

Stencil Skeleton Evaluation

Francesco Piccinno <stack.box@gmail.com>

July 12, 2011

Contents

1	Introduction	2
2	Abstract implementation	2
3	Cost model: completion time	4
4	Performance measurements and tests	5
4.1	Statistics for min function	6
4.2	Statistics for linsolve function	6
5	Concrete implementation	8
6	Build instructions	25

List of Figures

1	Data parallel partition scheme	2
2	Different partition schemes	3
3	Two ways of partitioning the same matrix	4
4	min function: Completion time	6
5	min function: Efficiency, Speedup and Scalability	7
6	linsolve function: Completion time	8
7	linsolve function: Efficiency, Speedup and Scalability	9
8	Bogus data derivation	11

List of Tables

1	Calculation time timings	5
2	Average of standard deviation between workers calculation time	8

List of listings

1	Spawner: <code>run.py</code>	9
2	Worker logic: <code>slave.py</code>	10
3	Master logic: <code>main.py</code>	10
4	Stencil skeleton implementation: <code>stencil.py</code>	11
5	Stencil skeleton implementation: <code>puzzle.py</code>	17
6	Communication: <code>communicator/mpi.py</code>	19
7	Communication constants: <code>communicator/enum.py</code>	20
8	Matrix class: <code>matrix.py</code>	21
9	Functional code: <code>functional.py</code>	24

Abstract

The aim of the report is to briefly introduce the concept of data-parallel paradigm and to study and evaluate the implementation of a stencil skeleton, a particular parallel computation in which data dependence pattern is implemented by the mean of inter-worker communications.

1 Introduction

A data parallel computation, which can be defined both on streams and or single data value, consists in the partition of possibly large data structures to different processing entities which are able to compute a function on the assigned partition in parallel. This parallel pattern bases its functionalities on function replication. As a consequence, the knowledge of the sequential computation form is necessary both for function replication and data partitioning. This is the main reason for the parallel program design complexity. Latency and memory capacity are optimized compared to other stream parallel paradigms like farms, though a potential load unbalance exists.

To better understand, let us consider a data parallel computation operating on a bi-dimensional matrix as depicted in Figure 1. The matrix is partitioned in blocks of $g = \frac{M^2}{n}$ elements. Successively each block is transmitted to each worker through a *scatter* operation. Then all workers apply the same function F to the corresponding partition in parallel. At the end, the result may be collected through a *gather* operation.

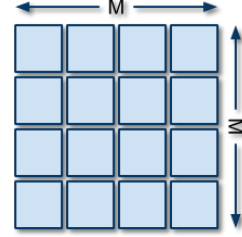


Figure 1: Data parallel partition scheme

Map

When all the workers are *fully independent* and operate on their own partition, without any cooperation

Stencil-based

When a workers operate in parallel and *cooperate* through data exchanges. In this case data dependencies are implemented by inter-worker communications.

The form of a stencil may be *statically* predictable or *dynamically* exploited according to the current data values. In case of static stencil we have two other possibilities: a *fixed* stencil where the communication scheme is fixed throughout the computation, and a *variable* stencil where the communication scheme varies from a computation step to the next one, although the overall behavior is statically predictable.

In this report we will concentrate our attention to *static fixed stencils* applied on bi-dimensional matrix.

2 Abstract implementation

In this section we briefly discuss the general idea of our implementation leaving the implementation details out of this section. The interested reader can find a more detailed analysis of the real implementation in “Concrete implementation” section.

The main idea we followed is very simple. Since we are dealing with a data-parallel pattern, some kind of partition of the input problem is required to properly solve the problem. The schema can be synthesized in four steps:

1. The input matrix is loaded by the collector and equally split among the available workers
2. Each worker communicate its contour segments needed to the other neighbors
3. After having received the necessary data from the neighbors, the computation on the given partition is started in parallel on each worker
4. The resulting partition is sent back toward a collector, which is in charge of collating partial results, thus providing the final result.

As you might imagine, the central point in this scheme which can introduce substantial optimization is the partition phase. This can be implemented in different ways on the basis on the input offsets.

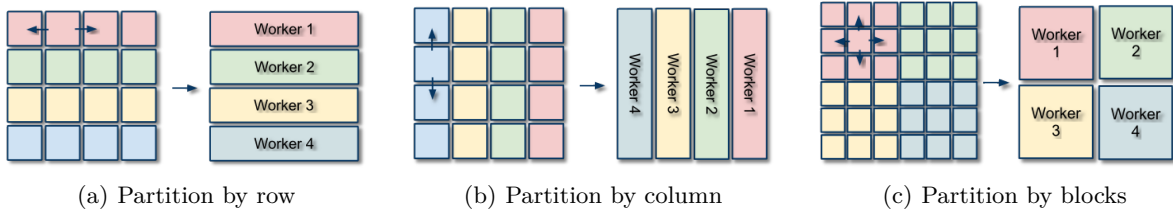


Figure 2: Different partition schemes

In our implementation, the partition phase tries to derive an optimal partition by trying to minimize the communication between workers. We have isolated three different situations depicted, which are depicted in Figure 2:

Row partition Whenever the input offsets set contains dependencies that are row dependent only, a row partition scheme is derived. This means that all the elements in the input offset set are of the type $(0, x)$. In this case no communication at all is required between workers, therefore a full *map* parallelism is possible if each worker has an entire row assigned.

Column partition This is the mirror of the *row partition* case and it is verified whenever a given element depends only on elements which are on the same column. In this situation the input offset set contains only elements of the type $(x, 0)$. No communications is required, on the condition that at each worker is assigned an entire column of the input matrix.

Block partition Whenever a *row* neither a *column partition* is applicable, a block partition is derived. In this case the input offset set contains “heterogeneous” points, so no kind of trivial derivation is possible. In our case a minimum rectangle called “Elementary Rectangle” is derived. This rectangle is evaluated by taking as reference point a single element and deriving the smallest rectangle that encompass all data dependencies needed to evaluate the function in that specific point.

If we are in presence of a *Block partition* situation another optimization is still possible. Indeed, we can try to still partition the matrix by “row-blocks” or “column-blocks” by vertically or horizontally stacking single *elementary rectangles* thus achieving a more optimized partition. In fact this situation is able to avoid communications between “elementary” workers that are still present in a situation like the one depicted in Figure 2(c).

3 Cost model: completion time

Since the cost model specifically depends on the offsets passed as input to the stencil skeleton, we have decided to derive a general approximated cost model which tries to collect the most configurations still maintaining a reasonable precision in the evaluation. Before reaching the final formulation, we explain the steps that brought us to the successful derivation.

In our case the starting matrix is read from the disk and then partitioned in some way among n workers through the use of a remote communication which has a cost of $T_{scatter}(n, g)$, where g is the block size assigned to each worker ($g = \frac{M^2}{n}$, assuming “balanced” partitions).

Then each worker after having received its own partition, is in charge to communicate to all the neighbors the necessary contour partition. Of course, this communication has an associated cost, which has to be evaluated but varies according to the specific input offsets passed to the skeleton. To model this part of the formula we have to introduce a parameter which is rc (real communications), which represents the number of communications we have to make with the neighbors. To better understand the situation, let's assume to have as input offsets: $(0, -1)$, $(-1, 0)$, $(0, 1)$ $(1, 0)$

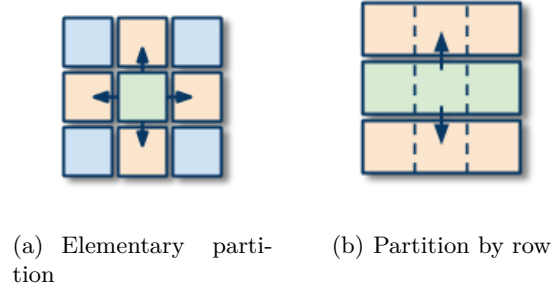


Figure 3: Two ways of partitioning the same matrix

The situation is the one depicted in Figure 3(a). In this case we are assuming the “elementary” partition 1 by 1 where each worker has responsibility for just 1 matrix element. In this case before starting the computation the worker has to communicate respectively to the upper, next, lower and previous neighbor its data partition, assuming the “owners-compute-rule” to be satisfied. At the end the parameter rc will be 4. This has to be multiplied by the $T_{comm}(P)$ where P represent the size of the data to be transmitted to the neighbors. In this specific case this is equal to 1, but in other cases this may vary based on the size of the partition and the offsets, as you will see. In fact, by keeping the same input offsets and by only changing the size of the sub-matrix partition assigned to each worker the things changes, as presented in Figure 3(b). In this specific case, not only the P parameter changes and becomes 3 but also the rc changes to 2.

Then a calculation phase is started in parallel by each worker on its specific partition which takes a total time of $g \cdot T_f$, where T_f is the time needed to calculate the function among all inputs (in our case 5 considering the central value) while g is the size of the partition assigned to each worker or better the number of items a given worker has assigned inside the partition.

At the end of the calculation, a second phase is started which consists in transferring back the computed partition to the root node ($T_{gather}(n, g)$). This can be evaluated by simply taking into account the double communication at the first phase. Thus the resulting formula is:

$$2T_{scatter}(n, g) + rc \cdot T_{comm}(P) + g \cdot T_f$$

Function	M	T_{calc} (sec)	T_f (sec)
min	3000	61.923	$6.880 \cdot 10^{-06}$
	4000	95.880	$5.992 \cdot 10^{-06}$
	5000	146.599	$5.864 \cdot 10^{-06}$
	6000	206.726	$5.742 \cdot 10^{-06}$
linsolve	3000	711.255	$7.903 \cdot 10^{-05}$
	4000	1274.833	$7.968 \cdot 10^{-05}$

Table 1: Calculation time timings

As a measure of comparison, the sequential algorithm cost model is given by the last part of the previous formulation, considering that the partition in this case is the entire input matrix: $M^2 \cdot T_f$.

4 Performance measurements and tests

In this section we present performance measurements we have collected during our experiments. For this scope we have decided to test two different functions with different complexity. The first one is **min** which applied to various elements returns the minimum among them. The latter is **linsolve** a simple function which tries to solve a system of linear algebra equations. The only assumption we made during the experiment was to have as input offsets the following set: $(0, -1)$, $(-1, 0)$, $(0, 1)$ $(1, 0)$.

To better understand the final goal of the **linsolve** function take in consideration the scheme below, which assumes “ e ” as the target element to be computed:

$$\begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} \rightarrow \begin{cases} bx + hy = e \\ dx + fy = e \end{cases} \rightarrow result = \frac{x + y}{2}$$

Regarding the ideal statistics we have used some instrumentation inside the code to evaluate the calculation time $T_{calc} = M^2 \cdot T_f$ and used it to derive the correct T_f value. Results collected and derived are presented in Table 1.

Since the implementation we have provided makes use of MPI message passing library which does not guarantee any constant behaviour for the setup “socket” procedure, we preferred to avoid to take in consideration the T_{setup} parameter which should be paid once for every communication as $T_{comm}(L) = T_{setup} + L \cdot T_{transm}$ formula suggests. By the way this is not a wrong assumption, because in MPI after the setup procedure all peers involved in the MPI “world” are free to communicate each other without paying any more the setup procedure. So the setup is made once and for all, thus ignoring this parameter, by placing instrumentation code in proper position, does not impact in any way our study case.

The evaluation of T_{transm} parameter was simply evaluated by taking a statistic mean on top of the scattering time of the matrix registered by our instrumentation code. The resulting value results to be equal to $3.437 \cdot 10^{-07} \frac{item}{sec}$.

After having set up all the parameters we collected all the needed timings to derive the proper statistics. Since we were dealing with a data-parallel skeleton, we have focused our attention on the following parameters:

Completion time	Time needed to complete the computation on the given input	T_c
Efficiency	Provides a measure of how close is the effective performance with respect to the ideal one.	$\epsilon(n) = \frac{T_c^{(n)}}{T_c^{(1)}}$
Speedup	Provides a measure of the relative speed of the parallel computation with respect to the sequential one.	$s(n) = \frac{T_c^{(1)}}{T_c^{(n)}} = \epsilon(n) \cdot n$
Scalability	The ratio between parallel execution time with parallelism degree equal to 1 and the parallel execution with parallelism degree equal to n	$scalab(n) = \frac{T_c^{(1)}}{T_c^{(n)}}$

4.1 Statistics for min function

In this section we present the statistics we have collected regarding the `min` function. Since no interesting differences were reported for the mid-sized cases, we decided to just attach the statistics regarding the limiting cases. The former when the array size is equal to 3000 by 3000 elements. The latter where the array is composed of 6000 by 6000 elements.

A compare between completion time of the two limit cases is presented in Figure 4. As you can imagine, in this situation a solution of this kind does not offer any kind of advantages, since the time taken by the function to compute a given element is very low compared to the effort taken by the skeleton to distribute data and work among workers through the network. Therefore the transmission overhead highly impacts on the possibility to scale up.

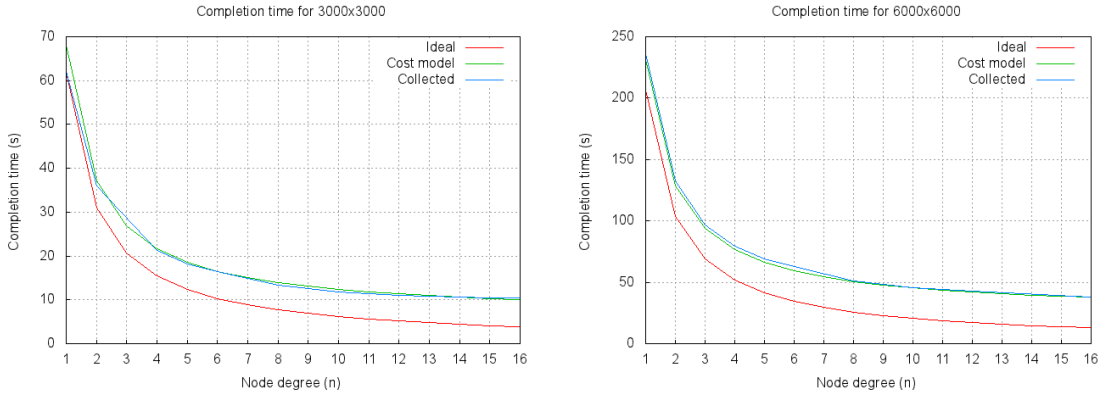
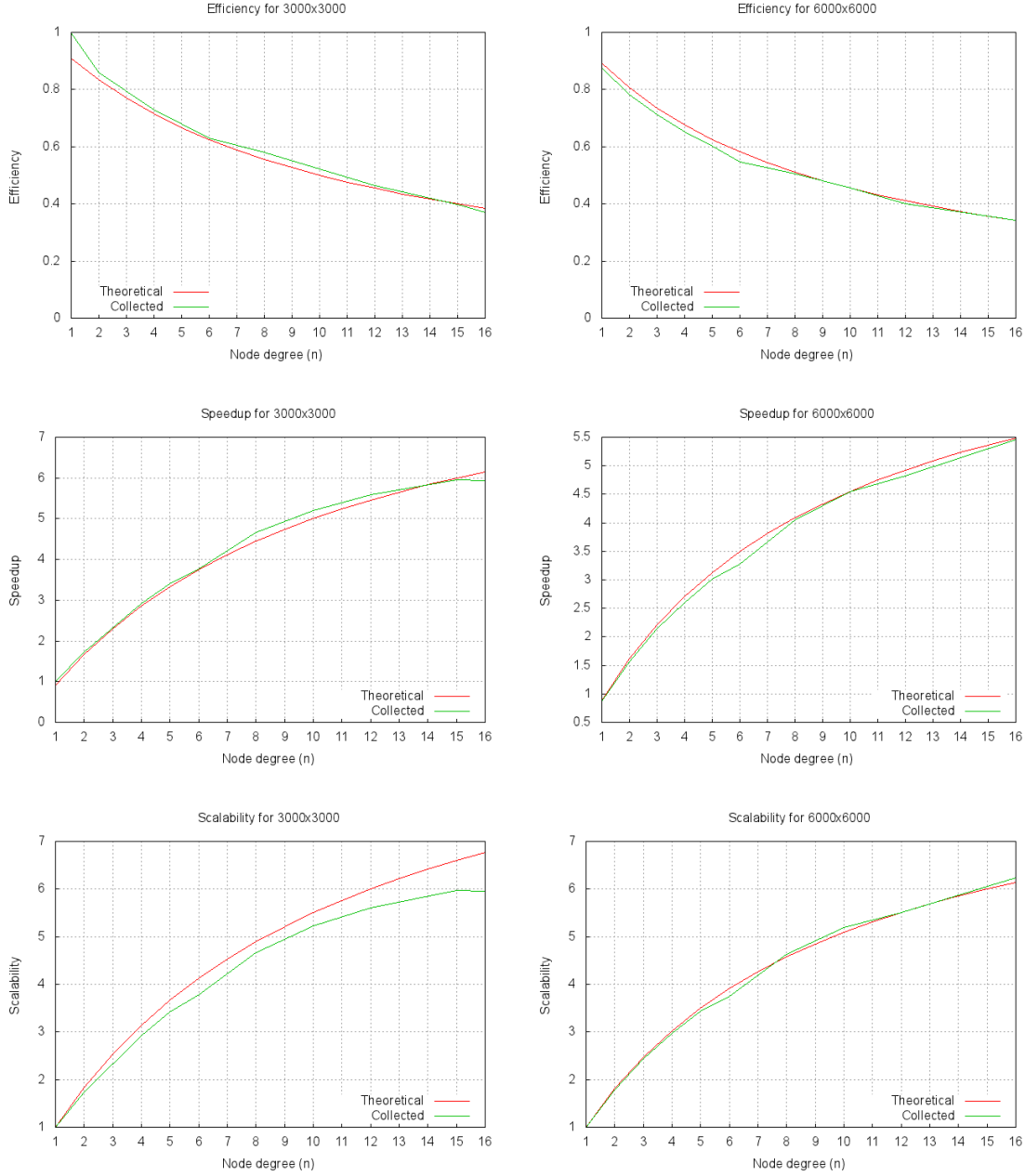


Figure 4: `min` function: Completion time

The thing is more evident if you take as reference Figure 5 which presents statistics regarding speedup, efficiency and scalability. As you can see, for the first case, with the maximum number of processing elements we can achieve a speedup equals to 6, which is quite low for the employed resources. By taking in consideration the second case, an efficiency of less than 35% produces a speedup of ≈ 5.5 for 16 processing elements.

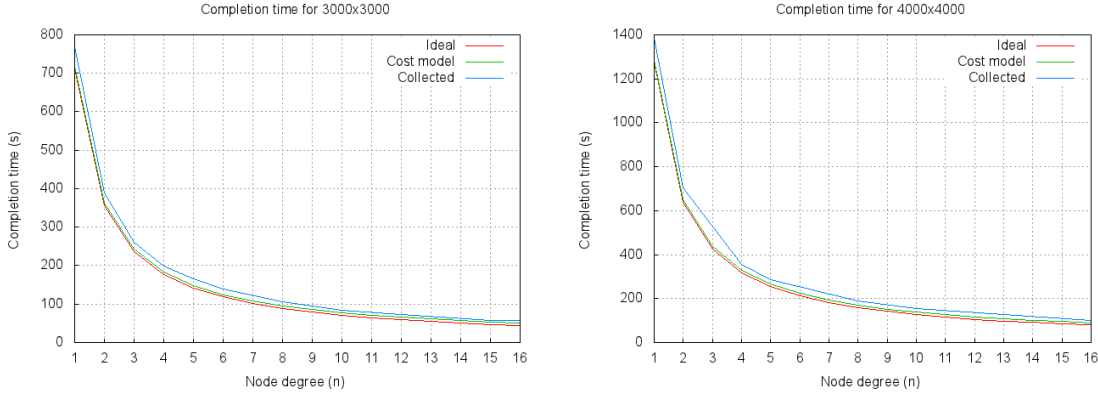
4.2 Statistics for linsolve function

Regarding the parallelization of `linsolve` function, we got interesting results which demonstrate the goodness of our implementation. Figure 6 shows the completion time graph for the two runs. As you can see, the cost model elaborated really resemble in a very precise way the behaviour

Figure 5: *min* function: Efficiency, Speedup and Scalability

of the program. Indeed all the three lines almost describe the same decay function. This is because our computation time is very large compared to the other overhead we incur on when we split the matrix through the network of multi-cores.

On the other hand, Figure 7 shows a comparison between *speedup* and *efficiency* statistics. As you may notice, there is the presence of differences between the real and the predicted situation. Indeed this is not strange at all. The behavior is mostly caused by the semantic of *linsolve* function itself. In fact, the function tends to generate load unbalancing among workers, thus resulting in a slight increase of completion time and therefore in fluctuations between theoretical efficiency and the real one. This is because some partitions might have a major number of solvable systems while the others might have less or at least simpler one.

Figure 6: `linsolve` function: Completion time

Inputs	min	<code>linsolve</code>
3000	0.2143	1.296
4000	0.1156	1.101

Table 2: Average of standard deviation between workers calculation time

Anyhow, the final results demonstrate how much the parameter T_f impacts on the speedup.

This factor was further investigated and studied, by taking the standard deviation on the registered time interval taken by each worker to compute the given partition thus revealing some kind of asymmetry between the `min` and `linsolve` functions. As reference take in consideration Table 2. From a rapid look you can notice a magnitude of order of difference between the two functions. This means that the T_f parameter we take as reference is not fixed but may vary depending on the input values and should have taken as an average case measure. By the way, although slightly imprecise, the curve still describe with a given approximation the real situation.

5 Concrete implementation

Regarding the concrete implementation part, we decided to use Python¹ programming language jointly with `mpi4py`² library which provides a wrapper to Message Passing Interface API specification.

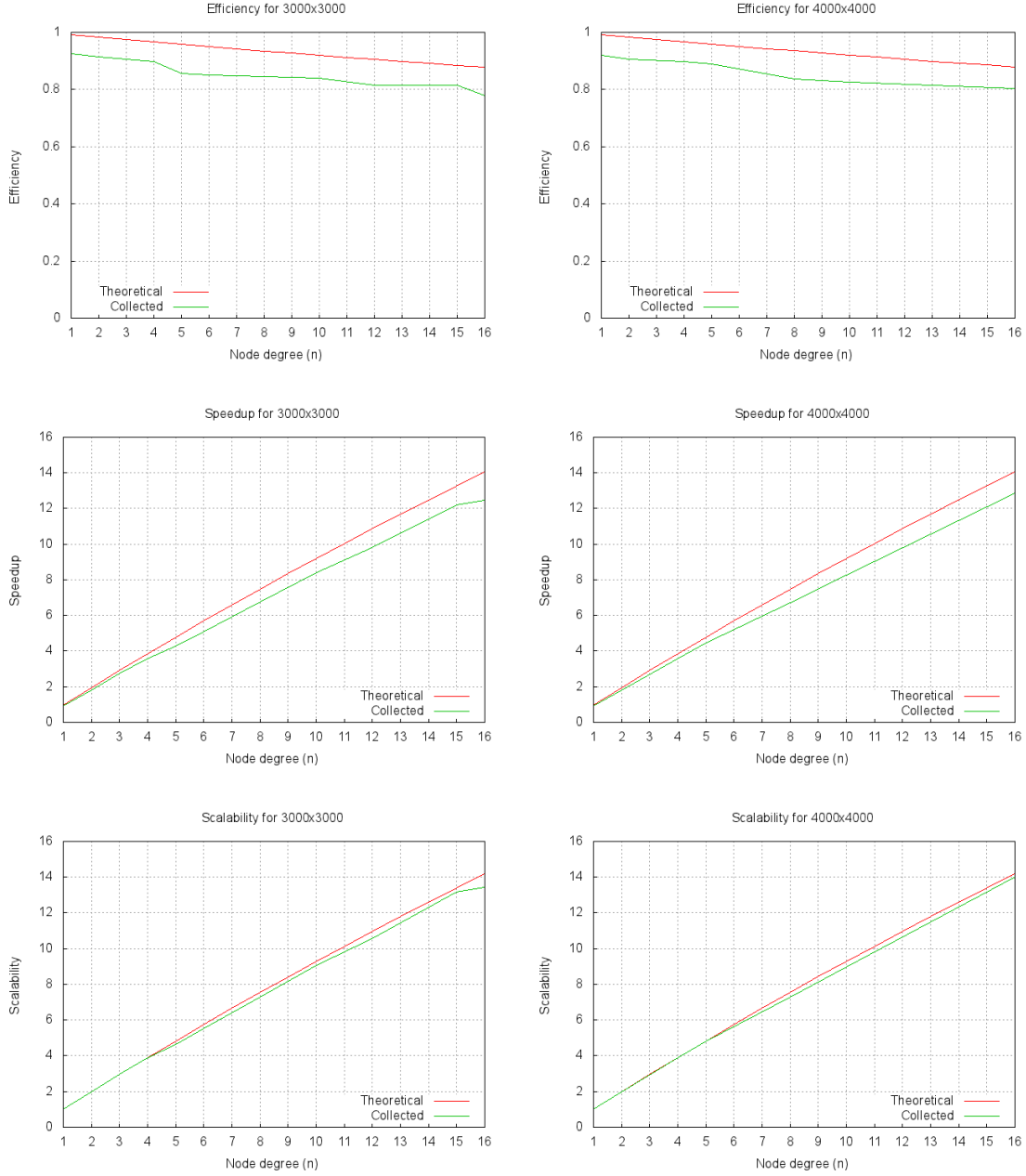
Since Python (actually CPython the official Python implementation) does not provide an efficient implementation of threads due to GIL³ (Global Interpreter Lock) we decided to just use process as method for exploiting cluster of multi-cores resources. Regarding communications everything is delegated to the library which is in charge of introducing optimization whenever two workers are allocated on the same machine and they need to exchange data.

We will start our analysis by taking a look to the `run.py` (Listing 1) which is responsible to set up everything and start the real computation. It accepts as parameter two parameters. The first one may be set to “seq” to spawn the sequential version of the program or to “par” for the parallel version. The second parameter instead just pass the matrix file where the function should be applied to. The program accepts an optional third parameter which gives indication

¹<http://python.org>

²<http://mpi4py.scipy.org>

³<http://wiki.python.org/moin/GlobalInterpreterLock>

Figure 7: `linsolve` function: Efficiency, Speedup and Scalability

about the partition scheme. If it is not present the software will try to do some calculation to derive an optimal partition scheme by following the procedure we have described in the “Abstract implementation” section.

As you can see in case of parallel version is required to be run, the `run` function is called which is in charge of setting up the master (`main` function) and the workers (`slave` function).

Listing 1: Spawner: `run.py`

```
1 import sys
2 from mpi4py import MPI
3 from main import main, sequential
4 from slave import slave
```

```

5
6 def run(matrix_file, rows=None, cols=None):
7     rank = MPI.COMM_WORLD.Get_rank()
8
9     if rank == 0:
10         main(matrix_file, rows, cols)
11     else:
12         slave()
13
14 if __name__ == "__main__":
15     if len(sys.argv) < 3:
16         print "Usage: %s <par|seq> <matrix-file> [rowsxcols]" % sys.argv[0]
17         sys.exit(0)
18
19     if sys.argv[1] == 'seq':
20         sequential(sys.argv[2])
21     elif sys.argv[1] == 'par':
22         if len(sys.argv) == 4:
23             params = map(int, sys.argv[3].split('x'))
24         else:
25             params = (None, None)
26
27     run(sys.argv[2], *params)

```

The worker logic is presented in Listing 2. The worker just waits to receive its partition and then create a `StencilWorker` instance and starts the computation.

Listing 2: Worker logic: `slave.py`

```

1 #!/usr/bin/env python
2 from mpi4py import MPI
3 from stencil import StencilWorker
4
5 def slave():
6     worker = StencilWorker(*MPI.COMM_WORLD.scatter(root=0))
7     worker.start()

```

Listing 3 shows the two function which are called by the `run.py` module. The first function is the one responsible for the parallel execution of the program, while the second is meant for sequential execution. Instrumentation code is also present to take into account loading time of the matrix, and completion time.

Listing 3: Master logic: `main.py`

```

1 import cPickle as pickle
2 import time
3
4 from mpi4py import MPI
5 from stencil import Stencil
6 from functional import offsets
7
8 def main(matrix_file, rows=None, cols=None):
9     start = time.time()
10    matrix = pickle.load(open(matrix_file, "r"))
11    print "%.10f seconds to load the matrix" % (time.time() - start)
12
13    start = time.time()
14    st = Stencil(offsets)
15    st.apply(matrix, rows, cols)
16    print "%.10f seconds to apply on %d processors" % \
17        (time.time() - start, MPI.COMM_WORLD.Get_size() - 1)
18
19 def sequential(matrix_file):
20    print "Trying sequential version"
21    start = time.time()
22    matrix = pickle.load(open(matrix_file, "r"))
23    print "%.10f seconds to load the matrix" % (time.time() - start)

```

```

24
25     st = Stencil(offsets)
26     start = time.time()
27     st.seq_apply(matrix)
28     print "%.10f seconds to apply sequentially" % \
29         (time.time() - start)
30     #matrix.dump()

```

Now we will go deeper in the analysis of the stencil skeleton by introducing Listing 4. The main class in this case is the `Stencil` whose constructor takes as parameter the input offset set. Immediately after the `analyze_offsets` method is called which is in charge of extracting communication data dependencies and to convert them into elementary vectors tuples. This is needed whenever a dependency of the type (x,y) with $x \neq 0$ and $y \neq 0$ is present. In this specific case all the input offsets are analyzed to derive a tuple like (max_x, max_y) , for the composed directions (eg. up-left, down-right, ...).

The `apply` function is the most important function in this class. It is in charge of deriving a proper partition scheme in case (if required, by calling the homonym method) and to scatter the matrix among the workers. Method `fix_data_dep` is called before the actual scattering phase takes place. This is necessary to avoid bogus function evaluation in various corner cases, such the one depicted in Figure 8.



Figure 8: Bogus data derivation

In this situation, a partition 2×2 is derived while the input offset set is equal to $(-1, -1)$. Therefore although not only the up-left neighbor is involved but partially also the left and the upper neighbors. This corner case is verified whenever the composed direction tuple (x, y) is lesser than the partition dimension ($x < height \vee y < width$). If this condition is met, `fix_data_dep` function is in charge of adding redundant communication to the pattern (in this case the communication with the up and left neighbors).

After this check a phase of set-up is made, to interconnect the workers by the mean of `autoconnect` method, which is called twice first to connect left-, right-, up- and down- neighbors each other and then for the remaining up-right, up-left, down-right and down-left neighbors.

Then the real scattering phase is done and after the gathering procedure followed by a reconstruction phase. A final barrier is then used to synchronize all the remote entities and conclude the execution.

Regarding the class `StencilWorker` which is in charge of modeling the behaviour of the workers, no further comment is needed. The class first sends the contour part of the partition to the other workers, then it waits to receive the parts which are sent by the other neighbors. After reception the calculation phase is executed by making use of the `Puzzle` class. The final step is taking part in the global gathering phase where each worker sends back to the master (node with rank 0) its own computed part.

Listing 4: Stencil skeleton implementation: `stencil.py`

```

1 import time

```

```

2 import numpy
3 import logging
4
5 from mpi4py import MPI
6
7 from matrix import Matrix
8 from puzzle import Puzzle
9 from functional import function
10
11 from communicator.enum import *
12 from communicator.mpi import Communicator
13
14 comm = MPI.COMM_WORLD
15 rank = comm.Get_rank()
16 name = MPI.Get_processor_name()
17
18 logging.basicConfig()
19
20 log = logging.getLogger("stencil")
21 log.setLevel(logging.INFO)
22
23 class StencilWorker(object):
24     _id = 0
25
26     def __init__(self, offsets, data_segments, partition, pwidth, pheight, \
27                 connections=None):
28         self.wid = StencilWorker._id
29         self.offsets = offsets
30         self.data_segments = data_segments
31         self.partition = partition
32
33         self.col, self.row = 0, 0
34         self.width, self.height = self.partition.cols, self.partition.rows
35         self.pwidth, self.pheight = pwidth, pheight
36
37         self.conns = connections
38
39         if rank == 0:
40             log.debug("Worker %d has %s" % (StencilWorker._id, self.partition))
41             StencilWorker._id += 1
42         else:
43             self.comm = Communicator(self)
44
45     def autoconnect(self, wdict, rows=None, cols=None):
46         if self.conns:
47             if self.conns.up:
48                 self.conns.up_left = self.conns.up.conns.left
49                 self.conns.up_right = self.conns.up.conns.right
50             if self.conns.down:
51                 self.conns.down_left = self.conns.down.conns.left
52                 self.conns.down_right = self.conns.down.conns.right
53
54             return
55
56         tot = int(rows * cols)
57
58         i_col = self.wid % cols
59         i_row = int(self.wid / cols)
60
61         start = cols * i_row
62
63         # Here we calculate our position in the general matrix
64         self.col, self.row = i_col * self.height, i_row * self.width
65
66         links = []
67
68         for wid in ((self.wid + 1) % cols + start,
69                   (self.wid - 1) % cols + start,
70                   (self.wid - cols) % tot,
71                   (self.wid + cols) % tot):

```

```
72
73     if wid == self.wid:
74         links.append(None)
75     else:
76         links.append(wdict[wid])
77
78     self.conns = Connections(*links)
79
80     def __repr__(self):
81         return 'Worker(%d [%d %d %d %d])' % (self.wid, self.row, self.col,
82             self.width, self.height)
83
84     def start(self):
85         log.info("Worker started on processor %d %s" % (rank, name))
86
87         # First we send all the required partitions to the neighbors and then
88         # receive all the segments and reconstruct a local matrix where the
89         # function will be evaluated.
90
91         puzzle = Puzzle(self.partition)
92
93         for lbl, enum in zip(LABELS, ORDERED):
94             self.comm.send(getattr(self.conns, lbl), enum)
95
96         for lbl, enum in zip(RLABELS, REVERSED):
97             m = self.comm.receive(getattr(self.conns, lbl), enum)
98             if m: puzzle.add_piece(m, enum)
99
100         start = time.time()
101         puzzle.apply(self.offsets, function)
102         log.info("Worker %d: %.10f seconds to compute the partition" % \
103             (rank - 1, time.time() - start))
104
105         log.info("Sending back the computed sub-partition from %d" % rank)
106
107         start = time.time()
108         comm.gather(self.partition.matrix, root=0)
109         log.info("Worker %d: %.10f seconds to send back the partition" % \
110             (rank - 1, time.time() - start))
111
112         comm.Barrier()
113
114 class Stencil(object):
115     def __init__(self, offsets):
116         self.offsets = offsets
117         self.data_segments = None
118         self.analyze_offsets()
119
120     def analyze_offsets(self):
121         """
122         This function is in charge of extracting the proper sub-partitions that
123         must be transmitted to the other workers in a stencil fashion
124         """
125         extract = lambda l, rev: sorted(filter(l, self.offsets), reverse=rev)
126
127         right = extract(lambda x: x[0] == 0 and x[1] > 0, True)
128         left = extract(lambda x: x[0] == 0 and x[1] < 0, False)
129         up = extract(lambda x: x[0] < 0 and x[1] == 0, False)
130         down = extract(lambda x: x[0] > 0 and x[1] == 0, True)
131
132         # Sort the first key which is the row so order is False since we are
133         # interested in far jumps which are < 0
134         up_left = extract(lambda x: x[0] < 0 and x[1] < 0, False)
135
136         if up_left:
137             max_row = up_left[0][0]
138             max_col = sorted(up_left, key=lambda x: x[1], reverse=False)[0][1]
139             tgt_up_left = (max_row, max_col)
140         else:
141             tgt_up_left = None
```

```

142
143     up_right = extract(lambda x: x[0] < 0 and x[1] > 0, False)
144
145     if up_right:
146         max_row = up_right[0][0]
147         max_col = sorted(up_right, key=lambda x: x[1], reverse=True)[0][1]
148         tgt_up_right = (max_row, max_col)
149     else:
150         tgt_up_right = None
151
152     # Sort the first key which is the row so order is True since we are
153     # interested in far jumps which are > 0
154     down_left = extract(lambda x: x[0] > 0 and x[1] < 0, True)
155
156     if down_left:
157         max_row = down_left[0][0]
158         max_col = sorted(down_left, key=lambda x: x[1], reverse=False)[0][1]
159         tgt_down_left = (max_row, max_col)
160     else:
161         tgt_down_left = None
162
163     down_right = extract(lambda x: x[0] > 0 and x[1] > 0, True)
164
165     if down_right:
166         max_row = down_right[0][0]
167         max_col = sorted(down_right, key=lambda x: x[1], reverse=True)[0][1]
168         tgt_down_right = (max_row, max_col)
169     else:
170         tgt_down_right = None
171
172     tgt_right = right and right[0] or None
173     tgt_left = left and left[0] or None
174     tgt_up = up and up[0] or None
175     tgt_down = down and down[0] or None
176
177     log.debug("List of possible targets:")
178     log.debug("Left: %s Right: %s" % (tgt_left, tgt_right))
179     log.debug("Up : %s Down : %s" % (tgt_up, tgt_down))
180
181     log.debug("Up-left : %s Up-right : %s" % \
182               (tgt_up_left, tgt_up_right))
183     log.debug("Down-left : %s Down-right : %s" % \
184               (tgt_down_left, tgt_down_right))
185
186     # Now we try to assign correct mapping. Please beware that if you
187     # change the enumeration order you also need to change the order of
188     # this assignment.
189
190     self.data_segments = [
191         tgt_left,
192         tgt_right,
193         tgt_down,
194         tgt_up,
195         tgt_down_right,
196         tgt_up_left,
197         tgt_down_left,
198         tgt_up_right,
199     ]
200
201     def fix_data_deps(self, prow, pcol):
202         """
203         Try to fix data dependencies to avoid bogus data
204         @param prow is the number of rows each worker has assigned to
205         @param pcol is the number of cols each worker has assigned to
206         """
207         def adjust(corner1, corner2, direction, check_idx, set_idx):
208             targets = filter(lambda x: x,
209                             [self.data_segments[REVERSED[corner1]],
210                              self.data_segments[REVERSED[corner2]]])
211

```

```
212         if not targets: return
213
214         targets.sort(reverse=True, key=lambda x: abs(x[check_idx]))
215         rect = targets[0]
216         checker = (prow, pcol)
217
218         # If our partition to transmit is smaller than the partition size
219         # assigned it means that we are loosing part of the block needed to
220         # the computation.
221         if rect and abs(rect[check_idx]) < checker[check_idx]:
222             target = self.data_segments[REVERSED[direction]]
223
224             # If it is so we have to evaluate the maximum number of items
225             # we have other dimension.
226             if (target and abs(target[set_idx]) < abs(rect[set_idx])) or \
227                 not target:
228
229                 if check_idx == 0: out = (0, rect[set_idx])
230                 else: out = (rect[set_idx], 0)
231
232                 self.data_segments[REVERSED[direction]] = out
233                 log.info("Fixing data dependencies by adding %s in %s " \
234                         "direction" % (str(out), LABELS[direction]))
235
236         adjust(UP_LEFT, DOWN_LEFT, LEFT, 0, 1)
237         adjust(UP_RIGHT, DOWN_RIGHT, RIGHT, 0, 1)
238
239         adjust(UP_LEFT, UP_RIGHT, UP, 1, 0)
240         adjust(DOWN_LEFT, DOWN_RIGHT, DOWN, 1, 0)
241
242     def apply(self, matrix, rows=None, cols=None):
243         nw = comm.Get_size() - 1
244
245         if rows is not None and cols is not None:
246             log.info("Skipping auto-derivation")
247             rw = (matrix.cols / cols) * (matrix.rows / rows)
248         else:
249             rows, cols, rw = matrix.derive_partition(self.offsets, nw)
250
251         self.fix_data_deps(rows, cols)
252
253         wdict = {}
254         workers = []
255
256         for _ in range(rw):
257             # Just assign an empty invalid partition the correct partition will
258             # be derived when needed.
259             worker = StencilWorker(self.offsets,
260                                   self.data_segments,
261                                   Matrix.from_list(rows, cols, None),
262                                   matrix.cols, matrix.rows)
263             wdict[worker.wid] = worker
264             workers.append(worker)
265
266         log.info("After partition => Rows: %d Cols: %d" % \
267                 (matrix.rows/rows, matrix.cols/cols))
268
269         for worker in workers:
270             worker.autoconnect(wdict, matrix.rows / rows, matrix.cols / cols)
271         for worker in workers:
272             worker.autoconnect(wdict)
273         for worker in workers:
274             worker.conns.convert_to_id()
275
276         data = [None,]
277
278         for worker in workers:
279             partition = matrix.partition(rows, cols, worker.wid)
280             data.append((worker.offsets, worker.data_segments, partition, \
281                         worker.pwidth, worker.pheight, worker.conns))
```

```

282     start = time.time()
283     comm.scatter(data, root=0)
284     log.info("%.10f seconds to scatter the matrix" % \
285             (time.time() - start))
286
287     matrix = self.reconstruct(matrix.cols/cols)
288
289     log.info("Terminated.")
290     comm.Barrier()
291
292     #print "Result is"
293     #print matrix
294
295 def reconstruct(self, col_stop):
296     start = time.time()
297     data = comm.gather(None, root=0)
298     log.info("%.10f seconds to gather the results" % \
299             (time.time() - start))
300
301     start = time.time()
302
303     rows = []
304     row, col = 0, 0
305     curr_row = None
306
307     for partition in data[1:]:
308         # Creates rows then merge back
309         if curr_row is None:
310             curr_row = partition
311         else:
312             curr_row = numpy.column_stack((curr_row, partition))
313
314     col += 1
315
316     if col == col_stop:
317         col = 0
318         row += 1
319
320     rows.append(curr_row)
321     curr_row = None
322
323     matrix = None
324
325     for row in rows:
326         if matrix is None:
327             matrix = row
328         else:
329             matrix = numpy.vstack((matrix, row))
330
331     log.info("%.10f seconds to reconstruct the result" % \
332             (time.time() - start))
333     return matrix
334
335 def seq_apply(self, matrix):
336     old = matrix.clone()
337     rows, cols = matrix.rows, matrix.cols
338
339     for i in range(matrix.rows):
340         for j in range(matrix.cols):
341             val = old.matrix[i][j]
342
343             for (x, y) in self.offsets:
344                 val = function(val,
345                               old.matrix[(i + x) % rows][(j + y) % cols]
346                               )
347
348             matrix.matrix[i][j] = val
349

```

5 CONCRETE IMPLEMENTATION

Listing 5 presents the class which is used by the `StencilWorker` to re-arrange all contour partitions received by the neighbors in a new sub-matrix, which will then be used to execute the evaluation of the functional code. We make use of utilities functions such as `vstack` and `hstack` which are provided by the `numpy`⁴ library, respectively to vertically or horizontally stack small matrices into bigger ones.

Listing 5: Stencil skeleton implementation: `puzzle.py`

```
1 import time
2 import logging
3
4 from numpy import zeros, vstack, hstack
5 from matrix import Matrix
6 from communicator.enum import UP, DOWN, LEFT, RIGHT, \
7                               DOWN_LEFT, DOWN_RIGHT, \
8                               UP_LEFT, UP_RIGHT
9
10 logging.basicConfig()
11
12 log = logging.getLogger("puzzle")
13 log.setLevel(logging.INFO)
14
15 class Puzzle(object):
16     def __init__(self, center):
17         self.center = center
18         self.pieces = [None, ] * 8
19
20         self.max_up = 0
21         self.max_down = 0
22         self.max_left = 0
23         self.max_right = 0
24
25     def add_piece(self, piece, position):
26         if position in (UP, UP_LEFT, UP_RIGHT):
27             self.max_up = max(self.max_up, piece.rows)
28         if position in (DOWN, DOWN_LEFT, DOWN_RIGHT):
29             self.max_down = max(self.max_down, piece.rows)
30         if position in (LEFT, DOWN_LEFT, UP_LEFT):
31             self.max_left = max(self.max_left, piece.cols)
32         if position in (RIGHT, DOWN_RIGHT, UP_RIGHT):
33             self.max_right = max(self.max_right, piece.cols)
34
35         self.pieces[position] = piece
36
37     def pad(self, part, where, dtype):
38         if part is None:
39             if where == RIGHT and self.max_right > 0:
40                 return zeros((1, self.max_right))
41             if where == LEFT and self.max_left > 0:
42                 return zeros((1, self.max_left))
43             if where == UP and self.max_up > 0:
44                 return zeros((self.max_up, 1))
45             if where == DOWN and self.max_down > 0:
46                 return zeros((self.max_down, 1))
47
48             return None
49
50         if isinstance(part, Matrix): m = part.matrix
51         else: m = part
52
53         if where == RIGHT and m.shape[1] < self.max_right:
54             pad = zeros((m.shape[0], self.max_right - m.shape[1]), dtype)
55             return hstack((m, pad))
56
57         elif where == LEFT and m.shape[1] < self.max_left:
58             pad = zeros((m.shape[0], self.max_left - m.shape[1]), dtype)
```

⁴<http://numpy.scipy.org>

```

59         return hstack((pad, m))
60
61     elif where == UP and m.shape[0] < self.max_up:
62         pad = zeros((self.max_up - m.shape[0], m.shape[1]), dtype)
63         return vstack((pad, m))
64
65     elif where == DOWN and m.shape[0] < self.max_down:
66         pad = zeros((self.max_down - m.shape[0], m.shape[1]), dtype)
67         return vstack((m, pad))
68
69     return m
70
71 def apply(self, offsets, function):
72     start = time.time()
73     matrix = self.center.matrix.copy()
74
75     up = self.pad(self.pieces[UP], UP, matrix.dtype)
76     down = self.pad(self.pieces[DOWN], DOWN, matrix.dtype)
77
78     is_empty = lambda x: x is not None
79
80     matrix = vstack(filter(is_empty, (up, matrix, down)))
81
82     left = self.pad(self.pieces[LEFT], LEFT, matrix.dtype)
83     right = self.pad(self.pieces[RIGHT], RIGHT, matrix.dtype)
84
85     upl = self.pad(self.pad(self.pieces[UP_LEFT], UP, matrix.dtype),
86                    LEFT, matrix.dtype)
87     downl = self.pad(self.pad(self.pieces[DOWN_LEFT], DOWN, matrix.dtype),
88                     LEFT, matrix.dtype)
89     upr = self.pad(self.pad(self.pieces[UP_RIGHT], UP, matrix.dtype),
90                   RIGHT, matrix.dtype)
91     downr = self.pad(self.pad(self.pieces[DOWN_RIGHT], DOWN, matrix.dtype),
92                    RIGHT, matrix.dtype)
93
94     # Hacky fix
95     if self.max_up == 0: upl, upr = None, None
96     if self.max_down == 0: downl, downr = None, None
97     if self.max_left == 0: downl, upl = None, None
98     if self.max_right == 0: downr, upr = None, None
99
100    coll = filter(is_empty, (upl, left, downl))
101
102    if coll: left = vstack(coll)
103    else: left = None
104
105    coll = filter(is_empty, (upr, right, downr))
106
107    if coll: right = vstack(coll)
108    else: right = None
109
110    matrix = hstack(filter(is_empty, (left, matrix, right)))
111
112    log.info("%.10f seconds to create auxiliary matrix" % \
113            (time.time() - start))
114
115    dest = self.center
116    rows, cols = matrix.shape
117    idisp, jdisp = self.max_up, self.max_left
118
119    for i in range(self.max_up, self.max_up + self.center.rows):
120        for j in range(self.max_left, self.max_left + self.center.cols):
121            val = matrix[i][j]
122
123            for (x, y) in offsets:
124                val = function(val, matrix[(i + x) % rows][(j + y) % cols])
125
126            dest.matrix[i - idisp][j - jdisp] = val

```

To finish our analysis we present the communicator class which is responsible to exchange data between workers whenever there is the necessity to do it. Listing 6 presents the `Communicator` class, which provides communication by means of MPI API specification. The only method which deserves an explanation is the `send` method. It is in charge, by looking in the `data_segment` structure, to literary extract the needed sub-matrix from the central one. This is done by the means of the `extract` method provided by the high-level `Matrix` class.

Listing 6: Communication: `communicator/mpi.py`

```
1 import logging
2 from mpi4py import MPI
3 from enum import REVERSED, LABELS
4
5 logging.basicConfig()
6
7 log = logging.getLogger("comm-mpi")
8 log.setLevel(logging.INFO)
9
10 comm = MPI.COMM_WORLD
11 rank = comm.Get_rank()
12
13 class Communicator(object):
14     def __init__(self, parent):
15         self.parent = parent
16
17     def send(self, remote, direction):
18         rect = self.parent.data_segments[direction]
19
20         if rect is None or remote == None or remote == rank - 1:
21             return
22
23         if rect[0] > 0:
24             row_stop = rect[0]
25             row_start = 0
26         elif rect[0] < 0:
27             row_stop = self.parent.height
28             row_start = self.parent.height - abs(rect[0])
29         else:
30             row_stop = self.parent.height
31             row_start = 0
32
33         if rect[1] > 0:
34             col_stop = rect[1]
35             col_start = 0
36         elif rect[1] < 0:
37             col_stop = self.parent.width
38             col_start = self.parent.width - abs(rect[1])
39         else:
40             col_stop = self.parent.width
41             col_start = 0
42
43         m = self.parent.partition.extract(row_start, row_stop, \
44                                         col_start, col_stop)
45
46         log.debug("%d --> send to : %s (direction %s)" % \
47                 (rank - 1, str(remote), LABELS[direction]))
48
49         comm.send(m, dest=remote + 1, tag=direction)
50         log.debug("%d --> send to : %s DONE" % (rank - 1, str(remote)))
51
52     def receive(self, remote, direction):
53         rect = self.parent.data_segments[REVERSED[direction]]
54
55         if rect is None or remote is None or remote == rank - 1:
56             return
57
58         log.debug("%d <-- recv from: %s (direction %s)" % \
59                 (rank - 1, str(remote), LABELS[direction]))
```

```

60
61     data = comm.recv(source=remote + 1, tag=REVERSED[direction])
62
63     log.debug("%d <-- recv from: %s (direction %s) DONE" % \
64               (rank - 1, str(remote), LABELS[direction]))
65
66     return data

```

Listing 7 is presented just to complete the global vision. It just contains various constants and enumerations which are used in various part of the code, and the `Connections` class which is just a class placeholder that is used to store information about the connections between workers (neighbor information).

Listing 7: Communication constants: `communicator/enum.py`

```

1 RIGHT,      \
2 LEFT,       \
3 UP,         \
4 DOWN,       \
5 UP_LEFT,    \
6 DOWN_RIGHT, \
7 UP_RIGHT,   \
8 DOWN_LEFT = range(8)
9
10 LABELS = 'right left up down up_left down_right up_right down_left'
11 RLABELS = 'left right down up down_right up_left down_left up_right'
12
13 LABELS = LABELS.split(' ')
14 RLABELS = RLABELS.split(' ')
15 ORDERED = range(8)
16 REVERSED = [LEFT, RIGHT, DOWN, UP, DOWN_RIGHT, UP_LEFT, DOWN_LEFT, UP_RIGHT]
17
18 class Connections(object):
19     def __init__(self, right, left, up, down):
20         self.right = right
21         self.left = left
22         self.up = up
23         self.down = down
24
25         self.up_left = None
26         self.up_right = None
27         self.down_left = None
28         self.down_right = None
29
30     def __repr__(self):
31         return "right=%s, left=%s, up=%s, down=%s," \
32               "ul=%s, ur=%s, dl=%s, dr=%s" % \
33               (self.right, self.left, self.up, self.down,
34                self.up_left, self.up_right, self.down_left, self.down_right)
35
36     def convert_to_id(self):
37         for lbl in LABELS:
38             obj = getattr(self, lbl)
39
40             if obj is None:
41                 setattr(self, lbl, None)
42                 continue
43
44             wid = obj.wid
45
46             setattr(self, lbl, wid)

```

The last file is the higher level encapsulation of the `numpy.array` data structure. The class `Matrix` is shown in Listing 8. Besides typical methods, such as `get` `set` or `random`, it provides also a `derive_partition` method which applies the same reasoning we have presented in Section 2 to derive an optimal partitioning scheme and `extract` which is used by the `Puzzle` class as

mentioned before.

Listing 8: Matrix class: matrix.py

```
1 import itertools
2 import logging
3 import numpy
4 import copy
5
6 logging.basicConfig()
7
8 log = logging.getLogger("matrix")
9 log.setLevel(logging.INFO)
10
11 class Matrix(object):
12     """
13     The text below is considered a test case.
14
15     >>> m = Matrix.from_string(6, 5, "2 3 4 5 6 5 6 8 9 1 4 4 6 7 3 6 6 3 2 3 0 4 6 3 0 3 5 1 6
16         99")
17     >>> row, col, proc = m.derive_partition([(0, 1), (0, 2)], 100)
18     >>> (row, col, proc)
19     (1, 2, 15)
20     >>> row, col, proc = m.derive_partition([(1, 0), (3, 0)], 100)
21     >>> (row, col, proc)
22     (3, 1, 10)
23     >>>
24
25     Test the limit
26     >>> row, col, proc = m.derive_partition([(1, 0), (3, 0)], 8)
27     >>> (row, col, proc)
28     (3, 2, 5)
29
30     >>> m = Matrix.from_string(50, 40, ' '.join([i for i in map(lambda x: str(x), range(50 * 40)
31         )]))
32     >>> m.derive_partition([(-3, -2), (2, 2)], 10)
33     (5, 40, 10)
34     >>> m.derive_partition([(-1, -1), (2, 2), (0, -1)], 111)
35     (5, 4, 100)
36     """
37
38     def __init__(self, rows, cols, contents=None):
39         self.rows, self.cols = rows, cols
40
41         if contents is not None:
42             self.matrix = contents
43         else:
44             self.matrix = numpy.empty((rows, cols))
45
46     def clone(self):
47         return copy.deepcopy(self)
48
49     def dump(self):
50         print str(self.matrix)
51
52     @staticmethod
53     def from_string(rows, cols, contents):
54         mtx = Matrix(rows, cols)
55         elems = itertools.imap(int,
56                                 contents.replace("\n", " ").strip().split(" "))
57         try:
58             for i in range(rows):
59                 for j in range(cols):
60                     mtx.matrix[i][j] = elems.next()
61         except StopIteration:
62             pass
63
64         return mtx
65
66     @staticmethod
```

```

65     def from_list(rows, cols, contents):
66         return Matrix(rows, cols, contents)
67
68     @staticmethod
69     def random(rows, cols):
70         contents = numpy.random.randint(0, 9, (rows, cols))
71         return Matrix(rows, cols, contents)
72
73     def derive_partition(self, offsets, nproc):
74         """
75         This method should derive a proper partition scheme starting from the
76         offsets you pass in input
77
78         @param offsets model the dependencies you need in order to properly
79             evaluate your function on the matrix.
80         @param nproc number of processors you have available
81         @return a tuple (rows, cols, rw) that can be used to derive the correct
82             partition
83         """
84
85     def get_min_step(lst, idx):
86         if not lst: return 1
87         return lst[0][idx]
88
89     row_dependent = sorted([x for x in offsets if x[0] != 0],
90                           key=lambda x: abs(x[0]), reverse=True)
91     col_dependent = sorted([x for x in offsets if x[1] != 0],
92                           key=lambda x: abs(x[1]), reverse=True)
93
94     # Now we need to get the maximum displacement element, so we
95     # sort by using abs()
96
97     def trivial_sol(depends, is_col=0):
98         aref, bref = self.cols, self.rows
99         if is_col == 0: aref, bref = bref, aref
100
101         astep = get_min_step(depends, is_col)
102         num = aref / astep
103
104         # Adaption in the other direction is not needed since we
105         # start from the barely minimum
106
107         while num > nproc:
108             if astep * 2 > self.cols:
109                 break
110
111             astep *= 2
112             num = int(aref / astep)
113
114         bstep = max(int(bref / (nproc - num)), 1)
115         eproc = int((aref * bref) / (bstep * astep))
116
117         while eproc > nproc:
118             bstep *= 2
119             eproc = int((aref * bref) / (bstep * astep))
120
121         return (is_col == 0) and (astep, bstep, eproc) \
122             or (bstep, astep, eproc)
123
124     # If we do not depend on rows in any way we can fully exploit
125     # the parallelism pattern and do a partition by rows.
126
127     if not row_dependent:
128         return trivial_sol(col_dependent, 1)
129     elif not col_dependent:
130         return trivial_sol(row_dependent, 0)
131
132     log.debug("Non-trivial case evaluation triggered")
133
134

```

```
135     # Now for non-trivial partition we derive the biggest square
136     # that our offsets can derive.
137
138     extract = lambda l, k, rev: sorted(filter(l, offsets), reverse=rev)
139
140     # Sort row descending
141     targets = extract(lambda x: x[0] > 0,
142                       lambda x: x[0], True)
143
144     height = targets and targets[0][0] or 0
145
146     targets = extract(lambda x: x[0] < 0,
147                       lambda x: x[0], False)
148
149     height += targets and abs(targets[0][0]) or 0
150
151     if height > 0:
152         height += 1
153
154     # Sort row descending
155     targets = extract(lambda x: x[1] > 0,
156                       lambda x: x[1], True)
157
158     width = targets and targets[0][1] or 0
159
160     targets = extract(lambda x: x[1] < 0,
161                       lambda x: x[1], False)
162
163     width += targets and abs(targets[0][1]) or 0
164
165     if width > 0:
166         width += 1
167
168     log.debug("Possible partition individuated %dx%d" % (height, width))
169
170     # Now we try to figure out how many workers we can spawn
171
172     if self.cols % width != 0:
173         width = max(1, self.cols / nproc)
174
175     if self.rows % height != 0:
176         height = max(1, self.rows / nproc)
177
178     def throttle(is_width):
179         aparam, bparam = width, height
180         if not is_width: aparam, bparam = bparam, aparam
181
182         astep = aparam
183         elems = self.rows * self.cols
184
185         while (elems / (aparam * bparam)) > nproc:
186             aparam += astep
187
188             if (elems / (aparam * bparam)) < nproc:
189                 break
190
191         return aparam
192
193     log.debug("Trying to fit %d processors" % nproc)
194
195     if width <= height:
196         width = throttle(True)
197         height = throttle(False)
198     else:
199         height = throttle(False)
200         width = throttle(True)
201
202     log.debug("Possible partition individuated %dx%d" % (height, width))
203
204     eproc = int((self.cols * self.rows) / (height * width))
```

```

205
206     return (height, width, eproc)
207
208     def partition(self, rows, cols, idx):
209         col_idx = (idx * cols)
210         row_idx = int(col_idx / self.cols) * rows
211         col_idx %= self.cols
212         return self.extract(row_idx, rows + row_idx, col_idx, cols + col_idx)
213
214     def extract(self, row_start, row_stop, col_start, col_stop):
215         return Matrix(row_stop - row_start, col_stop - col_start,
216                       self.matrix[row_start:row_stop, col_start:col_stop])
217
218     def get(self, i, j):         return self.matrix[i][j]
219     def set(self, i, j, val): self.matrix[i][j] = val
220
221 if __name__ == "__main__":
222     import doctest
223     doctest.testmod()

```

Regarding the functional code this is simply stored in the `functional.py` file presented in Listing 9. In the file the `offsets` variable and a function named `function` is exported to the other modules.

Listing 9: Functional code: `functional.py`

```

1 import numpy
2
3 offsets = ((-1, 0), (0, -1), (0, 1), (1, 0))
4
5 def variance(a, b):
6     if isinstance(a, int):
7         n, mean, m2 = 0, 0, 0
8     elif isinstance(a, list):
9         n, mean, m2 = a
10
11     n += 1
12     delta = b - mean
13     mean += delta / n
14     m2 = m2 + delta * (b - mean)
15
16     if n == 4:
17         return m2/(n - 1)
18     else:
19         return [n, mean, m2]
20
21 function = variance
22
23 def linsolve(a, b):
24     """
25     Here we try to interpret the matrix as follows
26         2
27     1 5 3
28         4
29
30     5 will be our a parameter. The aim of the function is to find a solution to
31     the given algebraic system:
32     1x + 3y = 5
33     2x + 4y = 5
34
35     Then the mean value is calculated and put in the center.
36     """
37     if isinstance(a, int):
38         return [a, b]
39     elif isinstance(a, list):
40         a.append(b)
41
42     if len(a) == 5:
43         equations = numpy.array(((a[1], a[4]),

```



```
44             (a[2], a[3]))
45         given = numpy.array((a[0], a[0]))
46         try:
47             return numpy.linalg.solve(equations, given).mean()
48         except Exception:
49             return 0
50     return a
51
52
53 #offsets = ((-2, -1), (-2, 0), (-2, -2))
54 function = min
```

6 Build instructions

In order to successfully run the source code you need at least CPython 2.5.x. For getting a usable build of `mpi4py` package you need a compiler and the development file of the CPython. You can easily install them by typing:

```
$ apt-get install python python-dev gcc
```

To avoid installing the extension system wide, we have used `virtualenv`⁵ to create an isolated Python environment.

```
$ wget http://pypi.python.org/packages/source/v/virtualenv/virtualenv-1.6.1.tar.gz
$ tar xzf virtualenv-1.6.1.tar.gz
$ python virtualenv-1.6.1/virtualenv.py mpi
New python executable in mpi/bin/python
Installing setuptools.....done.
Installing pip.....done.
$ . mpi/bin/activate
(mpi) $
```

After creating the environment you can install the MPI python wrapper. In our test-bed we were using machines with `MPICH-1`⁶ implementation of the 1.0 MPI API specification.

In our case since `MPiV1.0` is not well supported we have to do some tricks in order to get a working build. We have to thank Lisandro Dalcin the author of the library for his precious support.

```
(mpi) $ http://mpi4py.googlecode.com/files/mpi4py-1.2.2.tar.gz
(mpi) $ tar xzf mpi4py-1.2.2.tar.gz && cd mpi4py-1.2.2

# This is needed to fix bogus error message in the finalization state.
# Take a look to http://code.google.com/p/mpi4py/issues/detail?id=20 for
# more information
(mpi) $ wget http://mpi4py.googlecode.com/svn/trunk/src/python.c -O src/python.c
(mpi) $ export MPICC=/usr/bin/mpicc.mpich
(mpi) $ export MPICH_USE_SHLIB=yes
(mpi) $ python setup.py install
(mpi) $ python setup.py build_exe
(mpi) $ python setup.py install_exe
```

⁵<http://pypi.python.org/pypi/virtualenv>

⁶<http://www.mcs.anl.gov/research/projects/mpich2/>

At the end of this phase you should have a `python2.5-mpi` executable installed under `mpi/bin/` directory. Now to run a specific test you have first to create a random matrix. You have to change the current directory and make it point to `skippy/src`. Then simply spawn a python interpreter and type:

```
Python 2.5.2 (r252:60911, Jan  4 2009, 17:40:26)
[GCC 4.3.2] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> import matrix
>>> import cPickle # For serialization
>>> cPickle.dump(matrix.Matrix.random(3000, 3000), open("/tmp/matrix-3000-3000.mtx", "w"))
>>> quit()
```

This should create a randomly filled matrix 3000 by 3000 in your `/tmp` directory. Now to run the program on different machines you have to write down a *machinefile* which contains the specification about the machines to use in your network. For our tests we have used the following file:

```
axth1.cli.di.unipi.it:3
axth3.cli.di.unipi.it:2
axth4.cli.di.unipi.it:2
axth6.cli.di.unipi.it:2
axth7.cli.di.unipi.it:2
axth8.cli.di.unipi.it:2
axth10.cli.di.unipi.it:2
axth15.cli.di.unipi.it:2
axth29.cli.di.unipi.it:2
axth34.cli.di.unipi.it:2
```

The syntax is pretty straightforward, first the machine host-name or IP address, followed by colon and then by the number of slots to reserve on that specific machine. Than to run the test you can simply type:

```
(mpi) $ mpirun -machinefile multicores -np 3 `pwd`/mpi/bin/python2.5-mpi \
run.py par /tmp/matrix-3000-3000.mtx 3000x1500
3.8472688198 seconds to load the matrix
INFO:stencil:Skipping auto-derivation
INFO:stencil:After partition => Rows: 1 Cols: 2
INFO:stencil:Worker started on processor 1 axth15.cli.di.unipi.it
INFO:stencil:3.3783619404 seconds to scatter the matrix
INFO:stencil:Worker started on processor 2 axth27.cli.di.unipi.it
INFO:puzzle:0.0981070995 seconds to create auxiliary matrix
INFO:puzzle:0.1032311916 seconds to create auxiliary matrix
INFO:stencil:Worker 0: 27.8885238171 seconds to compute the partition
INFO:stencil:Sending back the computed sub-partition from 1
INFO:stencil:Worker 1: 29.1980512142 seconds to compute the partition
INFO:stencil:Sending back the computed sub-partition from 2
INFO:stencil:Worker 0: 2.8959200382 seconds to send back the partition
INFO:stencil:Worker 1: 3.1485359669 seconds to send back the partition
INFO:stencil:32.5039210320 seconds to gather the results
```

6 BUILD INSTRUCTIONS

INFO:stencil:0.1175770760 seconds to reconstruct the result
INFO:stencil:Terminated.
36.0049271584 seconds to apply on 2 processors