

19 congestion. It is generally believed that the transpositions seriously limit the scalability
20 of the spectral models. The experiments show that ~~although~~ the communication time of
21 the transposition is larger than those of the wide halo exchange for the Semi-Lagrangian
22 method and the allreduce in the GCR iterative solver for the Semi-Implicit method be-
23 low 2×10^5 MPI processes, ~~the~~. The transposition whose communication time decreases
24 quickly ~~as the~~ with increasing number of MPI processes ~~increases~~ demonstrates strong
25 scalability in the case of very large grids and moderate latencies; ~~the~~. The halo exchange
26 whose communication time decreases more slowly than that of transposition ~~as the~~ with
27 increasing number of MPI processes ~~increases~~ reveals its weak scalability; ~~in~~. In contrast,
28 the allreduce whose communication time increases ~~as the~~ with increasing number of MPI
29 processes ~~increases~~ does not scale well. From this point of view, the scalability of ~~the~~
30 spectral models could still be acceptable, ~~therefore~~. Therefore it seems to be premature
31 to conclude that the scalability of the grid-point models is better than that of spectral
32 models at exascale, unless innovative methods are exploited to mitigate the problem of
33 the scalability presented in the grid-point models.

34 **Keyword:** performance, scalability, MPI, communication, transposition, halo exchange,
35 all reduce, topology, routing, bandwidth, latency

36 1 Introduction

37 Current high performance computing (HPC) systems have thousands of nodes and millions
38 of cores. According to the 49th TOP500 list (www.top500.org) published on June 20, 2017,
39 the fastest machine (Sunway TaihuLight) had over than 10 million cores with a peak perfor-
40 mance approximately 125 PFlops (1 PFlops= 10^{15} floating-point operations per second), and
41 the second HPC (Tianhe-2) is made up of 16,000 nodes and has more than 3 million cores with
42 a peak performance approximately 55 PFlops. It is estimated that in the near future, HPC
43 systems will dramatically scale up in size. Next decade, it is envisaged that exascale HPC
44 system with millions of nodes and thousands of cores per node, whose peak performance ap-
45 proaches to or is beyond 1 EFlops (1 EFlops= 10^3 PFlops), will become available (Engelmann,
46 2014; Lagadapati et al., 2016). Exascale HPC poses several challenges in terms of power con-
47 sumption, performance, scalability, programmability, and resilience. The interconnect net-
48 work of exascale HPC system becomes larger and more complex, and its performance which
49 largely determines the overall performance of the HPC system is crucial to the performance

50 of distributed applications. Designing energy-efficient cost-scalable interconnect networks and
51 communication-efficient scalable distributed applications is an important component of HPC
52 hardware/software co-design to address these challenges. Thus, evaluating and predicting the
53 communication behaviour of distributed applications is obligatory; it is only feasible by mod-
54 elling the communications and the underlying interconnect network, especially for the future
55 supercomputer.

56 Investigating the performance of distributed applications on future architectures and the
57 impact of different architectures on the performance by simulation is a hardware/software
58 co-design approach for paving the way to exascale HPCs. Analytical interconnect network sim-
59 ulation based on an analytical conceptual model is fast and scalable, but comes at the cost of
60 accuracy owing to its unrealistic simplification (Hoeffler et al., 2010). Discrete event simulation
61 (DES) is often used to simulate the interconnect network, and it provides high fidelity since the
62 communication is simulated in more detailed level (e.g., flit, packet, or flow levels) to take into
63 account congestion (Janssen et al., 2010; Böhm and Engelmann, 2011; Dechev and Ahn, 2013;
64 Acun et al., 2015; Jain et al., 2016; Wolfe et al., 2016; Degomme et al., 2017; Mubarak et al.,
65 2017). Sequential DES lacks scalability owing to its large memory footprints and long exe-
66 cution time (Degomme et al., 2017). Parallel DES (PDES) is scalable since it can reduce the
67 memory required per node, but its parallel efficiency is not very good because of frequent
68 global synchronization of conservative PDES (Janssen et al., 2010) or high rollback overhead of
69 optimistic PDES (Acun et al., 2015; Jain et al., 2016; Wolfe et al., 2016). Generally, the simu-
70 lation of distributed applications can be divided into two complementary categories: offline and
71 online simulations. Offline simulation replays the communication traces from the application
72 running on a current HPC system. It is sufficient to understand the performance and dis-
73 cover the bottleneck of full distributed applications on the available HPC system (Tikir et al.,
74 2009; Noeth et al., 2009; Núñez et al., 2010; Dechev and Ahn, 2013; Casanova et al., 2015;
75 Acun et al., 2015; Jain et al., 2016; Lagadapati et al., 2016); however, is not very scalable be-
76 cause of the huge traces for numerous processes and limited extrapolation to future architecture
77 (Hoeffler et al., 2010; Núñez et al., 2010). Online simulation has full scalability to future system
78 by running the skeleton program on the top of simulators (Zheng et al., 2004; Janssen et al.,
79 2010; Engelmann, 2014; Degomme et al., 2017), but has the challenge of developing a skele-
80 ton program from a complex distributed application. Most simulations in the aforementioned
81 literatures have demonstrated the scalability of simulators. The simulator xSim (Engelmann,

2014) simulated a very simple MPI program, which only calls MPI_Init and MPI_Finalize without any communication and computation, up to 2^{27} processes. For collective MPI operations, Hoefer et al. (2010) obtained an MPI_Allreduce simulation of 8 million processes without consideration of congestion using LogGOPSim, Engelmann (2014) achieved an MPI_Reduce simulation of 2^{24} processes, and Degomme et al. (2017) demonstrated an MPI_Allreduce simulation of 65536 processes using SimGrid. For simulations at application level, Jain et al. (2016) used the TraceR simulator based on CODES and ROSS to replay 4.6×10^4 process traces of several communication patterns that are used in a wide range of applications. In addition, Mubarak et al. (2017) presented a 1.1×10^5 process simulations of two multigrid applications. However, to the best of our knowledge, there is no exascale simulation of complex communication patterns such as the MPI transposition (Multiple simultaneous MPI_Alltoallv) for the spectral method and the wide halo exchange (the width of a halo may be greater than the subdomain size of its direct neighbours) for the Semi-Lagrangian method used in atmospheric models.

With the rapid development of increasingly powerful supercomputers in recent years, numerical weather prediction (NWP) models have increasingly sophisticated physical and dynamical processes, and their resolution is getting higher and higher. Nowadays, the horizontal resolution of global NWP model is in the order of 10 kilometres. Many operational global spectral NWP models such as IFS at ECMWF, ARPEGE at METEO-FRANCE, and GFS at NCEP are based on the spherical harmonics transform method that includes Fourier transforms in the zonal direction and Legendre transforms in the meridional direction (Ehrendorfer, 2012). Moreover, some regional spectral models such as AROME at METEO-FRANCE (Seity et al., 2011) and RSM at NCEP (Juang et al., 1997) use the Bi-Fourier transform method. The Fourier transforms can be computed efficiently by fast Fourier transform (FFT) (Temperton, 1983). Even with the introduction of fast Legendre transform (FLT) to reduce the growing computational cost of increasing resolution of global spectral models (Wedi et al., 2013), it is believed that global spectral method is prohibitively expensive for very high resolution (Wedi, 2014).

A global (regional) spectral model performs FFT and FLT (FFT) in the zonal direction and the meridional direction, respectively. Because both transforms require all values in the corresponding directions, the parallelization of spectral method in global (regional) model is usually conducted to exploit the horizontal domain decomposition only in the zonal direction and meridional directions for FFT and FLT (FFT), respectively (Barros et al., 1995; Kanamitsu et al., 2005). Owing to the horizontal domain decomposition in a single horizontal direction for the

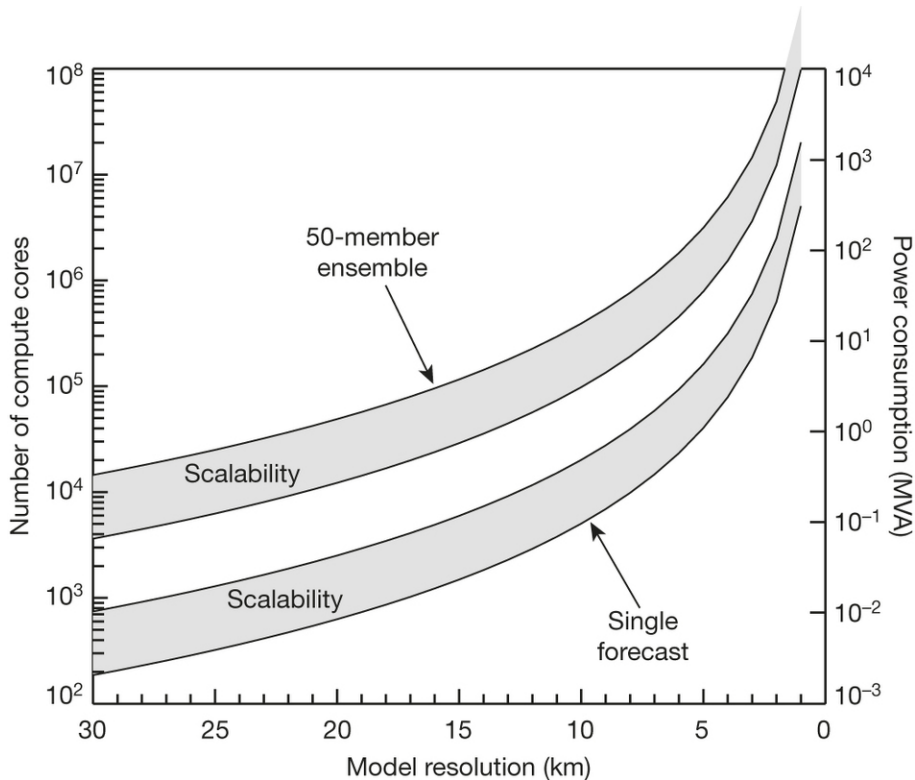


Fig. 1: CPU and power requirements as a function of NWP model resolution, adapted from Bauer et al. (2015). The left and right y axes are the number of cores and the power (in megavolt amps), respectively, required for a single 10-day model forecast (the lower shaded area including its bounds) and a 50-member ensemble forecast (the upper shaded area including its bounds) as a function of model resolution, respectively, based on current model code and compute technology. The lower and upper bounds of each shaded area indicate perfect scaling and inefficient scaling, respectively.

114 parallelization of spectral transforms, there is a transposition between the spectral transforms
 115 in the zonal direction and meridional directions. MPI (Message Passing Interface) transposition
 116 is an all-to-all personalized communication which can cause significant congestion over inter-
 117 connect network when the number of MPI tasks and the amount of exchanged data are large,
 118 and results in severe communication delay. Bauer et al. (2015) estimated that a global NWP
 119 model with a two-kilometre horizontal resolution requires one million compute cores for a single
 120 10-day forecast (Fig. 1). With one million compute cores, the performance and scalability of
 121 the MPI transposition become of paramount importance for a high resolution global spectral
 122 model. Thus, evaluating and predicting the performance and scalability of MPI transposition
 123 at exascale is one of the foremost subjects of this study.

124 The Semi-Lagrangian (SL) method is a highly efficient technique for the transport of mo-
 125 mentum, heat and mass in the NWP model because of its unconditional stability which permits
 126 a long time step (Staniforth and Côté, 1991; Hortal, 2002). However, it is known that the MPI

127 exchange of wide halo required for the interpolation at the departure point of high wind-speed
128 particles near the boundary of the subdomain causes significant communication overhead as
129 resolution increases towards kilometres scale and the HPC systems move towards exascale.
130 This communication overhead could reduce the efficiency of the SL method; thus, modelling
131 the performance and scalability of wide halo exchange at exascale is essential and is another
132 subject of this study.

133 With consideration of the efficiency of the Legendre transform and the scalability of MPI
134 transposition that may arise in the global spectral model on exascale HPC systems, a cou-
135 ple of global grid-point models have recently been developed (Lin, 2004; Satoh et al., 2008;
136 Qaddouri and Lee, 2011; Skamarock et al., 2012; Dubos et al., 2015; Zangl et al., 2015; Smolarkiewicz et al.
137 2016). Since spherical harmonics are eigenfunctions of the Helmholtz operator, the Semi-
138 Implicit (SI) method is usually adopted in order to implicitly handle the fast waves in the
139 global spectral model to allow stable integration with a large time step (Robert et al., 1972;
140 Hoskins and Simmons, 1975). However, for a grid-point model, the three-dimensional Helmholtz
141 equation is usually solved using Krylov subspace methods such as the generalized conjugate
142 residual (GCR) method (Eisenstat et al., 1983), and a global synchronization for the inner
143 product in Krylov subspace methods may become the bottleneck at exascale (Li et al., 2013;
144 Sanan et al., 2016). As it is not clear whether the three-dimensional Helmholtz equation can
145 be solved efficiently in a scalable manner, most of the aforementioned models use a horizontally
146 explicit vertically implicit (HEVI) scheme. The HEVI scheme typically requires some damping
147 for numerical stability (Satoh et al., 2008; Skamarock et al., 2012; Zangl et al., 2015), and its
148 time step is smaller than that of the SI method (Sandbach et al., 2015). Therefore, it is de-
149 sirable to know whether the SI method is viable or even advantageous for very high resolution
150 grid-point models running on exascale HPC systems. Thus, it is valuable to explore the per-
151 formance and scalability of global synchronization in solving the three-dimensional Helmholtz
152 equation using Krylov subspace methods; this forms the third subject of this study.

153 In this paper, we present the application of SST/macro 7.1, a coarse-grained parallel discrete
154 event simulator, to investigate the communication performance and scalability of atmospheric
155 models for future exascale supercomputers. The remainder of the paper is organized as fol-
156 lows. Section 2 introduces the simulation environment, the SST/macro simulator, and our
157 optimizations for reducing the memory footprint and accelerating the simulations. Section 3
158 reviews three key MPI operations used in the atmospheric models. Section 4 presents and

159 analyses the experimental results of the modelling communication of the atmospheric model
160 using SST/macro. Finally, we summarize the conclusions and discuss future work in section 5.

161 **2 Simulation Environment**

162 *2.1 Parallel Discrete Event Simulation*

163 Modelling application performance on exascale HPC systems with millions of nodes and a
164 complex interconnect network requires that the simulation can be decomposed into small tasks
165 that efficiently run in parallel to overcome the problem of large memory footprint and long
166 simulation time. PDES is such an approach for exascale simulation. Each worker in PDES is
167 a logical process (LP) that models a specific component such as a node, a switch, or an MPI
168 process of the simulated MPI application. These LPs are mapped to the physical processing
169 elements (PEs) that actually run the simulator. An event is an action such as sending an MPI
170 message or executing a computation between consecutive communications. Each event has its
171 start and stop times, so the events must be processed without violating their time ordering.
172 To model the performance of an application, PDES captures time duration and advances the
173 virtual time of the application by sending timestamped events between LPs.

174 PDES usually adopts conservative or optimistic parallelized strategies. The conservative
175 approach maintains the time ordering of events by synchronization to guarantee that no early
176 events arrive after the current event. Frequent synchronization is time-consuming so the effi-
177 ciency of the conservative approach is highly dependent on the look ahead time; a larger look
178 ahead time (that means less synchronization) allows a much greater parallelism. The optimistic
179 approach allows LPs to run events at the risk of time-ordering violations. Events must be rolled
180 back when time-ordering violations occurs. Rollback not only induces significant overhead, but
181 also requires extra storage for the event list. Rollback presents special challenges for online
182 simulation, so SST/macro adopts a conservative approach (Wike and Kenny, 2014).

183 *2.2 SST/macro Simulator*

184 Considering that the offline trace-driven simulation does not provide an easy way for extrap-
185 olating to future architectures, the online simulator SST/macro is selected here to model the
186 communications of the atmospheric models for future exascale HPC systems. SST/macro is a

187 coarse-grained parallel discrete event simulator which provides the best cost/accuracy trade-off
188 simulation for large-scale distributed applications (Janssen et al., 2010). SST/macro is driven
189 by either a trace file or a skeleton application. A skeleton application can be constructed from
190 scratch, or from an existing application manually or automatically by source-to-source trans-
191 lation tools. SST/macro intercepts the communications issued from the skeleton program to
192 estimate their time rather than actually execute it by linking the skeleton application to the
193 SST/macro library instead of the real MPI library. Since the purpose of this study is to investi-
194 gate the performance and scalability of communications in an atmospheric model, we construct
195 the communication-only skeleton program from scratch by identifying the key MPI operations
196 taking place in the atmospheric models.

197 Congestion is a significant factor that affects the performance and scalability of MPI appli-
198 cations running on exascale HPC systems. SST/macro has three network models: the analytical
199 model transfers the whole message over the network from point-to-point without packetizing
200 and estimates the time delay Δt predominantly based on the logP approximation

$$201 \quad \Delta t = \alpha + \beta N, \quad (1)$$

202 where α is the communication latency, β is the inverse bandwidth in second per byte, and N is
203 the message size in bytes; the packet-level model PISCES (Packet-flow Interconnect Simulation
204 for Congestion at Extreme Scale) divides the message into packets and transfers the packets
205 individually; the flow-level model will be deprecated in the future. Compared to the SimGrid
206 simulator, the packet-level model of SST/macro produces almost identical results (figure omit-
207 ted). Acun et al. (2015) also found that the SST/macro online simulation is very similar to
208 the TraceR simulation. Thus, we adopt the PISCES model with a cut-through mechanism
209 (SNL, 2017) to better account for the congestion. SST/macro provides three abstract machine
210 models for nodes: the AMM1 model is the simplest one which grants exclusive access to the
211 memory, the AMM2 model allows multiple CPUs or NICs (network interface controller) to
212 share the memory bandwidth by defining the maximum memory bandwidth allocated for each
213 component, the AMM3 model goes one further step to distinguish between the network link
214 bandwidth and the switch bandwidth. In this paper, the AMM1 model with one single-core
215 CPU per node is adopted since simulation of communications is the primary goal.

216 SST/macro provides several topologies of the interconnect network. In this study, three

217 types of topologies (Fig. 2) commonly used in current supercomputers, and their configurations
 218 are investigated. Torus topology has been used in many supercomputers (Ajima et al., 2009).
 219 In the torus network, messages hop along each dimension using the shortest path routing
 220 from the source to the destination (Fig. 2a), and its bisection bandwidth typically increases with
 221 increasing dimension size of the torus topology. The practical implementation of the fattree
 222 topology is an upside-down tree that typically employs all uniform commodity switches to
 223 provide high bandwidth at higher levels by grouping corresponding switches of the same colour
 224 (Fig. 2b). Fattree topology is widely adopted by many supercomputers for its scalability and
 225 high path diversity (Leiserson, 1985); it usually uses a D-mod-k routing algorithm (Zahavi et al.,
 226 2010) for desirable performance. A dragonfly network is a multi-level dense structure of which
 227 the high-radix routers are connected in a dense even all-to-all manner at each level (Kim et al.,
 228 2008). As shown in Fig. 2c, a typical dragonfly network consists of two levels: the routers at
 229 the first level are divided into groups and routers in each group form a two-dimension mesh
 230 of which each dimension is an all-to-all connected network; at the second level, the groups as
 231 virtual routers are connected in an all-to-all manner (Alverson et al., 2015). There are three
 232 available routing algorithms for dragonfly topology in SST/macro:

233 **minimal** transfers messages by the shortest path from the source to the destination. For
 234 example, messages travel from the blue router in group 0 to the red router in group 2 via
 235 the bottom-right corner in group 0 and the bottom-left corner in group 2 (Fig. 2c).

236 **valiant** randomly picks an intermediate router, and then uses a minimal routing algorithm to
 237 transfer messages from the source to the intermediate router and from the intermediate
 238 router to the destination. For example, the arrow path from the blue router in group 0
 239 to the red router in group 2 goes via the intermediate yellow node in group 1 in Fig. 2c.

240 **ugal** checks the congestion, and either switches to the valiant routing algorithm if congestion
 241 is too heavy, or otherwise uses the minimal routing algorithm.

242 Table 1 summaries the network topology configurations used in this paper. Torus-M (torus-
 243 L) configuration is a 3D torus of 25x25x25 (75x25x25) size. Fattree-M (fattree-L) configuration
 244 has 4 layers: the last layer consists of nodes while the other layers consist of switches with 25 (33)
 245 descendant ports per switch. We tested four configurations of dragonfly topology. Dragonfly-
 246 MM configuration has a medium size of a group of a 25x25 mesh with 25 nodes per switch

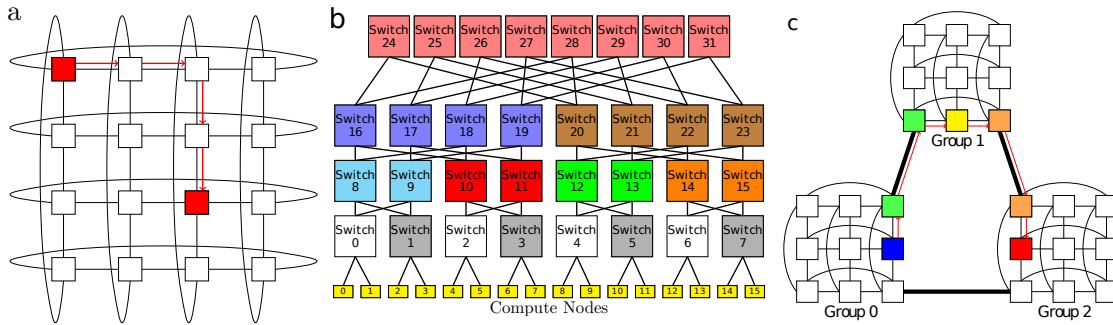


Fig. 2: Topology illustration: a, b, and c are the torus, fattree, and dragonfly topologies, respectively. Adapted from SNL (2017)

Table 1: Summary of the network topologies: the geometry of a torus topology specifies the size of each dimension; the first and second number in the geometry of a fattree topology are the number of layers and descendant ports per switch, respectively; the first two numbers and the last number in the geometry of a dragonfly topology indicate the group mesh size and the number of groups, respectively.

name	geometry	switches	nodes per switch	nodes	radix
torus-M	25,25,25	15625	25	390625	31
fattree-M	4,25	46875	25	390625	50
dragonfly-MM	25,25,25	15625	25	390625	97
dragonfly-SL	25,25,125	15625	5	390625	177
dragonfly-LS	125,125,5	15625	5	390625	257
torus-L	75,25,25	46875	25	1171875	31
fattree-L	4,33	107811	33	1185921	66
dragonfly-ML	25,25,75	46875	25	1171875	147

247 and medium number (=25) of groups. Dragonfly-SL configuration has a small size of a group
248 of a 25x25 mesh with 5 nodes per switch and large number (=125) of groups. Dragonfly-LS
249 configuration has a large size of a group of a 125x125 mesh with 5 nodes per switch and small
250 number (=5) of groups. Dragonfly-ML configuration has a medium size of a group of a 25x25
251 mesh with 25 nodes per switch and large number (=75) of groups. The fattree configuration
252 has a significant larger number of switches than other topologies for the same number of nodes,
253 which indicates that fattree is not cost- or energy-efficient. All the configurations with 390625
254 nodes are used for simulating transposition for the spectral transform method. Torus-L, fattree-
255 L, and dragonfly-ML with more than one million nodes are used for the cases of halo exchange
256 and allreduce communication since we cannot finish the simulation of transposition for the
257 spectral transform method (multiple simultaneous all-to-all personalized communications) on
258 such large configuration within 24 hours (see Section 3 for three key MPI communications in
259 the atmospheric model).

260 *2.3 Reduce the Memory Footprint and Accelerate the Simulation*

261 Although SST/macro is a parallel discrete event simulator that can reduce the memory foot-
262 print per node, its parallel efficiency degrades if more cores are used. Even with an MPI
263 transposition of 10^5 processes, this all-to-all personalized communication has almost 10^{10} dis-
264 crete events, which consumes a considerable amount of memory and takes a very long time
265 for simulation. Furthermore, almost every MPI program has a setup step to allocate memory
266 for storing the setup information such as the parameters and the domain decomposition of all
267 processes what each process must know in order to properly communicate with other processes,
268 therefore, it needs to broadcast the parameters to and synchronize with all processes before
269 actual communications and computation. Even if the setup information for a single process
270 needs only 10^2 bytes memory, a simulation of 10^5 processes MPI transposition will need one
271 terabyte ($10^2 \times 10^5 \times 10^5 = 10^{12}$ bytes) memory, which is not easily available on current com-
272 puters if the simulator runs on a single node. In addition, the MPI operations in the setup step
273 not only are time-consuming, but also affect subsequent communications. A common way to
274 eliminate this effect is to iterate many times to obtain a robust estimation of communication
275 time; however, one iteration is already very time-consuming for simulation. To circumvent the
276 issue of setup steps, we use an external auxiliary program to create a shared memory segment
277 on each node running SST/macro and initialize this memory with the setup information of all
278 the simulated MPI processes. Then, we modified SST/macro to create a global variable and
279 attach the shared memory to this global variable; this method not only reduces the memory
280 footprint and eliminates the side effect of communications in the setup step, but also avoids
281 the problem of filling up the memory address space if each simulated process attaches to the
282 shared memory.

283 Large-scale application needs a large amount of memory for computation; and in some
284 cases, such as spectral model, the whole memory for computation is exchanged between all the
285 processes. Even when computation is not considered, a large amount of memory for the message
286 buffers is usually required for MPI communications. Fortunately, the simulator only needs
287 message size, the source/destination, and the message tag to model the communication; thus,
288 it is not necessary to allocate actual memory. Since SST/macro can operate with null buffers,
289 the message buffer is set to null in the skeleton application, which significantly reduces the size
290 of memory required by the simulation of communication of the high resolution atmospheric

291 model.

292 **3 Key MPI Operations in Atmospheric Models**

293 **3.1 Transposition for the Spectral Transform Method**

294 A global spectral model generally uses spherical harmonics transform on the horizontal with
295 triangular truncation. The backward spherical harmonics transform is

$$296 \quad f(\theta, \lambda) = \sum_{m=-M}^M \left(e^{im\lambda} \sum_{n=|m|}^M f_n^m P_n^m(\cos \theta) \right), \quad (2)$$

297 where θ and λ are the colatitude and longitude, f_n^m is the spectral coefficients of the field f , and
298 P_n^m is the associated Legendre polynomials of degree m and order n . Moreover, the forward
299 spherical harmonics transform is

$$300 \quad f_n^m = \frac{1}{2} \int_{-1}^1 \left(P_n^m(\cos \theta) \frac{1}{2\pi} \int_0^{2\pi} f(\theta, \lambda) e^{-im\lambda} d\lambda \right) d \cos \theta, \quad (3)$$

301 In (2), the backward Legendre transform of each m can be computed independently; then,
302 the same is for the backward Fourier transform of each θ . Similar to (3), the forward Fourier
303 transform of each θ can be computed independently; then, the same is for the forward Legendre
304 transform of each m . This leads to a natural way to parallelize the spectral transforms. If
305 we start with the grid-point space (Fig. 3a), which is decomposed by cx/cy cores in the x/y
306 direction, cy simultaneous xz slab MPI transpositions lead to the partition (Fig. 3b) with cy/cx
307 cores in the y/z direction, and a spectral transform such as a forward FFT can be performed
308 in parallel since data w.r.t. λ are local to each core. Then, cx simultaneous xy slab MPI
309 transpositions lead to the partition (Fig. 3c) with cy/cx cores in the x/z direction, and a
310 spectral transform such as a forward FLT can be computed in parallel because data w.r.t. θ
311 are now local to each core. Finally, cy simultaneous yz slab MPI transpositions lead to the
312 spectral space (Fig. 3d) with cy/cx cores in the x/y direction, where the Semi-Implicit scheme
313 can be easily computed because spectral coefficients belonging to the same column are now
314 local to the same core. The backward transform is similar. It is of paramount importance that
315 the partition of the four stages described in Fig. 3 must be consistent so that multiple slab MPI
316 transpositions can be conducted simultaneously, which significantly reduces the communication

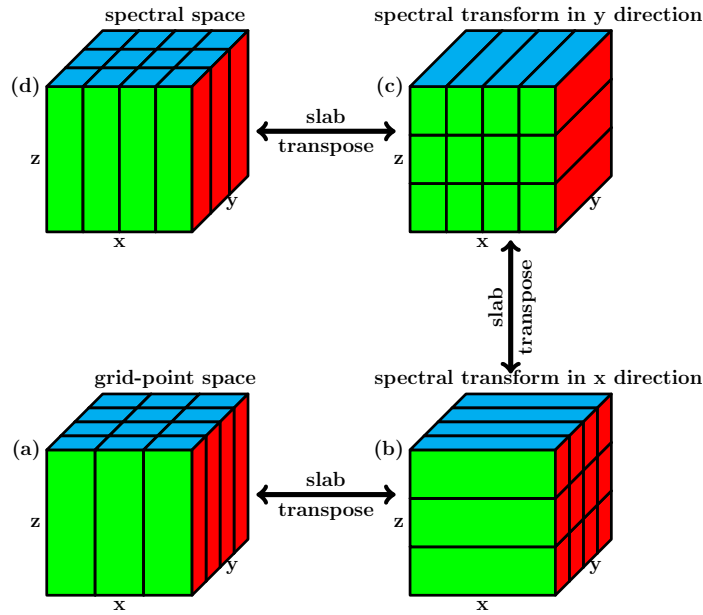


Fig. 3: Parallel scheme of regional spectral model: (a) 2D decomposition of 3D grid field with c_x/c_y cores in the x/y direction, (b) 2D decomposition of 3D grid field with c_y/c_x cores in the y/z direction, (c) 2D decomposition of 3D grid field with c_y/c_x cores in the x/z direction, and (d) 2D decomposition of 3D grid field with c_y/c_x cores in the x/y direction. Transposition between (a) and (b) can be conducted by c_y independent xz slab MPI transpositions, transposition between (b) and (c) can be conducted by c_x independent xy slab MPI transpositions, and transposition between (c) and (d) can be conducted by c_y independent yz slab MPI transpositions.

317 time of MPI transpositions from one stage to another. It is worth noting that the number of
 318 grid points in one direction is not always a multiple of the number of cores in the corresponding
 319 direction; thus, the partition shown in Fig. 3 can use as many as possible ~~computed~~ compute
 320 cores without any limit on c_x or c_y provided $c_x \times c_y = ncpu$, and c_x or c_y is not greater than
 321 the number of grid points in the corresponding direction. It is generally believed that the MPI
 322 transpositions from one stage to another poses a great challenge to the scalability of spectral
 323 models because each slab MPI transposition is an all-to-all personalized communications which
 324 is the most complex and time-consuming all-to-all communication.

325 There are different algorithms for all-to-all personalized communication. Table 2 lists the
 326 three algorithms for all-to-all personalized communication, whose performance and scalability
 327 are investigated in this study. Algorithm ring-k is our proposal algorithm for all-to-all per-
 328 sonalized communication which is a generalized ring alltoally algorithm. In algorithm ring-k,
 329 each process communicates with $2k$ processes to reduce the stages of communications and make
 330 efficient use of the available bandwidth, and thus reduces the total communication time.

Table 2: Three algorithms for all-to-all personalized communication.

name	description	stages
burst	Each process communicates with all other processes simultaneously by posting all non-block send and receive operations simultaneously. The burst messages cause significant congestion on the network. This algorithm is equivalent to the algorithm ring-k when k=n-1.	1
bruck	This algorithm is better for small message and a large latency since it has only $\lceil \log_2(n) \rceil$ stages of communications (Thakur et al., 2005). For k^{th} stage, each process sends the messages whose destination process id has one at the k^{th} bit (begin at Least Significant Bit) to process $i + 2^k$.	$\lceil \log_2(n) \rceil$
ring-k	In the first stage, process i sends to $i + 1, \dots, i + k$ and receive from $i - 1, \dots, i - k$ in a ring way (black arrows in Fig. 4a); in the second stage, process i sends to $i + 1 + k, \dots, i + 2k$ and receive from $i - 1 - k, \dots, i - 2k$ in a ring way (blue arrows in Fig. 4a); this continues until all partners have been communicated with. This algorithm is a generalization of the ring algorithm and efficiently uses the available bandwidth by proper selection of radix k .	$\lceil \frac{n-1}{k} \rceil$

3.2 Halo Exchange for Semi-Lagrangian Method

The SL method solves the transport equation:

$$\frac{D\phi}{Dt} = \frac{\partial\phi}{\partial t} + u\frac{\partial\phi}{\partial x} + v\frac{\partial\phi}{\partial y} + w\frac{\partial\phi}{\partial z} = 0, \quad (4)$$

where the scalar field ϕ is advected by the 3D wind $\mathbf{V} = (u, v, w)$. In the SL method, the grid-point value of the scalar field ϕ at next time step $t + \Delta t$ can be found by integrating (4) along the trajectory of the fluid parcel (Staniforth and Côté, 1991; Hortal, 2002)

$$\int_t^{t+\Delta t} \frac{D\phi}{Dt} dt = 0 \rightarrow \phi^{t+\Delta t} = \phi_d^t, \quad (5)$$

where $\phi^{t+\Delta t}$ is the value of the fluid parcel ϕ arriving at any grid point at $t + \Delta t$, and ϕ_d^t is the value of the same fluid parcel at its departure point d and departure time t . This means that the value of the scalar field ϕ at any grid point at $t + \Delta t$ is equal to its value at the departure point d and the departure time t . The departure point d usually does not coincide with any grid point, so the value of ϕ_d^t is obtained by interpolation using the surrounding grid-point values ϕ^t at time t . The departure point d is determined by iteratively solving the trajectory equation

344 (Staniforth and Côté, 1991; Hortal, 2002)

$$345 \quad \frac{D\mathbf{r}}{Dt} = \mathbf{V}(\mathbf{r}, t) \rightarrow \mathbf{r}^{t+\Delta} - \mathbf{r}_d^t = \int_t^{t+\Delta} \mathbf{V}(\mathbf{r}, t) dt, \quad (6)$$

346 where $\mathbf{r}^{t+\Delta}$ and \mathbf{r}_d^t are the position of the arrival and the departure point, respectively. From
347 (6), it is obvious that the departure point is far from its arrival point if the wind speed is large.
348 Thus, the departure point of one fluid parcel at the boundary of the subdomain of an MPI task
349 is far from its boundary if the wind speed is large and the wind blows from the outside. To
350 facilitate calculation of the departure point and its interpolation, MPI parallelization adopts
351 a “maximum wind” halo approach so that the halo is sufficiently large for each MPI task to
352 perform its SL calculations in parallel after exchanging the halo. This “maximum wind” halo
353 is named “wide halo” since its width is significantly larger than that of the thin halo of finite
354 difference methods whose stencils have compact support. With numerous MPI tasks, the width
355 of a wide halo may be larger than the subdomain size of its direct neighbour, which implies
356 that the process needs to exchange the halo with its neighbours and its neighbours’ neighbours,
357 which may result in a significant communication overhead which counteracts the efficiency of
358 the favourite SL method, and pose a great challenge to the scalability of the SL method.

359 Fig. 4b demonstrates the halo exchange algorithm adopted in this paper. First, the al-
360 gorithm posts the MPI non-block send and receive operations 1-4 simultaneously for the x-
361 direction sweep. After the x-direction sweep, a y-direction sweep is performed in a similar way
362 but the length of halo is extended to include the left and right ~~haloes~~halo in the x-direction
363 so that the four corners are exchanged properly. This algorithm needs two stages communi-
364 cations, but is simple to implement, especially for the wide halo exchange owing to its fixed
365 regular communication pattern (Fig. 9d). In Fig. 9d, the pixels (near purple colour) tightly
366 attached to the diagonal are due to the exchange in x-direction, the pixels of the same colour
367 but off diagonal are due because of the periodicity in x-direction; the pixels (near orange or red
368 colour) off diagonal are due to the exchange in y-direction, and the pixels of the same colour
369 but far off diagonal are because of the periodicity in y-direction. This algorithm also applies to
370 the thin halo exchange for finite difference methods which is extensively used in the grid-point
371 models. The study emphasizes on the wide halo exchange, but the thin halo exchange is also
372 investigated for comparison (see the red line in Fig. 9a).

3.3 Allreduce in Krylov Subspace Methods for the Semi-Implicit Method

The three-dimensional SI method leads to a large linear system which can be solved by Krylov subspace methods:

$$\mathbf{Ax} = \mathbf{b}, \quad (7)$$

where \mathbf{A} is a non-symmetric sparse matrix. Krylov subspace methods find the approximation \mathbf{x} iteratively in a k -dimensional Krylov subspace:

$$\mathcal{K} = \text{span}(\mathbf{r}, \mathbf{Ar}, \mathbf{A}^2\mathbf{r}, \dots, \mathbf{A}^{k-1}\mathbf{r}), \quad (8)$$

where $\mathbf{r} = \mathbf{b} - \mathbf{Ax}$. To accelerate the convergence, preconditioning is generally used:

$$\mathbf{M}^{-1}\mathbf{Ax} = \mathbf{M}^{-1}\mathbf{b} \quad (9)$$

where \mathbf{M} approximates \mathbf{A} well so that $\mathbf{M}^{-1}\mathbf{A}$ be conditioned better than \mathbf{A} and \mathbf{M}^{-1} can be computed cheaply. The GCR method is a Krylov subspace method of easy implementation and can be used with variable preconditioners. Algorithm 1 of GCR shows that there are two allreduces operations using the sum operation for the inner product in each iteration, thus, it has $2N$ allreduce operations if the GCR iterative solver reaches convergence in N iterations. Allreduce is an all-to-all communication and becomes expensive when the number of iterations becomes larger in GCR solver with numerous MPI processes.

Fig. 4c demonstrates the recursive-k algorithm for the allreduce operation, which is a generalization of the recursive doubling algorithm. The radix k is the number of processes in a group for the recursive-k algorithm. Let $p = \lfloor \log_k(ncpu) \rfloor$, this algorithm has ~~$2+p$~~ p stages of communications if the number of processes is ~~not a power of radix k , otherwise it has two~~ extra stages of communications in the beginning and ending of the algorithm. The following description of the recursive-k algorithm applies to any number of processes, that is, the first and last stage are not necessary when the number of processes is a power of radix k . In the first stage with stage id $j = 0$ (the first row in Fig. 4c), each remaining process whose id $i \notin [0, k^p - 1]$ sends its data to process $i - (ncpu - k^p)$ for the reduce operation. For the stage of stage id $j \in [1, p]$ (rows between the first row and second last row in Fig. 4c), ~~each process whose id~~ all the processes with the same value of $\text{mod}(i, k^{j-1})$ form a list of processes in ascending

Algorithm 1 Preconditioned GCR returns the solution \mathbf{x}_i when convergence occurs where \mathbf{x}_0 is the first guess solution and k is the number of iterations for restart.

```

1: procedure GCR( $\mathbf{A}, \mathbf{M}, \mathbf{b}, \mathbf{x}_0, k$ )
2:    $\mathbf{r}_0 \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}_0$ 
3:    $\mathbf{u}_0 \leftarrow \mathbf{M}^{-1}\mathbf{r}_0$ 
4:    $\mathbf{p}_0 \leftarrow \mathbf{u}_0$ 
5:    $\mathbf{s}_0 \leftarrow \mathbf{A}\mathbf{p}_0$ 
6:    $\gamma_0 \leftarrow \langle \mathbf{u}_0, \mathbf{s}_0 \rangle, \eta_0 \leftarrow \langle \mathbf{s}_0, \mathbf{s}_0 \rangle$   $\triangleright$  Allreduce(sum) of two doubles
7:    $\alpha_0 \leftarrow \frac{\gamma_0}{\eta_0}$ 
8:   for  $i = 1, \dots$ , until convergence do
9:      $\mathbf{x}_i \leftarrow \mathbf{x}_{i-1} + \alpha_{i-1}\mathbf{p}_{i-1}$ 
10:     $\mathbf{r}_i \leftarrow \mathbf{r}_{i-1} - \alpha_{i-1}\mathbf{s}_{i-1}$ 
11:     $\mathbf{u}_i \leftarrow \mathbf{M}^{-1}\mathbf{r}_i$ 
12:    for  $j = \max(0, i - k), \dots, i - 1$  do
13:       $\beta_{i,j} \leftarrow \frac{-1}{\eta_j} \langle \mathbf{A}\mathbf{u}_i, \mathbf{s}_j \rangle$   $\triangleright$  Allreduce(sum) of min(i,k) doubles
14:       $\mathbf{p}_i \leftarrow \mathbf{u}_i + \sum_{j=\max(0, i-k)}^{i-1} \beta_{i,j}\mathbf{p}_j$ 
15:       $\mathbf{s}_i = \mathbf{A}\mathbf{p}_i$ 
16:       $\gamma_i \leftarrow \langle \mathbf{u}_i, \mathbf{s}_i \rangle, \eta_i \leftarrow \langle \mathbf{s}_i, \mathbf{s}_i \rangle$   $\triangleright$  Allreduce(sum) of two doubles
17:       $\alpha_i \leftarrow \frac{\gamma_i}{\eta_i}$ 
18:   return  $\mathbf{x}_i$ 

```

400 order of i , where $i \in [0, k^p - 1]$ only reduces with the processes that are a distance of i is the
401 process id and $\text{mod}(i, k^{j-1})$ is the remainder of i divided by k^{j-1} apart from itself. Then,
402 every k processes in this ordered list form a group of processes, i.e., the first k processes form
403 the first group, the second k processes form the second group, Each group of processes
404 perform their allreduce operation independently. In the final stage with stage id $j = 1 + p$ (the
405 second last row in Fig. 4c), each process whose id $i \notin [0, k^p - 1]$ receives its final result from
406 process $i - (ncpu - k^p)$. The recursive-k algorithm uses large radix k to reduce the stages of
407 communications and the overall communication time.

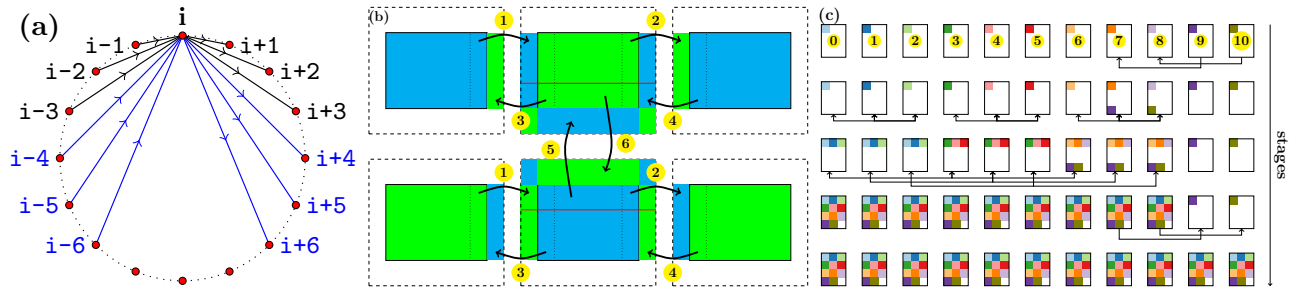


Fig. 4: Algorithms for three key MPI operations: (a) is the ring-k algorithm with k radix for all-to-all personalized communication generalized from ring alltoallv algorithm, (b) is the halo exchange algorithm, and (c) is the recursive-k algorithm with k radix generalized from the recursive doubling algorithm.

Table 3: A three-dimensional grid for assessing the communication of the atmospheric model. Δx and Δy are given as if this grid is a uniform global longitude-latitude grid. In fact, this grid resembles the grid of a regional spectral atmospheric model or the uniform longitude-latitude grid used by some global models.

nx	ny	nz	Δx	Δy	grid points
28800	14400	256	0.0125°	0.0125°	> 100 billion
memory size			max processes		
> 800 GB per double field			3686400 for a 2D partition		

4 Experimental Results

4.1 Experiment Design

In the next decade, it is estimated the resolution of global NWP model will approach kilometre-scale and the HPC will move towards exascale. What would the performance of a global NWP model with a very high resolution on exascale HPC be? In this paper, we are especially interested in the strong scaling of an atmospheric model, that is, how does the atmospheric model with fixed resolution (such as the one presented in Table 3) behave as the number of processes increases? In this study, these strong scalings of the three key MPI operations in the atmospheric model are assessed for $10^2, 2 \times 10^2, \dots, 9 \times 10^2, 10^3, 2 \times 10^3, \dots, 9 \times 10^3, 10^4, 2 \times 10^4, \dots, 9 \times 10^4, 10^5, 2 \times 10^5, \dots, 9 \times 10^5, 10^6$ MPI tasks; but the maximum number of processes is 2×10^5 for the MPI transposition owing to the hard time limitation in our cluster. Table 3 presents a summary of the three-dimensional grid for assessing the communication of the kilometre-scale atmospheric model. The number of grid points of this grid is beyond 100 billion, and one field of double precision variable for this grid requires more than 800 gigabytes of memory. Only with such a large grid, is it possible to perform a 2D domain decomposition for a spectral model with more than one million processes so that modelling the communication of the atmospheric model at exascale HPC become possible.

Besides the topology and its configuration, the routing algorithm, and the collective MPI algorithm; the bandwidth and the latency of the interconnect network of an HPC system have a great impact on the performance of communications. First, we simulate the transposition for the spectral transform method in the simulator for three topologies (torus-M, fattree-M, and dragonfly-MM in Table 1), three configurations of dragonfly topology (dragonfly-MM, dragonfly-SL, and dragonfly-LS in Table 1), three routing algorithms (minimal, valiant, and

431 ugal), and three alltoally algorithms (Table 2). In addition, we compare the simulations of the
 432 transposition for the spectral transform method between four interconnect bandwidths (10^0 ,
 433 10^1 , 10^2 , and 10^3 GB/s) and between four interconnect latencies (10^1 , 10^2 , 10^3 , and 10^4 ns).
 434 After a thorough investigation of the transposition for the spectral transform method, we test
 435 the halo exchange for the SL method with different halo widths (3, 10, 20, and 30 grid points),
 436 three topologies (torus-L, fattree-L, dragonfly-ML in Table 1), and three routing algorithms
 437 (minimal, valiant, and ugal). Finally, the allreduce operation in Krylov subspace methods for
 438 the SI method is evaluated on different topologies (torus-L, fattree-M, dragonfly-ML in Table
 439 1), and the statistics of the optimal radix of recursive-k algorithms for allreduce operations are
 440 presented.

441 **4.2 Transposition for the Spectral Transform Method**

442 Fig. 5a shows that the communication times for the burst, bruck, ring-1, and ring-4 algorithms
 443 decrease as the number of MPI processes increases. The ring-1 and ring-4 algorithms are
 444 almost identical for less than 5×10^4 MPI processes, but ring-4 performs better than ring-1 for
 445 more than 10^5 MPI processes. The burst and bruck algorithms perform worse than the ring-k
 446 algorithm. The SST/macro simulator cannot simulate the burst algorithm for more than 2×10^4
 447 MPI processes because the burst messages result in huge events and large memory footprint.
 448 The communication time of the bruck algorithm is significantly larger than that of the ring-k
 449 algorithm for less than 10^5 MPI processes; however, for a greater number of processes, it is
 450 better than the ring-1 algorithm since the bruck algorithm is targeted for small messages, and
 451 the more processes, the smaller message for a fixed sized problem. The performance of these
 452 alltoally algorithms is confirmed by actually running the skeleton program of transposition
 453 for the spectral transform method with 10^4 MPI processes on the research cluster of Météo
 454 France (Beaufix), which shows that the ring-4 algorithm is even better than the INTEL native
 455 MPI_Alltoally function (Fig. 6).

456 The differences in the communication times of the transpositions between the topology
 457 torus-M, fattree-M, and dragonfly-MM can be an order of magnitude (Fig. 5b). Messages have
 458 to travel a long distance in the topology torus-M which is a 3D torus, so its communication
 459 time is the largest. The best performance of the topology fattree-M can be attributed to its
 460 non-blocking D-mod-k routing algorithm, but its communication time gradually increases as

461 the number of MPI processes increases beyond 10^4 . The performance of topology dragonfly-
462 MM is between that of torus-M and fattree-M (Fig. 5b), it can achieve a better performance by
463 tuning the configuration of the dragonfly topology (Fig. 5c). By comparing Fig. 5b and Fig. 5c,
464 we can see that the topologies of dragonfly-SL and dragonfly-LS are still not as good as the
465 fattree-M, but their performance is very close to that of fattree-M and they lose less scalability
466 than fattree-M for more than 5×10^4 MPI processes.

467 The differences in communication time of the transpositions between the routing algorithms
468 of minimal, valiant and ugal are also an order of magnitude (Fig. 5d), which indicates that the
469 impact of routing algorithm on communication is significant. The valiant routing algorithm
470 performs the best, but the communication time begins to increase when the number of MPI
471 processes is larger than 3×10^4 . The ugal routing algorithm performs the worst, and the
472 performance of minimal routing algorithm is in between that of valiant and ugal routing al-
473 gorithms. The valiant routing algorithm has the longest path for messages from the source to
474 the destination with a randomly chosen intermediate node; thus, theoretically, its communica-
475 tion time is larger. On the contrary, the minimal routing algorithm that moves the messages
476 using the shortest path from the source to the destination has the smallest communication
477 time. The congestion between processes in Fig. 7 shows that the valiant routing algorithm for
478 the dragonfly-MM topology (Fig. 7b) and the minimal routing algorithm for the dragonfly-SL
479 topology (Fig. 7d) are less congested and have a more uniform congestion, the minimal routing
480 algorithm for the dragonfly-MM topology is moderately congested, but its congestion is not
481 uniform (Fig. 7a), the congestion of the ugal routing algorithm for the dragonfly-MM topology
482 is large and highly non-uniform (Fig. 7c). These congestions in Fig. 7 are consistent with the
483 communication times in Fig. 5c and Fig. 5d, that is, the more uniform congestion, the lower
484 communication time because the latter is determined by the longest delay event and uniform
485 congestion can avoid the hotspot of the congestion with the longest delay event. Fig. 8 con-
486 firms this that a high percentage of delay events has a delay time of less than 30 us using the
487 valiant routing algorithm for the dragonfly-MM topology and the minimal routing algorithm
488 for the dragonfly-SL topology; however the minimal routing algorithm for the dragonfly-MM
489 topology has a significant percentage of events that delays by more than 50 us, especially there
490 are a large number of events delayed by more than 100 us using the ugal routing algorithm
491 for the dragonfly-MM topology. Thus, the configuration of the interconnect network and the
492 design of its routing algorithm should make the congestion as uniform as possible if congestion

493 is inevitable.

494 Although the communication time with a bandwidth of 10^0 GB/s is apparently separated
495 from those with bandwidths of 10^1 , 10^2 , and 10^3 GB/s, the curves describing the communication
496 times with bandwidths of 10^1 , 10^2 , and 10^3 GB/s overlap (Fig. 5e). The communication times
497 with latencies of 10^1 and 10^2 ns are almost identical; that with a latency of 10^3 (10^4) ns is
498 slightly (apparently) different from those with latencies of 10^1 and 10^2 ns (Fig. 5f). Equation
499 (1) indicates that the communication time stops decreasing only when α (β) approaches zero and
500 β (α) is constant. Neither α in Fig. 5e nor β in Fig. 5f approaches zero, but the communication
501 time stops decreasing. The inability of the analytical model (1) to explain this suggests that
502 other dominant factors such as congestion contribute to the communication time. Latency
503 is the amount of time required to travel the path from one location to another. Bandwidth
504 determines how ~~many~~ much data per second can be moved in parallel along that path, and
505 limits the maximum number of packets travelling in parallel. Because both α and β are greater
506 than zero, congestion occurs when data arrives at a network interface at a rate faster than the
507 media can service; when this occurs, packets must be placed in a queue to wait until earlier
508 packets have been serviced. The longer the wait, the longer the delay and communication
509 time. Fig. 8b and Fig. 8c show the distributions of the delay caused by congestion for different
510 bandwidths and different latencies, respectively. In Fig. 8b, the distributions of the delay for
511 bandwidths of 10^1 , 10^2 , and 10^3 GB/s are almost identical, which explains their overlapped
512 communication times in Fig. 5e; and the distribution of the delay for a bandwidth of 10^0 GB/s
513 is distinct from the rest since near 20 percent of events are delayed by less than 10 us but a
514 significant percentage of events are delayed more than 100 us, which accounts for its largest
515 communication time in Fig. 5e. In Fig. 8c, the distributions of the delay for latencies of 10^1 and
516 10^2 ns are the same; the distributions of the delay for a latency of 10^3 ns is slightly different from
517 the formers; but the distributions of the delay for a latency of 10^4 ns has a large percentage of
518 events in the right tail which resulted in the longest communication time; these are consistent
519 with their communication times in Fig. 5f.

520 In summary, the alltoallv algorithm, the topology and its configuration, the routing al-
521 gorithm, the bandwidth, and the latency have great impacts on the communication time of
522 transpositions. In addition, the communication time of transpositions decreases as the number
523 of MPI processes increases in most cases; however, this strong scalability is not applicable for
524 the fattree-M topology (the red line in Fig. 5b), the dragonfly-SL and dragonfly-LS topologies

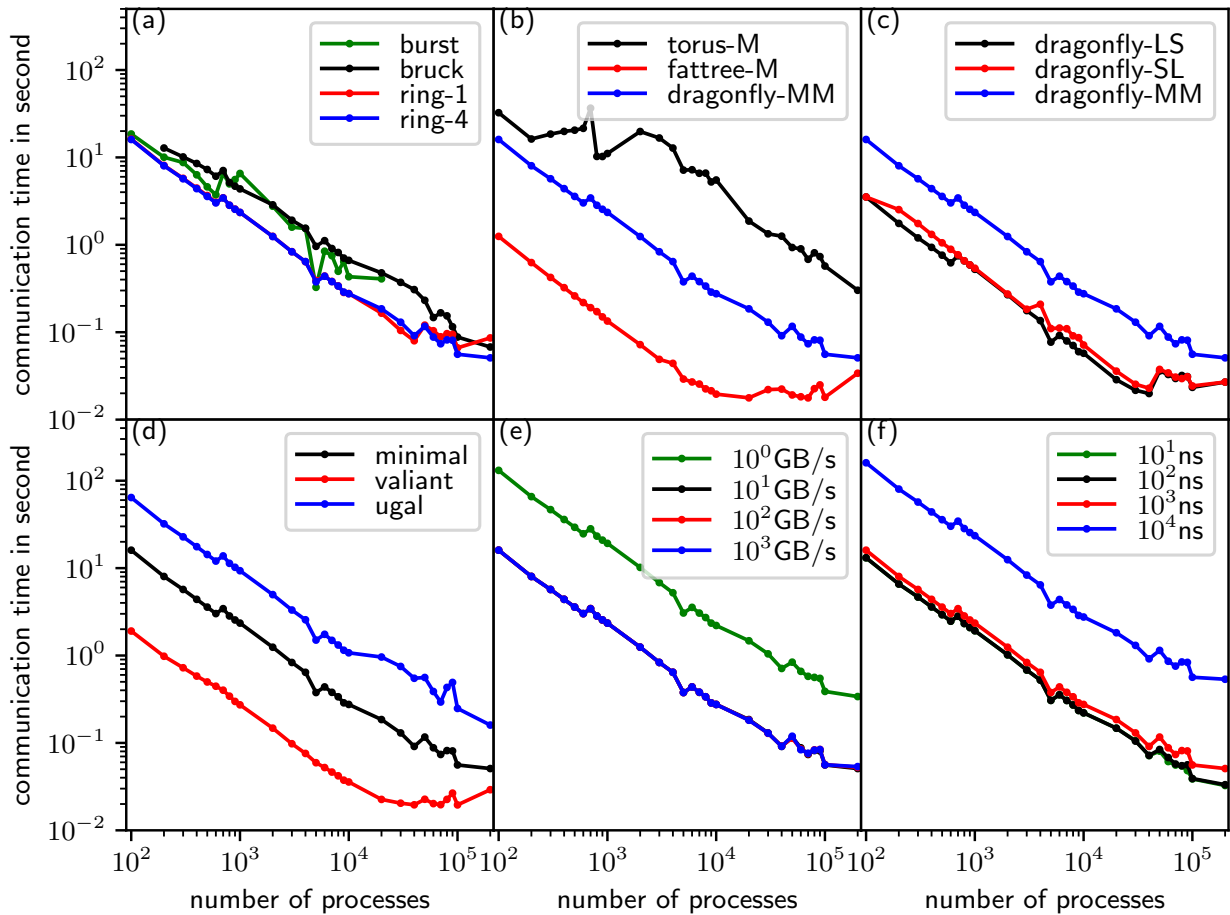


Fig. 5: Communication times of transposition for (a) alltoall algorithms, (b) topologies, (c) configurations of the dragonfly topology, (d) routing algorithms for the dragonfly topology, (e) bandwidth, and (f) latency. The circle markers indicate the numbers of processes of the corresponding simulations.

525 (red and black lines in Fig. 5c), and the valiant routing algorithm (the red line in Fig. 5d) when
 526 the number of MPI processes is large. Thus, the topology of the interconnect network and its
 527 routing algorithm have a great impact on the scalability of transpositions for the spectral trans-
 528 form method. Since the transposition for spectral transform method is a multiple simultaneous
 529 all-to-all personalized communication, congestion has a great impact on its performance.

530 4.3 Halo Exchange for the Semi-Lagrangian Method

531 The most common application of the wide halo exchange is the SL method. For the resolution
 532 of 0.0125° in Table 3 and a time step of 30 seconds, the departure is approximately 5 grid points
 533 away from its arrival if the maximum wind speed is 200 m/s; therefore, the width of the halo
 534 is at least 7 grid points using the ECMWF quasi-cubic scheme (Ritchie, 1995); there are more
 535 grid points if a higher order scheme such as the SLICE-3D (Zerroukat and Allen, 2012) is used.

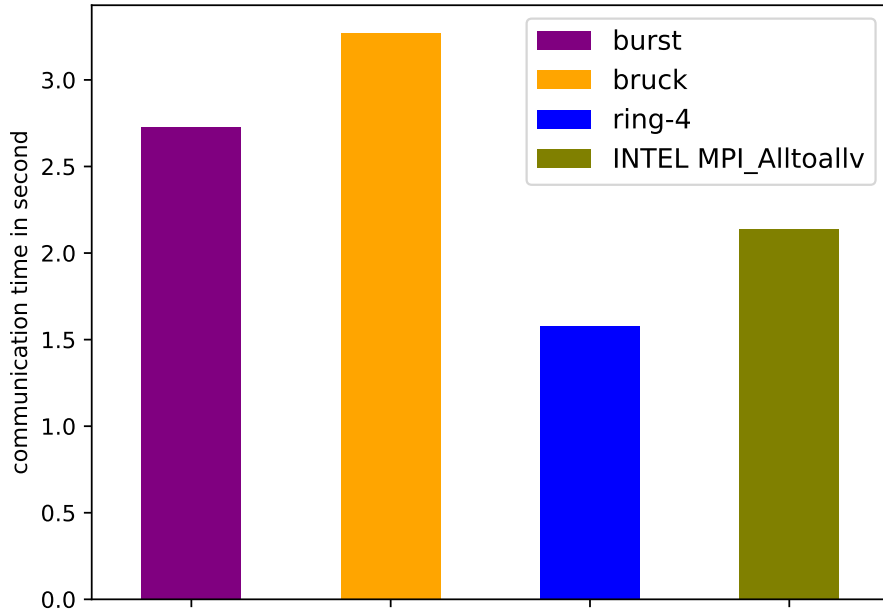


Fig. 6: Actual communication time of transposition for the spectral transform method with 10^4 MPI processes run on beaufix cluster in Météo France.

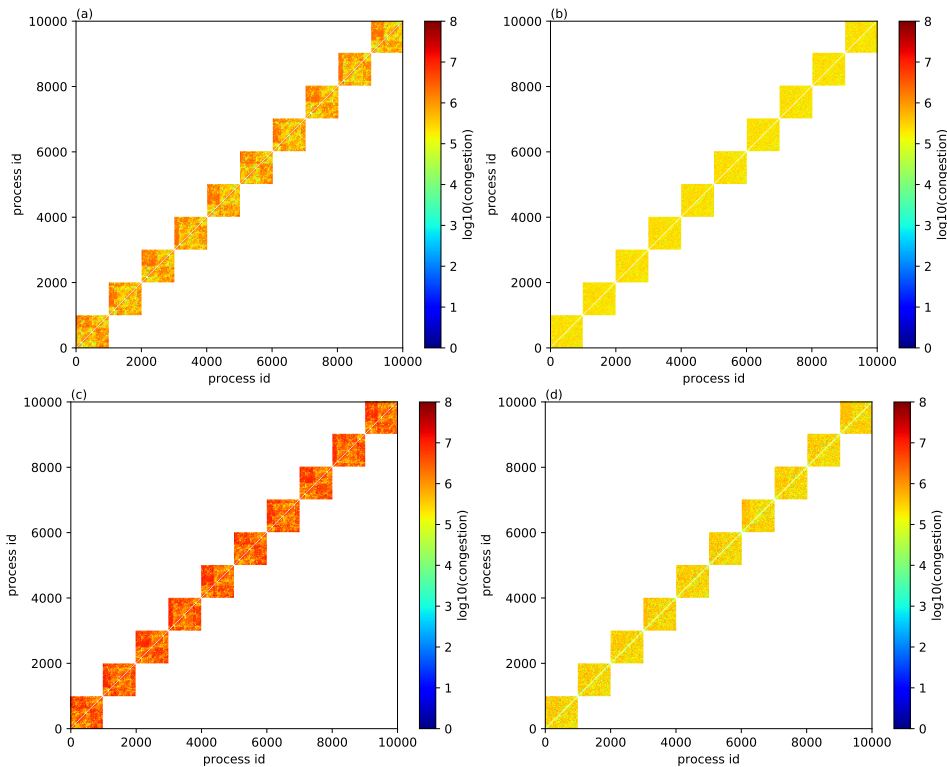


Fig. 7: Congestion of transposition using (a) minimal routing algorithm for the dragonfly-MM topology, (b) valiant routing algorithm for the dragonfly-MM topology, (c) ugal routing algorithm for the dragonfly-MM topology, and (d) minimal routing algorithm for the dragonfly-SL topology.

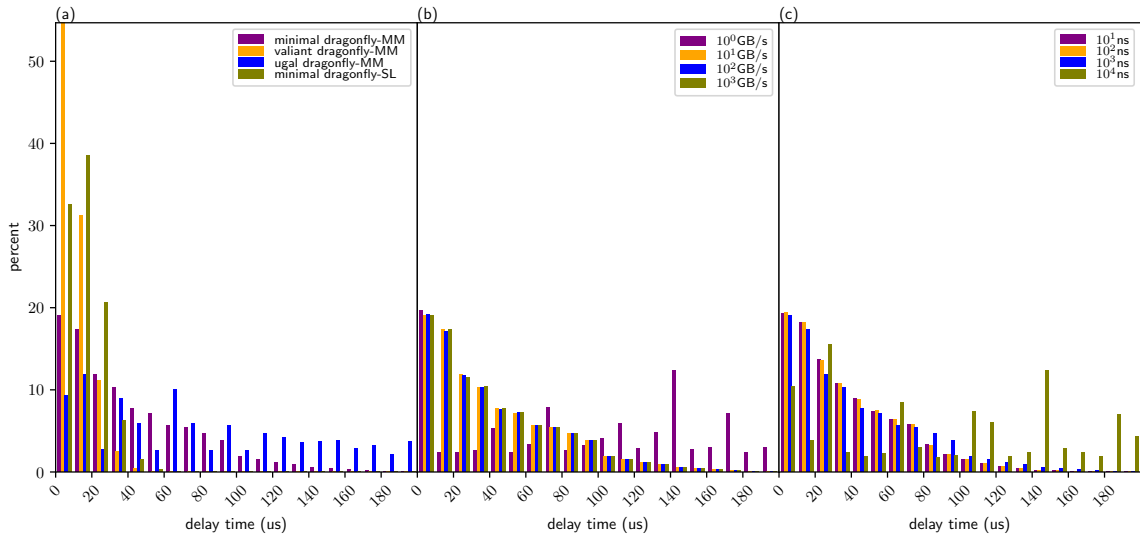


Fig. 8: Distribution of delayed events of transposition for the spectral transform method with 10^4 MPI processes using (a) different routing algorithms and topology configurations, (b) different bandwidths, and (c) different latencies, simulated by SST/macro.

536 In Fig. 9a, the communication time of the halo exchange decreases more slowly ~~as the with~~
537 increasing number of processes ~~increases~~ than that of transposition for the spectral transform
538 method. This is because the message size decreases more slowly than that of transposition owing
539 to the fixed width of the halo (figure omitted). If the communication time of the transposition
540 (halo exchange) continues its decreasing (increasing) trend in Fig. 9a, they meet at certain
541 number of MPI processes; then, the communication time of the halo exchange is larger than
542 that of the transposition. In addition, it can be seen that the wider the halo, the longer the
543 communication time. The halo exchange of a thin halo of 3 grid points, for such as the 6th
544 order central difference $F'_i = \frac{-F_{i-3}+9F_{i-2}-45F_{i-1}+45F_{i+1}-9F_{i+2}+F_{i+3}}{60\Delta}$ (the red line in Fig. 9a), is
545 significantly faster than that of wide halo for the SL method (green and blue lines in Fig. 9a).
546 Thus, the efficiency of the SL method is counteracted by the overhead of the wide halo exchange
547 where the width of the halo is determined by the maximum wind speed. Wide halo exchange
548 for the SL method is expensive at exascale, especially for the atmospheric chemistry models
549 where a large number of tracers need to be transported. On-demand exchange is a way to
550 reduce the communication of halo exchange for the SL method, and will be investigated in a
551 future study.

552 Significant differences in the communication times of the wide halo exchange of 20 grid
553 points for topology torus-L, fattree-L, and dragonfly-ML are shown in Fig. 9b. It can be
554 seen that topology torus-L performs the worst, fattree-L is the best, and the performance of
555 dragonfly-ML is between that of torus-L and fattree-L. The communication time of the wide

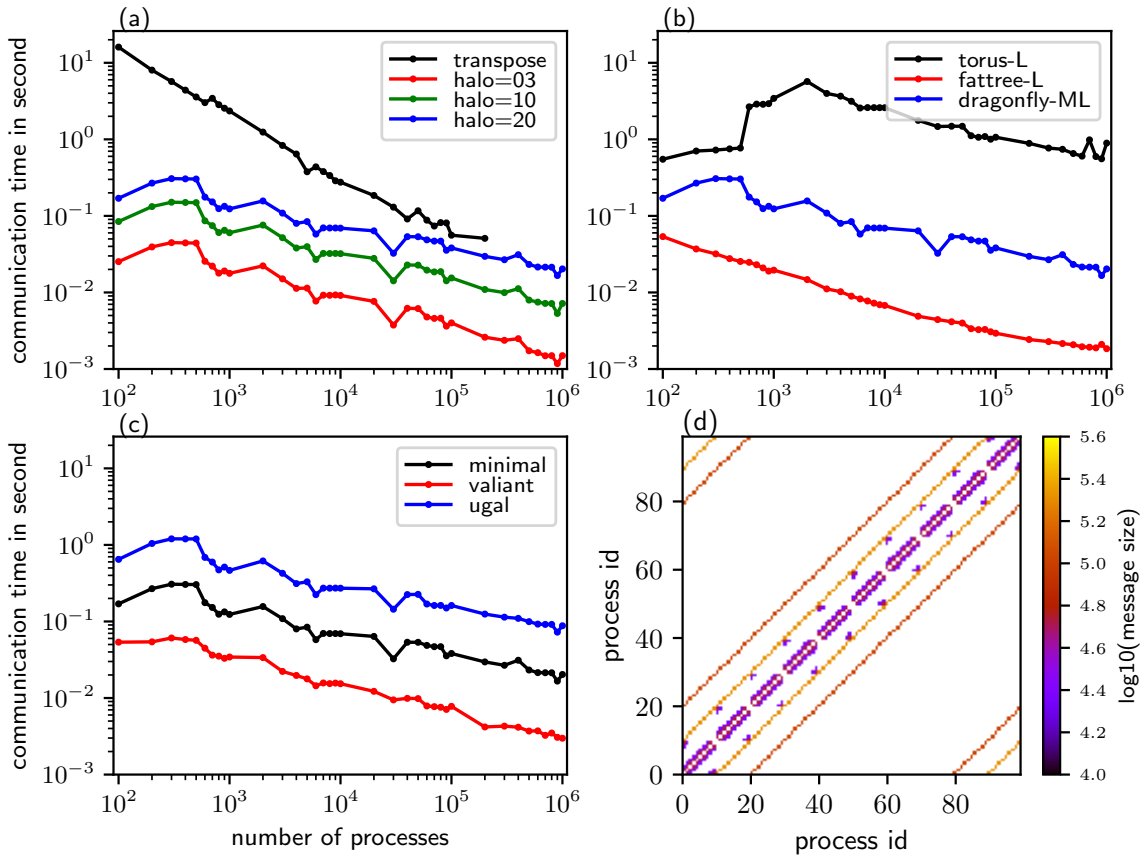


Fig. 9: (a) is the communication times of the halo exchange with a halo of 3 (red line), 10 (green line), and 20 (blue line) grid points, and the communication time of transposition for the spectral transform method is shown for comparison (black line). (b) is the communication times of the halo exchange with a halo of 20 grid points for the topology of torus-L (black line), fattree-L (red line), and dragonfly-ML (blue line). (c) is the communication times of the halo exchange with a halo of 20 grid points for the routing algorithm of minimal (black line), valiant (red line), and ugal (blue line). (d) illustrates the communication pattern of the halo exchange with a wide halo. The circle markers in (a)–(c) indicate the numbers of processes of the corresponding simulations.

556 halo exchange of 20 grid points for the topology tour-L abruptly increases at approximately 10^3
 557 MPI processes, and then gradually decreases when the number of MPI tasks becomes larger
 558 than 3×10^3 MPI processes. The impact of the routing algorithm on the communication time
 559 of the wide halo exchange of 20 grid points (Fig. 9c) is the same as on that of transposition
 560 (Fig. 5d): the routing algorithm valiant performs the best, the routing algorithm ugal performs
 561 the worst, and the routing algorithm minimal is between valiant and ugal.

562 4.4 Allreduce in Krylov Subspace Methods for the Semi-Implicit Method

563 If, in average, the GCR with a restart number $k = 3$ is convergent with $N = 25$ iterations, the
 564 number of allreduce calls is $2 \times N = 50$. The black and blue lines are the communication times

565 of 50 allreduce operations using MPI_Allreduce and the recursive-k algorithm, respectively;
566 that is, the estimated communication time of one single GCR call (Fig. 10a). Contrary to that
567 of transposition, the communication time of GCR increases as the number of MPI processes
568 increases. Following the trend, the communication of a single GCR call may be similar to or
569 even larger than that of a single transposition when the number of MPI processes approaches
570 to or is beyond one million. Although it is believed that the spectral method does not scale
571 well owing to its time-consuming transposition, it does not suffer from this expensive allreduce
572 operation for the SI method because of its mathematical advantage that spherical harmonics are
573 the eigenfunctions of Helmholtz operators. In this sense, a grid-point model with the SI method
574 in which the three-dimensional Helmholtz equation is solved by Krylov subspace methods may
575 also not scale well at exascale unless the overhead of allreduce communication can be mitigated
576 by overlapping it with computation (Sanan et al., 2016).

577 Fig. 10b shows the communication times of allreduce operations using the recursive-k algo-
578 rithm on the topologies of torus-L, fattree-L, and dragonfly-ML. The impact of topology on the
579 communication performance of allreduce operations is obvious. The topology of torus-L has the
580 best performance, but is similar to that of dragonfly-ML for more than 5×10^5 MPI processes;
581 and fattree-L has the worst performance. However, the impact of three routing algorithms
582 (minima, valiant, and ugal) for the dragonfly-ML topology has a negligible impact on the com-
583 munication performance of allreduce operations (figure omitted); this may be because of the
584 tiny messages (only 3 doubles for the restart number $k = 3$) communicated by the allreduce
585 operation.

586 One advantage of the recursive-k algorithm of the allreduce operation is that the radix k
587 can be selected to reduce the stages of communication by making full use of the bandwidth
588 of the underlying interconnect network. We repeat the experiment, whose configuration is
589 as that of the blue line in Fig. 10a, for the proper radix $k \in [2, 32]$, and the optimal radix
590 is that with the lowest communication time for a given number of MPI processes. For each
591 number of MPI processes, there is an optimal radix. The statistics of all the optimal radices are
592 shown in Fig. 10c. It can be seen that the minimum and maximum optimal radices are 5 and
593 32, respectively. Thus, the recursive doubling algorithm that is equivalent to the recursive-k
594 algorithm with radix $k=2$ is not efficient since the optimal radix is at least 5. The median
595 number of optimal radices is approximately 21, and the mean number is less than but very
596 close to the median number. We cannot derive an analytic formula for the optimal radix since

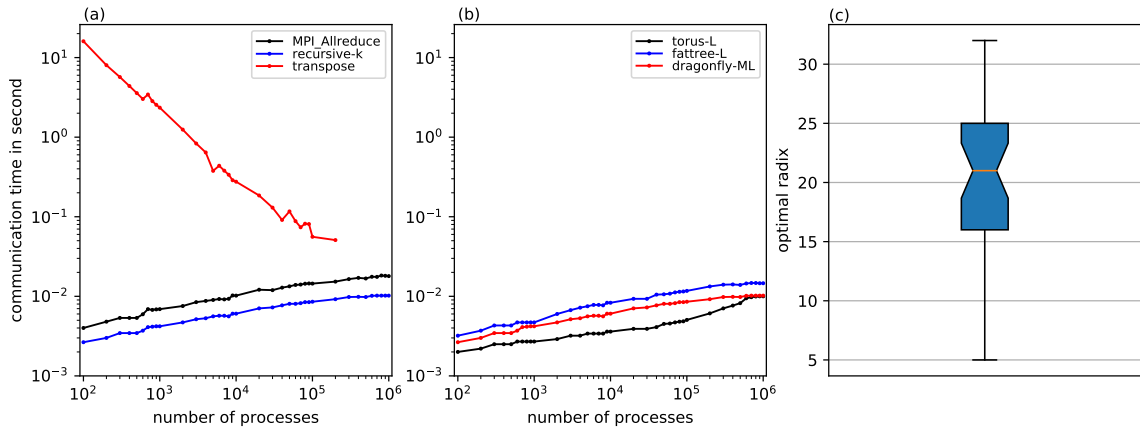


Fig. 10: (a) is the communication times of the allreduce operation using the MPI_Allreduce (black line) and the recursive-k algorithm (blue line), and the communication time of transposition for the spectral transform method is shown for comparison (red line). (b) is the communication times of the allreduce operation using the recursive-k algorithm for the topology torus-L (black line), fattree-L (blue line), and dragonfly-ML (red line). (c) is the statistics of the optimal radices for the recursive-k algorithm. The circle markers in (a)–(b) indicate the numbers of processes of the corresponding simulations.

597 modelling the congestion is difficult in an analytic model. However, for a given resolution of
 598 NWP model and a given HPC system, fortunately, the number of processes, bandwidth, and
 599 latency are fixed; thus, it is easy to perform experiments to obtain the optimal radix.

600 5 Conclusion and Discussion

601 This work shows that it is possible to make simulations of the MPI patterns commonly used in
 602 NWP models using very large numbers of MPI tasks. This enables the possibility to examine
 603 and compare the impact of different factors such as latency, bandwidth, routing and network
 604 topology on response time. We have provided an assessment of the performance and scalability
 605 of three key MPI operations in an atmospheric model at exascale by simulating their skeleton
 606 programs on an SST/macro simulator. After optimization of the memory and efficiency of
 607 the SST/macro simulator and construction of the skeleton programs, a series of experiments
 608 was carried out to investigate the impacts of the collective algorithm, the topology and its
 609 configuration, the routing algorithm, the bandwidth, and the latency on the performance and
 610 scalability of transposition, halo exchange, and allreduce operations. The experimental results
 611 show that:

- 612 1. The collective algorithm is extremely important for the performance and scalability of
 613 key MPI operations in the atmospheric model at exascale because a good algorithm can

614 make full use of the bandwidth and reduce the stages of communication. The generalized
615 ring-k algorithm for the alltoally operation and the generalized recursive-k algorithm for
616 the allreduce operation proposed herein perform the best.

- 617 2. Topology, its configuration, and the routing algorithm have a considerable impact on the
618 performance and scalability of communications. The fattree topology usually performs
619 the best, but its scalability becomes weak with a large number of MPI processes. The
620 dragonfly topology balances the performance and scalability well, and can maintain almost
621 the same scalability with a large number of MPI processes. The configurations of the
622 dragonfly topology indicate that a proper configuration can be used to avoid the hotspots
623 of congestion and lead to good performance. The minimal routing algorithm is intuitive
624 and performs well. However, the valiant routing algorithm (which randomly chooses an
625 intermediate node to uniformly disperse the communication over the network to avoid
626 the hotspot/bottleneck of congestion) performs much better for heavy congestion.
- 627 3. Although they have an important impact on communication, bandwidth and latency
628 cannot be infinitely grown and reduced owing to the limitation of hardware, respectively.
629 Thus, it is important to design innovative algorithms to make full use of the bandwidth
630 and to reduce the effect of latency.
- 631 4. It is generally believed that the transposition for the spectral transform method, which is
632 a multiple simultaneous all-to-all personalized communication, poses a great challenge to
633 the scalability of the spectral model. This work shows that the scalability of the spectral
634 model is still acceptable in terms of MPI transposition. However, the wide halo exchange
635 for the Semi-Lagrangian method and the allreduce operation in the GCR iterative solver
636 for the Semi-Implicit method, both of which are often adopted by the grid-point model,
637 also suffer the stringent challenge of scalability at exascale.

638 In summary, both software (algorithms) and hardware (characteristics and configuration)
639 are of great importance to the performance and scalability of the atmospheric model at exascale.
640 The software and hardware must be co-designed to address the challenge of the atmospheric
641 model for exascale computing.

642 As shown previously, the communications of the wide halo exchange for the Semi-Lagrangian
643 method and the allreduce operation in the GCR iterative solver for the Semi-Implicit method

644 are expensive at exascale. The on-demand halo exchange for the Semi-Lagrangian and the
645 pipeline technique to overlap the communication with the computation for the GCR itera-
646 tive solver are not researched in this study and should be investigated. All the ~~computed~~
647 compute nodes in this work only contain one single-core CPU, which is good for assessing the
648 communication of the interconnect network; however, the architectures of current and future
649 supercomputers are multi-core and multi-socket nodes, even non-CPU architectures. These
650 more complex hierarchies seem to complicate the inter-process communications. However, an
651 MPI rank can be bound to any core for multi-core and multi-socket nodes; ~~with INTEL MPI~~
652 ~~library.~~ For example, an MPI rank can be bound to any processor/co-processor for MIC ar-
653 chitectures such as Xeon Phi ; ~~with CUDA-aware MPI~~ ; using the INTEL MPI library, and
654 an MPI rank can be bound to a CPU core but can communicated with GPUs for GPU ar-
655 chitectures using a CUDA-aware MPI. Because a multi-core node behaves more or less like a
656 more powerful single core node when the OpenMP is used for the intra-node parallelization,
657 the conclusions in this study could be generalized to the complex hierarchical system. Multiple
658 MPI processes per node may be good for the local pattern communication such as thin halo
659 exchange since the shared memory communication mechanism is used, but may result in con-
660 gestion in the network interface controller for inter-node communication. The congestion can
661 be mitigated or even eliminated, if each node has more network interface controllers (NICs) or
662 a network interface controller with multi-ports (as a mini-switch). From this point of view, the
663 conclusions should still be valid for the complex hierarchical architectures, but the scalability
664 might be affected. The more MPI processes, the less computation per node ~~without limitation~~
665 if there is only one single-core CPU per node, thus, computation is not considered in this pa-
666 per. ~~However, the bandwidth of memory limits the performance and scalability of computation~~
667 ~~for~~ Because multi-core or many-core ~~systems~~ processors share a memory bus, it is possible for
668 a memory-intensive application (such as an atmospheric model) to saturate the memory bus
669 and result in degraded performances of all the computations running on that processor. The
670 assessment of ~~computation~~ computations is currently underway and a detailed paper will be
671 presented separately; the purpose of this subsequent study is to model the time response of a
672 time step of a model such as the regional model (AROME) used by Météo-France.

673 Code Availability

674 The code of the SST/macro simulator is publicly available at [https://github.com/sstsimulator/sst-](https://github.com/sstsimulator/sst-macro)
675 macro. The skeleton programs, scripts, and our modified version of SST/macro 7.1.0 for the
676 simulations presented the paper are available at <https://doi.org/10.5281/zenodo.1066934>.

677 Competing Interests

678 The authors declared no competing interests.

679 Acknowledgements

680 This work was supported by ~~Centre National de Recherches Météorologiques, Météo France~~
681 ~~within~~ the ESCAPE (~~Energy-efficient Scalable Algorithms for Weather Prediction at Exascale~~[Energy-efficient](#)
682 [Scalable Algorithms for Weather Prediction at Exascale](#)) project. [The ESCAPE project has](#)
683 [received funding from the European Union's Horizon 2020 research and innovation programme](#)
684 [under grant agreement No 671627.](#)

685 References

- 686 Acun, B., N. Jain, A. Bhatele, M. Mubarak, C. D. Carothers, and L. V. Kale. *Preliminary*
687 *Evaluation of a Parallel Trace Replay Tool for HPC Network Simulations*, pages 417–429.
688 Springer International Publishing, Cham, 2015. ISBN 978-3-319-27308-2.
- 689 Ajima, Y., S. Sumimoto, and T. Shimizu, Nov 2009: Tofu: A 6d mesh/torus interconnect for
690 exascale computers. *Computer*, **42**(11), 36–40.
- 691 Alverson, B., E. Froese, L. Kaplan, and D. Roweth. *Cray XC Series Network*. Cray Inc., 2015.
- 692 Barros, S. R. M., D. Dent, L. Isaksen, G. Robinson, G. Mozdzyński, and F. Wollenweber, 1995:
693 The IFS model: a parallel production weather code. *Parallel Comput.*, **21**, 1621–1638.
- 694 Bauer, P., A. Thorpe, and G. Brunet, 2015: The quiet revolution of numerical weather predic-
695 tion. *Nature*, **525**(7567), 47–55.

- 696 Böhm, S. and C. Engelmann. xsim: The extreme-scale simulator. In *2011 International Con-*
697 *ference on High Performance Computing Simulation*, pages 280–286, July 2011.
- 698 Casanova, H., A. Gupta, and F. Suter, 2015: Toward more scalable off-line simulations of mpi
699 applications. *Parallel Processing Letters*, **25**(03), 1541002.
- 700 Dechev, D. and T. H. Ahn, 2013: Using sst/macro for effective analysis of mpi-based applica-
701 tions: Evaluating large-scale genomic sequence search. *IEEE Access*, **1**, 428–435.
- 702 Degomme, A., A. Legrand, G. S. Markomanolis, M. Quinson, M. Stillwell, and F. Suter, 2017:
703 Simulating MPI Applications: The SMPI Approach. *IEEE Transactions on Parallel and*
704 *Distributed Systems*, **28**(8), 2387–2400.
- 705 Dubos, T., S. Dubey, M. Tort, R. Mittal, Y. Meurdesoif, and F. Hourdin, 2015: DYNAMICO-
706 1.0, an icosahedral hydrostatic dynamical core designed for consistency and versatility.
707 *Geosci. Model Dev.*, **8**, 3131–3150.
- 708 Ehrendorfer, M., 2012: *Spectral numerical weather prediction models*. SIAM.
- 709 Eisenstat, S. C., H. C. Elman, and M. H. Schultz, 1983: Variational iterative methods for
710 nonsymmetric systems of linear equations. *SIAM J. Numer. Anal.*, **20**(2), 345–357.
- 711 Engelmann, C., 2014: Scaling to a million cores and beyond: using light-weight simulation
712 to understand the challenges ahead on the road to exascale. *Future Generation Computer*
713 *Systems*, **30**(0), 59–65.
- 714 Hoeffler, T., T. Schneider, and A. Lumsdaine. LogGOPSim - Simulating Large-Scale Appli-
715 cations in the LogGOPS Model. In *Proceedings of the 19th ACM International Symposium*
716 *on High Performance Distributed Computing*, pages 597–604. ACM, Jun. 2010. ISBN 978-1-
717 60558-942-8.
- 718 Hortal, M., 2002: The development and testing of a new two-time-level semi-Lagrangian scheme
719 (SETTLS) in the ECMWF forecast model. *Q. J. R. Meteorol. Soc.*, **128**, 1671–1687.
- 720 Hoskins, B. J. and A. J. Simmons, 1975: A multi-layer spectral model and the semi-implicit
721 method. *Q. J. R. Meteorol. Soc.*, **101**, 637–655.

722 Jain, N., A. Bhatele, S. White, T. Gamblin, and L. V. Kale. Evaluating hpc networks via
723 simulation of parallel workloads. In *SC16: International Conference for High Performance*
724 *Computing, Networking, Storage and Analysis*, pages 154–165, Nov 2016.

725 Janssen, C. L., H. Adalsteinsson, S. Cranford, J. P. Kenny, A. Pinar, D. A. Evensky, and
726 J. Mayo, 2010: A Simulator for Large-Scale Parallel Computer Architectures. *International*
727 *Journal of Distributed Systems and Technologies*, **1**(2), 57–73.

728 Juang, H. H., S. Hong, , and M. Kanamitsu, 1997: The NCEP regional spectral model: an
729 update. *Bull. Am. Meteorol. Soc.*, **78**(10), 2125–2143.

730 Kanamitsu, M., H. Kanamaru, Y. Cui, and H. Juang. Parallel implementation of the regional
731 spectral atmospheric model. Technical report, Scripps Institution of Oceanography, Univer-
732 sity of California at San Diego, and National Oceanic and Atmospheric Administration for
733 the California Energy Commission, PIER Energy-Related Environmental Research, 2005.
734 CEC-500-2005-014.

735 Kim, J., W. J. Dally, S. Scott, and D. Abts. Technology-driven, highly-scalable dragonfly
736 topology. In *2008 International Symposium on Computer Architecture*, pages 77–88, June
737 2008.

738 Kuhnlein, C. and P. K. Smolarkiewicz, 2017: An unstructured-mesh finite-volume MPDATA
739 for compressible atmospheric dynamics. *J. Comput. Phys.*, **334**, 16–30.

740 Lagadapati, M., F. Mueller, and C. Engelmann. Benchmark generation and simulation at
741 extreme scale. In *2016 IEEE/ACM 20th International Symposium on Distributed Simulation*
742 *and Real Time Applications (DS-RT)*, pages 9–18, Sept 2016.

743 Leiserson, C. E., Oct 1985: Fat-trees: Universal networks for hardware-efficient supercomput-
744 ing. *IEEE Transactions on Computers*, **C-34**(10), 892–901.

745 Li, L., W. Xue, R. Ranjan, and Z. Jin, 2013: A scalable Helmholtz solver in GRAPES over
746 large-scale multicore cluster. *Concurrency Computat.: Pract. Exper.*, **25**, 1722–1737.

747 Lin, S.-J., 2004: A "vertically Lagrangian" finite-volume dynamical core for global models.
748 *Mon. Wea. Rev.*, **132**, 2293–2307.

749 Mubarak, M., C. D. Carothers, R. B. Ross, and P. Carns, January 2017: Enabling parallel
750 simulation of large-scale hpc network systems. *IEEE Trans. Parallel Distrib. Syst.*, **28**(1),
751 87–100.

752 Noeth, M., P. Ratn, F. Mueller, M. Schulz, and B. R. de Supinski, 2009: Scalatrace: Scalable
753 compression and replay of communication traces for high-performance computing. *Journal*
754 *of Parallel and Distributed Computing*, **69**(8), 696 – 710.

755 Núñez, A., J. Fernández, J. D. Garcia, F. Garcia, and J. Carretero, Jan 2010: New techniques
756 for simulating high performance mpi applications on large storage networks. *The Journal of*
757 *Supercomputing*, **51**(1), 40–57.

758 Qaddouri, A. and V. Lee, 2011: The Canadian global environmental multiscale model on the
759 Yin-Yang grid system. *Q. J. R. Meteorol. Soc.*, **137**, 1913–1926.

760 Ritchie, H., 1995: Implementation of the Semi-Lagrangian Method in a High-Resolution Version
761 of the ECMWF forecast model. *Mon. Wea. Rev.*, **123**, 489–514.

762 Robert, A., J. henderson, and C. Turnbull, 1972: An implicit time integration scheme for
763 baroclinic models of the atmosphere. *Mon. Wea. Rev.*, **100**, 329–335.

764 Sanan, P., S. M. Schnepp, and D. A. May, 2016: Pipelined, flexible krylov subspace methods.
765 *SIAM J. Sci. Comput.*, **38**(5), C441–C470.

766 Sandbach, S., J. Thuburn, D. Vassilev, and M. G. Duda, 2015: A Semi-Implicit version fo the
767 MPAS-atmosphere dynamical core. *Mon. Wea. Rev.*, **143**, 3838–3855.

768 Satoh, M., T. Matsuno, H. Tomita, H. Miura, T. Nasuno, and S. Iga, 2008: Nonhydrostatic
769 icosahedral atmospheric model (NICAM) for global cloud resolving simulations. *J. Comput.*
770 *Phys.*, **227**, 3486–3514.

771 Seity, Y., P. Brousseau, S. Malardel, G. Hello, P. Bénard, F. Bouttier, C. Lac, and V. Masson,
772 2011: The AROME-France convective-scale operational model. *Mon. Wea. Rev.*, **139**, 976–
773 991.

774 Skamarock, W. C., J. B. Klemp, M. G. Duda, L. D. Fowler, and S.-H. Park, 2012: A mul-
775 tiscale nonhydrostatic atmospheric model using centroidal voronoi tesselations and C-grid
776 staggering. *Mon. Wea. Rev.*, **140**, 3090–3105.

777 Smolarkiewicz, P. K., W. Deconinck, M. Hamrud, C. Kühnlein, G. Mozdzynski, J. Szmelter,
778 and N. P. Wedi, 2016: A finite-volume module for simulating global all-scale atmospheric
779 flows. *J. Comput. Phys.*, **314**, 287–304. doi: <https://doi.org/10.1016/j.jcp.2016.03.015>.

780 SNL, L. C. *SST/macro 7.1: User's Manual*. Sandia National Labs, Livermore, CA, Jun 2017.

781 Staniforth, A. and J. Côté, 1991: Semi-Lagrangian integration schemes for atmospheric models—
782 a review. *Mon. Wea. Rev.*, **119**, 2206–2223.

783 Temperton, C., 1983: Self-sorting mixed-radix fast Fourier transforms. *J. Comput. Phys.*, **52**,
784 1–23.

785 Thakur, R., R. Rabenseifner, and W. Gropp, February 2005: Optimization of collective com-
786 munication operations in mpich. *Int. J. High Perform. Comput. Appl.*, **19**(1), 49–66.

787 Tikir, M. M., M. A. Laurenzano, L. Carrington, and A. Snavely. *PSINS: An Open Source*
788 *Event Tracer and Execution Simulator for MPI Applications*, pages 135–148. Springer Berlin
789 Heidelberg, Berlin, Heidelberg, 2009.

790 Wedi, N. P., M. Hamrud, and G. Mozdzynski, 2013: A fast spherical harmonics transform for
791 global NWP and climate models. *Mon. Wea. Rev.*, **141**, 3450–3461.

792 Wedi, N. P., 2014: Increasing horizontal resolution in numerical weather prediction and climate
793 simulations: illusion or panacea? *Phil. Trans. R. Soc. A*, **372**, 20130289.

794 Wike, J. J. and J. P. Kenny. Using Discrete Event Simulation for Programming Model Ex-
795 ploration at Extreme-Scale: Macroscale Components for the Structural Simulation Toolkit
796 (SST). Technical report, Sandia National Laboratories, 2014. SAND2015-1027.

797 Wolfe, N., C. D. Carothers, M. Mubarak, R. Ross, and P. Carns. Modeling a million-node slim
798 fly network using parallel discrete-event simulation. In *Proceedings of the 2016 Annual ACM*
799 *Conference on SIGSIM Principles of Advanced Discrete Simulation*, pages 189–199. ACM,
800 2016.

801 Zahavi, E., G. Johnson, D. J. Kerbyson, and M. Lang, 2010: Optimized infiniband™ fat-tree
802 routing for shift all-to-all communication patterns. *Concurrency and Computation: Practice*
803 *and Experience*, **22**(2), 217–231.

- 804 Zangl, G., D. Reinert, P. Ripodas, and M. Baldauf, 2015: The ICON (icosahedral non-
805 hydrostatic) modelling framework of DWD and MPI-M: description of the non-hydrostatic
806 dynamical core. *Q. J. R. Meteorol. Soc.*, **141**, 563–579.
- 807 Zerroukat, M. and T. Allen, 2012: A three-dimensional monotone and conservative semi-
808 Lagrangian scheme (SLICE-3D) for transport problems. *Q. J. R. Meteorol. Soc.*, **138**, 1640–
809 1651.
- 810 Zheng, G., G. Kakulapati, and L. V. Kale. Bigsim: a parallel simulator for performance
811 prediction of extremely large parallel machines. In *18th International Parallel and Distributed*
812 *Processing Symposium, 2004. Proceedings.*, pages 78–, April 2004.