



# Basic Profiling Tools and Hardware Counter Usage

Get started with some standard tools on blizzard

Hendryk Bockelmann, DKRZ

# Access to performance analysis tools

- ▶ standard profiling available (like `-p` and `-pg` compiler flags)
- ▶ source `/usr/lpp/ppe.hpct/env_sh` for IBM tools (libmpitrace and HPM) to set `IHPCT_BASE` and paths
- ▶ use modules for vampirtrace, scalasca and vampir(server)

module environment:

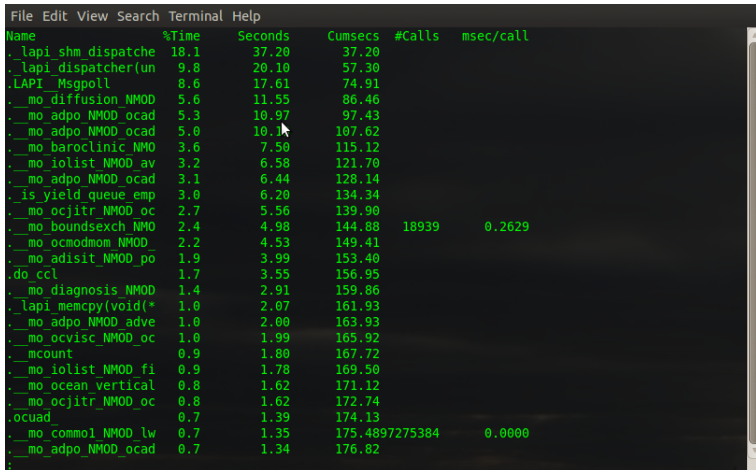
- ▶ `module av`: shows all modules available
- ▶ `module add/rm`: adds, removes modules
- ▶ on blizzard use: `scalasca/1.4.2rc3-aixpoe-ibm`, `vampirtrace/5.13rc1-aixpoe-ibm` (latest versions), maybe go back to stable versions ...
- ▶ on lizard use: `vampir/vampirclient-7.5.0`

# Sequential profiling with prof

- ▶ displays object file profile data
- ▶ for each text symbol in object file the percentage of execution time, number of times that function was called, and the average number of milliseconds per call is shown
- ▶ easy to use: compile and link with additional `-p` option, run your program as usual
- ▶ produce profiling information:  

```
prof ./foo.x -m mon.out > foo.mon.rpt
```
- ▶ output is flat profile of called functions, even MPI calls are profiled at lowest level (ie. `_lapi_*`)

# Seq profiling with prof (example)



Name	%Time	Seconds	Cumsecs	#Calls	msec/call
lapi_shm_dispatche	18.1	37.20	37.20		
lapi_dispatcher(un	9.8	20.10	57.30		
LAPI_Msgpoll	8.6	17.61	74.91		
mo_diffusion_NMOD	5.6	11.55	86.46		
mo_adpo_NMOD_ocad	5.3	10.97	97.43		
mo_adpo_NMOD_ocad	5.0	10.11	107.62		
mo_baroclinic_NMO	3.6	7.50	115.12		
mo_iolist_NMOD_av	3.2	6.58	121.70		
mo_adpo_NMOD_ocad	3.1	6.44	128.14		
is_yield_queue_emp	3.0	6.20	134.34		
mo_ocjitr_NMOD_oc	2.7	5.56	139.90		
mo_boundsexch_NMO	2.4	4.98	144.88	18939	0.2629
mo_ocmodmom_NMOD	2.2	4.53	149.41		
mo_adisit_NMOD_po	1.9	3.99	153.40		
do_ccl	1.7	3.55	156.95		
mo_diagnosis_NMOD	1.4	2.91	159.86		
lapi_memcpy(void(*	1.0	2.07	161.93		
mo_adpo_NMOD_adve	1.0	2.00	163.93		
mo_ocvisc_NMOD_oc	1.0	1.99	165.92		
mcount	0.9	1.80	167.72		
mo_iolist_NMOD_fi	0.9	1.78	169.50		
mo_ocean_vertical	0.8	1.62	171.12		
mo_ocjitr_NMOD_oc	0.8	1.62	172.74		
ocoad	0.7	1.39	174.13		
mo_commol_NMOD_lw	0.7	1.35	175.4897275384		0.0000
mo_adpo_NMOD_ocad	0.7	1.34	176.82		

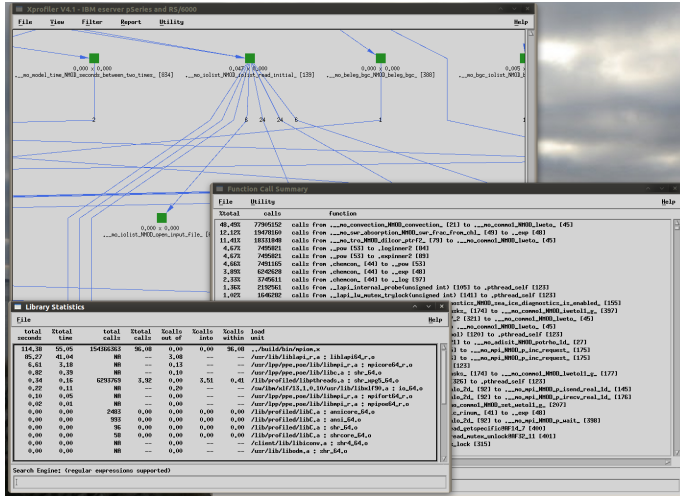
# Sequential profiling with gprof

- ▶ in addition to prof, gprof shows the call-graph data and not only flat profiles
- ▶ easy to use, available on most systems: compile and link with additional `-pg` option, run your program as usual
- ▶ produce profiling information:  

```
gprof ./foo.x gmon.out > foo.gmon.rpt
```
- ▶ drawback: profiling support is added by the compiler, so if you wish to obtain profiling information from any shared libraries, you need to also compile them with `-pg` (therefore MPI calls `_lapi_*` appear spontaneously)



# Use xprofiler for gmon.out



# Processor usage report with tprof

- ▶ tprof charges processor time to object files, processes, threads, subroutines (user mode, kernel mode and shared library) - using AIX trace utility started in background
- ▶ advantage: subroutine-level profiling without modifying executable programs (no recompile, relink needed)
- ▶ modify batch script:  

```
trcstop 2>/dev/null  
tprof -usz -p <binary> -x poe /path/to/binary
```
- ▶ profiles all processes on system (use node not\_shared)
- ▶ output of tprof is very flat, but shows all needed information
- ▶ only one node will be profiled (the one on which poe is started!)





# System profiling with tprof (example)

```
File Edit View Search Terminal Help
Total Ticks For All Processes (USER) = 497135
=====
User Process      PID      PPID      PGRP      Ticks    %    Address  Bytes
=====
profiling-tests/mpiom-gprof/build/bin/mpiom.x 497002  51.67  100000240 84ba74
PID-1          127    0.01  100079e0 f0c3d0d
/etc/pmdv5     5      0.00  1000001f8 ef4a8
/usr/bin/poe   1      0.00  1000001f8 e3a88

Profile: /work/k20200/k202082/profiling-tests/mpiom-gprof/build/bin/mpiom.x

Total Ticks For All Processes (/work/k20200/k202082/profiling-tests/mpiom-gprof/build/bin/mpiom.x) = 497002

Subroutine      Ticks    %    Source      Address  Bytes
=====
ffusion NMOD octdiff_trf_ 86799  9.02  src/mo_diffusion.f90  623c00  d3a0
o_adpo NMOD ocadpo_trf_v_ 45991  4.78  n../src/mo_adpo.f90   67f760  3780
o_adpo NMOD ocadpo_trf_u_ 43318  4.50  n../src/mo_adpo.f90   682e00  3f20
o_adpo NMOD ocadpo_trf_z_ 36246  3.77  n../src/mo_adpo.f90   680e00  4940
o_baroclinic NMOD occlit_ 24764  2.57  rc/mo_baroclinic.f90  64ad60  110a0
glist element accumulate 20318  2.11  ../src/mo_ololist.f90 56fde0  4700
o_ocjitr NMOD ocjitr_trf_ 19923  2.07  ../src/mo_ocjitr.f90  69fae0  5940
,carchm 18691  1.94  rc_hamocc/carchm.f90  713040  4900
NMOD bounds_exch halo_3d_ 15508  1.61  rc/mo_boundsexch.f90  1167800 9620
o_ocmodmom NMOD ocmodmom_ 14818  1.54  /src/mo_ocmodmom.f90  65be00  d400
,ocprod 13556  1.41  rc_hamocc/ocprod.f90  6ff240  fd80
mo_adisit NMOD potrho_1d_ 13028  1.35  ../src/mo_adisit.f90  5de5a0  760
fusion NMOD octdiff_base_ 12944  1.35  src/mo_diffusion.f90  61c720  74e0
```

# System profiling with tprof (example)

```
File Edit View Search Terminal Help

Total Ticks For All Processes (SH-LIBs) = 247344

Shared Object                               Ticks    %    Address  Bytes
-----
/usr/lib/liblapi_r.a[lblapi64_r.o]          218930   22.76 900000002380480 15f734
/usr/lpp/ppe.poe/lib/libmpi_r.a[mpicore64_r.o] 25859   2.69 9000000020a2c00 27d15c
/usr/lpp/ppe.poe/lib/libc.a[shr_64.o]       1666    0.17 900000000000900 3108fc
/usr/lib/libpthreads.a[shr_xpg5_64.o]      440     0.05 9000000005e0200 30a90
/usr/lpp/ppe.poe/lib/libmpi_r.a[mpifort64_r.o] 198     0.02 900000001b17fc0 fe24
/usr/lpp/ppe.poe/lib/libmpi_r.a[mpipoe64_r.o] 70      0.01 90000000206f280 271b7
/usr/lib/libtrace.a[shr.o]                 47      0.00 d1d5e280 2e17c
ibm/xf/13.1.0.10/usr/lib/libxf90.a[io_64.o] 28      0.00 90000000504cf00 67be0
/usr/lib/libc.a[shr_64.o]                   23      0.00 900000000000900 3108fc
/usr/lib/libpthreads.a[shr_xpg5.o]         21      0.00 d04ea180 3157b
/usr/lib/libc.a[shr.o]                      19      0.00 d0118680 31404c
/usr/lib/libptools.a[shr.o]                11      0.00 d42b2280 33228
/usr/lib/libc.a[shr_64.o]                   7       0.00 900000000000900 3108fc
/usr/lib/libBaseLib.a[shr_64.o]             5       0.00 9000000007f6c00 24740
/usr/lib/libodm.a[shr_64.o]                 5       0.00 9000000003ebe80 15e48
/usr/lib/libllrapi.a[shr.o]                 5       0.00 d2e09260 6d9dfc
/usr/lib/libhal_ib.a[hal_ib64.o]           3       0.00 90000000282d600 ef20
/usr/lpp/ppe.poe/lib/libc.a[shr.o]          2       0.00 d0118680 31404c
/usr/lib/libllrapi.a[shr_64.o]              2       0.00 900000000a51ce0 6e58f0
/usr/lib/libpnstd.a[shr_64.o]              1       0.00 90000000274cd80 c8030
/opt/freeware/lib/ganglia/modcpu.so        1       0.00 d05c6128 4b0f
/usr/lib/libc.a[ansicore_64.o]              1       0.00 9000000005a3400 f850
```

# Summary of \*prof

src change, output given:

- ▶ gprof: standard, supported profiling tool on many UNIX systems - no need to learn new stuff
- ▶ tprof: collects data with no impact on execution time, works on optimized binaries without any need for recompilation - prof, gprof might have overhead, might not work on optimized binary, need recompile, relink
- ▶ prof, gprof: provide subprogram profiling, exact counts of number of times every subprogram is called - tprof does not
- ▶ gprof: provides call graph - prof, tprof do not

# Summary of \*prof

what data is given:

- ▶ all \*prof obtain processor consumption estimates for each subprogram by sampling program counter of user program
- ▶ tprof collects processor usage information for whole system - prof, gprof get profiling information for single program & time in user mode

treatment of MPI data:

- ▶ all \*prof tools give flat information on CPU-time used
- ▶ all MPI-stuff is hidden, each process gives its own profile (no correlations can be drawn)

→ MPI profiling needed

# MPI profiling

What we know by now:

- ▶ time spent in user subroutines
- ▶ time spent in MPI-lib

What we do not know by now:

- ▶ which MPI routines take the time
- ▶ correlations between MPI-tasks (send/receive/wait)

→ tools are needed (this is why you are here)

- ▶ libmpitrace
- ▶ scalasca
- ▶ vampir(trace)

# HPCTlib: libmpitrace

- ▶ calls to MPI routines are intercepted by library functions
- ▶ MPI profile data of MPI routines can be collected during a programs execution
- ▶ on `MPI_Finalize`, data is gathered and profile data is written
- ▶ compile your application with `-g` to enable mapping of performance information back to the application source
- ▶ link your application with  
`-L/usr/lpp/ppe.hpct/lib64 -lmpitrace`
- ▶ profile data written to `mpi_profile.<rank>` and XML file `mpi_profile_<rank>.viz` for visualization

# libmpitrace: profile

- ▶ default: data only for rank 0 and min, max, avg rank - set envVar OUTPUT\_ALL\_RANKS=yes for more
- ▶ mpi\_profile.<rank> are human readable

```
File Edit View Search Terminal Help
-----
MPI Routine #calls MPI # #calls avg. bytes avg. time(sec) program execution
-----
MPI Comm_size 1 1 0.0 0.000 information about the rank
MPI Comm_rank 3 3 0.0 0.000
MPI Send 6417 88090.5 0.377 mapping of performance
MPI Isend 245155 8376.3 1.234
MPI Irecv 259470 12461.0 0.804
MPI Sendrecv 36600 1158.4 0.160
MPI Wait 490310 0.0 30.138 by environment variable
MPI Waitall 2045 0.0 0.297
MPI Bcast 231 101330.0 0.078
MPI Barrier 25 54.2 0.001
MPI Gather 2046 647699.5 0.440
MPI Reduce 25 54.2 0.001 MPI Trace MPI
MPI Allreduce 11555 431.5 2.707
-----
total communication time = 47.101 seconds.
total elapsed time = 116.732 seconds.
-----
Message size distributions:
-----
MPI Send #calls avg. bytes time(sec)
6386 83653.2 0.362
31 1003837.9 0.015
MPI Isend #calls avg. bytes time(sec)
158192 816.0 0.357
39548 1632.0 0.095
35724 32640.0 0.737
```

Hendryk Bockelmann, DKRZ



# libmpitrace: profile

view mpi\_profile\_<rank>.viz files with peekperf

The screenshot shows the peekperf application interface. The left pane, titled 'DATA VISUALIZATION WINDOW', displays a tree view of MPI data. The right pane, titled 'SOURCE CODE WINDOW', shows the corresponding source code with line numbers and comments.

Label	Count	WallClock	Transferred B
- MPI_Allreduce	11555	2.10527	4.98579e+06
- MPI_Barrier	2	0.008804	0
- MPI_Bcast	231	2.11603	2.34072e+07
- MPI_Comm_rank	2	2e-06	0
- MPI_Comm_size	1	1e-06	0
- MPI_free	196124	0.245181	1.5757e+09
- MPI_send	198169	1.10373	1.54424e+09
- MPI_recv	207	2.78786	1.63225e+07
- MPI_Reduce	25	0.000605	1904
- MPI_Sendrecv	23250	0.173451	2.42439e+07
- MPI_wait	302248	6.63888	0
- MPI_Waitall	2045	10.2707	0
- mo_boundseach.f90			
- mo_diagnose.f90			
- mo_forcing.f90			
- mo_forcing_NMCD_read_nameList_forcctl_(mo_forcing.f90)			
- MPI_Bcast_420	1	2.5e-05	20
- MPI_Bcast_421	1	6e-06	8
- MPI_Bcast_422	1	7e-06	4
- MPI_Bcast_423	1	7e-06	4
- MPI_Bcast_424	1	6e-06	4
- MPI_Bcast_425	1	5e-06	4
- MPI_Bcast_426	1	6e-06	4
- MPI_Bcast_427	1	2.0e-05	6
- MPI_Bcast_428	1	4.4e-05	4
- MPI_Bcast_429	1	6e-06	4
- mo_grid.f90			
- mo_klobst.f90			
- mo_model_time.f90			
- mo_mpi.f90			
- mo_parallel.f90			
- mo_parallel_debug.f90			
- mo_spm_gather.f90			
- mo_tracer.f90			
- mpiom.f90			
- read_nameList.f90			

```
386 !> Define and read nameList FORCCTL
387 SUBROUTINE read_nameList_forcctl(model_start_time, io_in_forcctl, lerr)
388
389 TYPE(time_desc), INTENT(IN) :: model_start_time
390 INTEGER, INTENT(IN) :: io_in_forcctl
391 INTEGER, INTENT(OUT) :: lerr
392
393 NAMELIST /forcctl/ cfortdata, forcing_frequency, lwrite_forcing, &
#endif NOCDI
394
395 time_interp_forcing, lpsat_interp_forcing, &
396 luse_model_time, forcing_periodicity, &
#endif
397
398 periodic_forcing, &
399 ldebug_forcing, kldt_runoff_grid, forcing_start_time
400
401 ! Set default values
402 cfortdata = 'COMF' ! forcing data set
403 forcing_frequency = 68400_dp ! forcing data period [sec]
404 lwrite_forcing = .FALSE. ! write interpolated forcing files
405 #endif NOCDI
406 time_interp_forcing = .FALSE. ! perform time interpolation
407 lpsat_interp_forcing = .FALSE. ! perform spatial interpolation
408 luse_model_time = .TRUE. ! use model time axis to assign leap years
409 forcing_periodicity = 3659 ! number of available forcing years
410 #endif
411 !periodic_forcing = .FALSE. ! reposition forcing files to initial point at turn of year
412 !kldt_runoff_grid = .FALSE. ! turn on debug mode
413 kldt_runoff_grid = .TRUE. ! offering original grid for runoff data
414 forcing_start_time = model_start_time ! forcing start year
415
416 ! Read nameList forcctl
417 IF (i_p_poe == p_io) READ(io_in_forcctl,forcctl,ostate=lerr)
418
419 ! Broadcast external nameList settings
420 CALL p_bcast(cfortdata, p_io)
421 CALL p_bcast(forcing_frequency, p_io)
422 CALL p_bcast(lwrite_forcing, p_io)
423 CALL p_bcast(time_interp_forcing, p_io)
424 CALL p_bcast(lpsat_interp_forcing, p_io)
425 CALL p_bcast(luse_model_time, p_io)
426 CALL p_bcast(forcing_periodicity, p_io)
427 CALL p_bcast(periodic_forcing, p_io)
428 CALL p_bcast(ldebug_forcing, p_io)
```

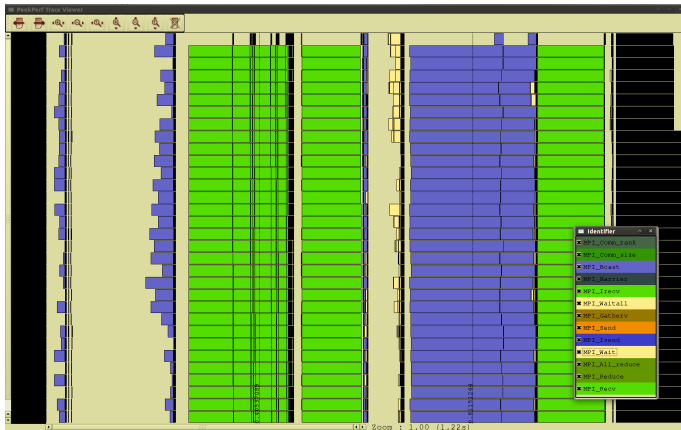
# libmpitrace: customization

- ▶ `OUTPUT_ALL_RANKS=[yes|no]`: show all results
- ▶ `TRACE_ALL_EVENTS=[yes|no]`: do tracing or profiling
- ▶ `MAX_TRACE_RANK=#`: maximum of traced ranks
- ▶ `TRACEBACK_LEVEL=#`: useful for nested MPI-calls within other functions/libraries
- ▶ manual tracing for selected portion of the program through API
- ▶ profiling data cannot be controlled by the API - always collected throughout the entire execution of the program

→ `single_trace_0` outputfile contains tracing information from MPI ranks for which tracing was enabled

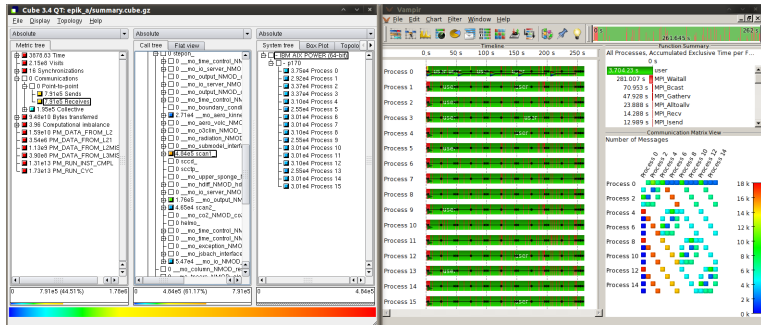
# libmpitrace: trace

view single\_trace\_0 file with peekview



# Summary on basic tools

We have profiles of user subroutines and MPI calls by at least 2 different tools - one would be better  
Additionally, the visualization could be better ...



# If CPU time is not enough ...

Knowing hotspots of the application is only the beginning:

- ▶ why does some routine take so much time?
- ▶ why does MPI take so much time?

Answers: not trivial but often related to one of these issues

- ▶ parallelization aspects (imbalance, race-conditions, etc.)  
→ covered by scalasca and vampir
- ▶ no suitable hardware usage (e.g. cache utilization)

# If CPU time is not enough ...

... getting hardware counter data from IBM HPC Toolkit:

i/ `hpccount` command provides

- ▶ execution wall clock time
- ▶ resource utilization statistics
- ▶ hardware performance counters information
- ▶ derived hardware metrics

for the whole application run:

```
poe hpccount -u -n -o <name> <prog>
```

# hpccount options

- g specifies the hardware counter group
- n suppresses output to stdout
- o writes output to file <name>
- u unique file names will be used

HPM\_ASC\_OUTPUT, HPM\_VIZ\_OUTPUT for ASCII or XML output  
HPM\_AGGREGATE:

- ▶ mirror.so: gets raw data from each MPI-task [default]
- ▶ average.so: counter groups distributed in round robin fashion!  
Aggregator takes avg over these subgroups

# hpcaccount example

```
hpcaccount (IBM HPC Toolkit v5.2.2.4) summary

##### Resource Usage Statistics #####

Total amount of time in user mode      : 119.653446 seconds
Total amount of time in system mode    : 0.053128 seconds
Maximum resident set size              : 401968 Kbytes
Average shared memory use in text segment : 7773 Kbytes*sec
Average unshared memory use in data segment : 34396194 Kbytes*sec
Number of page faults without I/O activity : 6925
Number of page faults with I/O activity  : 0
Number of times process was swapped out : 0
Number of times file system performed INPUT : 0
Number of times file system performed OUTPUT : 0
Number of IPC messages sent            : 0
Number of IPC messages received         : 0
Number of signals delivered             : 0
Number of voluntary context switches    : 493
Number of involuntary context switches   : 120

##### End of Resource Statistics #####

Execution time (wall clock time)       : 121.634124139789 seconds

PM_FPU_1FLOP (FPU executed one flop instruction)      : 76000179
PM_FPU_FMA (FPU executed multiply-add instruction)   : 16000000076
PM_FPU_FSQRT_FDIV (FPU executed FSQRT or FDIV instruction) : 8
PM_FPU_FLOP (FPU executed 1FLOP, FMA, FSQRT or FDIV instruction) : 16076000263
PM_RUN_INST_CMPL (Run instructions completed)        : 151953175375
PM_RUN_CYC (Run cycles)                              : 564868441806

Utilization rate                                     : 98.724 %
Instructions per run cycle                           : 0.269
Total scalar floating point operations               : 32076.000 M
Scalar flop rate (flops / WCT)                       : 263.709 Mflop/s
Scalar flops / user time                             : 267.116 Mflop/s
Algebraic floating point operations                  : 32076.000 M
Algebraic flop rate (flops / WCT)                    : 263.709 Mflop/s
Algebraic flops / user time                          : 267.116 Mflop/s
Scalar FMA percentage                               : 99.763 %
```



# If CPU time is not enough ...

... getting hardware counter data from IBM HPC Toolkit:  
ii/ instrumentation with libhpm

```
#include <libhpc.h>
int main(void) {
    ...
    MPI_Init(MPI_COMM_WORLD);
    MPI_Comm_rank(MPI_COMM_WORLD, &myrank);
    hpmInit(myrank, "my_program");
    ... maybe some initialization code
    hpmStart(1, "first section");
    ... some code you want to analyze
    hpmStop(1);
    ... more boring code
    hpmTerminate(myrank);
    MPI_Finalize();
}
```

```
PROGRAM HELLO_WORLD
    IMPLICIT NONE
#include "f_hpc.h"
    ...
    CALL MPI_INIT(ierrord)
    CALL MPI_COMM_RANK(MPI_COMM_WORLD, myrank, ierrord)
    CALL f_hpmInit(myrank, 'my_program')
    ... maybe some initialization code
    CALL f_hpmstart(1, 'first section')
    ... some code you want to analyze
    CALL f_hpmstop(1)
    ... more boring code
    CALL f_hpmterminate(myrank)
    CALL MPI_FINALIZE(ierrord)
END
```

don't forget `-I/usr/lpp/ppe.hpct/include` and  
`-L/usr/lpp/ppe.hpct/lib64 -lhpc -lpmapi`

# Choosing counters

on POWER6

- ▶ you can measure 6 counters simultaneously
- ▶ not all combinations allowed
- ▶ 202 performance counter groups (`pm1ist -g -1`)
- ▶ may need to sample multiple times for completeness
- ▶ or use multiplexing of counter groups

Use `hpccount/libhpc` to measure code efficiency by means of:

- ▶ Instructions per run cycle, Mflop/s (group 127)
- ▶ L1 cache usage (group 47)
- ▶ cache/memory access (group 7,11)

# HPM counter analysis

counters are difficult to understand

- ▶ limited documentation
- ▶ experience needed to see counter value indicating a problem
- ▶ use IPC (instructions per run cycle) for first estimate

$$\text{IPC} = \text{PM\_RUN\_INST\_CMPL} / \text{PM\_RUN\_CYC} = 1/\text{CPI}$$

category	IPC	description
1	< 0.4	Houston, we have a problem
2	0.4... 0.7	not tuned for POWER6
3	0.7... 0.9	acceptable
4	0.9... 1.3	very good (can be tough to get here)
5	> 1.3	wow ! (not always possible)
6	> 2.0	LINPACK, VMASS, ESSL, FFTW, ...

# HPM counter analysis

group 127:

- ▶ PM\_FPU\_1FLOP: FPU executed single FLOP operation
- ▶ PM\_FPU\_FMA: FPU executed a multiply-add
- ▶ PM\_FPU\_FSQRT\_FDIV
- ▶ PM\_FPU\_FLOP: 1flop, fma, sqrt, div operations for unit0 and 1

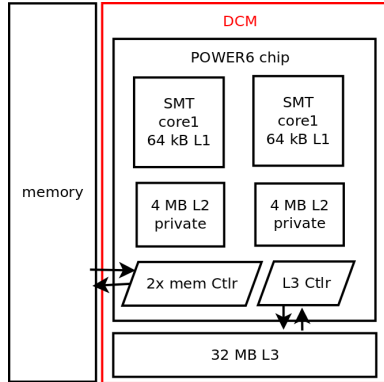
group 47:

- ▶ PM\_ST\_REF\_L1: L1 D cache store references
- ▶ PM\_LD\_REF\_L1: L1 D cache load references
- ▶ PM\_ST\_MISS\_L1: L1 D cache store misses
- ▶ PM\_LD\_MISS\_L1: L1 D cache load misses

# HPM counter analysis

group 7:

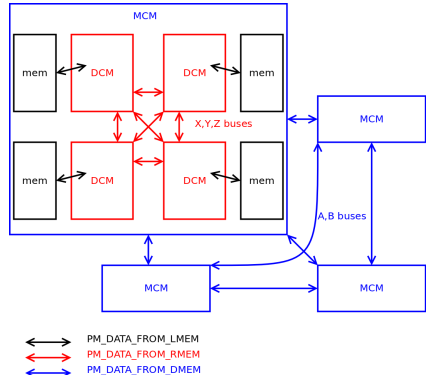
- ▶ PM\_DATA\_FROM\_L2: D Cache reloaded from local L2
- ▶ PM\_DATA\_FROM\_L21: D Cache reloaded from private L2 of other core on chip
- ▶ PM\_DATA\_FROM\_L2MISS: D Cache reloaded but not from local L2
- ▶ PM\_DATA\_FROM\_L3MISS: D Cache reloaded from beyond L3



# HPM counter analysis

group 11:

- ▶ **PM\_DATA\_FROM\_LMEM:** D Cache reloaded from local memory (attached to DCM)
- ▶ **PM\_DATA\_FROM\_RMEM:** D Cache reloaded from remote memory (attached to other DCM on same MCM)
- ▶ **PM\_DATA\_FROM\_DMEM:** D Cache reloaded from distant memory (attached to a different MCM)



# HPCC results out of the box

- i/ HPL: solves a dense linear system in double precision (linpack)  
 $N = 20000$ ,  $NB = 20 \vee 120$ ,  $P = 8$ ,  $Q = 4$

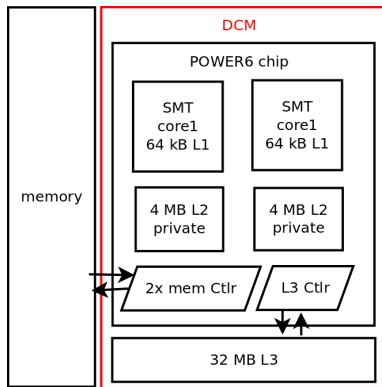
NB	HPL Gflop/s	libhpc: per MPI task			
		IPC	Gflop/s	% peak	L1 cache hit rate
20	2.39e+02	1.45 - 1.72	7.33 - 7.56	39.9 - 41.1 %	97.5 - 99.5 %
120	3.15e+02	1.82 - 2.09	9.75 - 9.8	52.3 - 56.0 %	98.5 - 99.4 %

NB	libhpc: per MPI task		
	PM_DATA_FROM_L2	PM_DATA_FROM_L2MISS	PM_DATA_FROM_LMEM
20	35131005 - 70195962	8913581 - 12180711	5921680 - 9110732
120	30105907 - 65539464	4992461 - 6045666	3332966 - 4227469

# POWER6 background

POWER6 cache/memory: size and latency

- ▶ L1: 64 kB, 2 cycles
- ▶ L2: 4 MB core-private, 25 cycles
- ▶ L3: 32 MB chip-shared, 150 cycles
- ▶ mem: 50 (100) GB, 500 cycles





# HPCC results out of the box

- ii/ DGEMM: measures floating point rate of double precision real matrix-matrix multiplication  
bench  $N = 1500$  using libessl or not

essl	HPL Gflop/s (min/max)	libhpc: per MPI task	
		IPC	Gflop/s
no	2.94 / 3.31	0.65 - 0.68	3.4 - 3.6
yes	10.84 / 15.36	2.49 - 2.59	14.2 - 14.8

essl	libhpc: per MPI task		
	% peak	PM_ST_MISS_L1	PM_LD_MISS_L1
no	18.5 - 19.6 %	1.23e6 - 1.47e6	215.68e6 - 221.91e6
yes	77.3 - 80.6 %	2.71e6 - 3.32e6	0.93e6 - 1.99e6