

Performance of Coarray Fortran vs MPI in a CFD Application

Mike Ashworth

Scientific Computing Department, STFC Daresbury Laboratory, Sci-Tech Daresbury, Warrington WA4 4AD, United Kingdom

Introduction

Most scientific simulation codes running on distributed memory computers use a high-level language with MPI.

MPI is

- portable (easy to implement)
- simple (easy to learn)
- flexible (can implement any parallel algorithm)
- efficient (easy to optimize)

However, MPI codes can have many lines of additional message passing code and be difficult to maintain.

Coarray Fortran has several advantages [3]:

- Easy to write code; the compiler looks after the communication
- References to local data are obvious as such
- Easy to maintain code esp. more concise than MPI
- Integrated with Fortran for type checking, etc.
- The compiler can optimize communication
- Local optimizations are still available
- Few demands on the compiler, e.g. for coherency

Shock-Boundary Layer Interaction (SBLI)

A finite difference formulation of Direct Numerical Simulation (DNS) of Turbulence from the University of Southampton [2].

A sophisticated DNS code that incorporates a number of advanced features:

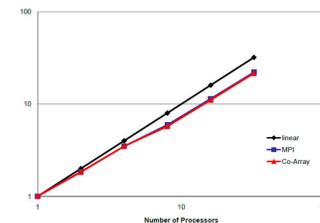
- High-order central differencing
- Shock-preserving advection scheme from the total variation diminishing (TVD) family
- Entropy splitting of the Euler terms
- Stable boundary scheme

Key computational features:

- Fortran 90 plus MPI
- 3D halo exchange
- With large domains scales to 168,000 cores

Previous Work

Figure 1. Speedup of the MPI and CAF codes relative to the MPI code on one processor of the Cray X1E



Ashby and Reid [1] produced a CAF implementation of SBLI and tested it on the Cray X1E achieving identical performance to MPI. They concluded:

- The main advantage of using coarrays is that the code is more readable and understandable, and thus more maintainable
- Our tests show that the performance of our code does not suffer from the change to coarrays
- With the ... Fortran 2008 standard, coarrays represent the first International Standard parallel programming language, ... their use will enhance productivity and further the exploitation of parallelism

The Coarray Implementation

Ashby and Reid handled variable array sizes by implementing arrays as allocatable components of coarrays of derived type.

Dynamic co-shaping, e.g. `q[nx,ny,*]`, replaces the use of Cartesian communicators in the MPI version.

The 'sync all' statement is used for synchronisation.

Three options are available in the code:

Option: Push vs. Pull

- push (put) stores local data into memory on a remote process: `a[k] = a`
- pull (get) read data from a remote process: `a = a[k]`

Option: Packing Data

- packed: small transfers are packed into a buffer as in the MPI code
- non-packed: all transfers are expressed as separate coarray assignments no matter how small

Option: Fortran 77 vs Fortran 90

- Fortran 77: coarray assignments are made using nested do loops
- Fortran90: coarray assignments are expressed using Fortran 90 array syntax

Performance

- Simple channel flow
- Domain of 120 x 120 x 120 run for 100 timesteps
- Run on ARCHER, the UK national supercomputer
- Cray XC30 node is 12-core Intel Ivybridge 2.7 GHz CPUs
- Cray PE 2.1.1, Cray Fortran 8.2.6
- Nodes linked using the Cray Aries interconnect
- Performance plotted is a constant (1000) divided by the execution time measured by system_clock

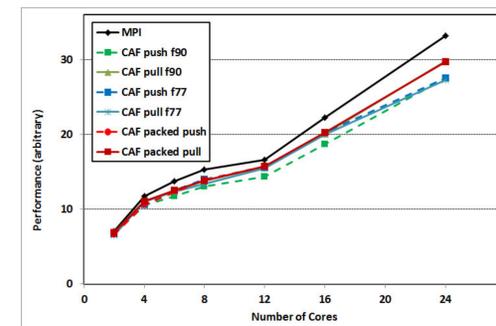


Figure 2. Performance of MPI and various forms of CAF for SBLI up to 24 cores (one node) on the Cray XC30

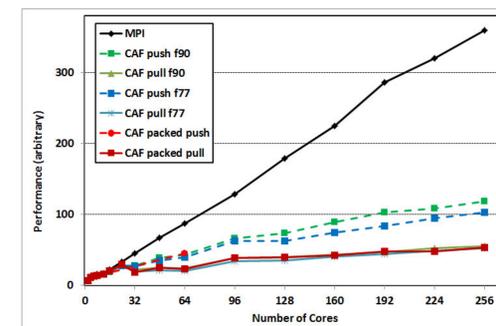


Figure 3. Performance of MPI and various forms of CAF for SBLI up to 256 cores on the Cray XC30

- CAF on a single node achieves up to 94% of MPI
- On multiple nodes this fraction ranges from 15%-33%
- Push (put) is better than pull (get)
- Packing into buffers and the Fortran 90 array syntax make little difference to performance

Profiling

Sample	Time	Group	Function	Sample	Time	Group	Function
100.0%	731.2		Total	100.0%	3045.3		Total
64.4%	471.2		USER	63.8%	1949.3		USER
29.6%	216.2	36.8%	lchm	29.6%	915.3	36.8%	lchm
8.3%	60.7	15.3%	lchm_euler	8.3%	257.0	15.3%	lchm_euler
6.6%	48.4	21.6%	ldia_2deriv	6.6%	207.0	21.6%	ldia_2deriv
6.4%	46.3	22.5%	ldia_3deriv	6.4%	200.2	22.5%	ldia_3deriv
5.6%	41.1	17.9%	ldia_2deriv	5.6%	172.1	17.9%	ldia_2deriv
2.8%	18.6	15.4%	ldia_3deriv	2.8%	105.7	15.4%	ldia_3deriv
1.8%	13.1	8.9%	ldia_3deriv	1.8%	57.1	8.9%	ldia_3deriv
1.3%	9.2	11.5%	ldia_3deriv	1.3%	41.1	11.5%	ldia_3deriv
29.3%	214.0		MPI	2.6%	79.2		MPI
11.3%	82.4	95.6%	lchm_comm_dup	1.7%	50.4	95.6%	lchm_comm_dup
9.3%	68.3	222.7%	lchm_comm	1.6%	48.0	222.7%	lchm_comm
4.4%	32.4	47.5%	lchm_comm	1.3%	39.2	47.5%	lchm_comm
3.0%	22.0	173.0%	lchm_comm	1.0%	31.0	173.0%	lchm_comm
6.3%	45.9		ETC	36.0%	1095.1		ETC
3.9%	28.7	16.3%	RTOR_M_00	29.3%	860.8	16.3%	RTOR_M_00
1.5%	10.6	9.4%	RTOR_M_00	4.1%	126.8	9.4%	RTOR_M_00

Figure 4. CrayPAT output for the SBLI MPI version (left) and for CAF push f90 code (right) on 128 cores on the Cray XC30

Profiling of both MPI and CAF code shows:

- We use non-blocking MPI calls so overheads are in MPI_WaitAny
- Equivalent wait times for CAF are in STOP2 in the ETC group
- Many of the user subroutines show a large increase in runtime
- Especially the main program `pdns3d_` which increases from 9 to 452 samples
- CrayPat shows that the high water mark for memory usage for MPI is 62 MB but for CAF it is 102 MB

Many of these features are not understood and investigations are ongoing.

References

- Ashby, J.V. and Reid, J.K., Migrating a Scientific Application from MPI to Coarrays, Proceedings of the Cray User Group Conference, 2008
- Redford, J.A., Sandham, N.D. and Roberts, G.T., Direct numerical simulation of transitional flow at high Mach number coupled with a thermal wall model, Computers & Fluids 45 pp. 37-46, 2011
- Reid, J.K., Parallel programming in Fortran with Coarrays, RAPS workshop, ECMWF, Reading, 5 November 2008

Acknowledgement

This work is part of the Software Outlook program funded through a Service Level Agreement between STFC and the Engineering and Physical Sciences Research Council.