

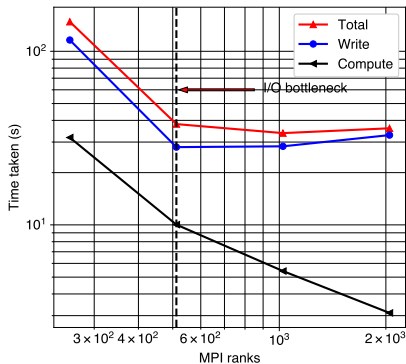
Hiding I/O using SMT on the ARCHER2 HPC Cray EX system

Shrey Bhardwaj
Paul Bartholomew
Mark Parsons

EPCC, University of Edinburgh

May 9, 2023

- 1 Overview
- 2 Methods
- 3 Results
- 4 Conclusions



xCompact-3D¹

¹<https://doi.org/10.5281/zenodo.7898885>

- Hypothesis: Can hyperthreads be used as an I/O server to gain additional “effective” bandwidth?

- Hypothesis: Can hyperthreads be used as an I/O server to gain additional “effective” bandwidth?
- Developed `iocomp` library¹
- Compared performance of mapping dedicated I/O server on SMT cores against other configurations
- Benchmarks used: STREAM² and HPCG³

¹<https://github.com/iocomp-org/iocomp.git>

²<https://github.com/jeffhammond/STREAM.git>

³<https://github.com/hpcg-benchmark/hpcg.git>

- HPE CRAY system
- 5860 nodes
- 2× AMD EPYC™ 7742,
2.25 GHz, 64-core
- Memory/node = 256GiB
- Memory/core = 2GiB
- 128 cores, 256 threads
(including SMT) per
node

¹<https://docs.archer2.ac.uk>

- HPE CRAY system
- 5860 nodes
- 2× AMD EPYC™ 7742, 2.25 GHz, 64-core
- Memory/node = 256GiB
- Memory/core = 2GiB
- 128 cores, 256 threads (including SMT) per node
- Rank of SMT = 128 + rank of core *within* a node

Package L#0

Group0

NUMANode L#0 P#0 (62GB)

L3 (16MB)				L3 (16MB)			
L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)
L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)
L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)
Core L#0	Core L#1	Core L#2	Core L#3	Core L#4	Core L#5	Core L#6	Core L#7
PU L#0 P#0	PU L#2 P#1	PU L#4 P#2	PU L#6 P#3	PU L#8 P#4	PU L#10 P#5	PU L#12 P#6	PU L#14 P#7
PU L#1 P#128	PU L#3 P#129	PU L#5 P#130	PU L#7 P#131	PU L#9 P#132	PU L#11 P#133	PU L#13 P#134	PU L#15 P#135
L3 (16MB)				L3 (16MB)			
L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)	L2 (512KB)
L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)	L1d (32KB)
L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)	L1i (32KB)
Core L#8	Core L#9	Core L#10	Core L#11	Core L#12	Core L#13	Core L#14	Core L#15
PU L#16 P#8	PU L#18 P#9	PU L#20 P#10	PU L#22 P#11	PU L#24 P#12	PU L#26 P#13	PU L#28 P#14	PU L#30 P#15
PU L#17 P#136	PU L#19 P#137	PU L#21 P#138	PU L#23 P#139	PU L#25 P#140	PU L#27 P#141	PU L#29 P#142	PU L#31 P#143

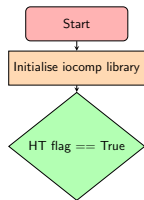
Istopo output clip for ARCHER2.

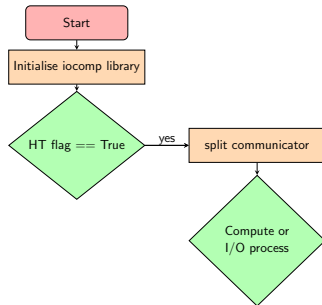
¹<https://docs.archer2.ac.uk>

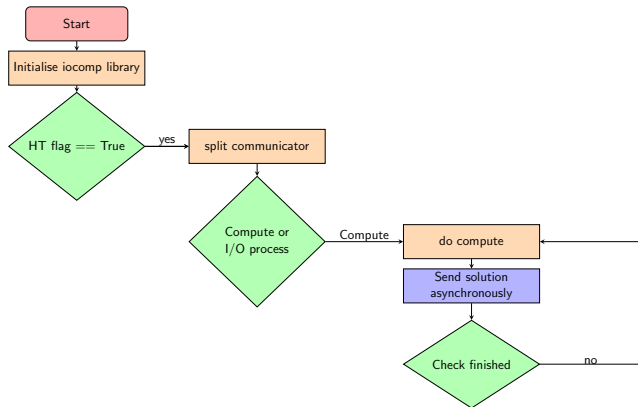
- 1 Overview
- 2 Methods**
- 3 Results
- 4 Conclusions

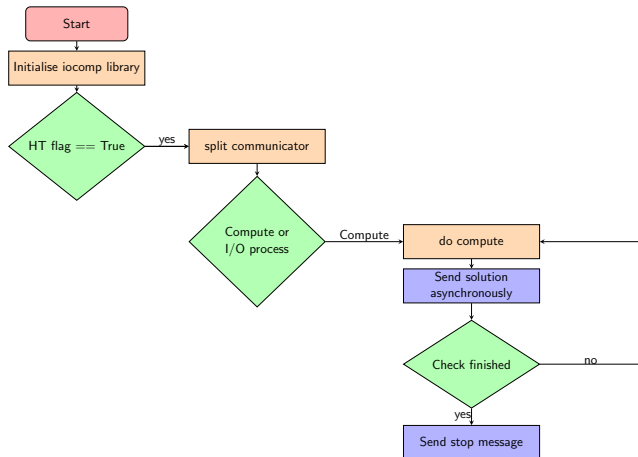
- Framework created to enable splitting of processes into I/O and compute servers¹
- After splitting processes, I/O server gathers data using MPI asynchronous sends from the client process
- Compute server is the client process
- Benchmarks used as client process include HPCG, HPL and STREAM

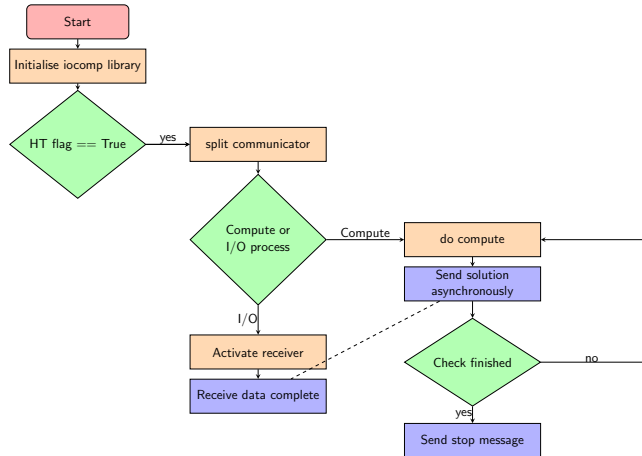
¹<https://github.com/iocomp-org/iocomp.git>

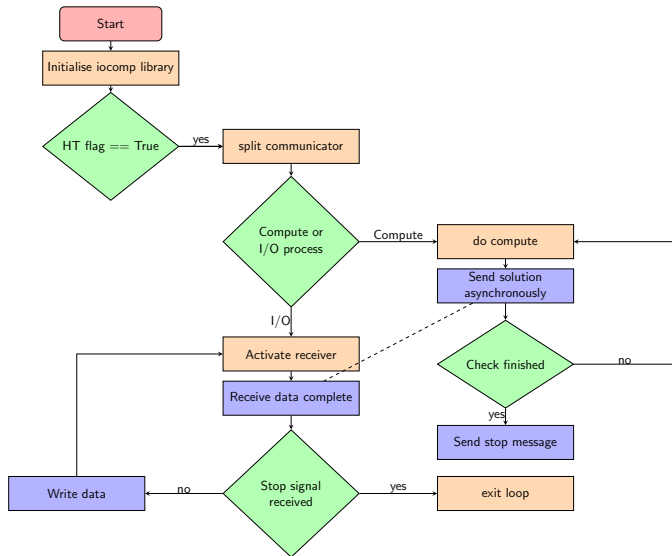


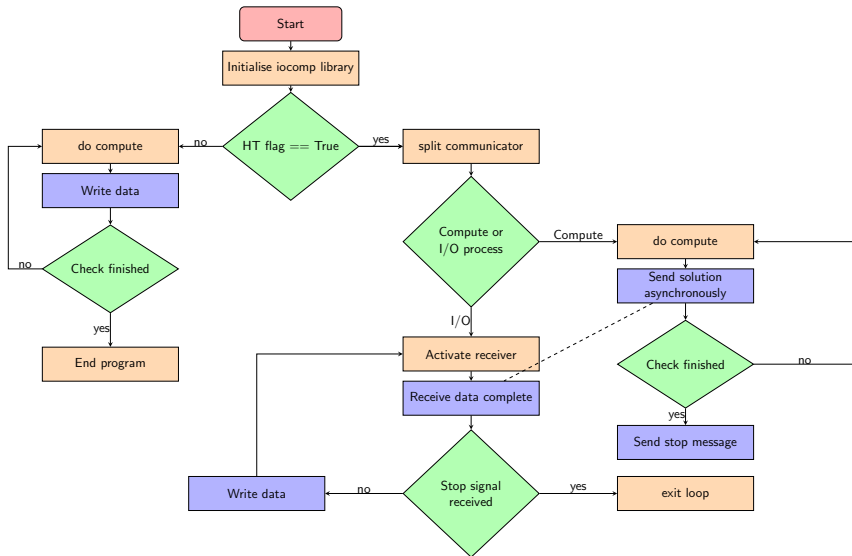












```
1 MPI_Comm iocompInit(struct iocomp_params *iocompParams, MPI_Comm comm, bool FLAG,
2                     int ioLibNum, int fullNode)
3 // iocompParams is structure for the library
4 // comm is the global communicator
5 // FLAG is used to switch between the direct synchronous and the asynchronous I/O
6 // ioLibNum is used to select the I/O library
7 // fullNode is used to specify number of ranks placed in 1 node
```

```
1  MPI_Comm iocompInit(struct iocomp_params *iocompParams, MPI_Comm comm, bool FLAG,  
2      int ioLibNum, int fullNode)  
3  // MPI_Comm will be the new "global" communicator for the client process
```

```
1  for(;;) {
2      MPI_Probe();
3      len = MPI_Get_count();
4      if (len > 0) {
5          MPI_Recv();
6          write_data();
7      }
8      else {
9          // Recieved ghost message, exit.
10         break();
11     }
12 }
```

```
1 // find closest sq root
2 root = sqrt(dataSize);
3 // if its a perfect square
4 if(root*root == dataSize) {
5     dim[0] = root;
6     dim[1] = root;
7 }
8 // if square root is a factor
9 else if(dataSize%root == 0) {
10    dim[0] = root;
11    dim[1] = dataSize/root;
12 }
13 // else search for closest factors
14 else {
15    for(int i = 1; i < root; i++) {
16        if(dataSize%(root-i) == 0) {
17            dim[0] = (root-i);
18            dim[1] = dataSize/(root-i);
19            break;
20        }
21    }
22 }
```

HT flag disabled

- Sequential
 - Default case, with sequential compute and I/O processing

HT flag disabled

- Sequential
 - Default case, with sequential compute and I/O processing

HT flag enabled

- Consecutive
 - Uses 2x number of cores as sequential, **without** SMT
 - Compute and I/O processes are placed on separate cores

HT flag disabled

- Sequential
 - Default case, with sequential compute and I/O processing

HT flag enabled

- Consecutive
 - Uses 2x number of cores as sequential, **without** SMT
 - Compute and I/O processes are placed on separate cores
- Hyperthread
 - Uses the same number of physical cores as sequential, **with** SMT
 - Corresponding SMT cores as I/O processes

HT flag disabled

- Sequential
 - Default case, with sequential compute and I/O processing

HT flag enabled

- Consecutive
 - Uses 2x number of cores as sequential, **without** SMT
 - Compute and I/O processes are placed on separate cores
- Hyperthread
 - Uses the same number of physical cores as sequential, **with** SMT
 - Corresponding SMT cores as I/O processes
- Oversubscribe
 - Uses the same number of cores as sequential **without** SMT
 - Compute and I/O processes are placed on the same cores

Number of nodes	MPI size	Division of MPI processes	
		Compute	I/O
1	128	0-63	64-127
2	256	0-127	128-255
3	384	0-127, 256-319	128-255, 320-383
4	512	0-127, 256-383	128-255, 384-511

Table: Division of MPI processes under different MPI sizes with a “fullNode” value of 128.

- Representing an extreme end of memory-bound computational kernel
- Contiguous memory access runs at limit of largest level in memory hierarchy used
- Implemented as a test case for the iocomp library under stream directory

¹<https://github.com/iocomp-org/iocomp.git>

```
1  for(i=0;i<LOOPCOUNT;i++){
2      copy(C) + test_triad(A)
3      if ((i+1)%WRITE_FREQ == 0){
4          wait_triad(A)
5          send_copy(C)
6      }
7      scale(B) + test_copy(C)
8      if (i%WRITE_FREQ == 0){
9          wait_copy(C)
10         send_scale(B)
11     }
12     add(C) + test_scale(B)
13     if (i%WRITE_FREQ == 0){
14         wait_scale(B)
15         send_add(C)
16     }
17     triad(A) + test_add(C)
18     if (i%WRITE_FREQ == 0){
19         wait_add(C)
20         send_triad(A)
21     }
22 }
```

STREAM overview.

```
1  add(c, a, b ...)
2  {
3      for(int i = 0; i<size; i++)
4          {
5              c[i] = a[i] + b[i];
6              if(i%WRITE_FREQ == 0)
7                  {
8                      dataSendTest(...);
9                  }
10         }
11 }
```

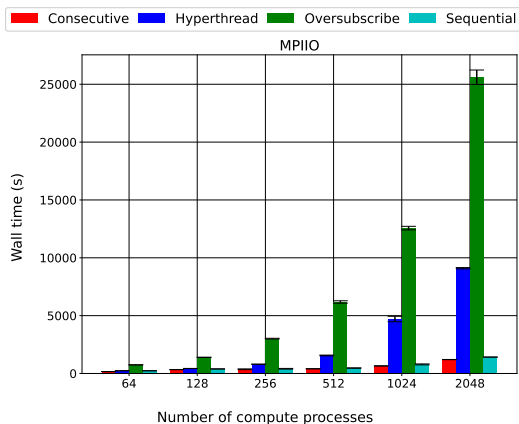
STREAM add kernel.

- More representative of typical HPC applications
 - combining multigrid preconditioner and Conjugate Gradient solver
 - adds indirection and communication
- Another example of a memory-bound program

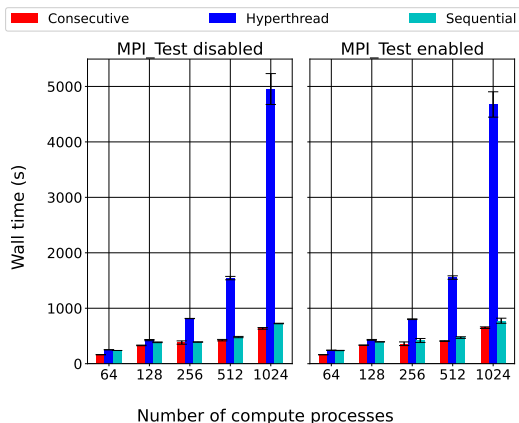
¹<https://github.com/iocomp-org/iocomp-hpcg.git>

```
1 initialise_matrix();
2 for(i=0;i<numberOfCgSets ;i++) // numberOfCgSets=10
3 {
4     dataSend(matrix...); } Send time
5
6     ZeroVector(x)
7     dataSendTest(matrix...); } Compute time
8     CG(A,x,...);
9     testnorms_data.values;
10
11     dataWait(matrix...); } Wait time
12
13     testnorms();
14 }
15
16
17 }
```

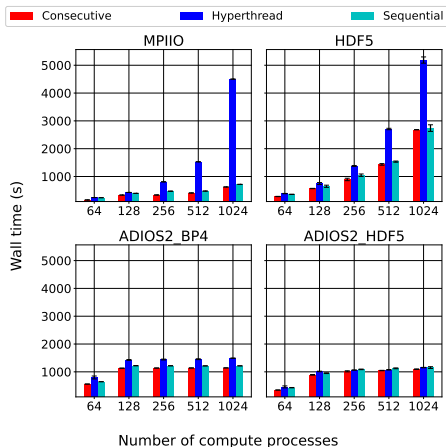
- 1 Overview
- 2 Methods
- 3 Results**
- 4 Conclusions



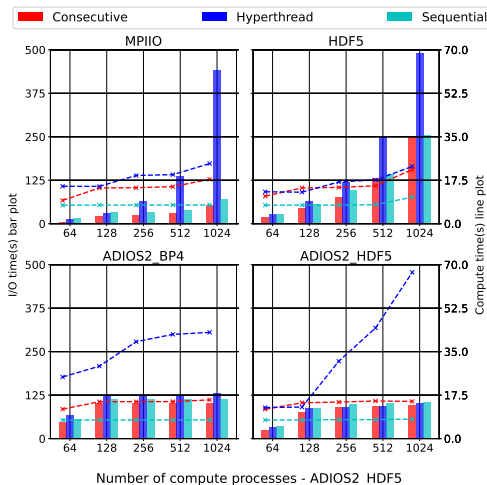
STREAM wall time comparison using MPIIO comparing oversubscribe mapping to the other mappings with local size 0.125GiB.



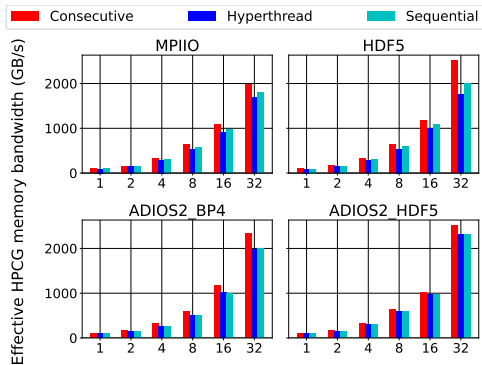
Wall time comparison between STREAM runs using MPIIO I/O backend with “MPI_Test” enabled and disabled for a local size of 0.125GiB.



STREAM wall time comparison using different I/O libraries for local size of 0.125GiB.

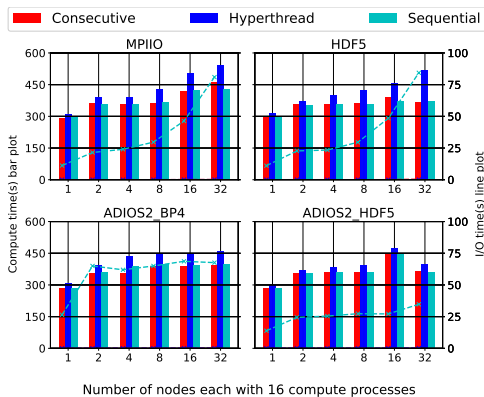


Breakdown of compute time vs I/O time for local size 0.125GiB.

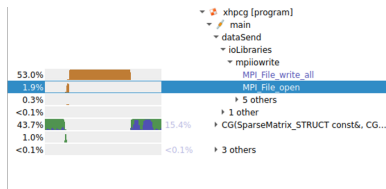


Number of nodes each with 16 compute processes

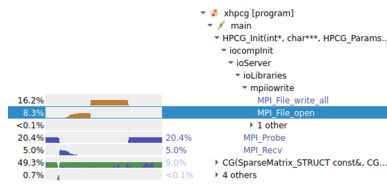
HPCG effective bandwidth 1.69GiB local size.



Breakdown of compute time vs I/O time for HPCG with 1.69GiB local size.



MAP breakdown of Sequential run using HPCG with local size 0.125GiB.



MAP breakdown of Hyperthread run using HPCG with local size 0.125GiB¹.

¹To note: these timing proportions should be multiplied by 2 to get the total time of writing due to the averaging being used per node by MAP

- 1 Overview
- 2 Methods
- 3 Results
- 4 Conclusions**

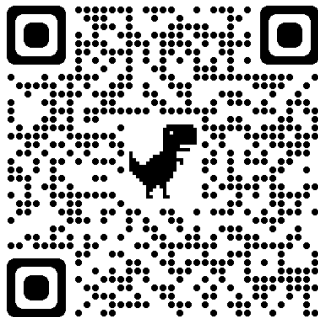
- iocomp was created to compare different mappings for an I/O server
 - Different cases considered; Hyperthread, Consecutive, Oversubscribe and Sequential
 - Different I/O backends were also tested with these mappings
 - HPCG and STREAM benchmarks were tested
 - Consecutive was the best performer, and hyperthreads performance was dependent on the I/O backend used

- iocomp was created to compare different mappings for an I/O server
 - Different cases considered; Hyperthread, Consecutive, Oversubscribe and Sequential
 - Different I/O backends were also tested with these mappings
 - HPCG and STREAM benchmarks were tested
 - Consecutive was the best performer, and hyperthreads performance was dependent on the I/O backend used
- Future work
 - Hyperthread implementation would be more optimised with a shared memory access
 - FEniCSx, a PDE solver will be integrated and tested using iocomp
 - Checkpointing simulations will be tested with the library

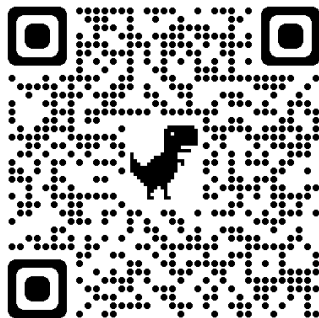
I would like to thank Dr. Paul Bartholomew, Prof. Mark Parsons and EPCC for guidance and suggestions.

This work used the ARCHER2 UK National Supercomputing Service ¹ and was supported by an EPCC funded studentship as part of the ASiMoV project (EP/S005072/1).

¹<https://www.archer2.ac.uk>



QR code for iocomp-org

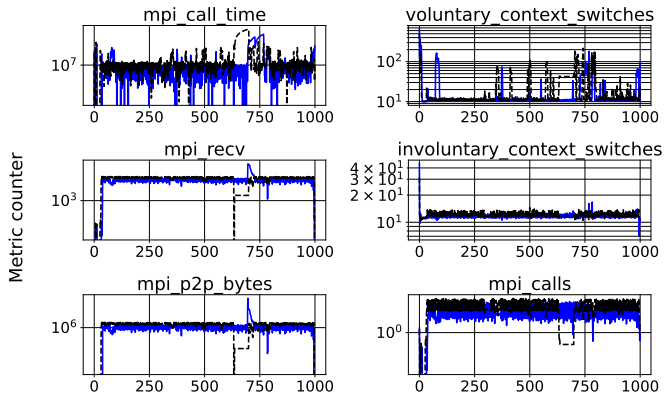


QR code for iocomp-org

→ email: shrey.bhardwaj@ed.ac.uk

→ iocomp: <https://github.com/iocomp-org/iocomp.git>

5 Appendix



Time steps, dashed line for sequential and solid line for hyperthread