>, <34,5,0,2,4>, <42,5,2,2,5>, <50,5,4,2,6>, <58,5,6,2,7> |
| 3 | <3,1,1,3,0>, <11,1,3,3,1>, <19,1,5,3,2>, <27,1,7,3,3>, <35,5,1,3,4>, <43,5,3,3,5>, <51,5,5,3,6>, <59,5,7,3,7> |
| 4 | <4,2,0,4,0>, <12,2,2,4,1>, <20,2,4,4,2>, <28,2,6,4,3>, <36,6,0,4,4>, <44,6,2,4,5>, <52,6,4,4,6>, <60,6,6,4,7> |
| 5 | <5,2,1,5,0>, <13,2,3,5,1>, <21,2,5,5,2>, <29,2,7,5,3>, <37,6,1,5,4>, <45,6,3,5,5>, <53,6,5,5,6>, <61,6,7,5,7> |
| 6 | <6,3,0,6,0>, <14,3,2,6,1>, <22,3,4,6,2>, <30,3,6,6,3>, <38,7,0,6,4>, <46,7,2,6,5>, <54,7,4,6,6>, <62,7,6,6,7> |
| 7 | <7,3,1,7,0>, <15,3,3,7,1>, <23,3,5,7,2>, <31,3,7,7,3>, <39,7,1,7,4>, <47,7,3,7,5>, <55,7,5,7,6>, <63,7,7,7,7> |

400

**Table 7. Performance of ~~DaRong~~DiRong1.0 and the comparison with the original global routing network generation (Global) when concurrently increasing the grid size and core number.**

| Core number of each toy component model | Grid size | Execution time (s) of ~~DaRong~~DiRong1.0 | Execution time (s) of Global | Global/~~DaRongDi~~Rong1.0 |
|---|---|---|---|---|
| 200~~250~~ | 500,000~~500,000~~ | 0.029~~0.032~~ | 0.214 ~~0.262~~ | 7.38 ~~8.19~~ |
| 400~~450~~ | 1,000,000~~1,000,000~~ | 0.033~~0.034~~ | 0.453 ~~0.492~~ | 13.73 ~~14.47~~ |
| 800~~900~~ | 2,000,000~~2,000,000~~ | 0.039~~0.041~~ | 1.008 ~~1.158~~ | 25.85 ~~28.24~~ |
| 1600~~1600~~ | 4,000,000~~4,000,000~~ | 0.045~~0.045~~ | 1.949 ~~1.949~~ | 43.31 ~~43.31~~ |

Formatted Table