

PFLOTRAN Case Study

Brian Wylie & Markus Geimer, Jülich Supercomputing Centre Tobias Hilbrich & Matthias Lieber, Technische Universität Dresden Chee Wai Lee & Wyatt Spear, University of Oregon

supported by the PFLOTRAN & VI-HPS development teams







RHEINISCH-WESTFÄLISCHE TECHNISCHE HOCHSCHULE





UNIVERSITY OF OREGON









- PFLOTRAN application
 - characteristics, scalability (BG/P & XT5)
- Scalasca
 - summary, summary+PAPI & trace analyses
- Vampir
 - timeline visualization, clustering, communication matrix, etc.
- TAU
 - PDT selective instrumentation
 - ParaProf experiment manager, PerfExplorer,
 2D/3D profiles, histograms, callgraph, etc.

PFLOTRAN

- 3D reservoir simulator developed by LANL/ORNL/PNNL
 - http://ees.lanl.gov/source/orgs/ees/pflotran/
 - approx. 80 thousand lines of Fortran90, combining
 - PFLOW non-isothermal, multi-phase groundwater flow solver
 - PTRAN reactive, multi-component contaminant transport solver
 - employs PETSc, LAPACK, BLAS & HDF5 I/O libraries
 - 87 PFLOTRAN source files (72 modules) + 789 PETSc
 - run with "2B" input dataset for 10 timesteps
 - most of run time in initialization phase, typically amortized
 - uses 3-dimensional (non-MPI) PETSc Cartesian grid
 - alternating FLOW & TRAN(sport) steps in each timestep
 - scaled on Jugene BG/P to 64k, Jaguar XT5 to 32k
 - TRAN(sport) step scaled much better than flow step
 - FLOW step generally faster, but crossover at larger scale

PFLOTRAN (2B.10ts) scalability





- Demonstrate performance measurement and analysis of PFLOTRAN using "2B" problem dataset and more than 10 thousand MPI processes
 - IBM BG/P (jugene.fz-juelich.de)
 - Cray XT5 (jaguar.nccs.ornl.gov)
- Challenge issued for Dagstuhl seminar 10181 on "Program development for extreme-scale computing" (3-7 May 2010)
 - two months notice with download/build/run instructions
 - http://www.dagstuhl.de/10181

PFLOTRAN "2B" test case





Simulation of U(VI) Migration

- Hanford 300 area
- 1-year simulation:
 - 900x1300x20m
 - $-\Delta x/\Delta y = 5 m$
 - 1.87M grid cells
 - 15 chemical species
 - 28M DoF total
- 1-year simulation:
 - $-\Delta t = 1$ hour
 - 5-10 hour runtime on Cray XT5 (4k cores)



- Automatic application instrumentation
 - both PFLOTRAN application (Fortran) & PETSc library (C)
 - USR routines instrumented by IBM XL & Cray/PGI compilers
 - MPI routines via interposition of instrumented library (PMPI)
- Initial (small-scale) summary measurements used to define filter files specifying all purely computational routines
 - distinct filters required for IBM XL and PGI compilers
- Summary & trace experiments collected using filters
- Post-processing of analysis reports
 - cut to extract timestep loop; incorporation of 3D application topology
- Analysis report examination in GUI

- Determined by scoring summary expt using fully-instrumented executable
 - single timestep measurement sufficient



- 1146 PFLOTRAN+PETSc routines executed
 - plus 29 MPI library routines
 - 856 not on a callpath to MPI, purely local calculation (USR)
 - 291 on callpaths to MPI, mixed calculation & comm. (COM)
- Using measurement filter listing all USR routines
 - maximum callpath depth 22 frames
 - 1732 unique callpaths (399 in FLOW, 375 in TRAN)
 - 633 MPI callpaths (121 in FLOW, 114 in TRAN)
- Of 399 FLOW callpaths, 40 missing from TRAN, 20 only in TRAN, 14 with "similar" names
 - richardsjacobian vs rtjacobian, _MPIAIJ vs _MPIBAIJ, etc.



Scalasca summary analysis: entire program



Scalasca summary analysis: stepperrun



Scalasca analysis: excl. Execution time



Scalasca analysis: floating-point operations



Scalasca trace analysis: MPI waiting time



11.

Scalasca scalability analysis



MPI collective communication becomes a bottleneck

Scalasca scalability analysis

1.3

1.2

1.1

1.0

0.9

0.8

0.7

0.6

0.5

0.4

0.3

0.2

0.1

0.0

2048

4096

8192

32768

16384 Processes 65536

131072

Proportion of phase time



- Proportion of phase time for calculation diminishes, as cost of communication grows
- Only collective communication cost becomes significant, exceeding calculation times

Scalasca scalability analysis



- Collect traces with up to 16,384 MPI processes on BG/P for parallel analysis with Vampir7
- Use of manual instrumentation API for selective tracing
 - disable measurement during initialization phase and only enable trace of one timestep
 - reduces traces to a managable size
- Interactive exploratory analysis of traces
 - cluster similar processes
 - zoom into time interval of interest
 - scroll/zoom timelines for processes
 - investigate communication patterns
 - examine message distributions

Full execution time-line





Timestep 7 zoom





PFLOW phase zoom





End of init phase + 1st timestep





Clustering of processes





Comparing process timelines





1st timestep: imbalance





1st timestep: imbalance (zoom)



Comm matrix + message hist







- Use PDT for selective instrumentation
 - Select user source routines taking more than 1% of exclusive execution time
 - 4 from PFLOTRAN, 19 from PETSc [+44 MPI]
- Collect profiles on Jaguar Cray XT5 from runs with varying numbers of MPI processes
 - Include PAPI preset hardware counters
- Use ParaProf Manager to browse experiments
 - Examine 2D & 3D graphical profile, histogram and callgraph presentations

ParaProf Manager (16380)



TAU: ParaPro	_ + X	
File Options Help		
Applications	TrialField	Value
- T Standard Applications	Name	pflotran_xt5_16380_reduced.xml
• 🗖 Default App	Application ID	0
• • • Persuit Evn	Experiment ID	0
Genetic Exp Second se	Trial ID	0
	CPU Cores	6
	CPU MHz	2600.000
PARI_FF_OFS	CPU Type	6-Core AMD Opteron(tm) Processor 23
	CPU Vendor	AuthenticAMD
- PAPI L1 DCM	CWD	/lustre/widow1/scratch/wspear/pflotran
- PAPI_RES_STL	Cache Size	512 KB
- 🕒 PAPI_TOT_CYC	Executable	/var/spool/alps/1887061/pflotran
- 🕒 PAPI_L2_TCA	File Type Index	10
- PAPI_L2_TCM	File Type Name	TAU Snapshot
🔶 🚍 Default (jdbc:derby:/home/users/wspear/tau2/x86_64/lib/perfdmf)	Hostname	nid10413
Fill peri_s3d (idbc:postgresgl://apollo.cs.uoregon.edu:5432/peri_s3d)	Local Time	2010-05-03T04:49:09-04:00
Control scott (idbc:derby/home/users/scotth/ ParaProf/nerfdmf)	MPI Processor Name	nid10413
C neri pflotran (idbc:nostgresci) (canollo cs. uoregan edu;5422 (neri pflotran))	Memory Size	16385788 kB
pen_photran (dubc.postgresql.//apolio.cs.doregon.edd.5452/pen_photran)	Node Name	nid10413
	OS Machine	x86_64
	OS Name	Linux
	OS Release	2.6.16.60-0.39_1.0102.4787.2.2.41
	OS Version	#1 SMP Thu Nov 12 17:58:04 CST 2009
	Starting Timestamp	1272875793196482
	TAU Architecture	crayoni
	TAU Config	-pdt=/ccs/home/wspear/bin/pdtoolkit
	TAU Makefile	/ccs/home/wspear/bin/tau2/craycnl/lib
	TAU MetaData Merge Time	0.000305 seconds
	TAU Profile Merge Time	1.208 seconds
	TAU Unification Time	0.001299 seconds
	TAU Version	2.18.2-cvs
	TAU_CALLPATH	off
	TAU_CALLPATH_DEPTH	2
	TAU_COMM_MATRIX	off
	TAU_COMPENSATE	off
	TAU_PROFILE	on
	TAU_PROFILE_FORMAT	merged



ParaProf full profile (16380)





ParaProf 3D full profile (16380)



3D profile (w/o MPI_Allreduce)



ParaProf histogram (16380)





ParaProf call graph





ParaProf mean profile (8184)



Metric: TIME Value: Exclusive Units: seconds

273.947	MPL Alreduce()
85.232	Petc_ErrorCode MatSolve_SegBAIJ_N(Mat, Vec, Vec) {{baijfact2.c} {2055,1}-{2108,1}]
81.019	PetscErrorCode oursnesjacobian(SNES, Vec, Mat *, Mat *, MatStructure *, void *) {{zsnesf.c} {86,1}-{91,1}}
79.55	PetscErrorCode MatLUFactorNumeric_SegBAI_N(Mat, Mat, const MatFactorInfo *) [{baijfact4.c} {12,1}-{87,1}]
69.588	REACTION_MODULE::RTOTAL [{reaction.pp.F90} {4058,12}]
62.672	PetscErrorCode MatMult SedBAIL N(Mat. Vec. Vec) [{bail2.c} {601.1}-{652.1}]
58,906	REACTION MODULE: RMULTIRATESORPTION [{reaction.pp.F90} {45 15 .12}]
43.953	PetscErrorCode oursnesfunction(SNES, Vec. Vec. void *) {zsnesf.c} {62, 1}-{67, 1}
41.82	PetscErrorCode PetscMallocValidate(int. const char *, const char *, const char *) C [(mtr.c) {125,1}-{159,1}]
29.113	PetscErrorCode MatSolve SegAll NaturalOrdering(Mat. Vec. Vec) [{aiifact.c} {1134.1}-{1186.1}]
25.803	PetscErrorCode MatDiagonalScale SegBAI(Mat. Vec. Vec) [(bail2.c) {1640.1)-{1694.1}]
24.545 🥅	PetscErrorCode MatMult_SegAJ(Mat, Vec, Vec) [{aij.c} {973, 1}-{1029, 1}]
23.035 🥅	MPI_Start()
18.161	MPI_Waitany()
17.19 🗖	PetscErrorCode MatZeroEntries_SegBAI((Mat) [{baij2.c} {1726,1}-{1734,1}]
7.631 🚪	INPUT_MODULE::INPUTCREATE [{input.pp.F90} {1705,10}]
6.168	PetscErrorCode VecAXPBYPCZ_Seq(Vec, PetscScalar, PetscScalar, PetscScalar, Vec, Vec) [{bvec1.c} {194,1}-{225,1}]
6.019	PetscErrorCode VecWAXPY_Seq(Vec, PetscScalar, Vec, Vec) [{dvec2.c} {718,1}-{750,1}]
5.898 🛽	MPI_Waitall()
4.386	TIMESTEPPER_MODULE::STEPPERUPDATETRANSPORTSOLUTION [{timestepper.pp.F90} {6840,12}]
4.341	MPI_Bcast()
3.262	PetscErrorCode VecNorm_MPI(Vec, NormType, PetscReal *) [{pvec2.c} {49,1}-{105,1}]
3.256	MPI_Startall()
2.982	PetscErrorCode MatMultAdd_SeqBAJJ_N(Mat, Vec, Vec, Vec) [{baij2.c} {1116,1}-{1169,1}]
2.921	PetscErrorCode VecDot_Seq(Vec, Vec, PetscScalar *) [{bvec1.c} {13,1}-{45,1}]
2.761	PetscErrorCode VecSet_Seq(Vec, PetscScalar) [{dvec2.c} {560,1}-{574,1}]
2.635	PetscErrorCode MatMultAdd_SeqAIJ(Mat, Vec, Vec, Vec) [{aij.c} {1034,1}-{1096,1}]
2.602	MPI_Barrier()
2.319	PetscErrorCode PetscTrFreeDefault(void *, int, const char *, const char *, const char *) [{mtr.c} {259,1}-{328,1}]
2.29	PetscErrorCode VecScatterBegin(VecScatter, Vec, Vec, InsertMode, ScatterMode) C [{vscat.c} {1509,1}-(1553,1}]
2.203	TIMESTEPPER_MODULE::STEPPERSTEPTRANSPORTDT [{timestepper.pp.F90} {3761,12}]
1.844	PetscErrorCode KSPSolve(KSP, Vec, Vec) C [{itfunc.c} {298,1}-{594,1}]
1.807	MPI_Comm_dup()
1.71	PetscErrorCode VecScatterEnd(VecScatter, Vec, Vec, InsertMode, ScatterMode) C [{vscat.c} {1581,1}–{1599,1}]
1.387	MPI_scatterv()
1.279	PetscErrorCode VecCopy_Seq(Vec, Vec) {{bvec1.c}{104,1}-{11/,1}}
1.14	PetscErrorCode VecDotNorm2(Vec, Vec, PetscScalar *, PetscScalar *) C [{vinv.c} {1197,1}-(1227,1]
0.699	PetscErrorCode PetschreeAlign(vold ", Int, const char ", const char ", const char ") [[mai.c] {58,1]-[81,1]]
0.619	Void Scatter_I(Perscint, const Perscint ", const Perscicaar ", const Perscicaar ", insertMode) [{vpscat.c} {523,1}-{548,1}]
0.609	Petscerror.ode MatDiagonalscaleLocal_MPIBAJ(Mat, Vec) C [mmbal].C} {506,1}-{340,1}
0.568	REACTION_MODULE: KKINETICMINERAL [reaction, pp. F90] (4904, 12)]
0.558	INPUT_MODULE::INPUT_READFLOT_RANSTRINGSLAVE [[Input_pp./s0][2039.12]]
0.555	
0.538	GLOBAL_MODOLEGLOBALUPDATEDENANDSATRATCH [gjj00ai.pp.F90}[2350,12]]
0.492	reisterror code maimmiskajoragonaistajetodaisetopi(mai, Vec) ((mindai).c) (222),1)+(286,1))
0.438	Timesterrer, module, sterrersterredmet (timestepper.pp.1903/2546,12))
0.426 1	volu nauki tuneusuhu, tuhsu metsuhuri, tuhsu metsusualahiri, metsusualahiri) kVDscat.c)(489,1+495,1))

ParaProf 3D correlation cube



PerfExplorer charts







Total TIME Breakdown for pflotran_2b:xt5



Number of Processors

- INPUT_MODULE::INPUTCREATE INPUT_MODULE::INPUTDESTROY MPI_Allreduce() MPI_Bcast() MPI_Start()
- MPI_Waitany() PetscErrorCode MatDiagonalScale_SeqBAIJ() PetscErrorCode MatLUFactorNumeric_SeqBAIJ_N()
- PetscErrorCode MatMult_SeqAIJ() PetscErrorCode MatMult_SeqBAIJ_N()
- PetscErrorCode MatSolve_SeqAJ_NaturalOrdering() PetscErrorCode MatSolve_SeqBAJ_N()
- PetscErrorCode MatZeroEntries_SeqBAIJ() PetscErrorCode oursnesfunction() PetscErrorCode oursnesjacobian()
- PetscErrorCode PetscMallocValidate() REACTION_MODULE::RMULTIRATESORPTION
- REACTION_MODULE::RTOTAL _ other

PerfExplorer aligned bar chart





PetscErrorCode MatMult_SeqBAIJ_N() = PetscErrorCode MatSolve_SeqAIJ_NaturalOrdering() = PetscErrorCode MatSolve_SeqBAIJ_N()

PetscErrorCode MatZeroEntries_SeqBAIJ() = PetscErrorCode PetscMallocValidate() = PetscErrorCode oursnesfunction()

PetscErrorCode oursnesjacobian() = REACTION_MODULE::RMULTIRATESORPTION = REACTION_MODULE::RTOTAL = other

ParaProf Manager (131040)



TAU:	ParaProf Manager	_ + ×
File Options Help		
Applications	TrialField	Value
👇 🔚 Standard Applications	Name	pflotran_xt5_131040_reduced.xml
🔶 🚍 Default App	Application ID	0
- C Default Exp	Experiment ID	0
Image: a contract of the second s	Trial ID	0
	CPU Cores	6
	CPU MHz	2600.000
- PAPL TOT INS	CPU Type	6-Core AMD Opteron(tm) Processor 23 (D0)
- PAPI L1 DCA	CPU Vendor	AuthenticAMD
- PAPI_L1_DCM	CWD	/lustre/widow1/scratch/wspear/pflotran_2b_hug
- PAPI_RES_STL	Cache Size	512 KB
- PAPI_TOT_CYC	Executable	/var/spool/alps/1891721/pflotran
- PAPI_L2_TCA	File Type Index	10
- PAPI_L2_TCM	File Type Name	TAU Snapshot
🔶 🔚 Default (jdbc:derby:/home/users/wspear/tau	Hostname	nid11823
🔶 🔚 peri_s3d (jdbc:postgresgl://apollo.cs.uoregon	Local Time	2010-05-03T20:45:44-04:00
Image: Scott (idbc:derby/home/users/scottb/ ParaPrise)	MPI Processor Name	nid11823
- T peri pflotran (idbc:postgresgl: (/apollo cs.upre	Memory Size	16385788 kB
	Node Name	nid11823
	OS Machine	x86_64
	OS Name	Linux
	OS Release	2.6.16.60-0.39_1.0102.4787.2.2.41-cnl
	OS Version	#1 SMP Thu Nov 12 17:58:04 CST 2009
	Starting Timestamp	1272932052236170
	TAU Architecture	craycni
	TAU Config	-pdt=/ccs/home/wspear/bin/pdtoolkit/ -pdt_c
	TAU Makefile	/ccs/home/wspear/bin/tau2/craycnl/lib/Makefil
	TAU MetaData Merge Time	0.00041 seconds
	TAU Profile Merge Time	12.96 seconds
	TAU Unification Time	0.02815 seconds
	TAU Version	2.18.2-cvs
	TAU_CALLPATH	off
	TAU_CALLPATH_DEPTH	2
	TAU_COMM_MATRIX	off
	TAU_COMPENSATE	off
	TAU_PROFILE	on
	TAU_PROFILE_FORMAT	merged
	TAU_SAMPLING	off
	TAU_THROTTLE	on
	TAU_THROTTLE_NUMCALLS	100000
	TAU_THROTTLE_PERCALL	10
	TAU_TRACE	off
	TAU_TRACK_HEADROOM	off
	TAU_TRACK_HEAP	off
	TAU_TRACK_MESSAGE	off
	Timestamp	1272933945055505
	UTC Time	2010-05-04T00:45:44Z
	pid	32472
	username	wspear

ParaProf full profile (131040)



ParaProf histogram (131040)





- Duplication of MPI Communicators by HDF5 dominates initialization at scale (particularly on XT5)
- MPI_Allreduce collective communication remains a severe bottleneck in the timestep loop
 - for both FLOW & TRAN phases
 - particularly MatZeroRowsLocal/PetscMaxSum
- Computational imbalance from inactive grid cells evident for processes at one end of 3D process grid
 - widespread in the computational routines
 - since only a minority of processes, little to gain

- Performance analysis of complex applications at largescale requires care
 - full automatic instrumentation is convenient, but may produce more detail than desirable
 - selective measurement and/or instrumentation may be used to reduce overhead and size of event traces
 - even basic reduced execution profiles for many thousands of processes rapidly become awkwardly large
 - powerful analyses and interactive customizable visualizations required for an effective initial overview leading to in-depth refinement of performance issues
- VI-HPS experts are available to assist

- Applications can be prepared for measurement using manual, selective and automatic instrumentation and PAPI counters
- Scalasca runtime summary & automatic event trace analyses quantify and isolate locations of MPI communication & synchronization overheads
- Vampir visual trace analysis supports interactive exploration of detailed process interaction and clustering of similar traces
- TAU experiment management, graphical profile displays and extensive instrumentation options facilitate customization of measurements and analysis presentations
- VI-HPS tools can be used independently, however, maximum benefit offered by exploiting their complementary integrated capabilities



- PFLOTRAN developers for providing an interesting large-scale test case and assistance building/running their code
- Dagstuhl seminar organizers for setting the tools challenge and hosting the associated discussions
- NIC/JSC & NCCS/ORNL for use of compute-time allocations on their large-scale facilities
- Support and development teams for the various tools for their prompt support