



EZTRACE - ADT EZPerf easy performance trace analyser

HASTARAN Matias - SYLVAND Guillaume - RUE François

INTRODUCTION

Modern HPC applications are complex:

- Complex hardware: NUMA architecture, hierarchical caches, accelerators;
- Hybrid programming models: MPI/OpenMP/Pthread.
 - Understand the performance of those application is difficult.

Generating program behavior trace with a low overhead

- Code instrumentation

INTRODUCTION

EZPerf environment

- Project based on EZTrace performance tool
- EPI HiePACS

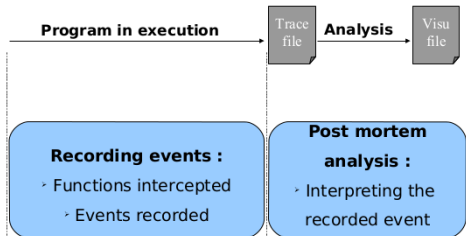
EZPerf Team

- Guillaume SYLVAND
Scientist project manager (INRIA)
- François Trahay
EZTrace's architect (Télécom SudParis)
- François Rué
Technical project manager (INRIA)
- Matias HASTARAN
Engineer(INRIA)

1

EZTRACE

Execution in two phases



- `mpirun -np 16 eztrace myProg`
- `eztrace_convert /tmp/<username>_eztrace_log_rank_*`

User module

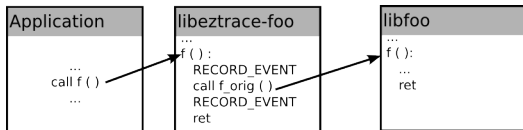
```
BEGIN_MODULE
NAME myModule
DESC "description _for _my _module"
ID 42
void do_something ()
END_MODULE
```

- eztrace_create_plugin plugin.tpl
- export EZTRACE_LIBRARY_PATH=module_path
- Run eztrace

2

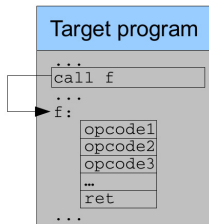
Interception mechanism in eztrace

LD_PRELOAD mechanism



- Pros:
 - Easy to use
 - Easy to run
 - Easy to do
- Cons:
 - Only C functions
 - Only in dynamic libraries

EZTrace Instrumentation

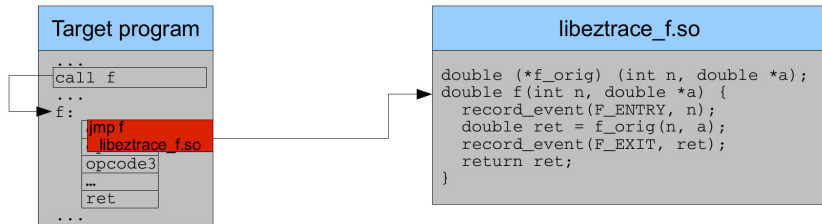


libeztrace_f.so

```
double (*f_orig) (int n, double *a);  
double f(int n, double *a) {  
    record_event(F_ENTRY, n);  
    double ret = f_orig(n, a);  
    record_event(F_EXIT, ret);  
    return ret;  
}
```

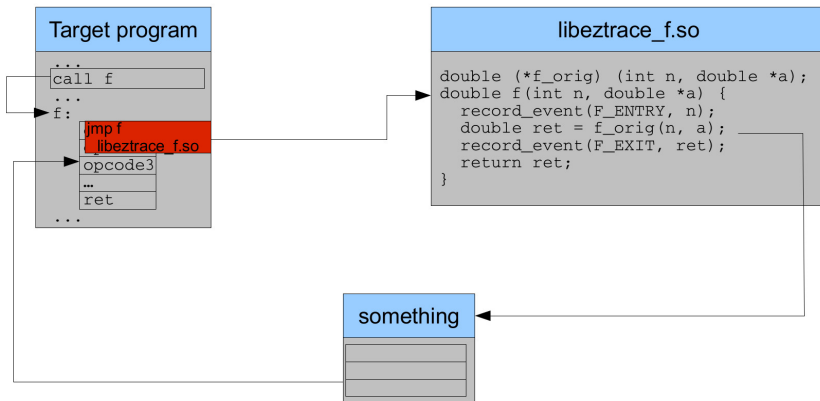
- Coarse-grain instrumentation
- Interception code is compiled into a shared library
- Mechanism similar to LD_PRELOAD

EZTrace Instrumentation



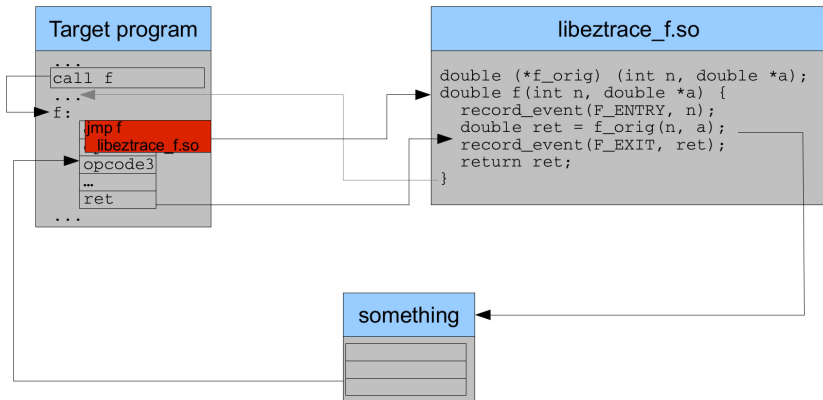
- f_orig contains the address of the original function f
- First f function opcodes are overwritten by a jump

EZTrace Instrumentation



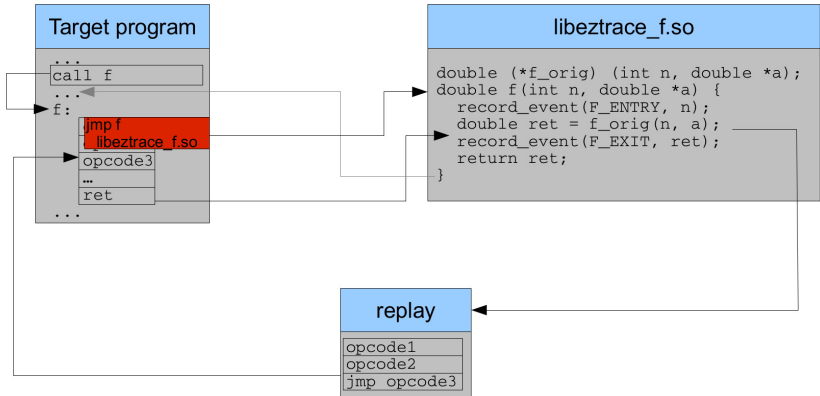
- A trampoline is created somewhere in memory
- `f_orig` will jump to this trampoline
- The trampoline will jump to the next opcode

EZTrace Instrumentation



- The return of original `f` function will jump to the next event of new `f` function
- The return of the new `f` function will jump to the next event

EZTrace Instrumentation



- The trampoline is fill with the overwritten opcodes
- During the execution it will play the start of the original f function

EZTrace instrumentation

- Function are instrumented using plugins at the application startup
- Functions are instrumented by modifying the target program in memory
- Pros:
 - Efficient and simple
 - Easy to use
 - Easy to extends: little architecture-dependent code
- Cons:
 - Limited to functions entries and exits

Evaluation

Interception method	no instrumentation	PIN	DYNINST	EZTRACE
Statically linked library	4.7 ns	1287 ns	1294 ns	245.2 ns
Shared library	5.3 ns	1293 ns	-	227.4 ns

Table: Function duration when recording events during a function call

Interception method	no instrumentation	PIN	DYNINST	EZTRACE
Execution time (s)	0.45	3.16	3.28	2.26
Overhead (ns / iteration)	+0	+341	+413	+227

Table: Execution time of a molecular dynamics simulation using 4 OpenMP threads

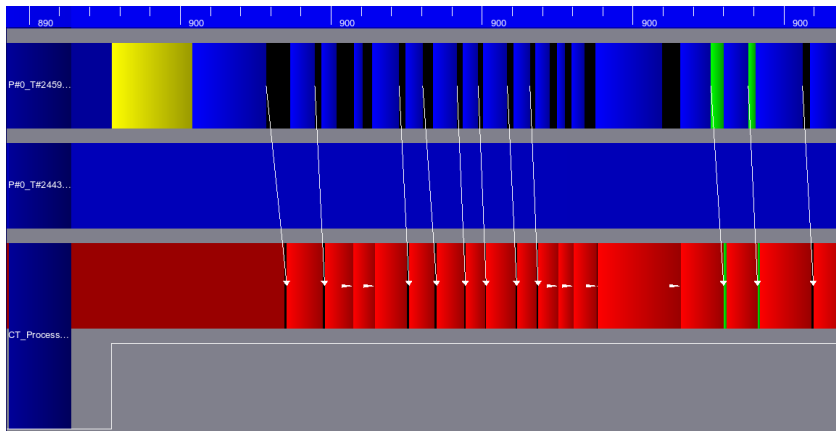
3

Recent Developments

Module

- MPI Module: Stabilization and new testsuite
- OpenMP : Module and eztrace_cc stabilization
- CUDA : New module !!

Cuda



Library replacement

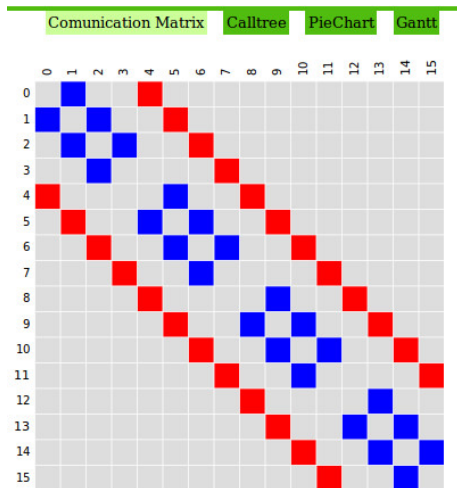
Goodbye FxT. Welcome LiTL.

- Easiest to install
- Trace file are lighter

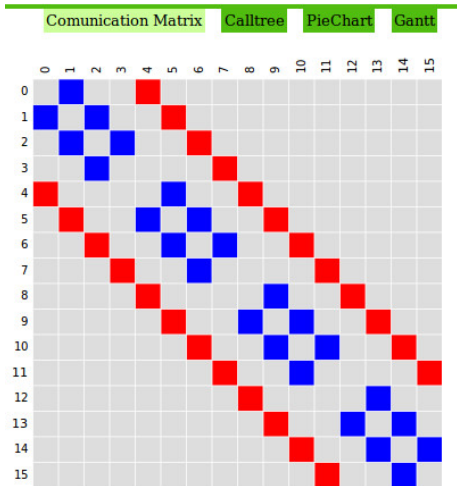
EZTrace statistic tool

```
MPI:
---
MPI_SEND           :6720 calls
MPI_IRecv          :6720 calls
MPI_WAIT           :6720 calls
MPI_BARRIER       :4 calls
MPI_REDUCE         :4 calls
CT_Process #0 --   1680 messages sent
Size of messages (byte):   min: 8 max: 56000 average: 27737 total: 46598912
Time spent sending messages (ms): min: 0.000908 max: 0.267776 average: 0.013695 total: 23.007082
Time spent computing while sending messages (ms): min: 0.000000 max: 0.000000 average: 0.000000 total: 0.000000
0
Time spent in MPI_Send or waiting for a Isend to complete (ms): min: 0.000908 max: 0.267776 average: 0.013695 total:
23.007082
Time spent receiving messages (ms): min: 0.003911 max: 0.315338 average: 0.019409 total: 32.006908
Time spent computing while receiving messages (ms): min: 0.002304 max: 0.313383 average: 0.013882 total: 23.3220
53
Time spent in MPI_Recv or waiting for a Irecv to complete (ms): min: 0.000628 max: 0.088700 average: 0.005527 total:
9.284855
CT_Process #1 --   1680 messages sent
Size of messages (byte):   min: 8 max: 56000 average: 27737 total: 46598912
Time spent sending messages (ms): min: 0.000908 max: 0.578365 average: 0.014614 total: 24.551360
Time spent computing while sending messages (ms): min: 0.000000 max: 0.000000 average: 0.000000 total: 0.000000
0
Time spent in MPI_Send or waiting for a Isend to complete (ms): min: 0.000908 max: 0.578365 average: 0.014614 total:
24.551360
Time spent receiving messages (ms): min: 0.003492 max: 3.989680 average: 0.024444 total: 41.065636
Time spent computing while receiving messages (ms): min: 0.002235 max: 0.269941 average: 0.019668 total: 33.0427
60
Time spent in MPI_Recv or waiting for a Irecv to complete (ms): min: 0.000628 max: 3.970682 average: 0.004776 total:
8.022876
```

Web interface

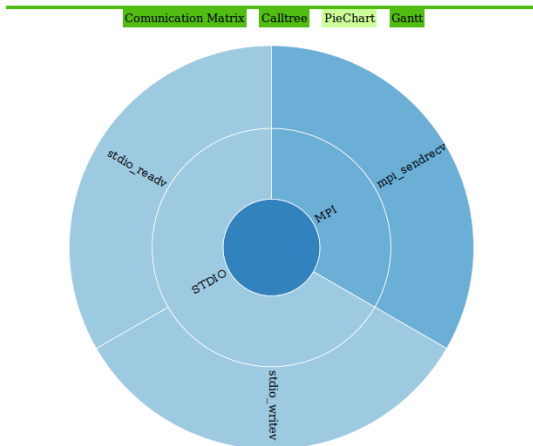


Web interface

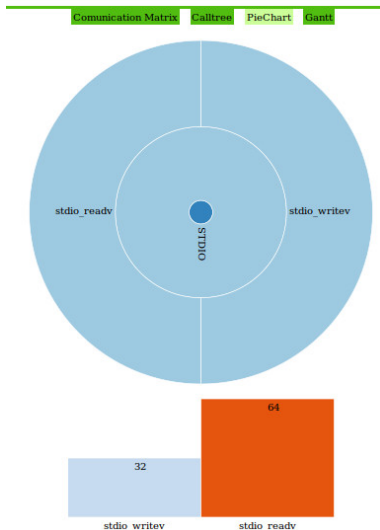


0	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

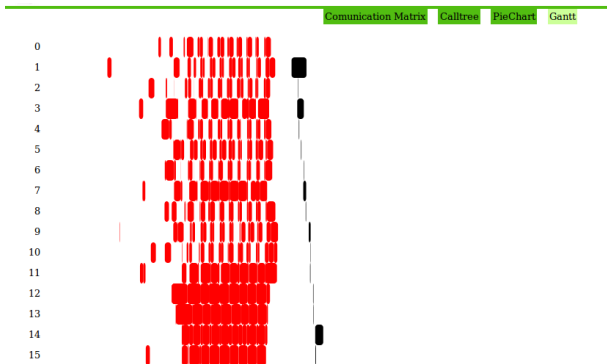
Web interface



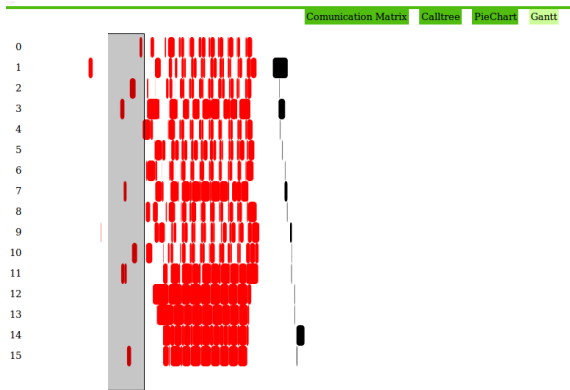
Web interface



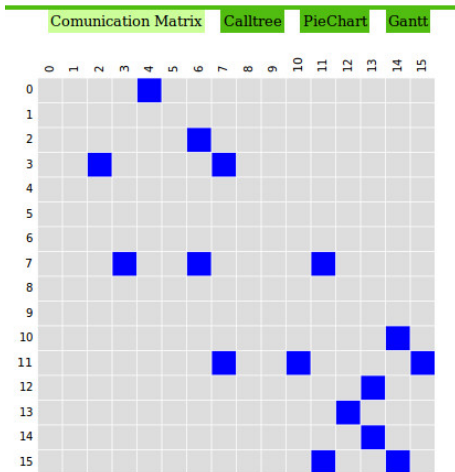
Web interface



Web interface



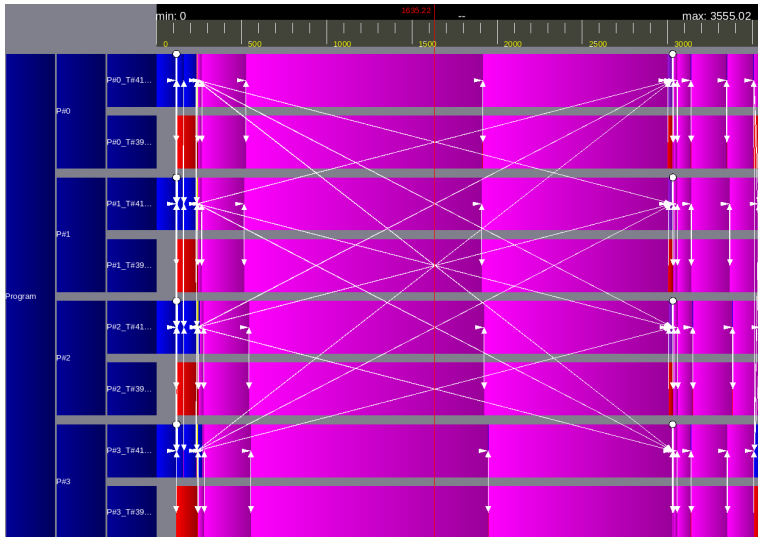
Web interface



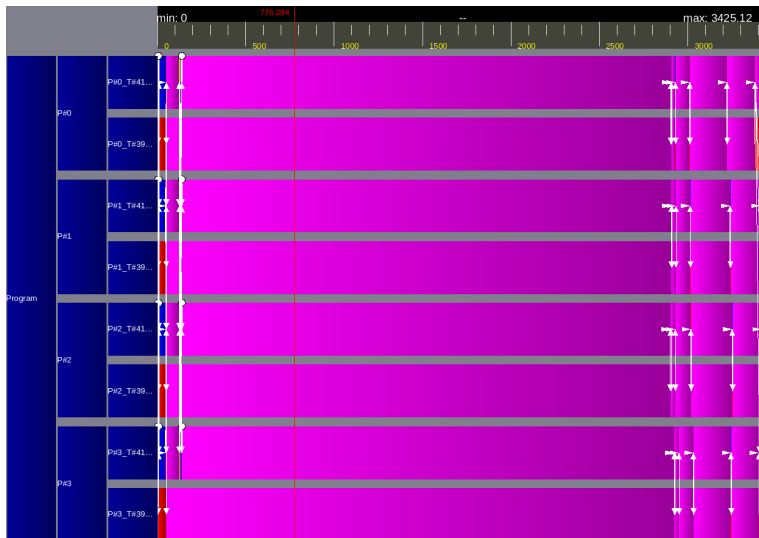
4

User Case : ScaIFMM

ScaIFMM

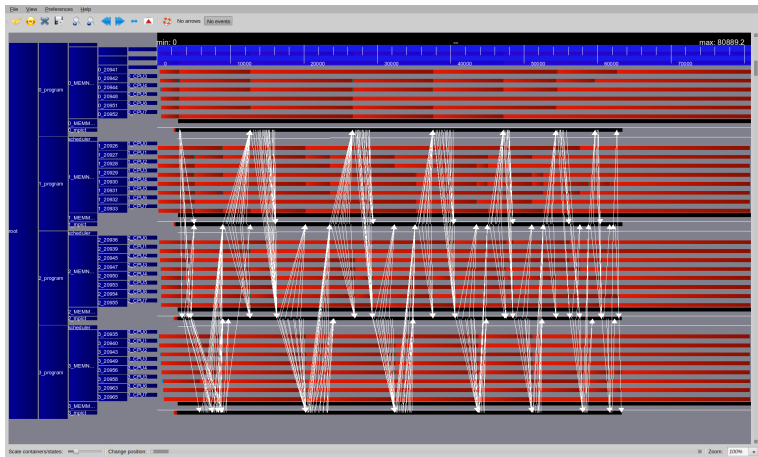


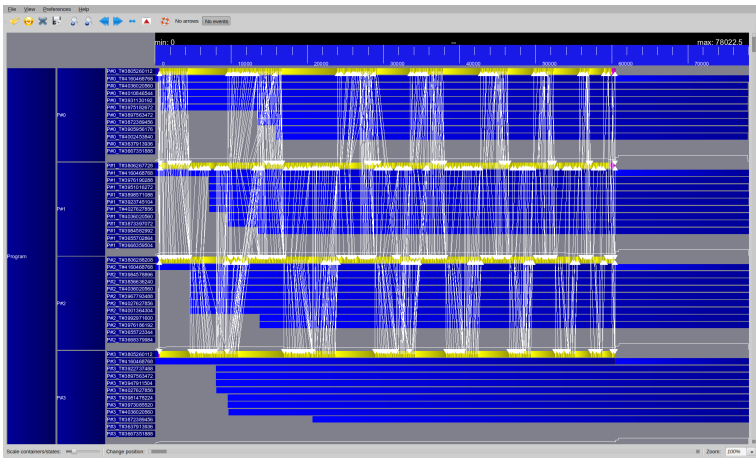
ScaIFMM



5

User Case : StarPU

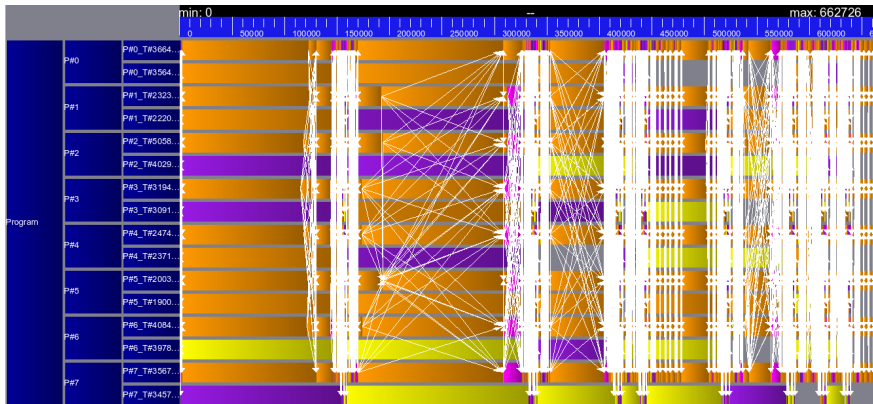




6

User Case : PaMPA

PaMPA



7

Thank you!

Contact

- EZTrace is freely available under the CeCILL-B license
- Website : <http://eztrace.gforge.inria.fr/>
- Contact : eztrace-devel@lists.gforge.inria.fr