

Determinação de Hotspots de execução de kernels recorrendo a PERF (linux-tools)

Email: a57816@alunos.uminho.pt

3. Services and Advanced Research Computing with HTC/HPC clusters

- Matrizes contidas na LLC:
Sabemos que as 3 matrizes terão de ter uma tamanho de coluna/linha inferior a:

$$\text{lado matriz} \leq \frac{\sqrt{\left(\frac{12288KB}{4 \text{ Bytes p/float}}\right)}}{3 \text{ matrizes}} \leq 1011$$

para que os datasets estejam completamente contidos em cache. Foi escolhido o lado da matriz de 512 elementos.

- Matrizes contidas na memória principal:
Sabemos que as 3 matrizes terão de ter uma tamanho de coluna/linha superior a 1011 elemento por linha/coluna. Foi escolhido o lado da matriz de 2048 elementos.

Doravante, todos os comandos apresentados terão por standard o dataset menor (matrizes de 512*512), e o kernel naive. Em caso de serem utilizados outro kernel ou dataset será feita a referência antes da apresentação dos resultados.

4. Parte 1

4.1. Passo 0 – uma análise aos software e hardware events disponíveis na máquina

Recorrendo ao seguinte comando:

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf list | wc -l
```

Podemos ter uma noção em termos de eventos de hardware e software disponíveis na máquina de teste:

1 784

4.2. 1º Passo – o encontrar dos tempos base sem overhead de medição para o menor dataset (dataset em LLCACHE)

Por forma a podermos estabelecer um limite dentro do aceitável das possíveis alterações à performance das versões do kernel necessitamos de primeiramente ter uma medição com o mínimo de overhead possível, para o menor dataset em profiling (matriz 512*512), para o caso **naive**. Recorrendo ao seguinte comando:

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf stat -e \
```

```
cpu-clock ./naive
```

```
1 Performance counter stats for './naive':
2
3      201.696023      cpu-clock (msec)
4
5      0.206153405 seconds time elapsed
```

Teremos uma base de referência futura, para podermos validar as nossas medições com overhead.

Para além do tempo de execução é ainda demonstrado o número de cpu-clocks necessários para a execução.

Como seria de esperar, é possível num única medição registar vários eventos. Recorrendo ao seguinte comando:

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf stat -e \
```

```
cpu-clock,faults ./naive
```

```
1 [a57816@compute-431-l ESC_FLAME_GRAPH]$ perf stat -e \
2   cpu-clock,faults ./naive
3
4 Performance counter stats for './naive':
5
6      196.395754      cpu-clock (msec)
7      844            faults
8
9      0.200874285 seconds time elapsed
```

4.3. O processo de procura de hotspots de uma aplicação

O processo de procura de hotspots de uma aplicação pode ser traduzido em 2 passos simples:

- Passo 1: Recolha de dados para profiling da aplicação
- Passo 2 : Tratamento e apresentação da informação recolhida.

4.3.1. Passo 1. Podemos recolher informação de profiling simplesmente correndo o seguinte comando:

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf record -e \
```

```
cpu-clock,faults ./naive
```

Neste caso específico o PERF corre a aplicação naive e recolhe dados relativos aos eventos: cpu-clock e page-faults.

```
1 [a57816@compute-431-l ESC_FLAME_GRAPH]$ perf record -e \
2   cpu-clock,faults ./naive
3 [ perf record: Woken up 1 times to write data ]
4 [ perf record: Captured and wrote 0.046 MB perf.data \
5   (830 samples) ]
```

4.3.2. Passo 2. Relativamente ao passo 2 - Tratamento e apresentação da informação recolhida - o seguinte comando permite visualizar a informação presente no ficheiro perf.data (ficheiro default no qual é gravada a informação):

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf report -<
```

```
stdio
```

```
1 [a57816@compute-431-l ESC_FLAME_GRAPH]$ perf report -<
2 stdio
3
4 # To display the perf.data header info, please use --header=
5 # header/--header-only options.
6
7 #
8 # Samples: 808 of event 'cpu-clock'
9 # Event count (approx.): 808
10
11 # Overhead Command Shared Object Symbol
12 # .....
13
14 94.68% naive naive [.]
15 multiply_matrices
16 2.10% naive naive [.]
17 initialize_matrices
18 1.36% naive libc-2.12.so [.] __random
19 0.87% naive libc-2.12.so [.] __random_r
20 0.25% naive [kernel.kallsyms] [k]
21 __mem_cgroup_commit_charge
```

```

15 0.25% naive naive [.] rand@plt
16 0.12% naive [kernel.kallsyms] [k] __call_rcu
17 0.12% naive [kernel.kallsyms] [k] ←
    __mem_cgroup_uncharge_common
18 0.12% naive libc-2.12.so [.] malloc
19 0.12% naive libc-2.12.so [.] rand
20
21
22 # Samples: 22 of event 'faults'
23 # Event count (approx.): 1157
24 #
25 # Overhead Command Shared Object Symbol
26 # .....
27 #
28 63.61% naive naive [.] ←
    initialize_matrices
29 29.56% naive libc-2.12.so [.] ←
    __init_cpu_features
30 5.62% naive ld-2.12.so [.] dl_main
31 0.52% naive ld-2.12.so [.] ←
    _dl_setup_hash
32 0.26% naive ld-2.12.so [.] _dl_start
33 0.17% naive [kernel.kallsyms] [k] ←
    __clear_user
34 0.17% naive libc-2.12.so [.] ←
    __strchr_sse2
35 0.09% naive ld-2.12.so [.] _start

```

```

21 #
22 63.61% naive naive
23 29.73% naive libc-2.12.so
24 6.48% naive ld-2.12.so
25 0.17% naive [kernel.kallsyms]

```

Como pode confirmar pelo output, as duas opções no comando permitem agregar os resultados por comando e shared object.

Por forma a simplificar o processo de catalogação e posterior tratamento de dados é possível aceder a metadados das medições, através da junção da opção **-header** ao comando anterior. Desta forma temos acesso a por exemplo, os dados de hardware do ambiente de teste, o comando perf utilizado e informação relativa aos eventos medidos. Esta funcionalidade torna-se bastante útil dada a tremenda facilidade que temos hoje em dia de gerar dados de teste. A consequência da facilidade de tal geração é a posterior necessidade de tratamento correcto dos mesmos.

```

[a57816@search6 ESC_FLAME_GRAPH]$ perf report --stdio ←
--sort comm,dso --header

```

Existem outras alternativas a esse mesmo comando sendo esta:

```

[a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf report --tui

```

que apresenta o Terminal-based User Interface (TUI).

Voltemos ao output apresentado pelo comando com a opção **stdio**. Denote no seguinte detalhe:

```

1 Samples: 808 of event 'cpu-clock'
2 Event count (approx.): 808

```

Como poderá confirmar foram registadas 808 medições do evento "cpu-clock".

Adicionando a opção **-sort** conseguimos uma visualização agregada. Em conjugação com a opção **comm** e **dso** obtemos o seguinte resultado.

```

[a57816@search6 ESC_FLAME_GRAPH]$ perf report --stdio ←
--sort comm,dso

```

```

1 [a57816@search6 ESC_FLAME_GRAPH]$ perf report --stdio ←
2 --sort comm,dso
3 [kernel.kallsyms] with build id 09←
4 b53ce5dfe09c7b1ad9473e1356d7d922512e27 not found, ←
5 continuing without symbols
6 # To display the perf.data header info, please use ←
7 header/--header-only options.
8 #
9 # Samples: 808 of event 'cpu-clock'
10 # Event count (approx.): 808
11 #
12 # Overhead Command Shared Object
13 # .....
14 #
15 97.03% naive naive
16 2.48% naive libc-2.12.so
17 0.50% naive [kernel.kallsyms]
18
19 # Samples: 22 of event 'faults'
20 # Event count (approx.): 1157
21 #
22 # Overhead Command Shared Object
23 # .....
24 #
25 63.61% naive naive
26 29.73% naive libc-2.12.so

```

```

1 [a57816@search6 ESC_FLAME_GRAPH]$ perf report --stdio ←
2 --sort comm,dso --header
3 [kernel.kallsyms] with build id 09←
4 b53ce5dfe09c7b1ad9473e1356d7d922512e27 not found, ←
5 continuing without symbols
6 # =====
7 # captured on: Mon May 23 23:00:56 2016
8 # hostname : compute-431-1.local
9 # os release : 2.6.32-279.14.1.el6.x86_64
10 # perf version : 4.0.0
11 # arch : x86_64
12 # nrcpus online : 24
13 # nrcpus avail : 24
14 # cpudesc : Intel(R) Xeon(R) CPU X5650 @ 2.67GHz
15 # cpuid : GenuineIntel,6,44,2
16 # total memory : 12319324 kB
17 # cmdline : /share/jade/SOFT/perf/perf record -e cpu←
18 clock,faults ./naive
19 # event : name = cpu-clock, type = 1, config = 0x0, ←
20 config1 = 0x0, config2 = 0x0, excl_usr = 0, ←
21 excl_kern = 0, excl_host = 0, excl_guest = 1, ←
22 precise_ip = 0, attr_mmap2 = 0, attr
23 # event : name = faults, type = 1, config = 0x2, ←
24 config1 = 0x0, config2 = 0x0, excl_usr = 0, ←
25 excl_kern = 0, excl_host = 0, excl_guest = 1, ←
26 precise_ip = 0, attr_mmap2 = 0, attr_mm
27 # HEADER_CPU_TOPOLOGY info available, use -I to display
28 # HEADER_NUMA_TOPOLOGY info available, use -I to ←
29 display
30 # pmu mappings: cpu = 4, tracepoint = 2, software = 1
31 # =====
32 # Samples: 808 of event 'cpu-clock'
33 # Event count (approx.): 808
34 #
35 # Overhead Command Shared Object
36 # .....
37 #
38 97.03% naive naive
39 2.48% naive libc-2.12.so
40 0.50% naive [kernel.kallsyms]
41
42 # Samples: 22 of event 'faults'
43 # Event count (approx.): 1157
44 #
45 # Overhead Command Shared Object
46 # .....
47 #
48 63.61% naive naive
49 29.73% naive libc-2.12.so

```

```

41 6.48% naive ld-2.12.so
42 0.17% naive [kernel.kallsyms]

```

4.3.3. Aumentar a granularidade do profiling – fase I. O exemplo dado anteriormente apresenta os dados agregados relativos a toda a aplicação. Ora, daí podemos observar que **97.03%** do tempo de execução é despendido em métodos da própria aplicação, enquanto que os restantes **2.48%** são despendidos na biblioteca do sistema **libc-2.12.so**.

O próximo passo será portanto aumentar a granularidade do profiling nestes dois "shared objects".

```

[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf report --
stdio --dsos=naive,libc-2.12.so

```

```

1 [a57816@compute-431-l ESC_FLAME_GRAPH]$ perf report --
2 stdio --dsos=naive,libc-2.12.so
3 # To display the perf.data header info, please use --
4 header/--header-only options.
5 #
6 # Samples: 808 of event 'cpu-clock'
7 # Event count (approx.): 808
8 #
9 # Overhead Command Shared Object Symbol
10 # .....
11 94.68% naive naive [...]
12 multiply_matrices
13 2.10% naive naive [...]
14 initialize_matrices
15 1.36% naive libc-2.12.so [...] __random
16 0.87% naive libc-2.12.so [...] __random_r
17 0.25% naive naive [...] rand@plt
18 0.12% naive libc-2.12.so [...] malloc
19 0.12% naive libc-2.12.so [...] rand
20 #
21 # Samples: 22 of event 'faults'
22 # Event count (approx.): 1157
23 #
24 # Overhead Command Shared Object Symbol
25 # .....
26 63.61% naive naive [...]
27 initialize_matrices
28 29.56% naive libc-2.12.so [...]
29 __init_cpu_features
30 0.17% naive libc-2.12.so [...] __strchr_sse2

```

O método **multiply_matrices()** é o responsável pela maioria do tempo de computação, sendo aproximadamente **94.86%**.

A segunda função, em termos de tempo de computação é a **initialize_matrices** com 2.10%.

Podemos ainda analisar as **page faults** da aplicação, sendo que a maioria ocorre no método **initialize_matrices** – 63.61%, como seria de esperar. Este comportamento é o esperado, e não é responsável por qualquer problema de performance da aplicação, dado que esse mesmo método apenas ocupa 2.10% do tempo total da aplicação, não afectando o processo de multiplicação de matrizes, que é o responsável pela maior porção do mesmo.

Podemos ainda "validar" o número de page faults ocorrido, associando-o ao tamanho da matriz em uso. Ora, **512*512 floats** representa um total de **512*512*(4 bytes) = 1024KB**, multiplicado por 3 (3 matrizes), representa um

total de 3072 KB, que representa um total de páginas $\frac{3072KB}{4KB \text{ por pag.}} = 768$. Ora, 63.61% do total de **page faults** contabilizadas representa $0.6361 * 1157 = 735,9677 = 736 \text{ pags.}$, como esperado.

4.3.4. Aumentar a granularidade do profiling – fase II. Agora que temos a completa percepção de qual o método que está a consumir maioritariamente o tempo de execução, podemos realizar um análise não por biblioteca, mas por método, associando então ao método **multiply_matrices()** o respectivo código máquina e posição relativa ao início do método, recorrendo ao comando:

```

[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf annotate --
stdio --dsos=naive --symbol=multiply_matrices

```

Produzindo o seguinte output:

```

1 [a57816@compute-431-l ESC_FLAME_GRAPH]$ perf annotate --
2 stdio --dsos=naive --symbol=multiply_matrices
3 Percent | Source code & Disassembly of naive for cpu-
4 clock
5
6 :
7 :
8 : Disassembly of section .text:
9 :
10 : 0000000000400810 <multiply_matrices>:
11 : multiply_matrices():
12 : }
13 :
14 : }
15 :
16 : void multiply_matrices()
17 : {
18 0.00 : 400810: pxor %xmm2,%xmm2
19 0.00 : 400814: mov $0x7e97c0,%edi
20 0.00 : 400819: mov %rdi,%r8
21 0.00 : 40081c: xor %esi,%esi
22 0.00 : 40081e: sub $0x7e97c0,%r8
23 0.00 : 400825: nopl (%rax)
24 0.13 : 400828: lea 0x6f5580(%rsi),%rax
25 0.00 : 40082f: lea 0x7e97c0(%rsi),%rcx
26 0.00 : 400836: mov %rdi,%rdx
27 0.00 : 400839: movaps %xmm2,%xmm1
28 0.00 : 40083c: nopl 0x0(%rax)
29 :
30 :
31 :
32 :
33 :
34 :
35 :
36 :
37 :
38 :
39 :
40 :
41 :
42 :
43 :
44 :
45 :
46 :
47 :
48 :
49 :
50 :
51 :
52 :
53 :
54 :
55 :
56 :
57 :
58 :
59 :
60 :
61 :
62 :
63 :
64 :
65 :
66 :
67 :
68 :
69 :
70 :
71 :
72 :
73 :
74 :
75 :
76 :
77 :
78 :
79 :
80 :
81 :
82 :
83 :
84 :
85 :
86 :
87 :
88 :
89 :
90 :
91 :
92 :
93 :
94 :
95 :
96 :
97 :
98 :
99 :
100 :
101 :
102 :
103 :
104 :
105 :
106 :
107 :
108 :
109 :
110 :
111 :
112 :
113 :
114 :
115 :
116 :
117 :
118 :
119 :
120 :
121 :
122 :
123 :
124 :
125 :
126 :
127 :
128 :
129 :
130 :
131 :
132 :
133 :
134 :
135 :
136 :
137 :
138 :
139 :
140 :
141 :
142 :
143 :
144 :
145 :
146 :
147 :
148 :
149 :
150 :
151 :
152 :
153 :
154 :
155 :
156 :
157 :
158 :
159 :
160 :
161 :
162 :
163 :
164 :
165 :
166 :
167 :
168 :
169 :
170 :
171 :
172 :
173 :
174 :
175 :
176 :
177 :
178 :
179 :
180 :
181 :
182 :
183 :
184 :
185 :
186 :
187 :
188 :
189 :
190 :
191 :
192 :
193 :
194 :
195 :
196 :
197 :
198 :
199 :
200 :
201 :
202 :
203 :
204 :
205 :
206 :
207 :
208 :
209 :
210 :
211 :
212 :
213 :
214 :
215 :
216 :
217 :
218 :
219 :
220 :
221 :
222 :
223 :
224 :
225 :
226 :
227 :
228 :
229 :
230 :
231 :
232 :
233 :
234 :
235 :
236 :
237 :
238 :
239 :
240 :
241 :
242 :
243 :
244 :
245 :
246 :
247 :
248 :
249 :
250 :
251 :
252 :
253 :
254 :
255 :
256 :
257 :
258 :
259 :
260 :
261 :
262 :
263 :
264 :
265 :
266 :
267 :
268 :
269 :
270 :
271 :
272 :
273 :
274 :
275 :
276 :
277 :
278 :
279 :
280 :
281 :
282 :
283 :
284 :
285 :
286 :
287 :
288 :
289 :
290 :
291 :
292 :
293 :
294 :
295 :
296 :
297 :
298 :
299 :
300 :
301 :
302 :
303 :
304 :
305 :
306 :
307 :
308 :
309 :
310 :
311 :
312 :
313 :
314 :
315 :
316 :
317 :
318 :
319 :
320 :
321 :
322 :
323 :
324 :
325 :
326 :
327 :
328 :
329 :
330 :
331 :
332 :
333 :
334 :
335 :
336 :
337 :
338 :
339 :
340 :
341 :
342 :
343 :
344 :
345 :
346 :
347 :
348 :
349 :
350 :
351 :
352 :
353 :
354 :
355 :
356 :
357 :
358 :
359 :
360 :
361 :
362 :
363 :
364 :
365 :
366 :
367 :
368 :
369 :
370 :
371 :
372 :
373 :
374 :
375 :
376 :
377 :
378 :
379 :
380 :
381 :
382 :
383 :
384 :
385 :
386 :
387 :
388 :
389 :
390 :
391 :
392 :
393 :
394 :
395 :
396 :
397 :
398 :
399 :
400 :
401 :
402 :
403 :
404 :
405 :
406 :
407 :
408 :
409 :
410 :
411 :
412 :
413 :
414 :
415 :
416 :
417 :
418 :
419 :
420 :
421 :
422 :
423 :
424 :
425 :
426 :
427 :
428 :
429 :
430 :
431 :
432 :
433 :
434 :
435 :
436 :
437 :
438 :
439 :
440 :
441 :
442 :
443 :
444 :
445 :
446 :
447 :
448 :
449 :
450 :
451 :
452 :
453 :
454 :
455 :
456 :
457 :
458 :
459 :
460 :
461 :
462 :
463 :
464 :
465 :
466 :
467 :
468 :
469 :
470 :
471 :
472 :
473 :
474 :
475 :
476 :
477 :
478 :
479 :
480 :
481 :
482 :
483 :
484 :
485 :
486 :
487 :
488 :
489 :
490 :
491 :
492 :
493 :
494 :
495 :
496 :
497 :
498 :
499 :
500 :
501 :
502 :
503 :
504 :
505 :
506 :
507 :
508 :
509 :
510 :
511 :
512 :
513 :
514 :
515 :
516 :
517 :
518 :
519 :
520 :
521 :
522 :
523 :
524 :
525 :
526 :
527 :
528 :
529 :
530 :
531 :
532 :
533 :
534 :
535 :
536 :
537 :
538 :
539 :
540 :
541 :
542 :
543 :
544 :
545 :
546 :
547 :
548 :
549 :
550 :
551 :
552 :
553 :
554 :
555 :
556 :
557 :
558 :
559 :
560 :
561 :
562 :
563 :
564 :
565 :
566 :
567 :
568 :
569 :
570 :
571 :
572 :
573 :
574 :
575 :
576 :
577 :
578 :
579 :
580 :
581 :
582 :
583 :
584 :
585 :
586 :
587 :
588 :
589 :
590 :
591 :
592 :
593 :
594 :
595 :
596 :
597 :
598 :
599 :
600 :
601 :
602 :
603 :
604 :
605 :
606 :
607 :
608 :
609 :
610 :
611 :
612 :
613 :
614 :
615 :
616 :
617 :
618 :
619 :
620 :
621 :
622 :
623 :
624 :
625 :
626 :
627 :
628 :
629 :
630 :
631 :
632 :
633 :
634 :
635 :
636 :
637 :
638 :
639 :
640 :
641 :
642 :
643 :
644 :
645 :
646 :
647 :
648 :
649 :
650 :
651 :
652 :
653 :
654 :
655 :
656 :
657 :
658 :
659 :
660 :
661 :
662 :
663 :
664 :
665 :
666 :
667 :
668 :
669 :
670 :
671 :
672 :
673 :
674 :
675 :
676 :
677 :
678 :
679 :
680 :
681 :
682 :
683 :
684 :
685 :
686 :
687 :
688 :
689 :
690 :
691 :
692 :
693 :
694 :
695 :
696 :
697 :
698 :
699 :
700 :
701 :
702 :
703 :
704 :
705 :
706 :
707 :
708 :
709 :
710 :
711 :
712 :
713 :
714 :
715 :
716 :
717 :
718 :
719 :
720 :
721 :
722 :
723 :
724 :
725 :
726 :
727 :
728 :
729 :
730 :
731 :
732 :
733 :
734 :
735 :
736 :
737 :
738 :
739 :
740 :
741 :
742 :
743 :
744 :
745 :
746 :
747 :
748 :
749 :
750 :
751 :
752 :
753 :
754 :
755 :
756 :
757 :
758 :
759 :
760 :
761 :
762 :
763 :
764 :
765 :
766 :
767 :
768 :
769 :
770 :
771 :
772 :
773 :
774 :
775 :
776 :
777 :
778 :
779 :
780 :
781 :
782 :
783 :
784 :
785 :
786 :
787 :
788 :
789 :
790 :
791 :
792 :
793 :
794 :
795 :
796 :
797 :
798 :
799 :
800 :
801 :
802 :
803 :
804 :
805 :
806 :
807 :
808 :
809 :
810 :
811 :
812 :
813 :
814 :
815 :
816 :
817 :
818 :
819 :
820 :
821 :
822 :
823 :
824 :
825 :
826 :
827 :
828 :
829 :
830 :
831 :
832 :
833 :
834 :
835 :
836 :
837 :
838 :
839 :
840 :
841 :
842 :
843 :
844 :
845 :
846 :
847 :
848 :
849 :
850 :
851 :
852 :
853 :
854 :
855 :
856 :
857 :
858 :
859 :
860 :
861 :
862 :
863 :
864 :
865 :
866 :
867 :
868 :
869 :
870 :
871 :
872 :
873 :
874 :
875 :
876 :
877 :
878 :
879 :
880 :
881 :
882 :
883 :
884 :
885 :
886 :
887 :
888 :
889 :
890 :
891 :
892 :
893 :
894 :
895 :
896 :
897 :
898 :
899 :
900 :
901 :
902 :
903 :
904 :
905 :
906 :
907 :
908 :
909 :
910 :
911 :
912 :
913 :
914 :
915 :
916 :
917 :
918 :
919 :
920 :
921 :
922 :
923 :
924 :
925 :
926 :
927 :
928 :
929 :
930 :
931 :
932 :
933 :
934 :
935 :
936 :
937 :
938 :
939 :
940 :
941 :
942 :
943 :
944 :
945 :
946 :
947 :
948 :
949 :
950 :
951 :
952 :
953 :
954 :
955 :
956 :
957 :
958 :
959 :
960 :
961 :
962 :
963 :
964 :
965 :
966 :
967 :
968 :
969 :
970 :
971 :
972 :
973 :
974 :
975 :
976 :
977 :
978 :
979 :
980 :
981 :
982 :
983 :
984 :
985 :
986 :
987 :
988 :
989 :
990 :
991 :
992 :
993 :
994 :
995 :
996 :
997 :
998 :
999 :

```

```

50 :         for (i = 0 ; i < MSIZE ; i++) {
51 :             for (j = 0 ; j < MSIZE ; j<←
52 :                 ++ ) {
53 :                     float sum = 0.0 ;
54 :                     for (k = 0 ; k <←
55 :                         MSIZE ; k++) {
56 :                             33.86 : 40085d: jne 400840 <←
57 :                             multiply_matrices+0x30>
58 :                             sum = sum + <←
59 :                             (matrix_a[i][k] * matrix_b[k][j]) ;
60 :                             }
61 :                             matrix_r[i][j] = sum<←
62 :                             ;
63 :                             0.00 : 40085f: movss %xmm1,0x601340(%r8,%<←
64 :                             rsi,1)
65 :                             0.13 : 400869: add $0x4,%rsi
66 :                             void multiply_matrices()
67 :                             {
68 :                                 int i, j, k ;
69 :                                 for (i = 0 ; i < MSIZE ; i++) {
70 :                                     for (j = 0 ; j < MSIZE ; j<←
71 :                                         ++ ) {
72 :                                             0.00 : 40086d: cmp $0x7d0,%rsi
73 :                                             0.00 : 400874: jne 400828 <←
74 :                                             multiply_matrices+0x18>
75 :                                             0.00 : 400876: add $0x7d0,%rdi
76 :                                             void multiply_matrices()
77 :                                             {
78 :                                                 int i, j, k ;
79 :                                                 for (i = 0 ; i < MSIZE ; i++) {
80 :                                                     0.00 : 40087d: cmp $0x8dda00,%rdi
81 :                                                     0.00 : 400884: jne 400819 <←
82 :                                                     multiply_matrices+0x9>
83 :                                                     0.00 : 400886: repz retq

```

Caso pretendamos retirar o código na linguagem de alto nível, mantendo apenas a associação a código máquina com percentagem de tempo de execução, podemos recorrer ao seguinte comando:

```
Percent | Source code & Disassembly of naive for cpu<←
clock
```

```

1 Percent | Source code & Disassembly of naive for cpu<←
2 clock
3 :
4 :
5 :
6 : Disassembly of section .text:
7 :
8 : 0000000000400810 <multiply_matrices>:
9 : multiply_matrices():
10 0.00 : 400810: pxor %xmm2,%xmm2
11 0.00 : 400814: mov $0x7e97c0,%edi
12 0.00 : 400819: mov %rdi,%r8
13 0.00 : 40081c: xor %esi,%esi
14 0.00 : 40081e: sub $0x7e97c0,%r8
15 0.00 : 400825: nopl (%rax)
16 0.13 : 400828: lea 0x6f5580(%rsi),%rax
17 0.00 : 40082f: lea 0x7e97c0(%rsi),%rcx
18 0.00 : 400836: mov %rdi,%rdx
19 0.00 : 400839: movaps %xmm2,%xmm1
20 0.00 : 40083c: nopl 0x0(%rax)
21 24.97 : 400840: movss (%rdx),%xmm0
22 0.00 : 400844: add $0x7d0,%rax
23 0.00 : 40084a: mulss -0x7d0(%rax),%xmm0
24 18.56 : 400852: add $0x4,%rdx
25 22.35 : 400856: cmp %rcx,%rax
26 0.00 : 400859: addss %xmm0,%xmm1
27 33.86 : 40085d: jne 400840 <←
28 multiply_matrices+0x30>
29 0.00 : 40085f: movss %xmm1,0x601340(%r8,%<←
30 rsi,1)

```

```

29 0.13 : 400869: add $0x4,%rsi
30 0.00 : 40086d: cmp $0x7d0,%rsi
31 0.00 : 400874: jne 400828 <←
32 multiply_matrices+0x18>
33 0.00 : 400876: add $0x7d0,%rdi
34 0.00 : 40087d: cmp $0x8dda00,%rdi
35 0.00 : 400884: jne 400819 <←
36 multiply_matrices+0x9>
37 0.00 : 400886: repz retq

```

Pela análise da execução dos dois comandos anteriores, temos então a resposta há muito aguardada. O nosso primeiro trabalho de detective dá-se por terminado com descoberta da "principal responsável":

```
sum = sum + (matrix_a[i][k] * matrix_b[k][j]) ;
```

4.4. Método de profiling

Recolhidos os dados e descoberto o "culpado" resta-nos saber se o "modus operandi" da recursão ao comando PERF foi o mais indicado.

O PERF recorre a sampling estatístico para recolher os dados de profiling. Cada amostra é recolhida com base num tempo constante, tendo por base o **Linux time clock** para medir também a passagem do mesmo. Quando termina o intervalo de tempo constante é gerada uma interrupção e o PERF determina o trabalho que estava a ser realizado pelo CPU no momento da interrupção, registando por exemplo metadados do CPU, program counter, e software, produzindo aquilo que é chamado de um **sample**, sendo este escrito num buffer, posteriormente escrito no ficheiro **perf.data**.

Dado que nos iremos basear no sampling estatístico para inferir conclusões acerca do comportamento de um kernel necessitamos de ser cautelosos quanto ao número de samples a recolher. **Se por um lado o excesso de amostras levará a uma interferência nos dados a serem medidos, a falta delas levará à incerteza.** Poderemos estar a concluir erradamente que uma porção de código é a responsável pela degradação de performance devido à falta de amostras que nos apontem noutro sentido.

Uma forma de aumentar o número de amostras é especificando a frequência sobre a qual as mesmas serão recolhidas, através do seguinte comando:

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf record -e <←
cpu-clock --freq=8000 ./naive
```

```

1 [a57816@compute-431-l ESC_FLAME_GRAPH]$ perf record -e <←
2 cpu-clock --freq=8000 ./naive
3 [ perf record: Woken up 1 times to write data ]
4 [ perf record: Captured and wrote 0.070 MB perf.data <←
5 (1636 samples) ]

```

A execução do comando anterior resultou na recolha do dobro de amostras (1636 amostras vs as anteriores 830) em comparação com a medição efectuada em 4.3.1, devido também ao aumento em dobro da frequência – 8000 vs 4000, tal como pode ser confirmado pela execução do comando seguinte:

```
[a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf report --stdio --show-nr-samples --dsos=naive
```

```
1 [a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf report --stdio --show-nr-samples --dsos=naive
2 # To display the perf.data header info, please use --header/--header-only options.
3 #
4 # dso: naive
5 # Samples: 1K of event 'cpu-clock'
6 # Event count (approx.): 1636
7 #
8 # Overhead      Samples  Command      Symbol
9 # .....
10 #
11 94.74%          1550  naive        [.]
12 1.47%           24    naive        [.]
13 0.06%           1     naive        [.] rand@plt
```

Confirmamos portanto que se mantém o ponto crítico da nossa aplicação. 94.74% do tempo da mesma foi despendido no método **multiply_matrices**, tal como pretendíamos confirmar.

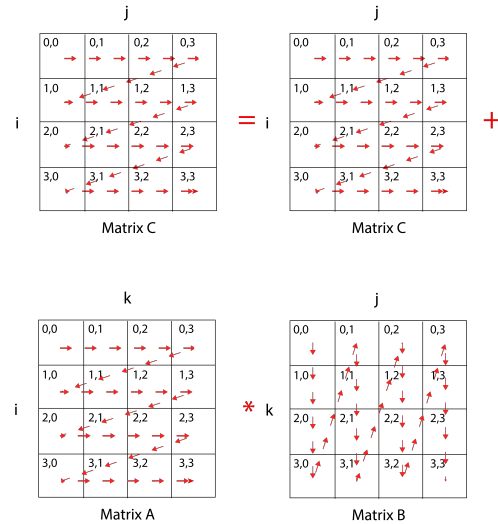
5. Parte 2

Encontrado o ponto alvo na "investigação" anterior, restamos resolver o problema. Atente nas seguintes transcrições de código, pertencentes nomeadamente à versão naive e interchange:

- Naive:

```
1 void multiply_matrices()
2 {
3     int i, j, k ;
4
5     // Textbook algorithm
6     for (i = 0 ; i < MSIZE ; i++) {
7         for (j = 0 ; j < MSIZE ; j++) {
8             float sum = 0.0 ;
9             for (k = 0 ; k < MSIZE ; k++) {
10                 sum = sum + (matrix_a[i][k] * matrix_b[k][j]) ;
11             }
12             matrix_r[i][j] = sum ;
13         }
14     }
15 }
```

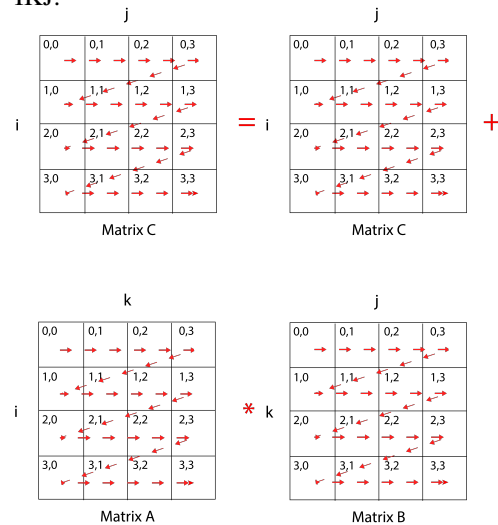
Ilustração do acesso às posições da matriz via ordem IJK:



- Interchange:

```
1 void multiply_matrices()
2 {
3     int i, j, k ;
4
5     // Loop nest interchange algorithm
6     for (i = 0 ; i < MSIZE ; i++) {
7         for (k = 0 ; k < MSIZE ; k++) {
8             for (j = 0 ; j < MSIZE ; j++) {
9                 matrix_r[i][j] = matrix_r[i][j] +
10                     (matrix_a[i][k] * matrix_b[k][j]) ;
11             }
12         }
13     }
14 }
```

Ilustração do acesso às posições da matriz via ordem IKJ:



A ordem dos ciclos I-J-K na versão naive resulta na constante invalidação das linhas de cache que correspondam a dados da matriz B. Ora, este problema irá acentuar-se em caso de necessidade de acesso à memória principal, aumento a latência da recolha dos mesmos e provocando um atraso ainda mais visível quando comparada com a versão I-K-J.

A seguinte tabela resume vários contadores de performance para as duas versões do kernel – naive e interchange,

para o dataset menor em estudo:

Tabela 2: Performance events (naive vs. interchange) para o nó compute-431

# EVENT NAME	NAIVE	INTERCHANGE
cpu-cycles	535187277	399561216
instructions	1044692763	1152237507
cache-references	8196140	429971
cache-misses	36522	43034
branch-instructions	126101720	132065934
branch-misses	258384	249858
bus-cycles	0	0
L1-dcache-loads	246027409	253077242
L1-dcache-load-misses	56436199	7577858
L1-dcache-stores	9973628	128034804
L1-dcache-store-misses	322982	106020
LLC-loads	7391770	262810
LLC-load-misses	2671	1001
LLC-stores	218407	69369
LLC-store-misses	18512	0
dTLB-load-misses	2239	950
iTLB-load-misses	446	9
branch-loads	0	0
branch-load-misses	129163483	129898962
	5688441	5560030

Da análise da tabela 2 podemos confirmar que:

- O número de CPU cycles é menor para a versão interchange, reflectindo-se num menor tempo de solução.
- O número de instruções é aproximadamente o mesmo dado para este dataset a invalidação dos dados não ser reflectida em leitura da memória principal mas apenas da LLCACHE (mais à frente neste relatório iremos analisar a influência da leitura da memória principal no tempo de execução).
- O número de "cache-references" é muito maior para o caso naive (tal como esperado pela invalidação dos dados de matrizes).
- O número de LLC loads é também muito maior para o caso naive (pela razão enumerada anteriormente).

Podemos então concluir que é este problema de padrão de acesso aos dados que resulta na degradação de performance da versão naive quando comparada com a versão interchange.

Analisemos agora as métricas compostas, obtidas com base nos valores apresentados:

Tabela 3: Performance rates (naive vs. interchange) para o nó compute-431

RATIO OR RATE	NAIVE	INTERCHANGE
Elapsed time (seconds)	0.2041	0.1597
Instructions per cycle	1.95 IPC	2.88 IPC
L1 cache miss ratio	22,9389 %	2,9942 %
L1 cache miss rate PTI	54,0218	6,5766
L3 cache miss ratio	0,0361	0,3808 %
Data TLB miss ratio	0,00027	0,0022
Data TLB miss rate PTI	0,0021	0,0008
Branch mispredict ratio	0,002	0,0019
Branch mispredict rate PTI	0,2473	0,2168

A degradação de performance das versões, para o caso do dataset menor, não é traduzida numa exponenciação dos tempos totais para a solução, muito devido aos dados apesar de estarem a ser acedidos de forma ineficiente, manterem-se contidos na LLCACHE. Focaremos a nossa análise do-ravante no dataset maior.

6. Parte 3

6.1. Período de amostragem e frequência de amostragem

Dado que pretendemos encontrar os hotspots de ambas as versões, iremos usar a amostragem baseada em cpu-cycles. Desta forma, porções da aplicação que consumam mais tempo terão um maior número de amostras registadas. No entanto, este tipo de amostragem tem um preço – pode ser demasiado pesada na avaliação de desempenho da aplicação.

Temos que ter em conta que cada máquina tem um número de contadores máximos e fixados a certo tipo de eventos específicos. Ora, quando requeremos um grande número de eventos no profiling, é necessária a recorrência a multiplexagem, que gasta tempo de computação. Essa mesma multiplexagem poderá ter um peso demasiado elevado na avaliação de desempenho.

Analisemos o tempo que demoram as versões large_naive e large_interchange sem qualquer tipo de ferramenta de amostragem por forma a podermos validar os resultados da avaliação de desempenho futura. Denote que por forma a validar os resultados foram realizadas 50 medições, sendo o valor apresentado o K Best (sendo K = 3) :

EVENT NAME	L. NAIVE	L. INTERCHANGE
Elapsed time (seconds)	65.61	10.43

Temos portanto agora uma base de referência também para estas versões que recorrem ao maior dataset em estudo. Voltemos portanto ao profiling das versões com os maiores datasets.

6.2. Análise comparativa dos evento de hardware para as versões large_naive e large_interchange, para o maior dataset em teste

Da análise de execução para as versões large_naive e large_interchange, foi produzida a tabela 4, que resume vários contadores de performance para as duas versões do kernel, para o dataset maior em estudo.

Denote que o overhead do processo de profiling não foi abusivo, dado que em comparação com os tempos sem qualquer tipo de profiling o aumento do tempo total da aplicação para ambos os casos rondou os 5%. Esta validação era necessária por forma a garantirmos uma comparação justa entre ambas as versões.

Tabela 4: Performance events (naive vs. interchange) para o nó compute-431

# EVENT NAME	NAIVE	INTERCHANGE
Elapsed time	68.119634290	11.053004702
Measurement Overhead	$\frac{68.12}{65.61} - 1 = 4\%$	$\frac{11.05}{10.43} - 1 = 6\%$
instructions	38 376 000 000	41 776 400 000
cycles	78390700000	12 384 800 000
cache-references	4 744 200 000	14 400 000
cache-misses	4 008 300 000	11 200 000
LLC-loads	4 815 000 000	14 800 000
LLC-load-misses	4 073 600 000	14400000
dTLB-load-misses	1100000	100000
branches	3923900000	3847500000
branch-misses	2200000	1900000

Ambas as versões executam aproximadamente o mesmo número de instruções, sendo que é para a versão mais otimizada que registamos um ligeiro acréscimo no valor medido. Dado que o kernel **large_interchange** executa em metade do tempo de o valor do número de ciclos para esta versão é também 1/2 do número de ciclos necessários para a execução da versão **large_naive**.

É com base nos valores de cache miss e cache reference que podemos inferir o principal obstáculo da versão **large_naive**. A forma pouco eficiente de acesso aos dados da matriz B reduz drasticamente a performance do mesmo. Em comparação, a versão **large_interchange** acede de uma forma correcta aos dados, possibilitando ganhos de performance relativamente à versão original.

A tabela 5 apresenta as métricas compostas com base nos valores apresentados na tabela 4.

Podemos comprovar que o número de instruções por segundo para a versão **large_interchange** é sete vezes maior relativamente à versão original e pouco eficiente – **large_naive**.

Tal como poderá constatar a percentagem de cache misses para ambas as versões é extremamente alta, o que é compreensível dado o tamanho dos datasets. Contudo, é no número de **cache-references** que a versão **large_naive** tem valores **330 vezes superiores** aos registados para a versão **large_interchange**.

Tabela 5: Performance rates (large_naive vs. large_interchange) para o nó compute-431

RATIO OR RATE	NAIVE	INTERCHANGE
Elapsed time (seconds)	0.2041	0.1597
Instructions per cycle	1.95 IPC	2.88 IPC
L3 cache miss ratio	84.6023 %	97.2973 %
Data TLB miss ratio	0.0232	0.6944
Data TLB miss rate PTI	0.0287	0.0024
Branch mispredict ratio	0.0561	0.0494
Branch mispredict rate PTI	0.0573	0.0455

6.2.1. Confirmação de metodologia de medição via metadados. Tal como referido na seção 4.3.2 a ferramenta PERF guarda os meta-dados de medição juntamente com o ficheiro de output "perf.data". O seguinte comando imprime a informação relativa às métricas recolhidas:

```
perf evlist --input=large_interchange_output
```

```
1 cpu-cycles
2 instructions
```

Com a adição da opção -F podemos visualizar a frequência de sampling:

```
perf evlist -F --input=large_interchange_output
```

```
1 cpu-cycles: sample_freq=100000
2 instructions: sample_freq=100000
```

Na sequência do já demonstrado na secção 4.3.2 podemos visualizar apenas as características nas quais a medição foi realizada através do comando:

```
perf report --header-only --input=large_interchange_output
```

```
1 # =====
2 # captured on: Sun Jun 5 20:57:12 2016
3 # hostname : compute-431-l.local
4 # os release : 2.6.32-279.14.1.el6.x86_64
5 # perf version : 4.0.0
6 # arch : x86_64
7 # nrcpus online : 24
8 # nrcpus avail : 24
9 # cpudesc : Intel(R) Xeon(R) CPU X5650 @ 2.67GHz
10 # cpuid : GenuineIntel,6,44,2
11 # total memory : 12319324 kB
12 # cmdline : /share/jade/SOFT/perf/perf record -e cpu-cycles, instructions -c 100000 --output=large_interchange_output ./large_interchange
13 # event : name = cpu-cycles, type = 0, config = 0x0, config1 = 0x0, config2 = 0x0, excl_usr = 0, excl_kern = 0, excl_host = 0, excl_guest = 1, precise_ip = 0, attr_mmap2 = 0, attr_mmap = 1, attr_mmap_data = 0, id = { 3517, 3518, 3519, 3520, 3521, 3522, 3523, 3524, 3525, 3526, 3527, 3528, 3529, 3530, 3531, 3532, 3533, 3534, 3535, 3536, 3537, 3538, 3539, 3540 }
14 # event : name = instructions, type = 0, config = 0x1, config1 = 0x0, config2 = 0x0, excl_usr = 0, excl_kern = 0, excl_host = 0, excl_guest = 1, precise_ip = 0, attr_mmap2 = 0, attr_mmap = 0, attr_mmap_data = 0, id = { 3541, 3542, 3543, 3544, 3545, 3546, 3547, 3548, 3549, 3550, 3551, 3552, 3553, 3554, 3555, 3556, 3557, 3558, 3559, 3560, 3561, 3562, 3563, 3564 }
15 # HEADER_CPU_TOPOLOGY info available, use -I to display
16 # HEADER_NUMA_TOPOLOGY info available, use -I to display
17 # pmu mappings: cpu = 4, tracepoint = 2, software = 1
18 # =====
19 #
```

6.3. Análise comparativa do número de amostras por evento de hardware para as versões large_naive e large_interchange

Dado que estamos a recolher amostras estatísticas da execução das duas versões da aplicação, inferindo a partir dessas mesmas amostras conclusões relativas à sua performance, é necessário garantir que o processo de amostragem e recolha é válido para o tamanho do dataset.

Sabemos que o intervalo de confiança dos valores apresentados será tão maior quanto o número de amostras recolhido. Dessa mesma premissa apresenta-se a tabela 6 que relaciona precisamente o número de amostras por métrica.

Tabela 6: Sampling mode: large_naive vs. large_interchange para o nó compute-431

# EVENT NAME	L. NAIVE	L. INTERCHANGE
instructions	383K samples	417K samples
cycles	783K samples	123K samples
cache-references	47K samples	144 samples
cache-misses	40K samples	112 samples
LLC-loads	48K samples	148 samples
LLC-load-misses	40K samples	144 samples
dTLB-load-misses	11 samples	1 samples
branches	39K samples	38K samples
branch-misses	22 samples	19 samples

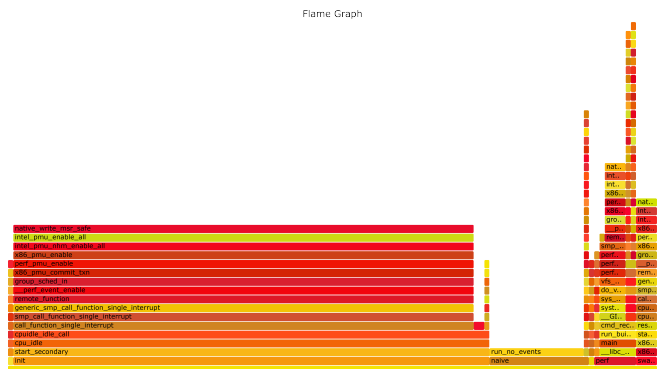
As métricas de maior peso para a nossa análise são precisamente instructions, cycles, cache-references, cache-misses, LLC-loads, e LLC-load-misses, sendo que é nestas onde se registam precisamente um maior número de amostras.

7. Análise gráfica com recurso a Flame Graphs

Os Flame Graphs podem ser produzidos através dos dados recolhidas pela ferramenta PERF, possibilitando a análise gráfica dos mesmos. Cada Rectângulo representa um método, e a largura desse mesmo rectângulo em comparação com a da imagem, demonstra o tempo de computação despendido na mesma. Podemos confirmar a pilha das chamadas também visualmente recorrendo a esta ferramenta. De seguida apresentam-se os Flame Graphs obtidos para as 2 versões da aplicação para os dois datasets em estudo, sendo os mesmos obtidos recorrendo à seguinte sequência de comandos, apenas alterando o executável e nome do ficheiro respectivamente:

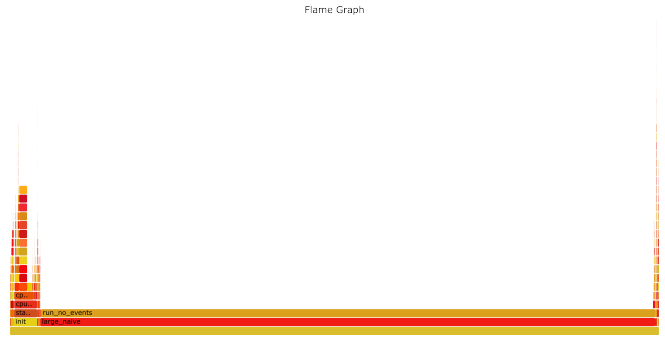
```
perf record -ag -F 99 ./large_naive
perf script | /share/jade/soft/FlameGraph/stackcollapse-
-perf.pl > out.perf-folded
cat out.perf-folded | /share/jade/soft/FlameGraph/
flamegraph.pl > perf_kernel_large_naive.svg
```

Figura 1: Flame Graph kernel **naive**



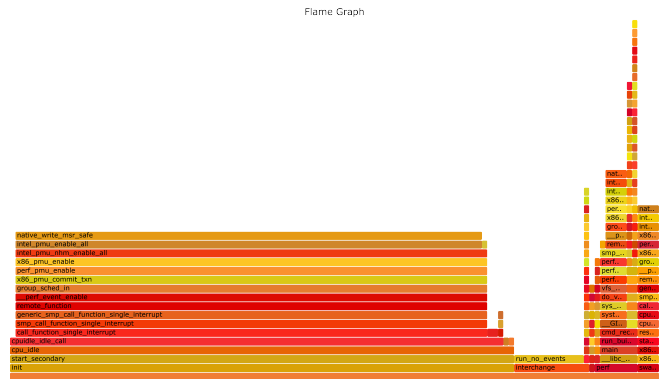
7.0.1. Flame Graph kernel naive.

Figura 2: Flame Graph kernel **large_naive**



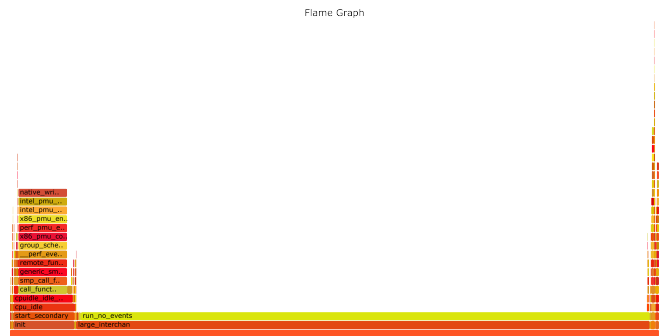
7.0.2. Flame Graph kernel large_naive.

Figura 3: Flame Graph kernel **interchange**



7.0.3. Flame Graph kernel interchange.

Figura 4: Flame Graph kernel **large_interchange**



7.0.4. Flame Graph kernel large_interchange. subsecçãoAnálise dos resultados gráficos O primeiro comando produziu o profiling de CPU stacks, a uma frequência de 99 Hertz. Ora, o comportamento visual vem confirmar o analisado durante todo este relatório. Dado que abordamos de uma forma introdutória os Flame Graphs será interessante recorrer a este tipo de ferramenta para outros kernels desenvolvidos em outras unidades curriculares.

8. Conclusão

Tal como mencionado no início do presente caso de estudo a ferramenta PERF mostra-se bastante útil e complementar a outras ferramentas já estudadas em termos de funcionalidades quando necessitamos de analisar a performance relativa de processos. Ou seja, no contexto da computação paralela será extremamente interessante recorrer a esta ferramenta na execução de algoritmos paralelos, sendo possível recorrendo a scripts como por exemplo o Flame Graph para adicionar funcionalidades à mesma.

Este foi apenas um trabalho introdutório mas permitiu demonstrar a capacidade de recolher e ao mesmo tempo tratar dados de todo um sistema extremamente complexo e vasto com apenas uma ferramenta. O caso de estudo ultrapassa portanto os resultados obtidos pelos comandos e aplicações utilizadas, prendendo-se uma vez mais com o desenvolvimento de capacidade prática no uso da ferramenta, e envolvimento com métodos de tratamento de grandes volumes de dados, e análise de métricas de sistemas de computação de alta performance.

Retrata sobretudo a capacidade analisar funcionalidades disponibilizadas e a sua correta aplicação na resolução de problemas de computação tendo sempre em conta o mínimo de alteração possível na performance dos kernels/sistemas a analisar.