

# SMPI Collective Operations Analysis

Arnaud Legrand

March 17, 2015

## Contents

|          |                                         |          |
|----------|-----------------------------------------|----------|
| <b>1</b> | <b>Parsing</b>                          | <b>1</b> |
| <b>2</b> | <b>Analysis</b>                         | <b>2</b> |
| 2.1      | Global overview . . . . .               | 2        |
| 2.2      | Linearity . . . . .                     | 5        |
| 2.3      | Large deployments . . . . .             | 6        |
| 2.4      | Small Deployments . . . . .             | 6        |
| 2.5      | Is there a "best" algorithm ? . . . . . | 7        |
| 2.6      | Conclusion . . . . .                    | 10       |

## 1 Parsing

Let's parse the data collected by Martin

```
1 grep -i precious collectives/log/* > res.log
2
3 use Data::Dumper;
4
5 open(INPUT,"res.log");
6 open(OUTPUT,"> res.csv");
7
8 my(@fields)=qw(date host algo numproc msgsize hostTime hostMem simTime);
9 print OUTPUT (join(',',@fields)."\n");
10
11 while(defined($line=<INPUT>)) {
12     chomp $line;
13     if(!($line =~ /PRECIOUS_RESULT/)) { next; }
14     my(%res)=();
```

```

13 ($head,$tail)=split(/PRECIOUS_RESULT /,$line);
14 $head =~ s/.*log\\//;
15 $head =~ s/.org:$//;
16 ($date,$host) = split(/_/, $head);
17 #      print "$tail\n";
18
19 my(%res) = split(/[: ]/, $tail);
20 $res{date}=$date;
21 $res{host}=$host;
22 #      print Dumper(\%res);
23
24 my(@out)=();
25 foreach(@fields) { push @out, $res{$_}; }
26 #      print Dumper(\@out);
27 print OUTPUT (join(',',@out))."\n";
28 }
29 close(OUTPUT);

```

## 2 Analysis

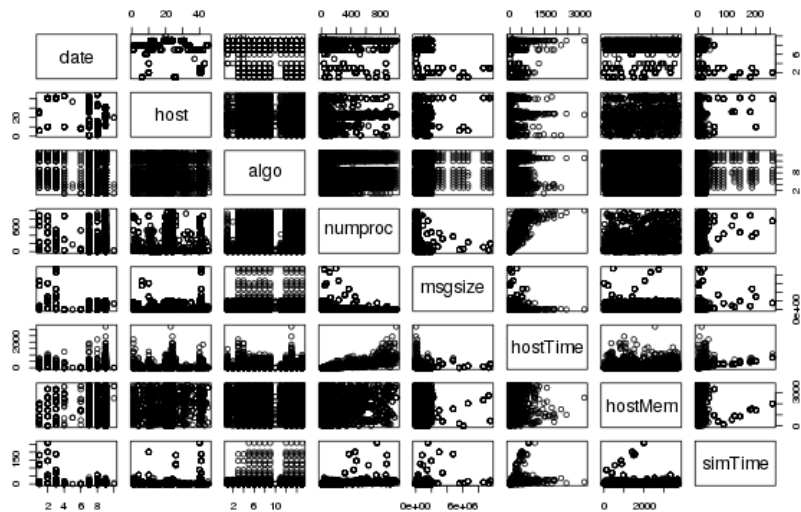
### 2.1 Global overview

```

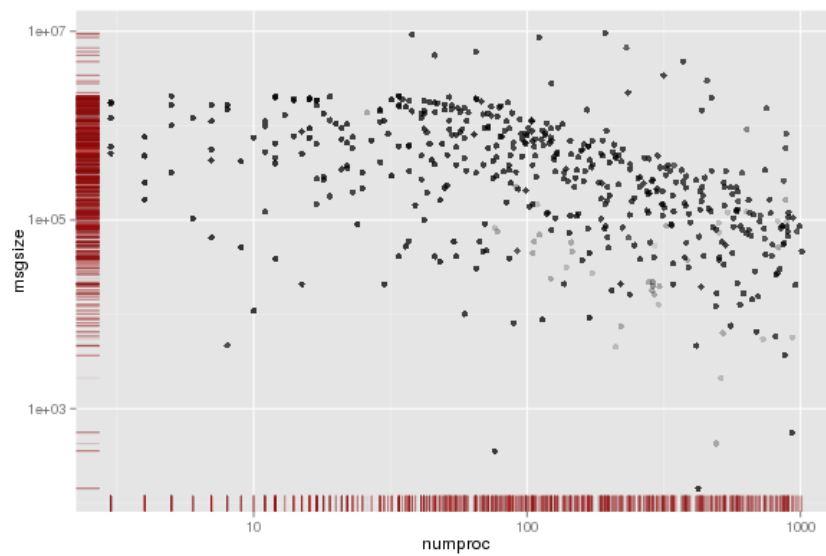
1 library(ggplot2)
2 library(plyr)
3 df=read.csv("res.csv",header=T,strip.white=T);

1 plot(df)

```



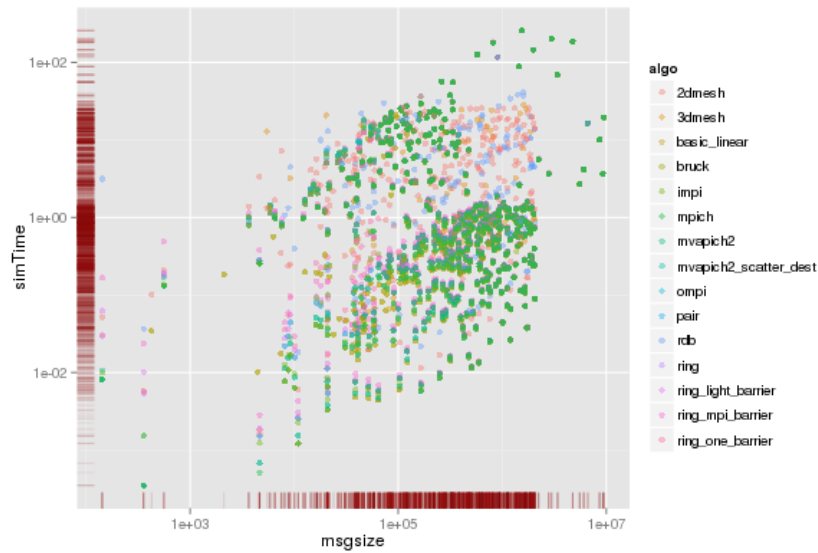
```
1 ggplot(data=df, aes(x=numproc, y=msgsize)) + geom_point(alpha=.1) +
2   geom_rug(alpha=.05, color="dark red") + scale_x_log10() + scale_y_log10()
```



```
1 ggplot(data=df, aes(x=numproc, y=simTime, color=algo)) + geom_point(alpha=.4) +
2   geom_rug(alpha=.05, color="dark red") + scale_x_log10() + scale_y_log10()
```



```
1 ggplot(data=df, aes(x=msgsize, y=simTime, color=algo)) + geom_point(alpha=.4) +
2   geom_rug(alpha=.05, color="dark red") + scale_x_log10() + scale_y_log10()
```

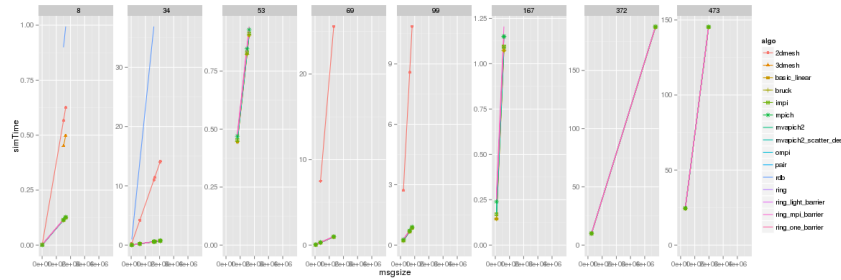


It's hard to really tell from the previous figures because there is too much data with different orders of magnitude but as can be seen in the next graph, `simTime` is roughly linear in `msgsize` and in `numproc`.

```

1 #sizes = sample(df[df$numproc>256,]$numproc , size=8, replace=F)
2 sizes = c(8, 34, 53, 69, 99, 167, 372, 473)
3 ggplot(data=df[df$numproc %in% sizes,],
4       aes(x=msgsize, y=simTime, color=algo, shape=algo)) +
5   geom_line() + geom_point() + ylim(0,NA) +
6   facet_wrap(~numproc,scales="free_y",nrow=1)
7 # + scale_colour_brewer(palette="Set1")

```



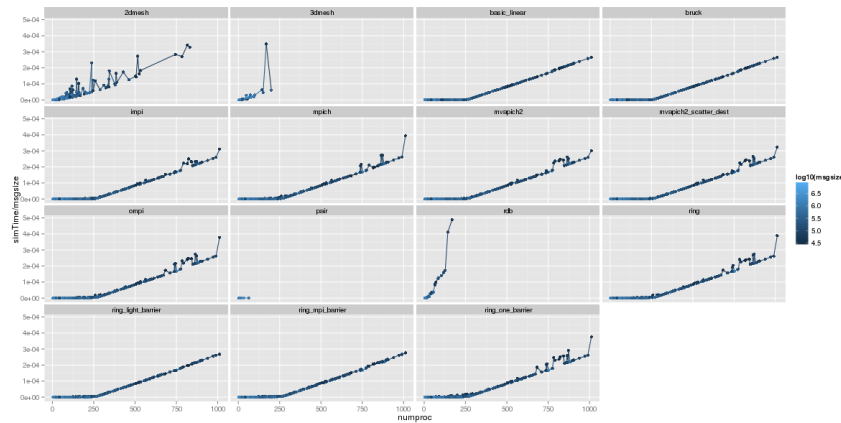
Therefore, in the following, I focus on `simTime/msgsize` and to avoid weird effects, I arbitrarily decide to focus on `msgsize > 25K`. This allows me to get `msgsize` "out of the picture".

## 2.2 Linearity

```

1 ggplot(data=df[df$msgsize > 25000,],
2       aes(x=numproc, y=simTime/msgsize, color=log10(msgsize))) +
3   geom_line() + geom_point() +
4   facet_wrap(~algo) # + scale_colour_brewer(palette="Set1")

```

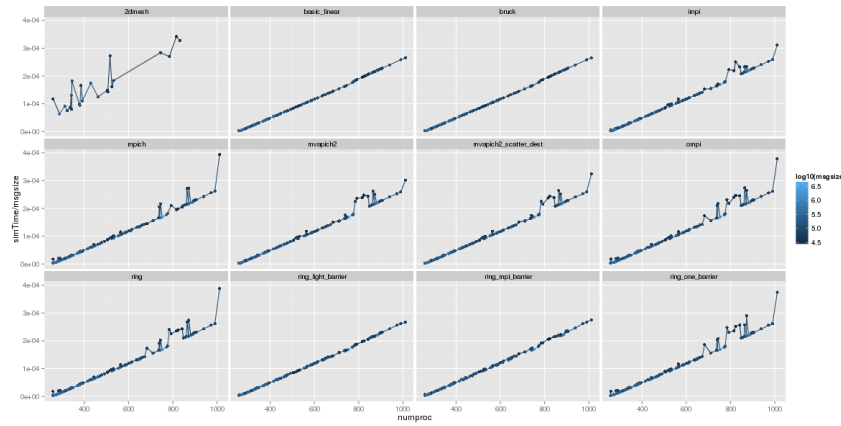


Interesting. It is piece-wise linear. This can easily be explained by the

fact that everything fits in a cabinet for `numproc <= 256`. Let's look at both settings separately.

## 2.3 Large deployments

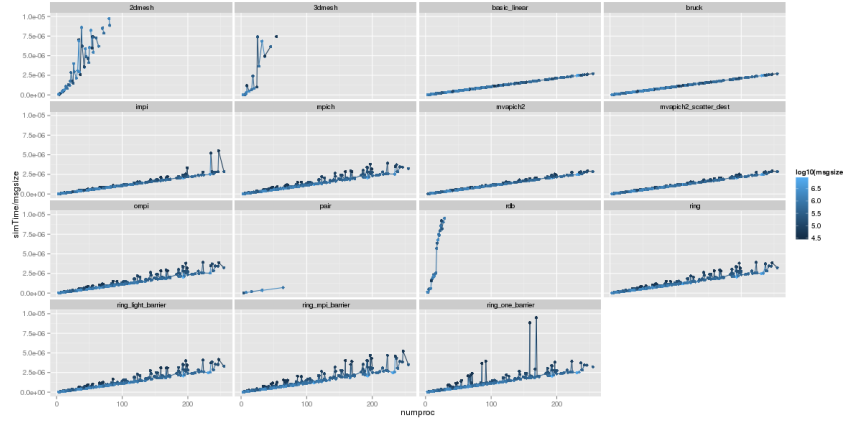
```
1 ggplot(data=df[df$msgsize > 25000 & df$numproc>256,],
2       aes(x=numproc, y=simTime/msgsize, color=log10(msgsize))) +
3       geom_line() + geom_point() +
4       facet_wrap(~algo) # + scale_colour_brewer(palette="Set1")
```



Some algorithms (`basic_linear` and `bruck` `ring_light_barrier` and `ring_mpi_barrier`) perform quite nicely. They are very stable and seem to all have the same best performance. Surprisingly `mpich`, `ompi`, `mvapich2` or `impi` have much more irregular behavior and sometime perform quite bad decision. It is surprising that this is solely due to the chosen algorithm and not to the instability of the platform...

## 2.4 Small Deployments

```
1 ggplot(data=df[df$msgsize > 25000 & df$numproc<256,],
2       aes(x=numproc, y=simTime/msgsize, color=log10(msgsize))) +
3       geom_line() + geom_point() + ylim(0,1E-5) +
4       facet_wrap(~algo) # + scale_colour_brewer(palette="Set1")
```



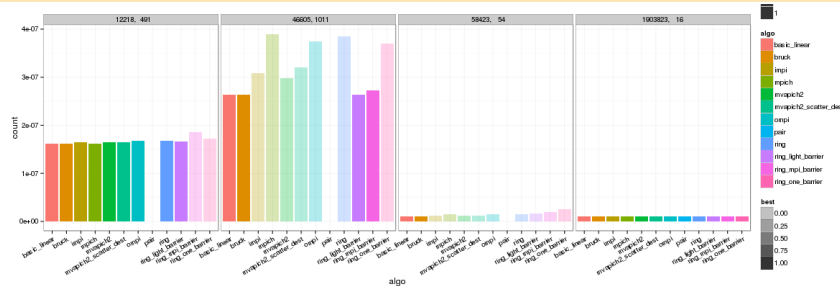
Again, some algorithms (`basic_linear` and `bruck`, and `mvapich2`) behave very well and are good regardless of the number of nodes involved. Again `mpich`, `impi` and `ompi` exhibit some surprising variability. The `ring_light_barrier` and `ring_mpi_barrier` are now less stable than in the large deployment setting.

## 2.5 Is there a "best" algorithm ?

Note that in the previous graphs, it is clear that some algorithms are more stable than others but they may be actually slightly less good in term of performance than more unstable ones. Let's find out by computing whenever for a given configuration (`msgsize`, `numproc`) an algorithm is within 5% of the best one.

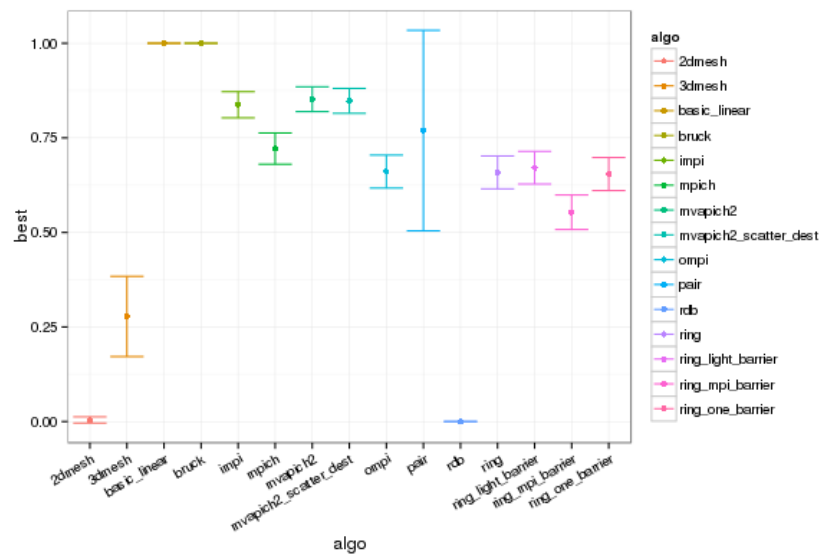
```
1 df = ddply(df,c("msgsize","numproc"), transform, min = min(simTime))
2 df = df[!is.na(df$min),]
3 df$best = 0
4 df[df$simTime <= 1.05 * df$min,]$best = 1
```

```
1 algo_excluded= c("2dmesh","3dmesh", "rdb")
2 msgsize = sample(df$msgsize , size=4, replace=F)
3 d = df[df$msgsize %in% msgsize & !(df$algo %in% algo_excluded),]
4
5 ggplot(data=d,
6         aes(x=algo, weight=simTime/msgsize/numproc, fill=algo,
7             alpha=best, size=factor(best))) +
8   geom_bar() + theme_bw() + scale_alpha(range = c(0.3, 1)) +
9   theme(axis.text.x = element_text(angle = 30, hjust=1)) +
10  facet_wrap(~msgsize*numproc, nrow=1)
```



From such graphs, it really looks again that **bruck** and **basic\_linear** always perform better than the other algorithms. Let's compute some statistics...

```
1 d = summarySE(df, c("best"), c("algo"))
2
3 ggplot(data=d, aes(x = algo, y=best, ymin=best-ci,
4                     ymax=best+ci, color=algo)) +
5   geom_point() + geom_errorbar() + theme_bw() +
6   theme(axis.text.x = element_text(angle = 30, hjust=1))
```



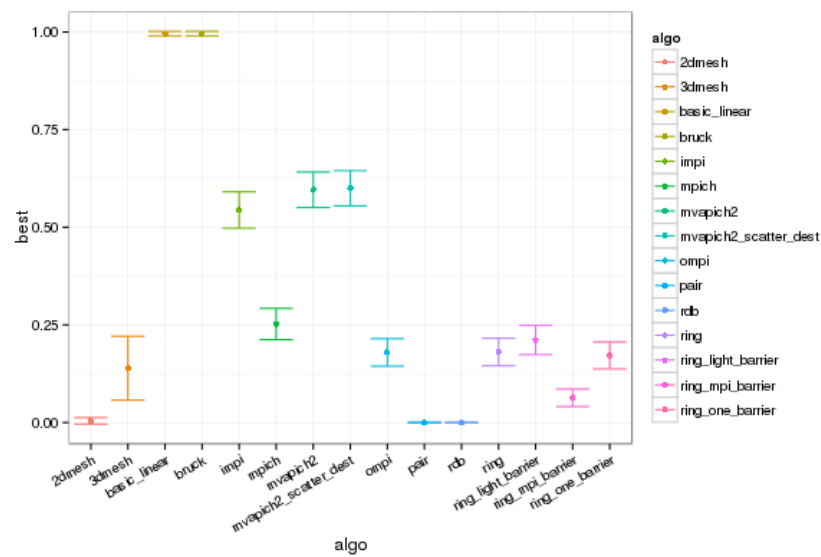
OK. Now, it is clear that **bruck** and **basic\_linear** are always within 5% of the optimal whereas **impi**, **mpich**, **mvapich2** and **ompi** are definitely not. The same keeps holding true with 2%



```

1 df$best = 0
2 df[df$simTime <= 1.02 * df$min,]$best = 1
3
4 d = summarySE(df,c("best"),c("algo"))
5
6 ggplot(data=d, aes(x = algo, y=best, ymin=best-ci,
7                     ymax=best+ci, color=algo)) +
8   geom_point() + geom_errorbar() + theme_bw() +
9   theme(axis.text.x = element_text(angle = 30, hjust=1))

```

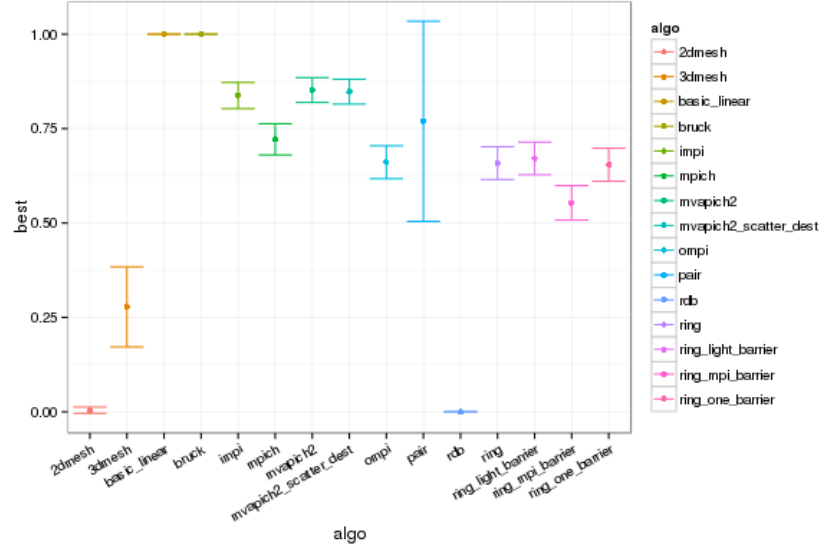


The same keeps holding true with 10%

```

1 df$best = 0
2 df[df$simTime <= 1.1 * df$min,]$best = 1
3
4 d = summarySE(df,c("best"),c("algo"))
5
6 ggplot(data=d, aes(x = algo, y=best, ymin=best-ci,
7                     ymax=best+ci, color=algo)) +
8   geom_point() + geom_errorbar() + theme_bw() +
9   theme(axis.text.x = element_text(angle = 30, hjust=1))

```



## 2.6 Conclusion

In this particular setup, **bruck** and **basic\_linear** perform much better than other more elaborate solutions and the smart optimizations of classical MPI implementations seems not only ineffective but even harmful.

A similar analysis should probably be conducted with other topologies.