

# Debuggers and Performance Tools

February 2013 | Markus Geimer

#### Outline





#### UNITE



- UNiform Integrated Tool Environment
- Standardizes tool access and documentation
  - Currently in use at JSC, RWTH, ZIH
- Based on "module" command
  - Standardized tool and version identification
    - <tool>/<version>-<special>
    - <special>: optional indicator if tool is specific for a MPI library, compiler, or 32/64 bit mode
- Tools only visible after
  - module load UNITE

**# once per session** 

- Basic usage and pointer to tool documentation via
  - module help <tool>

#### Example



```
% module load UNITE
UNITE loaded
% module help scalasca
Module Specific Help for scalasca/1.4.2:
Scalasca: Scalable Performance Analysis of Large-Scale
          Parallel Applications
Version 1.4.2
Basic usage:
1. Instrument application with skin
2. Collect & analyze execution measurement with scan
3. Examine analysis results with square
For more information:
- See ${SCALASCA_ROOT}/doc/manuals/QuickReference.pdf
  or type "scalasca -h"
- http://www.scalasca.org
- mailto:scalasca@fz-juelich.de
```

#### **Documentation**



- Use "module avail" to check latest status
- Websites
  - http://www.fz-juelich.de/ias/jsc/juqueen/
    - User Info
      - Debugging
      - Performance Analysis (<sup>(</sup>A))
  - http://www.vi-hps.org/training/material/
    - Performance Tools LiveDVD image
    - Links to tool websites and documentation
    - Tutorial slides





# **Debugging on JUQUEEN**

February 2013 | Alexandre Strube

**TOTALVIEW** Parallel Debugger



- UNIX Symbolic Debugger for C, C++, f77, f90, PGI HPF, assembler programs
- "Standard" debugger
- Special, non-traditional features
  - Multi-process and multi-threaded
  - C++ support (templates, inheritance, inline functions)
  - F90 support (user types, pointers, modules)
  - ID + 2D Array Data visualization
  - Support for parallel debugging (MPI: automatic attach, message queues, OpenMP, pthreads)
  - Scripting and batch debugging
  - Memory Debugging
- http://www.roguewave.com

#### **Debugger Setup**



- Compile and link your program with debug option: -g
- Use absolute paths for source code info: -qfullpath
- In case of optimized codes (XL), keep function call parameters: -qkeepparm
- Load modules

```
% ssh <u>-X</u> user@juqueen
[...]
juqueen% module load UNITE totalview
UNITE loaded
Totalview/8.11.0-0 loaded
juqueen% mpixlcxx hello.cpp -qfullpath -qkeepparm -g -o helloworld
juqueen%
```

#### **Debugger Startup**



- Interactively: call the lltv script
  - Creates a LoadLeveler batch script with required TotalView parameters
- If user cancels the script, it cancels the debugging job (does not eat your computing quota)
- DON'T attach to all ranks! This will be VERY slow.

#### **Launch Script**



juqueen% lltv -n <nodes> : -default\_parallel\_attach\_subset= <rank-range> runjob -a --exe <program> -p <num>

- Starts <program> with <nodes> and <num> processes, attaches to <rank-range>:
  - Rank: that rank only
  - RankX-RankZ: all ranks, both inclusive
  - RankX-RankZ:stride every strideth between RankX and RankZ

Example:

```
juqueen% lltv -n 2 : -default_parallel_attach_subset= \
2-6 runjob -a --exe helloworld -p 64
```

```
Creating LoadLeveler Job
Submitting LoadLeveler Interactive Job for Totalview
Wait for job juqueen1c1.32768.0 to be started:.....
```

#### JÜLICH FORSCHUNGSZENTRUM

#### Execution

- Totalview tries to debug "runjob" and shows no source code
  - Ignore it and press "GO"
- After some seconds, TotalView will detect parallel execution and ask if it should stop. Yes, it should stop.
- To find the correct point file/function to debug, use the "File-Open" command.
- Set your breakpoints, and press "GO" again. Debugging session will then start.
- To see a variable's contents, double click on it in the source.

#### **TotalView: Main Window**





#### **Totalview: Tools Menu**



Call Graph



#### Data visualization





# Performance Analysis Tools on JUQUEEN

February 2013 | Markus Geimer

#### **Typical Performance Analysis Procedure**



- Do I have a performance problem at all?
  - Time / speedup / scalability measurements
- What is the key bottleneck (computation / communication)?
  - MPI / OpenMP / flat profiling
- Where is the key bottleneck?
  - Call-path profiling, detailed basic block profiling
- Why is it there?
  - Hardware counter analysis
  - Trace selected parts (to keep trace size manageable)
- Does the code have scalability problems?
  - Load imbalance analysis, compare profiles at various sizes function-by-function

### **Remark: No Single Solution is Sufficient!**





A combination of different methods, tools and techniques is typically needed!

- Analysis
  - Statistics, visualization, automatic analysis, data mining, ...
- Measurement
  - Sampling / instrumentation, profiling / tracing, ...
- Instrumentation
  - Source code / binary, manual / automatic, ...

#### **Critical Issues**



- Accuracy
  - Intrusion overhead
    - Measurement itself needs time and thus lowers performance
  - Perturbation
    - Measurement alters program behavior
    - E.g., memory access pattern
  - Accuracy of timers & counters
- Granularity
  - How many measurements?
  - How much information / processing during each measurement?

Tradeoff: Accuracy vs. Expressiveness of data

# scalasca 🗖



- Scalable Analysis of Large Scale Applications
- Approach
  - Instrument C, C++, and Fortran parallel applications
    - Based on MPI, OpenMP, SHMEM, or hybrid
  - Option 1: scalable call-path profiling
  - Option 2: scalable event trace analysis
    - Collect event traces
    - Search trace for event patterns representing inefficiencies
    - Categorize and rank inefficiencies found
- http://www.scalasca.org

#### What is the Key Bottleneck?



- Generate flat MPI profile using Scalasca
  - Only requires re-linking
  - Low runtime overhead
- Provides detailed information on MPI usage
  - How much time is spent in which operation?
  - How often is each operation called?
  - How much data was transferred?
- Limitations:
  - Computation on non-master threads and outside of MPI\_Init/MPI\_Finalize scope ignored

#### Flat MPI Profile: Recipe



- Prefix your *link command* with "scalasca -instrument -comp=none"
- 2. Prefix your MPI *launch command* with "scalasca -analyze"
- 3. After execution, examine analysis results using "scalasca -examine epik\_<*title*>"

#### **Flat MPI Profile: Example**



module load UNITE scalasca
scalasca -analyze \
 runjob --ranks-per-node p --np n [...] --exe ./myprog

juqueen% scalasca -examine epik\_myprog\_*n*x*t*\_sum

## Flat MPI Profile: Example (cont.)





#### Where is the Key Bottleneck?



- Generate call-path profile using Scalasca
  - Requires re-compilation
  - Runtime overhead depends on application characteristics
  - Typically needs some care setting up a good measurement configuration
    - Filtering
    - Selective instrumentation
- Option 1 (recommended): Automatic compiler-based instrumentation
- Option 2: Manual instrumentation of interesting phases, routines, loops

#### **Call-path Profile: Recipe**



- Prefix your *compile & link commands* with "scalasca -instrument"
- 2. Prefix your MPI *launch command* with "scalasca -analyze"
- 3. After execution, compare overall runtime with uninstrumented run to determine overhead
- 4. If overhead is too high
  - Score measurement using "scalasca -examine -s epik\_<title>"
  - 2. Prepare filter file
  - 3. Re-run measurement with filter applied using prefix "scalasca -analyze -f <filter\_file>"
- 5. After execution, examine analysis results using "scalasca -examine epik\_<*title*>"

#### **Call-path Profile: Example**



```
juqueen% module load UNITE scalasca
juqueen% scalasca -instrument mpix1f90 -03 -c foo.f90
juqueen% scalasca -instrument mpix1f90 -03 -c bar.f90
juqueen% scalasca -instrument \
          mpixlf90 -03 -o myprog foo.o bar.o
In the job script: ##
##
module load UNITE scalasca
scalasca -analyze \setminus
   runjob --ranks-per-node p --np n [...] --exe ./myprog
```



juqueen% scalasca -examine -s epik\_myprog\_nxt\_sum cube3\_score -r ./epik\_myprog\_nxt\_sum/summary.cube Reading ./epik\_myprog\_nxt\_sum/summary.cube... done. Est. aggregate size of event trace (total\_tbc): 160,338,400,040 bytes Est. size of largest thread trace (max\_tbc): 133,910,372 bytes (When tracing set ELG\_BUFFER\_SIZE to avoid intermediate flushes or reduce requirements using filter file listing names of USR regions.)

INFO: Score report written to ./epik\_myprog\_nxt\_sum/epik.score

- Region/call-path classification
  - MPI (pure MPI library functions)
  - OMP (pure OpenMP functions/regions)
  - USR (user-level source local computation)
  - COM ("combined" USR + OpeMP/MPI)
  - ANY/ALL (aggregate of all region types)





juqueen%	less epik_m	iyprog_ <i>n</i> x <i>t</i> _	_sum/epik	.score	
flt type	max_tbc	time	%	region	
ANY	133910372	49656.65	100.00	(summary) ALL	
MPI	483104	6575.16	13.24	(summary) MPI	
OMP	1620780	27541.42	55.46	(summary) OMP	
COM	343320	4442.51	8.95	(summary) COM	
USR	131511360	10628.25	21.40	(summary) USR	
USR	42219648	4664.19	9.39	binvcrhs	
USR	42219648	2994.22	6.03	matmul_sub	
USR	42219648	2519.26	5.07	matvec_sub	
USR	1891008	156.64	0.32	binvrhs	
USR	1891008	210.20	0.42	lhsinit	
USR	1033416	76.88	0.15	exact_solution	
MPI	201000	37.20	0.07	MPI_Isend	
MPI	184920	33.19	0.07	MPI_Irecv	
MPI	96480	6362.37	6.25	MPI_Waitall	
COM	96480	1330.77	2.68	copy_x_face	
COM	96480	1331.61	2.68	copy_y_face	
OMP	88440	533.11	1.07	!\$omp parallel @foo.f90	$\square$
[]					

#### **Call-path Profile: Filtering**



- In this example, the 6 most fequently called routines are of type USR (max\_tbc is proportional to visit count)
- These routines contribute around 21% of total time
  - However, much of that is most likely measurement overhead for a few frequently-executed small routines
- Avoid measurements to reduce the overhead
- List routines to be filtered in simple text file

```
juqueen% cat filter.txt
binvcrhs
matmul_sub
matvec_sub
binvrhs
lhsinit
exact_solution
```



```
## To verify effect of filter:
juqueen% scalasca -examine -s -f filter.txt \
          epik_myprog_nxt_sum
In the job script:
##
                  ##
module load UNITE scalasca
scalasca -analyze -f filter.txt \
   runjob --ranks-per-node p --np n [...] --exe ./myprog
## After job finished:
                  ##
```

juqueen% scalasca -examine epik\_myprog\_nxt\_sum



Cube 3.0 QT: epik_bt-mz_C_2p64x32_sum/summary.cube.gz (on juqueen1.zam.kfa-juelich.de)										
<u>F</u> ile <u>D</u> isplay <u>T</u> opology <u>H</u> elp										
Absolute	Absolute -			Peer percent						
Metric tree	Call tree Fla	at view		Q Hardv	ware Topology	Threads x I	Processes			
<ul> <li>2.85e4 Time</li> <li>8.77e7 Visits</li> <li>128 Synchronizations</li> <li>3.67e5 Communications</li> <li>9.57e9 Bytes transferred</li> <li>3706.28 Computational imbalance</li> </ul>		bt								
	•		• • •	•		***				
0.00 2.85e4 (100.00%) 2.85e4	0.00	2.85e4 (100.00%)	2.85e4	<b>0.00</b> 0.00	<b>10</b> 2.85e4(	<b>0.00</b> 100.00%)	<b>100.00</b> 2.85e4			









#### Why is the Bottleneck There?



- This is highly application dependent!
- Might require additional measurements
  - Hardware-counter analysis
    - CPU utilization
    - Cache behavior
  - Selective instrumentation
  - Manual/automatic event trace analysis

#### **HW Counter Measurements w/ Scalasca**



- Scalasca supports both PAPI and native counters
- Available counters:

```
juqueen% module load UNITE papi
juqueen% less $PAPI_ROOT/doc/papi-5.0.1-avail.txt
juqueen% less $PAPI_ROOT/doc/papi-5.0.1-native_avail.txt
juqueen% less $PAPI_ROOT/doc/papi-5.0.1-avail-detail.txt
```

Specify using "-m" option of "scalasca -analyze":

```
module load UNITE scalasca
scalasca -analyze -f filter.txt \
    -m PAPI_FP_OPS:PAPI_TOT_INS:PAPI_TOT_CYC \
    runjob --ranks-per-node p --np n [...] --exe ./myprog_
```

#### **Manual Instrumentation w/ Scalasca**



- Can be used to mark initialization, solver & other phases
  - Annotation macros ignored by default
  - Enabled with "-user" flag of "scalasca -instrument"
- Appear as additional regions in analyses
  - Distinguishes performance of important phase from rest

```
#include "epik_user.inc"
subroutine foo(...)
! Declarations
EPIK_USER_REG(solve, "<solver>")
! Some code...
EPIK_USER_START(solve)
do i=1,100
[...]
end do
EPIK_USER_END(solve)
! Some more code...
end subroutine
```

```
void foo(...) {
   /* Declarations */
   EPIK USER REG(solve, "<solver>");
```

```
/* Some code... */
EPIK_USER_START(solve);
for (i=0; i<100; i++) {
   [...]
}
EPIK_USER_END(solve);
/* Some more code... */</pre>
```

#include "epik user.h"

#### Fortran (requires C preprocessor)

```
C / C++
```

#### **Selective Measurement w/ Scalasca**



- Can be used to temporarily disable measurement for certain intervals
  - Annotation macros ignored by default
  - Enabled with "-user" flag of "scalasca -instrument"
- Appear as synthetic PAUSED region in analyses

<pre>#include "epik_user.inc"</pre>
subroutine foo() ! Some code
EPIK_PAUSE_START()
! Loop will not be measured
do i=1,100
[]
end do
EPIK PAUSE END()
! Some more code
end subroutine

```
#include "epik_user.h"
void foo(...) {
   /* Some code... */
   EPIK_PAUSE_START();
   /* Loop will not be measured */
   for (i=0; i<100; i++) {
      [...]
   }
   EPIK_PAUSE_END();
   /* Some more code... */
}</pre>
```

Fortran (requires C preprocessor)

C / C++

#### **Scalasca API Limitations**



- START and END calls must be correctly nested
  - That is, you cannot start a region in one routine and end it in another
- Measurement control API calls are not allowed within OpenMP parallel regions
- Behavior of automatic trace analysis is undefined when PAUSING skips recording MPI events on subsets of processes
  - Examples:
    - A collective operation is excluded only on some ranks involved
    - A send/receive on one rank is excluded w/o the corresponding send/receive on the matching rank
  - Will generally lead to a deadlock

#### **Trace Generation & Analysis w/ Scalasca**



- Identifies inefficiency patterns in communication and synchronization operations
- Specify using "-t" option of "scalasca -analyze":

```
module load UNITE scalasca
scalasca -analyze -f filter.txt -t \
    runjob --ranks-per-node p --np n [...] --exe ./myprog
```

#### ATTENTION:

- Traces can quickly become extremely large!
- Remember to use proper filtering & selective instrumentation
- Ask us for assistance

## TAU

- Very portable tool set for instrumentation, measurementand analysis of parallel multi-threaded applications
- http://tau.uoregon.edu/
- Supports
  - Various profiling modes and tracing
  - Various forms of code instrumentation
  - C, C++, Fortran, Java, Python
  - MPI, multi-threading (OpenMP, Pthreads, ...)





#### **TAU Instrumentation**



- Flexible instrumentation mechanisms at multiple levels
  - Source code
    - manual
    - automatic
      - C, C++, F77/90/95 (Program Database Toolkit (PDT))
      - OpenMP (directive rewriting with Opari)
  - Object code
    - pre-instrumented libraries (e.g., MPI using PMPI)
    - statically-linked and dynamically-loaded (e.g., Python)
  - Executable code
    - dynamic instrumentation (pre-execution) (DynInst)
    - virtual machine instrumentation (e.g., Java using JVMPI)
- Support for performance mapping
- Support for object-oriented and generic programming

#### **TAU: Recipe**



- 1. module load UNITE tau # once per session
- 2. Specify programming model by setting TAU\_MAKEFILE to one of \$TAU\_MF\_DIR/Makefile.tau-\*
  - MPI: Makefile.tau-bgqtimers-papi-mpi-pdt
  - OpenMP/MPI: Makefile.tau-bgqtimers-papi-mpi-pdt-openmp-opari
- 3. Compile and link with
  - tau\_cc.sh file.c ...
  - tau\_cxx.sh file.cxx...
  - tau\_f90.sh file.f90 ...
- 4. Execute with real input data Environment variables control measurement mode
  - TAU\_PROFILE, TAU\_TRACE, TAU\_CALLPATH, ...
- 5. Examine results with paraprof

#### **TAU: Basic Profile View**



0	00		ParaProf: p	prof.dat/ParaProf/Desktop/bertie/Users/
Fil	Options	Windows	Help	
n.c.t, 0,0,0 - 512pro File Options Windows Help COUNTER NAME: P_WALL_CLOCK_TIME ( 345.5474 116.4951 103.2566 59.009 37.4 32.8 21. 11. 11. 11. 11. 11. 11. 11.	/samrai/taudata/neutr aconds) MPI_Allreduce algs::Hyperbo algs::Hyperbo algs::Hyperbo algs::Hyperbo algs::Hyperbo algs::Hyperbo billion algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition algs::Hyperbo condition condition algs::Hyperbo condition algs::Hyperbo condition condition algs::Hyperbo condition algs::Hyperbo condition condition algs::Hyperbo condition	onbackup/rs/sameer/Us () licLevelIntegrator3::a licLevelIntegrator3::a licLevelIntegrator3::a licLevelIntegrator3::a licLevelIntegrator3::c licLevelIntegrator3::c licLevelIntegrator3::e gAlgorithm3::find_bo gAlgorithm3::bdry_fil licLevelIntegrator3::ei licLevelIntegrator3::ei licLevelIntegrator3::ei licLevelIntegrator3::ei licLevelIntegrator3::find_bo	ers/ dvance_bdry_fill_create dvance_bdry_fill_create dvance_bvel() l_new_level_create lance_boxes dvance_bdry_fill_comm nementBoxes() barsen_fluxsum_create barsen_fluxsum_create tes_containing_tags _tags_create ror_bdry_fill_comm intersections_regrid_all l_new_level_comm	

#### **TAU: Callgraph Profile View**







#### **Vampir Event Trace Visualizer**



- Offline trace visualization for
  - VampirTrace OTF trace files
  - Scalasca EPILOG trace files
- Visualization of MPI, OpenMP and application events:



- All diagrams highly customizable (through context menus)
- Large variety of displays for ANY part of the trace
- http://www.vampir.eu
- Advantage:
  - Detailed view of dynamic application behavior
- Disadvantage:
  - Requires event traces (huge amount of data)
  - Completely manual analysis

#### **Vampir Displays**





## **Vampir: Timeline Diagram**



- Functions organized into groups
- coloring by group
- Message lines can be colored by tag or size



 Information about states, messages, collective and I/O operations available through clicking on the representation

February 2013

#### **Vampir: Process and Counter Timelines**



 Process timeline show call stack nesting

 Counter timelines for hardware or software counters



#### **Vampir: Execution Statistics**



- Aggregated profiling information: execution time, number of calls, inclusive/exclusive
- Available for all / any group (activity) or all routines (symbols)



Available for any part of the trace
 selectable through time line diagram

#### **Vampir: Process Summary**



- Execution statistics over all processes for comparison
- Clustering mode available for large process counts

🗮 🗽 🔍 🌕 🔄 🔚 🔚 🖏 🚱 💡 🖋

⊻iew <u>C</u> h	art <u>F</u> ilter			
🗮 😽	1 📶 🌀 🇧	7 🔝 📷 🖏 🤹 🕓		<b>A U</b> I
			Process Summary	
	0 s			
0	s 3	3s 6s !	9s 12s 15s 18s 21s 24s 27s 3	10 s
Pros 0	solve_em_	moduledriver_ MPI_	Wait moder_mbdld_mop_open_m extose_read	
Pros I	MPI_Bcast	solve_em_	moduledriver_MPI_Waitmoder_modldmopopenm	
Pros 2	MPI_Bcast	solve_en	m moduledriver_MPI_Wait moder_modld_m:op_m	
Pros 3	MPI_Bcast	solve_em_	moduledriver_MPI_Wait moder_modld_mop_open_m	
Pros 4	MPI_Bcast	solve_em_	moduledriver_MPL_Waitmoder_modIdmopopenm	
Pros 5	MPI_Bcast	scive_em_	moduledriver_MPI_waitmoder_modld_mop_open_m	1
Pros 6	MPI_Bcast	solve_em_	module_i.driver_MPI_Wait moder_modld_mop_open_m.i	
Pros /	MPI_Bcast	solve_em_	moduleriverMPI_Waitmoder_modIdmopopenm	4
Pros 8	MPI_Bcast	solve_em_	moduledriver_ MPI_wait moder_ modid_ mop_ cpen m	
Pros 9	MPI_Bcast	solve_em_	moduleariver_MPI_wait mods.er_modId_mop_open_m	
Pr10	MPI_Bcast	solve_em_	moduledriver_MPI_Wait moder_modld_mop_open_m	
Pr11	MPI_Bcast	scive_em_	module_idriver_MPI_Wait moder_modld_mop_ open m	
Pr12	MPI_Bcast	sove_em_	moduledriver_MPI_wait moder_modic_mop_open m	
Pr13	MPI_Bcast	solve_em_	moduledriver_MPI_wait moder_modid_mop_open m	
Pr14	MPI_Bcast	solve_em_	module driver_MPI; wait moder_ modid_ mop_ open_m	
Pr15	MPI_Bcast	solve_em_	moduleriver MPI_wait moder mcdld m open m	
Pr16	MPI_Bcast	solve_em_	moduledriver_ MPI_wait moder_ modld_ mop_ open_m	
Pr17	MPI_Bcast	solve_em_	moduledriver_MPI_Wait moder_modld_mop_ open_m	
Pr18	MPI_BCast	solve_em_	module ariver_ MPI wait mgaer_ modld_ mop_ open m	
Pr19	MPI_BCast	solve_em_	moduledriver_impt_waitmodermodidmopopen_m	
Pr20	MPI_Bcast	solve_em_	module driver MPI wait mode modld mop open m	
Pr21	MPI_Bcast	solve_em_	module driver MP wait moder modld mop open m	
Pr22	MPI_Bcast	solv	/e_emmoduledriver_MPI_waitmodermodldmopm	
Je/wi1.11	nocij		_emmoduleriverMPI_waitmodermodIdmopm	
		e_em	moduledriver_ MPI_waitmoder_ modidmopopenm	
		e_em_	module:anver_IMPt wait moder_modia_mop_cpen_m	
an tana		ve_em_	moduledriver_IMP_wait moder_ modid_ mop_ open m	
		ve_em_	moduledriver_MP_Waltmodermoddmopopen_m	1

	(	) s	з	s	6	s	9 s	12	2 s	1	s	18	s	21	s	24	4 s		27 s		30 9	5
3		MPI_	Bcast		sc	lve_em_		module_	driv	er_MPI	Wait	, i	noder_	mod	ld_	mop	_ op	en	m			
8		MPI_	Bcast			_solve_em		module	.river_	L MPI_	Wait			modi	Jve	r_ mod	ld_	mop	o_ op	en 🚬	m	
2	$\square$	MPI_	Bcast		٦ s	olve_em_	n	nodulec	river_	MPI_W	ait	lmo	der_ <mark>m</mark>	nodl	d_ n	10p_	opei	n T <mark>n</mark>	n			
1		solve	e_em_	mo	dule	driver_ MPI	Wait		mo	.r_ m	odld	mo	open	'n	n	e)	tc	ose_ <mark>re</mark>	ad .			
1		MPI_	Bcast		sol	ve_em_		module	. drive	r_ MPI_	Vait	mod.	.er_ <mark>mo</mark>	dld_	mo	p_ or	en	m	_			
4		MPI_	Bcast		so	lve_em_		module	.river	MPI_	Wait	٦	noder_	n <mark>mo</mark>	dld	mo	p_lop	ben	m	_		
6		MPI_	Bcast		1 s	olve_em_		module	rive	_ MP	_Wait	ևու	der_l	mod	ld_	mo	op op	en	m	-		
4	$\square$	MPI_	Bcast			solve_em		module	drive	r_ MPI_	Wait		L mod	er_	mod.	ld_ m	юр_	ope	n Trr			
6		MPI_	Bcast			solve_em	-	mod	uler	iver_ \	MPI_Wa	ait	L mod	er_ u	mod	ld_ r	no…p	ope	n m.			
7		MPI_	Bcast		- 1	solve_em_		n modul	erive	er_ M	I_Wait	-	mode	er_ m	odl	d_ mo	p_	open	1m.,	-		
9[		MPI_	Bcast		_	solve_em		l modu	leriv	er_ M	PI_Wait		<b>h_</b> mod	uve	r_ m	dld_	mo	.p_ op	en 🗖	m		
3		MPI_	Bcast			sol	ve_er	n_ 👘	L mod	duleri	rer	MPI_Wa	t	-	nod	.er_ˈu <mark>r</mark>	nod	ld_ m	op_	m		
1		MPI_	Bcast		sc	lve_em_		module	river	MPI_\	Vait	modu.	ver_ <mark>m</mark> o	odld	_ m.	<mark>o</mark> p	en	m				
5		MPI_	Bcast		- L	solve_em_		modul	erive	er_ M	PI_Wait	ւ հու	der_h	mod.	ld_	mop	ор	en <mark>n</mark>	n			
4		MPI_	Bcast			<b>L</b> solve_em	-	mod	uleri	ver_	1PI_Wa	it L <mark>re</mark>	oduve	r_1 <mark>m</mark>	odlo	d_ mo.	.p_ c	pen 1	m			

Process Summan

Vampir - [Trace View - /home/dolescha/tracefiles/feature-traces/wrf-p64-io-mer

¥ File View Help ⊻iew <u>C</u>hart <u>F</u>ilter

0.5

## **Vampir: Communication Statistics**



- Byte and message count, min/max/avg message length and min/max/avg bandwidth for each process pair
- Message length statistics





าท

• Available for any part of the trace

#### Vampir: Recipe (JUQUEEN)



- 1. module load UNITE vampirserver
- 2. Start Vampir server component using "vampirserver start smp"
  - Check output for port and pid
- 3. Connect to server from remote machine (see next slide) and analyze the trace
- 4. vampirserver stop <pid>
  - See above (2.)

#### Vampir: Recipe (local system)



- Open SSH tunnel to JUQUEEN using "ssh -L30000:localhost:<port>"
- 2. Start Vampir client component using "/usr/local/zam/unite/bin/vampir"
- 3. Select
  - 1. "Open other..."
  - 2. "Remote file"
  - 3. "Connect" (keep defaults)
  - 4. File "epik.esd" from Scalasca trace measurement directory

# HPCToolkit (Rice University) UJULICH

- Multi-platform sampling-based call-path profiler
- Works on unmodified, optimized executables
- http://hpctoolkit.org
- Advantages:
  - Overhead can be easily controlled via sampling interval
  - Advantageous for complex C++ codes with many small functions
  - Loop-level analysis (sometimes even individual source lines)
  - Supports POSIX threads
- Disadvantages:
  - Statistical approach that might miss details
  - MPI/OpenMP time displayed as low-level system calls

#### **HPCToolkit: Recipe**



- 1. Compile your code with "-g -qnoipa"
  - For MPI, also make sure your application calls MPI\_Comm\_rank first on MPI\_COMM\_WORLD
- 2. Prefix your *link command* with "hpclink"
  - Ignore linker warnings ;-)
- Run your application as usual, specifying requested metrics with sampling intervals in environment variable "HPCRUN\_EVENT\_LIST"
- 4. Perform static binary analysis with "hpcstruct --loop-fwd-subst=no <app>"
- 5. Combine measurements with "hpcprof -S <struct file> \ -I "<path\_to\_src>/\*" <measurement\_dir>"
- 6. View results with "hpcviewer <hpct\_database>"

#### **HPCToolkit: Metric Specification**



- General format: "name@interval [;name@interval ...]"
- Possible sample sources:
  - WALLCLOCK
  - PAPI counters
  - IO (use w/o interval spec)
  - MEMLEAK (use w/o interval spec)
- Interval: given in microseconds
  - E.g.,  $10000 \rightarrow 100$  samples per second

#### **Example hpcviewer**



	편 hpcviewer Eile Debug H @ sor.c 없	: sor <@jj28103> 🕥 <u>H</u> elp	associ source	ated code								
	670 fl 671 rl 672 al 673 fc 674 { 675 676 677 678 679 680 681 682 683 1 682 1	<pre><jm1 =="" field[k][j-1];<br=""><j =="" rhs[k][j];<br=""><j =="" ans[k][j];<br="">or (i=1+mod; i&lt;=nxl; i+=2) delta = omega*( fkj[i+1] +</j></j></jm1></pre>	fkj[i-1] +fkjp1[i] + fk - rkj[i] ); j[i]);	xjm1[i]								
	Calling Context View Callers View Flat View											
	Scope		WALLCLOCK (us).[0,0] (l)	WALLCLOCK (us).[0,0] (E)	WALLCLOCK (us).[1,0] (l)	WALLCLOCK (u						
	Experime	ent Aggregate Metrics	4.79e+06 100 %	4.79e+06 100 %	4.76e+06 100 %	4.76						
	∽ main		4.79e+06 100 %		4.76e+06 100 %	=						
	∽ 🖾 sor_it	er	4.68e+06 97.7%	4.01e+06 83.7%	4.66e+06 97.7%	3.95						
	✓ loop	p at sor.c: 344	2.67e+06 55.7%	2.00e+06 41.7%	2.71e+06 56.8%	2.00						
	⊽ ji	nlined from sor.c: 658	2.00e+06 41.7%	2.00e+06 41.7%	2.00e+06 42.0%	2.00						
	~	loop at sor.c: 662	2.00e+06 41.7%		2.00e+06 42.0%							
		loop at sor.c: 673	2.00e+06 41.7%	3.55e+04 0.7%	2.00e+06 42.0%	2.00						
		b inlined from sor c: 321	1.960+06 41.0%	8 380+05 17 5%	7 06e+05 14 8%	7.06						
		sor c: 675	6.59e+05 13 8%	6.59e+05 13 8%	6.23e+05 13 1%	6.23						
		sor.c: 682	2.40e+05 5.0%	2.40e+05 5.0%	3.96e+05 8.3%	3.96						
Callpath to		sor.c: 678	1.08e+05 2.2%	1.08e+05 2.2%	8.39e+04 1.8%	8.39						
hotepot						>						
ΠΟΙΒΡΟΙ				191M of 400M	Ū							