

Parallel Application Scaling, Performance, and Efficiency

NUG 2010 Supplemental Material (David Skinner and Katie Antypas)









Overview

- Review Some Basic MPI
- Domain Decomposition
- Load Balancing
- Case Study: FLASH Scaling
 Performance
- Performance Monitoring with IPM
- Many-core chips and the future of parallel programming







Overview and History of MPI

- Library (not language) specification
- Goals
 - Portability
 - Efficiency
 - Functionality (small and large)
- Most basic communications are 2 sided
- Pros
 - Programmer has control at low level
 - Performance model understood
 - Can be very high performing
- Cons
 - Programmer has control at low level
 - Error prone



Questions about memory usage as cores/node increase





Generic Message Passing

send(void* sendbuffer, int num_elements, int destination_rank)

receive(void* recvbuffer, int num_elements, int source_rank)

Processor 0						
$\mathbf{x} = 5$						
send(&x,	1,	1)				
$\mathbf{x} = 7$						

Processor 1		
receive(&x, print x	1,	0)

What rules are needed so that processor 1 receives "5" and not "7"?







Ways to Send Data

Non-Buffered

	Blocking	Non-blocking
d	Sending processor "blocks" or "waits" for receive	Sending processor completes call, but must be careful not to overwrite send buffer until receive operation has completed
	Sending process completes call after sendbuf has been copied to another buffer	
Office	of	ri



Buffered





Be careful with buffering ...

send(void* sendbuffer, int num_elements, int destination_rank)

receive(void* recvbuffer, int num_elements, int source_rank)

Processor 0
do i=1, 1000
 produce_data(&x)
 send(&x, 1, 1)
end do

Processor 1
do i=1, 1000
receive(&x, 1, 0)
consume(&x)
end do

What could go wrong with buffered send?





Example – John Mellor-Crummey, Rice University



Basic Point to Point

- Blocking Non-buffered
 - MPI_Send()
 - MPI_Recv()

```
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
if (rank == 0) {
    MPI_Send(&work, 1, MPI_INT, dest, TAG, MPI_COMM_WORLD);
} else {
    MPI_Recv(&result, 1, MPI_INT, src, TAG, MPI_COMM_WORLD,
    &status);
}
```







Non-Blocking Operations

- MPI_lsend()
- MPI_Irecv()
- "I" is for immediate
- Paired with MPI_Test()/MPI_Wait()







Non-Blocking Operations

```
MPI Comm rank(comm,&rank);
if (rank == 0) {
   MPI Isend(sendbuf,count,MPI REAL,1,tag,comm,&request);
   /* Do some computation */
   MPI Wait(&request,&status);
} else {
   MPI Irecv(recvbuf,count,MPI REAL,0,tag,comm,&request);
   /* Do some computation */
   MPI Wait(&request,&status);
}
```







Collective Operations

- May be layered on point to point
- May use tree communication patterns
 for efficiency
- Synchronization! (No non-blocking collectives)







Collective Operations

MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0, comm);







O(log P)





Quiz: MPI_Send()

- After I call MPI_Send()
 - The recipient has received the message
 - I have sent the message
 - I can write to the message buffer without corrupting the message
- I can write to the message buffer







Quiz: MPI_Isend()

- After I call MPI_lsend()
 - The recipient has started to receive the message
 - I have started to send the message
 - I can write to the message buffer without corrupting the message
- None of the above (I must call MPI_Test() or MPI_Wait())







Minimizing Latency

- Collect small messages together (if you can)
 - One 1024-byte message instead of 1024 one-byte messages
- Minimize other overhead (e.g., copying)
- Overlap with computation (if you can)







Example: Domain Decomposition









Naïve Approach

.....

BERKELEY

```
while (!done) {
  exchange(D, neighbors, myrank);
  dored(D);
  exchange(D, neighbors, myrank);
  doblack(D);
}
void exchange(Array D, int *neighbors, int myrank) {
  for (i = 0; i < 4; i++)
    MPI_send(...);
  for (i = 0; i < 4; i++)
    MPI recv(...);
}
```





Naïve Approach

- Deadlock! (Maybe)
- Can fix with careful coordination of receiving versus sending on alternate processes
- But this can still serialize









CCCC

BERKELEY

```
while (!done) {
  exchange(D, neighbors, myrank);
  dored(D);
  exchange(D, neighbors, myrank);
  doblack(D);
}
void exchange(Array D, int *neighbors, int myrank) {
  for (i = 0; i < 4; i++) {
    MPI_Sendrecv(...);
  }
}
```



Science

Immediate Operations

```
while (!done) {
  exchange(D, neighbors, myrank);
  dored(D);
  exchange(D, neighbors, myrank);
  doblack(D);
}
void exchange(Array D, int *neighbors, int myrank) {
  for (i = 0; i < 4; i++) {
    MPI Isend(...);
    MPI Irecv(...);
  }
  MPI Waitall(...);
                                                         CCCC
        Office o
IFRG
```



Basic Functions

MPI_Init

Initializes MPI

MPI_Comm_size MPI_Comm_rank

MPI_Send MPI_Recv

sends data receives data

MPI_Reduce MPI_Allreduce MPI_Bcast

MPI_Finalize



reduce data to single processor reduce all procs data to all procs broadcasts to all procs

Returns # tasks in communicator

Returns ID of current proc

Closes MPI





Load Balancing







Load Balance : cartoon

The Universal Parallel Science App



~All apps come down to the same basic pattern. Ok, Maybe there is no I/IO.

Unbalanced:



Balanced:









Load Balance: real code





MPI Rank →

Time \rightarrow







Communication Time: 64 tasks show 200s, 960 tasks show 230s





Load Balance: ~code

while(1) {
 do_flops(N_i);
 MPI_Alltoall();
 MPI_Allreduce();
}









Load Balance : analysis

- The 64 slow tasks (with more compute work) cause 30 seconds more "communication" in 960 tasks
- This leads to 28800 CPU*seconds (8 CPU*hours) of unproductive computing
- All imbalance requires is one slow task and a synchronizing collective!
- Parallel computers allow you to scale both your computation and your load imbalance.





Load Balance : FFT



Science









Dynamical Load Balance: Motivation











Imbalance most often a byproduct of data decomposition
Must be addressed **before** further MPI tuning can happen
Good software exists for graph partitioning / remeshing



•Dynamical load balance may be required for adaptive codes







Scaling Study







Scaling: definitions

- Scaling studies involve changing the degree of parallelism.
- Strong scaling
 Fixed problem size, more computer
- Weak scaling
 - Problem size grows with concurrency







Parallel Performance Measurements

Speed up = T_{serial} /T_{parallel}(n)

 T_{serial} = 100 secs
 T_{parallel}(2) = 80 secs
 25% speed up

- Efficiency = $T_{serial}/(n^*T_{parallel}(n))$ - 100/(2*80) =
 - -100/(200) -
 - 62% efficiency
- Perfect Scaling?



Be aware there are multiple definitions for these terms













Scaling: Analysis

- What is happening in the 8192 case?
 - Compute per core decreasing
 - Synchronization rate increasing
 - Surface to volume ratio increasing
- What else could happen?
 - Algorithmic scaling may change
 - Maybe we hit an architectural boundary in the machine (switch level, mid-plane, queue, etc.)
 - Maybe depleted some buffer space resource
 - Many more things...performance debugging at scale is detective work in the application + architecture space







Parallel programs are easier to mess up than serial ones. Here are a couple common pitfalls.







What's wrong here?



Communication Event Statistics (100.00% detail)

	Buffer Size	Ncalls	Total Time	Min Time	Max Time	%MPI	%Wall	
MPI_Allreduce	8	3278848	124132.547	0.000	114.920	59.35	16.88	
MPI_Comm_rank	0	35173439489	43439.102	0.000	41.961	20.77	5.91	
MPI_Wait	98304	13221888	15710.953	0.000	3.586	7.51	2.14	
MPI_Wait	196608	13221888	5331.236	0.000	5.716	2.55	0.72	
MPI_Wait	589824	206848	5166.272	0.000	7.265	2.47	0.70	







MPI_Barrier

Function	Total calls	Total ti	me (sec)	Total buffer size (MB)	Avg. Buffer Size/call (Bytes)
MPI_Barrier	6.02e+05	3.48e+05	44.23%	0	0
MPI_Allreduce	3.18e+07	2.31e+05	29.33%	3. <mark>61</mark> e+05	11,936
MPI_Send	1.29e+08	1.29e+05	16.36%	5.24e+04	426
MPI_Bcast	5.73e+07	6.08e+04	7.73%	5.39e+04	987
MPI_Reduce	1.08e+08	1.24e+04	1.58%	1.66e+05	1,620
MPI_Recv	1.29e+08	6.11e+03	0.78%	5.24e+04	426
MPI_Comm_rank	1.14e+03	5.92e-01	7.52e-05%	0	0
MPI_Comm_size	6.66e+02	0	0%	0	0



Is MPI_Barrier time bad? Probably. Is it avoidable? ~three cases:

- 1) The stray / unknown / debug barrier
- 2) The barrier which is masking compute imbalance
- 3) Barriers used for I/O ordering

ENERGY Office of Science





How to use IPM : basics

Many of the graphs in this talk were generated with a tool called IPM – Integrated Performance Monitoring: free and easy to install http://ipm-hpc.sourceforge.net/

On galera
1) >> mpicc test.c -lipm
2) Run job as usual
3) Appended to your output file

Science





How to use IPM : basics

# : #	#IPMv0.982	#########	############	*	*############	##########	############	####
π #	command .	unknown	(comploted)					
π #	host ·	n9-1-6/v	86 64 Linux		mni task	s•4 on 1 i	nodes	
" #	start ·	07/19/10	/1/.3/.2/		wallcloc	$\mathbf{v} \cdot \mathbf{n} = \mathbf{n} \mathbf{n} \mathbf{n} \mathbf{n}$		
π #	start .	07/19/10	/1/.3/.2/		² COmm	• 11 3 <i>1</i>	1 560	
π #	abytes .		$+00 \pm 0\pm 1$		oflon/se	~ 11.34	-+00 total	
" #	gbytes .	0.000000	100 cocar		grrob/sc			
π #:	#########	* # # # # # # # #	############	+++++++++++++++++++++++++++++++++++++++	***	##########	############	###
" #	region ·	*	[ntasks] =	- н н н н н Д				
#	regron .		[IICGDKD]	1				
#			ſtota	11	<avg></avg>	miı	n ma	ах
#	entries		L	4	1		1	1
#	wallclock		0.0453	3982	0.0113496	0.0091850	0.0131	011
#	user		0.161	974	0.0404935	0.0359	94 0.045	993
#	system		0.123	3979	0.0309947	0.0239	96 0.034	994
#	mpi		0.00594	119	0.0014853	0.00076183	37 0.00176	443
#	°- °⊂omm				11.3372	5.825	77 19.20	098
#	gflop/sec			0	0		0	0
#	gbytes			0	0		0	0
#								
#								
#			[tin	ne]	[calls]	<%mpi	> <%wall2	>
#	MPI_Allred	duce	0.00343	8561	372	57.8	83 7.5	57
#	MPI_Recv		0.00194	1091	558	32.	67 4.2	28
#	MPI_Send		0.000562	2498	558	9.4	47 1.2	24
#	MPI_Comm_s	size	1.21177€	e-06	4	0.0	0.0	0 0
#	MPI_Comm_	rank	9.65199€	e-07	4	0.0	0.0	00 🦯
U.S.	<u>╇</u> ₩₽₩₽₽₽₽₽₽₽₽	Offfide# # f# # #	###########	+++++++++++++++++++++++++++++++++++++++	+++++++++++++++++++++++++++++++++++++++	##########	# # # # # # # # # # # # # # #	####
En	IERGY	Science						BERKE











An Unbalanced Code









A Balanced Code













- 1872.82608032 MB
- 1123.69564819 MB
- 749.13043213 MB
- 374.56521606 MB
- □ 0.00000000 MB





Low Degree Regular Mesh **Communication Patterns** ERSC Cactus Point-to-Point Communication (bytes) MADbenchSG Point-to-Point Communication (byte LBMHD2D Point-to-Point Communication (bytes) 250 250 250 2.5e+08 200 200 200 CONTRACT OF 2.0e+08 STREET, ST 150 150 150 Processor Processor 1.5e+08 100 100 100 California and 1.0e+08 50 50 5.0e+07 0.0e+00 50 100 150 200 25 250 50 100 150 200 50 100 150 200 2 Processor Processor Processor GTC Point-to-Point Communication (bytes) FVCAM1D Point-to-Point Communication (bytes) 250 60 1.4e+09 200 50 1.2e+09 1.0e+09 40 150 Processor 8 Processor 8.0e+08 100 6.0e+08 20 4.0e+08 50 10 2.0e+08 0.0e+00 20 30 40 50 60 50 100 150 200 250 10

Processor



Processor

Processor





Parallel Programming in the future with many-core architectures.

Can MPI everywhere survive?

Slides for John Shalf's Future Architecture's Talk







Basic Multi-core Compute Node Architecture





U.S. DEPARTMENT OF

Coming Soon



Trend #3: Multicore / Manycore

- Power density limit single
 processor clock speeds
- Cores per chip is growing
- Simple doubling of cores is not enough to reach exascale
 - Also a problem in data centers, laptops, etc.
- Two paths to exascale:
 - Accelerators (GPUs)
 - Low power embedded cores
 - (Not x86 clusters)





NATIONAL ENERGY RESEARCH SCIENTIFIC COMPUTING CENTER

What's Wrong with MPI Everywhere

- We can run 1 MPI process per core (flat model for parallelism)
 - This works now and will work for a while
 - But this is wasteful of intra-chip latency and bandwidth (100x lower latency and 100x higher bandwidth on chip than off-chip)
 - Model has diverged from reality (the machine is NOT flat)
- How long will it continue working?
 - 4 8 cores? Probably. 128 1024 cores? Probably not.
 - Depends on performance expectations
- What is the problem?
 - Latency: some copying required by semantics
 - Memory utilization: partitioning data for separate address space requires some replication
 - How big is your per core subgrid? At 10x10x10, over 1/2 of the points are surface points, probably replicated
 - Memory bandwidth: extra state means extra bandwidth
 - Weak scaling: success model for the "cluster era;" will not be for the many core era -- not enough memory per core
 - Heterogeneity: MPI per CUDA thread-block?



SCIENTIFIC COMPUTING CENTER





However: MPI will likely persist

- Obviously MPI will not disappear in five years
- By 2014 there will be 20 years of legacy software in MPI
- Thus far, new systems are not sufficiently different to lead to new programming model
- MPI can evolve (like Fortran, the Fortran from 50 years ago is very different from the Fortran used today)





Why use Hierarchical Merse Hierarchical Matter (hybrid) model for parallelism?

- The machine is not flat
 - We lose a lot of performance by lying to ourselves
- Target: Get Strong scaling on-chip and weak-scaling
 off-chip
 - 100x higher bandwidth between cores on chip
 - 100x lower latency between cores on chip
 - If you pretend that every core is a peer (each is just a generic MPI rank) you are leaving a lot of performance on the table
 - You cannot domain-decompose things forever (cannot weak-scale forever)
- Potentially MPI between nodes and X within node
 - Where X could be OpenMP, UPC, OpenCL, CUDA, etc...







What is X?

• X is it OpenMP?

- Lots of synchronization
- Poor expression of locality (will not scale)
- X might be UPC or PGAS language
 - Explicit definition of local vs. remote
 - Very lightweight communication
- X might be CUDA or OpenCL
 - OpenCL is very CUDA-like cross-platform extension to C language
 - CUDA is also being extended to also taret multicore

• For all X

- Define better way to express fine-grained parallelism on-chip
- must rigorously determine semantics for interoperation with MPI

Scaling



MPI+X: Requirements for X

- Must be able to write once and run everywhere
 - Cannot develop architecture-specific code
 - Don't want to write code for each target! (just once please)
- Needs to be ubiquitous
 - Most people start a new code on a laptop and graduate to HPC systems
 - The complete development environment must be in both places (freely available)
- Must emphasize ability to deliver strong-scaling on-chip to replace clock-frequency scaling
 - Data parallelism might not be sufficient
 - We cannot rely on domain-decomposition for speed-up ad-infinitum (nothing to take up slack for CFL)

— Functional partitioning (Happening at macro-scale with ENERCY angeworks At micro-scale, requires bounded side-effects! its not magic)





Summary

Strong scaling on chip

- Memory is shrinking per chip and clocks stalled
- Solutions: UPC on-chip, OpenCL, domain-specific codegeneration
- Not-solutions: CUDA, OpenMP (not locality aware)
- Weak scaling between chips
 - Memory size is staying same per node
 - Probably MPI, but could be UPC, PGAS or other distributed memory locality aware models
- Frameworks for managing big programming teams
 - Should focus on modularity and agreement on interfaces
 - Benefits from functional semantics
- Languages for fine-grained parallelism + correctness
 - Defining exec model for fine-grained explicit parallelism is the challenge of our decade







More Info

- The Berkeley View/Parlab
 - <u>http://view.eecs.berkeley.edu</u>
 - <u>http://parlab.eecs.berkeley.edu/</u>
- NERSC System Architecture Group
 - <u>http://www.nersc.gov/projects/SDSA</u>
- LBNL Future Technologies Group
 <u>http://crd.lbl.gov/ftg</u>















Cover Stories from NERSC Research



NERSC is enabling new high quality science across disciplines, with over *1,600 refereed publications* last year



