

ParSA: High-throughput Scientific-data Analysis Framework with Distributed File System

Tao Zhang^{a,*}, XiangZheng Sun^b, Wei Xue^a

^a*Department of Computer Science and Technology, Tsinghua University, Beijing, China*

^b*Intel Corporation, Beijing, China*

Abstract

Scientific-data analysis and visualization has become a key component in nowadays large-scale simulations. Due to the rapidly increasing data volume and awkward I/O pattern among high-structured files, known serial methods/tools cannot scale well and usually lead to poor performance over traditional architectures. In this paper, we propose a new framework: ParSA (parallel scientific-data analysis) for high-throughput and scalable scientific analysis, with distributed file system. ParSA presents the optimization strategies for grouping and splitting logical units to maximize distributed I/O property of distributed file system as well as to maximize overlapping the data reading, processing and transferring during computation. Besides, ParSA provides the similar interfaces as the NCO (NetCDF Operator), which is used in most of climate data diagnostic packages, making it easy to port this framework. We utilize ParSA to accelerate well-known analysis methods for climate models on Hadoop Distributed File System (HDFS). Experimental results demonstrate the high efficiency and scalability of ParSA, getting the maximum 1.2GB/s throughput.

Keywords:

Data intensive, Scientific data analysis, Distributed file system

*Corresponding author. Tel: +86 18210192568

Email address: t-zhang11@mails.tsinghua.edu.cn (Tao Zhang)

1. Introduction

In most of modern scientific applications, huge amounts of data are produced. Large-scale simulations, such as climate modeling, high-energy physics simulation and genome mapping, generate hundreds of terabytes data volumes[1][2]. Additionally, it still increases as the high resolution model developing. As a consequence, analysis of scientific-data is data-intensive.

In fact, almost all of scientific-data are stored in high-structured files, some of which provide parallel I/O interface, such as Network Common Data Format version 4 (NetCDF4)[3], Hierarchical Data Format 5(HDF5)[4] and ADIOS BP data format (BP) [5], and some of which only support serial I/O interface, like Network Common Data Format version 3 (NetCDF3) [6]. All of them are self-describing, machine-independent data format.

In scientific-data analysis, large-scale scientific-data are stored in RAID-5/6 or parallel file system. Yet the analysis methods/ tools are always centralized approaches, such as NCO[7] and NCL[8], which are the most used in climate applications for crossing NetCDF files, leading to very poor scalability and time-consuming performance.

Inspired by big data solution in Internet Port Data Center (IPDC), numerous frameworks with distributed strategy have been developed. MapReduce is a program framework for processing and generating large data sets, providing automatic parallel mechanism and build-in fault-tolerance on a cluster[9]. However, this solution with MapReduce requires data firstly to be transformed into a text-based format[10]. SciHadoop[11] is a Hadoop plugin allowing scientists to specify logical queries over array-based data models. It executes queries as map/reduce programs defined over the logical data model. It shows remarkable improvements for holistic functions of NetCDF data sets for the following optimization goals: reduce total data transfers, reduce remote reads and reduce unnecessary reads. Nevertheless, SciHadoop using java language leads to the compatibility problem to the existing climate data analysis tools, which is written by C shell scripts, NCL and NCO commands. The SWAMP project [12] has provided the parallel NCO operations, but the reading performance is still bottleneck.

In this paper, we propose a new framework — Parallel Scientific-data Analysis (ParSA). We utilize the distributed I/O property with Hadoop Distributed File System (HDFS) to improved data reading throughput. What is more, ParSA optimizes the data layout schedule stored in the distributed file system to overlap the data reading, processing and transferring. Besides,

38 it provides parallel NCO operations, cooperating with HDFS, making it easy
39 to use the efficient tool, without changing a lot for current climate analysis
40 package.

41 **2. HDFS and scientific-data analysis**

42 In this section, we will present the property of distributed file system
43 HDFS, replicas and scheduler, which can be taken advantage of to optimize
44 distributed I/O performance. We also present the character of scientific-
45 data analysis, and discuss about the probability of analysis transportation
46 onto HDFS.

47 *2.1. HDFS*

48 HDFS[13] is an open source project, driven by Google File System (GFS)[14].
49 As a distributed, scalable and portable file system, HDFS is inherent for
50 large-scale data-intensive process.

51 In HDFS, there are two types of node, Namenode and Datanode. Namenode
52 maintains file system tree and metadata for all files or directories stored
53 in HDFS, and Datanodes are where the data are actually stored. When a
54 file is stored into HDFS, it will be split into file blocks as the storage unit
55 of HDFS. For achieving fault-tolerance, HDFS stores three replicas, by de-
56 fault, for each file blocks in different datanodes. Therefore, even individual
57 node halts down, all data, which are stored in the halted node, can be ac-
58 cessed from other replicas. All of I/O operations can be manipulated through
59 Clientnode.

60 Each Datanode can mount several hard disks and it manages these hard
61 disks by itself. By default, Datanode will store each block into the hard
62 disks in a round-robin way. For example, when a file “a” will be stored into
63 HDFS via Clientnode, Namenode will add “a” into the file system tree. Then
64 Clientnode begins writing the content into HDFS. Once Clientnode detects
65 that current writing size exceeds the block size, 64MB by default, it will
66 ask Namenode for a new block with unique block number. Simultaneously,
67 Namenode needs to re-code mapping relation between “a” and block numbers.
68 Since three replicas are used in HDFS for fault tolerance, Namenode will
69 select three Datanodes to store a block in the file “a”. In Datanode, HDFS
70 should choose a hard disk to each replica of a block. As shown in Figure 1,
71 the last block “k” is stored in disk “2”, then the new one “j” will locate at
72 disk “0” according to the round-robin rule.

73 Blocks of files are distributed in the two-tiered storage architecture, Datanode and Datanode's disk. It will make full use of the collective bandwidth
74 of HDFS if each replica of blocks is appropriately chosen and scheduled to
75 reduce remote reading among Datanodes.

76 MapReduce is usually introduced as a computation model cooperating
77 with HDFS. It can process the data with good locality. However, due to
78 the traits of scientific data and the operation pattern, MapReduce cannot be
79 directly utilized on scientific-data analysis with perfect performance.
80

81 2.2. *Scientific-data Analysis*

82 Scientific-data are usually stored in high-structured files, a kind of special
83 binary format, which can not be read directly by MapReduce. Each file
84 contains multiple logical units, and each unit corresponds to its own physical
85 meaning. The size of each logical unit varies a lot in one file and only the same
86 logical units can be manipulated in most of analysis operations.

87 In this paper, climate model data is used, the total size of which is extremely
88 huge. The high resolution ocean model, a sub-component model in climate
89 model, has 48TB for 100 simulated years, with monthly output. These data are
90 organized by thousands of files with timestamps named in NetCDF format, called
91 history results. Each file contains the same dozens of physical variables. Analysis
92 operations need to process the whole or part of the history results. These operations
93 include computing average, combining the same variables in different files, producing
94 the remapping file and fast Fourier transform (FFT), etc. They require to handle
95 whole or part of variables in multi-files.
96

97 2.3. *Problems and Challenges*

98 Problem/Challenge 1: How to define the operation unit considering the
99 storage unit? Logical unit size varies a lot from a few bytes to gigabytes. Then
100 it can not be used as the operation unit. Otherwise, the workflow of each
101 task will be imbalanced.

102 Problem/Challenge 2: How to schedule tasks among many files to increase
103 the locality and reduce the data network transfer? If all logical units with
104 the same meanings locate on the same Datanode, all data access will proceed
105 locally. However, the default block placement policy of HDFS cannot ensure that,
106 which leads to a huge amount of data network transfer, as mentioned above.

107 Problem/Challenge 3: How to best utilize the disk I/O?

108 3. Parallel scientific-data analysis

109 In order to solve the problems mentioned in previous section to take full
110 use of distributed I/O performance, we propose a new framework - Parallel
111 Scientific-data Analysis (ParSA).

112 3.1. Logical Unit Split and Group among File Blocks

113 No matter how a file is split, all data in one block are stored contiguously
114 in a physical location, since file block is the storage unit of HDFS. Thus the
115 block of HDFS can be used as the basic operation unit to keep continuity in
116 disk and improve I/O performance by proper scheduling to reduce network
117 read. In a block, we can group the small size of logical units or split the big
118 size one.

119 The logical unit split and group approach is shown in Figure 2. All logical
120 units located in one block are distributed into one group. If a logical unit
121 spans two or more blocks, this logical unit will split to several parts. For
122 example, LU1 spans three blocks. It will split into three split parts - LU1
123 (1), LU1 (2) and LU1 (3), unique identifier assigned to each part. Then LU1
124 (1) and LU0 locate in file block 0 and form a operation unit group. LU1 (2)
125 itself is distributed as a group. The remaining LU1 (3), LU2 and LU3 are
126 distributed to the third group. The mapping relationship between groups
127 and file blocks should be records, as shown in Figure 3.

128 Variables in NetCDF, used in climate data, can be viewed as multi-
129 dimension array. An array section is defined as a contiguous rectangular
130 block specified by two vectors - index vector and count vector. The index
131 vector indicates the start offset of the element in the corner closest to the
132 origin. The count vector gives the lengths of the edges along each of the
133 variable's dimension.

134 For instance, the variable given in Figure 4 (a) is a two-dimension slab.
135 If we access the whole variable, index vector is (0, 0) and count vector is (3,
136 5). However, if this variable locates on two file blocks stored with HDFS,
137 and variable are split between value (1, 1) and value (1, 2), it cannot be split
138 into these two pieces. The reason is that each of part is not a rectangular
139 shape, which can not use an index vector and count vector to represent. In
140 this situation, we should split the variable into four pieces, as Figure 4(b)
141 shows.

142 3.2. Scheduler

143 As discussed in section 2.3, quite a few operations are only executed on
 144 the same logical units among multi-files. In principle, blocks containing these
 145 units should be scheduled in one Datanode as many as possible to increase
 146 the access locality.

147 At first, information about file distribution should be gathered. As dis-
 148 cussion in section 2.1, file block has three replicas by default. In this paper,
 149 we define the tuple (N_i, D_j) as the location of a file block, which means
 150 Datanode number and hard disk number, respectively. Assume we get the
 151 following block location information for the first file block of all three files.

152 For the nine replicas for B0 of F1, F2 and F3, Datanode N0 covers 2
 153 replicas. N1 contains 2 replicas and other Datanode only contain 1 for each.
 154 Then Datanode N0 and N1 contain the same number of replicas. Therefore,
 155 Datanode N0 is chosen to process the B0 block of F1 and F2, instead of N1
 156 to process B0 of F0 and F2. The reason is that B0 of F0 and F2 locates
 157 on the same hard disk in N1 Datanode, which is more likely to impact the
 158 performance (explained in section 3.3). Although N1, N6 and N7 all have
 159 the B0 replica of F2, N1 is selected. Because when disk 0 of N0 is broken,
 160 the replica B0 of F0 can be accessed in N1. According to this principle, other
 161 blocks can be scheduled as the Figure 5(b) shown.

162 3.3. Workflow of ParSA

163 Although we utilize the file block location in HDFS to schedule which
 164 block replica are selected among Datanodes, I/O bandwidth of disks may not
 165 be efficiently utilized intra-Datanode. For example, the file block B0s of F0
 166 and F1 are assigned to N0 in Figure 5. Even though B0 of F0 and B0 of F1
 167 locate on different disks - D0 and D1 respectively, the two disks are accessed
 168 sequentially instead of parallelly, if the operations are implemented from logical
 169 level at Datanode N0 as following:

Algorithm 1 Sequential operation on disk in the same Datanode

```

for each logical-unit LU in B0 do
  for each file Fx whose block B0 locates on N0 do
    read LU from B0 of Fx
    operate on LU
  end for
end for

```

170 In this paper, we propose a new workflow to efficiently access the data
 171 on HDFS and maximize utilizing all disk bandwidth in one Datanode.

172 In ParSA, the framework of hybrid of MPI and pthread is used. MPI is
 173 responsible for parallel and communication among Datanodes. There is one
 174 MPI process in each Datanode. Pthead is responsible to implement the detail
 175 operations within each Datanode, constituting of three major components:
 176 reading-thread, processing-thread and receiving-thread.

177 **Reading-thread** If more than one thread dedicatedly access one disk,
 178 it cannot get a higher performance than only one thread. Thus we invoke
 179 one thread for each disk. This is defined as reading-thread and use $Tread_i$
 180 to present the i_{th} reading-thread. For each block, $Tread_i$ reads the queue
 181 of logical units, shown in Figure 3. After reading the logical unit, it is
 182 marked as LOCAL_FLAG, distinguishing with receiving data. The workflow
 183 of reading-thread is given below:

Algorithm 2 reading-thread

inputs:

datanode_number:the Datanode where the reading-thread locates.

disk_number:the disk number to be accessed.

scheduler_result:task scheduler result, given in Figure 5(b).

LU_groups_map:mapping relations between logical unit groups and file
 blocks in Figure 3.

```

for each tuple (node_no, disk_no) in scheduler_result do
  if node_no = datanode_number and disk_no = disk_number then
    get the file-index and block-index according to the index of tuple,
    shown in Figure 5(a).
    for each lu in LU_groups_map[block_index] do
      read values of lu of file_index file
      store file-index and values into lu
      mark lu as LOCAL_TAG
      insert (lu) into data queue
      notify processing-thread
    end for
  end if
end for
send ending-flag to processing-thread

```

184 **Processing-thread** After reading-threads read logical units into queue,
185 the processing-thread will be waked up. The processing- thread takes in
186 charge of calling user-defined function `local_data_process` to process these
187 logical units. Users only focus on implementing operations on logical unit.
188 Since each Datanode only covers part of entire dataset, the major process
189 should reduce the intermediate results, computed with each process in each
190 Datanode. This operation is another user-defined `recv_data_process` function.

191 **Receiving-thread** Receiving-thread is a component in major process. It
192 receives logical inits from worker processes, and puts the logical units into the
193 receiving queue. Then the processing-thread is notified to fetch these data.

194 The workflow of ParSA is shown in Figure 6. Major process locates in
195 Datanode 0, and a worker process is presented as Datanode2. Each read-
196 ing thread reads the logical units which is assigned to this thread. After
197 $Tread_0$ in Datanode 2 reading the logical LU1 (3) of F0, it reads LU2 of
198 F0, inserts the values into queue and notifies processing-thread to wake up.
199 The processing-thread processes all the values in queue sequentially. When
200 processing-thread detects that LU1 (3) are all processed in Datanode 2, it will
201 send the intermediate results to the major process. Once Receiving-thread
202 in major process gets the intermediate results, it inserts the results into the
203 logical unit queue. Since the logical unit queue is the shared resource, all
204 threads must access exclusively.

205 ParSA is the framework for scientific-data access on HDFS. Users can
206 only put their focus on logical operations. ParSA is implemented using MPI
207 and Pthread, providing more flexible configuration with regard to performance
208 than MapReduce. ParSA can achieve load-balance based on current block
209 replicas layout information. However, if the layout itself is unbalance, the
210 performance is also impacted. In the next section, how to further improve
211 the performance by optimizing layout will be discussed.

212 4. Block layout optimization

213 In section 3.2, the scheduler of ParSA utilizes the location information of
214 file blocks. In principle, file blocks containing the same logical units in the
215 one Datanode should be scheduled as many as possible for the operations
216 involving multi-files, while the default block layout strategy of HDFS cannot
217 guarantee this.

218 In HDFS, if Clientnode is one of the Datanodes, the first replica of each
219 file block will be stored in Clientnode, and other replicas are chosen random

220 Datanodes expect for Clientnode. If the Clientnode is not Datanode, the
221 first replica will be randomly stored in a Datanode.

222 Although the scheduler can assign tasks based on current file block repli-
223 cas distribution, the distribution should be control according to the principle
224 for task assignment in order to achieve load balancing. Since blocks of files
225 are stored in Datanode and Datanode’s disk, it needs to optimize the layout
226 in these the two tiered architecture.

227 4.1. Data layout inter-Datanodes

228 Taken into account the characteristic of scientific-data analysis, tightly
229 coupled data as above described should be assigned in one Datanode to
230 increase locality of processing data. Simultaneously, each Datanode should
231 have equal workload, with a pretty load-balance performance.

232 However, it may indeed exists a few analysis operations different from the
233 characteristic. ParSA expose an interface for users to determine where the
234 first replicas of blocks shoule locate. What’s more, the second replicas are
235 stored with uniform distribution automaticlly. The third abide by the HDFS
236 default strategy.

237 4.2. Data layout intra-Datanodes

238 As discussed in section 3.3, the workflow of ParSA invokes one thread
239 for one disk. If blocks are imbalance among disks in one Datanode, it will
240 waste the bandwidth of other disks. The load balance is not only the equal
241 number of file blocks in one Datanode, but also the distribution of blocks
242 containing the same logical units. The former can be guaranteed by HDFS
243 default disk allocation strategy with a round-robin way, while the later can
244 not be satisfied all the time. ParSA provides a round-robin mechanism for
245 these blocks distribution, as shown in Figure 7. The disk will be used for a
246 block, if it does not store one with the same logical units.

247 5. Evaluation

248 5.1. Experimental setting

249 To demonstrate the performance portability of ParSA framework, we
250 choose two HDFS clusters with different configuration, as shown in Table
251 1 and Table 2.

252 In order to demonstrate the efficiency of ParSA, we use climate data
253 generated by a ocean model, POP, a component of climate model. The

254 dataset has 100 NetCDF files, and each of file contains 60 two-dimension
255 variables (logical units) and 18 three-dimension variables, with 320x384 and
256 320x384x40 resolution, which has 309 MB datasize. The total size of this
257 dataset is 30.9 GB.

258 The major analysis operation is calculating time average using ncea com-
259 mand of NCO package, which is the basic preprocessing operation for almost
260 all of climate data analysis.

261 To evaluate the scalability of ParSA, we compare the speedup and through-
262 put performance of different number of Datanodes and different number of
263 disks per node. Cluster 1 is used for the former experiment since the number
264 nodes of cluster 1 is more than cluster 2. While cluser2 is used for the later
265 experiment since each node has more disks in cluster 2.

266 5.2. Scalability with different disks per Datanode

267 ParSA can get about 1.3 GB/s aggregation throuphout if per Datan-
268 ode uses 5 disks, as shown Figure 8(a). The throughput has a remarkable
269 scalability if per Datanode from 1 to 4 disks. However, it has little perfor-
270 mance improvement with 5 disks. The reason is that there are some serial
271 operations, such as writing results and communication, which impact on the
272 performance. Since a series of strategies propused by ParSA is aimed at I/O,
273 Figure 8(b) shows the speedup with only reading part. It's scalability is close
274 to the ideal speedup.

275 5.3. Scalability with different Datanodes

276 Two disks per node are used in this evaluation with cluster 1. It almost
277 keeps the linear speedup until 9 Datanodes as shown Figure 9. However, the
278 performance improvment is limited when the datanodes are greater than 9
279 due to the increasing overhead of communication and the critical section con-
280 trolled by mutex, as shown Figure 6. Although data processing is overlapped
281 with data accessing, the time non-overlapped occupies almost 38.89% on 16
282 datanodes, and this proportions are only 0.02% and 7.96% when datanodes
283 are 9 and 12 respectively.

284 In Figure 9(a), it can get about 800MB/s maximum throughput on 12
285 datanodes, with 24 disks totally. While the same number of disks in Figure
286 8(a), 4 disks per datanode, it can get about 1.2GB/s. The same performance
287 only need 2 disks per datanode, with 12 disks totally. The reason is that, on
288 one hand, the disks and network are difference between the two cluters; on

289 the other hand, cluster 1 requires more network access due to the few disks
290 per node.

291 The speedup of only reading is presented as Figure 9(b), which shows an
292 excellent scalability with different datanodes. The reason why the speedup is
293 higher than the ideal may result from performance variability.

294 5.4. Throughput comparison with different task optimization strategies

295 HDFS block is the basic operation unit in ParSA, splitting and grouping
296 the logical unit. We select the block replicas in different datanodes to reduce
297 network reading, optimize the workflow to support overlapping and paral-
298 lelism of reading and processing, as well as the optimization of blocks layout
299 for both inter-datanodes and intra-datanodes.

300 To present the improving performance, two additional strategies are com-
301 pared to ParSA. The first one takes logical unit as the operation unit. The
302 disk utilization is shown as Figure 10(a). Since the size of logical unit varies
303 quite a lot, a serious load imbalance exists among tasks. The second strat-
304 egy takes all logical units located in one file block as an operation unit, while
305 those logical unit which locate on more than one block do not be split and
306 treat the big size logical units as other operation unit. In this situation, it
307 can not guarantee that an entire operation unit stores in a disk. When a
308 thread try to read this unit, it requires read a part of value from the other
309 disk. In the mean time, the other thread is reading the same disk, leading
310 to disk resource competition. This utilization performance is presented as
311 Figure 10(b).

312 Compared with throughput to evaluate the speedup, ParSA outperforms
313 the above two strategies by a factor of 3.34X and 1.14X respectively.

314 6. Conclusion and future work

315 A high-throughput and scalable scientific framework combined with HDFS
316 distributed file system and NCO analysis tool has been designed, imple-
317 mented and evaluated. Maximize distributed I/O performance has been
318 gained by scheduling block replicas in different datanodes, optimizing the
319 blocks layout for both inter-datanodes and intra-datanodes, as well as over-
320 lapping the data reading, processing and transferring. Instead of Hadoop
321 framework, we use multi-process with MPI to parallel the NCO basic oper-
322 ations, without leading to compatible problem. In future work, we plan on

323 testing the framework on some complex climate data analysis, such as FFT,
324 principal component analysis and Wavelet transform.

325 References

- 326 [1] T. Kosar, M. Balman, A new paradigm: Data-aware scheduling in grid
327 computing, *Future Generation Computer Systems* 25 (4) (2009) 406–
328 413.
- 329 [2] A. J. Hey, A. E. Trefethen, The data deluge: An e-science perspective.
- 330 [3] R. Rew, E. Hartnett, J. Caron, et al., Netcdf-4: software implementing
331 an enhanced data model for the geosciences, in: 22nd International Con-
332 ference on Interactive Information Processing Systems for Meteorology,
333 Oceanograph, and Hydrology, 2006.
- 334 [4] M. Folk, A. Cheng, K. Yates, Hdf5: A file format and i/o library for
335 high performance computing applications, in: *Proceedings of Supercom-
336 puting*, Vol. 99, 1999.
- 337 [5] J. Lofstead, F. Zheng, S. Klasky, K. Schwan, Input/output apis and data
338 organization for high performance scientific computing, in: *Petascale
339 Data Storage Workshop*, 2008. PDSW'08. 3rd, IEEE, 2008, pp. 1–6.
- 340 [6] R. Rew, G. Davis, Netcdf: an interface for scientific data access, *Com-
341 puter Graphics and Applications*, IEEE 10 (4) (1990) 76–82.
- 342 [7] C. S. Zender, Analysis of self-describing gridded geoscience data with
343 netcdf operators (nco), *Environmental Modelling & Software* 23 (10)
344 (2008) 1338–1342.
- 345 [8] NCAR, Ncar command language, <http://www.ncl.ucar.edu/>.
- 346 [9] J. Dean, S. Ghemawat, Mapreduce: simplified data processing on large
347 clusters, *Communications of the ACM* 51 (1) (2008) 107–113.
- 348 [10] H. Zhao, S. Ai, Z. Lv, B. Li, Parallel accessing massive netcdf data based
349 on mapreduce, in: *Web Information Systems and Mining*, Springer,
350 2010, pp. 425–431.

- 351 [11] J. B. Buck, N. Watkins, J. LeFevre, K. Ioannidou, C. Maltzahn, N. Poly-
352 zotis, S. Brandt, Scihadoop: Array-based query processing in hadoop,
353 in: Proceedings of 2011 International Conference for High Performance
354 Computing, Networking, Storage and Analysis, ACM, 2011, p. 66.
- 355 [12] D. L. Wang, C. S. Zender, S. F. Jenks, Clustered workflow execution of
356 retargeted data analysis scripts, in: Cluster Computing and the Grid,
357 2008. CCGRID'08. 8th IEEE International Symposium on, IEEE, 2008,
358 pp. 449–458.
- 359 [13] D. Borthakur, Hdfs architecture guide, HADOOP APACHE PROJECT
360 [http://hadoop. apache. org/common/docs/current/hdfs design. pdf](http://hadoop.apache.org/common/docs/current/hdfs_design.pdf).
- 361 [14] S. Ghemawat, H. Gobioff, S.-T. Leung, The google file system, in: ACM
362 SIGOPS Operating Systems Review, Vol. 37, ACM, 2003, pp. 29–43.

Algorithm 3 processing-thread

inputs:

num_tasks:total number of **ending-flag** to wait

intermediate_results:store the intermediate results

cur_num_tasks = 0

while **cur_num_tasks** < **num_tasks** **do**

 wait until receiving notification

if receive **ending-flag** **then**

cur_num_tasks++

 fetch the logical unit queue

for each **lu** in logical unit queue **do**

 get **interm_res** according to **lu identifier** from **intermediate_results**

if **lu's flag** == **LOCAL_FLAG** **then**

 call **local_data_process** (**lu**, **interm_res**)

end if

if **lu's flag** == **RECV_FLAG** **then**

 call **recv_data_process** (**lu**, **interm_res**)

end if

if is not major process **then**

if all logical unit are processed in this Datanode **then**

 store **interm_res** into **lu**

 mark **lu** as **RECV_TAG**

 send (**lu**) to major process

end if

end if

end for

end if

end while
