

Introdução ao PERF

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

Filipe Oliveira

Departamento de Informática

Universidade do Minho

Email: a57816@alunos.uminho.pt

1. Introdução – Contextualização da utilização da ferramenta PERF

O presente trabalho prático surge como um “diário de bordo” relativo ao acompanhamento de um conjunto de 3 partes do tutorial: “**PERF tutorial: Finding execution hot spots**”.¹

A parte 1 descreve a utilização da ferramenta PERF como meio de identificação e análise de hotspots de execução em kernels. Cobre maioritariamente os comandos PERF mais básico, e respectivas opções e eventos de medição de performance relativa.

A parte 2 introduz eventos que permitem a medição de contadores de hardware e demonstra a utilização do PERF nesse sentido. Discute ainda a importância de certas métricas de performance e a sua influência relativa na performance geral de uma aplicação.

A parte 3 recorre aos eventos de contadores de hardware para, uma vez mais e desta vez de uma forma mais aprofundada, analisar os hotspots de um kernel/aplicação.

Todo o tutorial assenta em duas versões de uma aplicação de multiplicação de matrizes: **beginitemize**

Naive: Algoritmo de multiplicação de matrizes com acessos a memória “column-wise”.

Interchange: Algoritmo de multiplicação de matrizes com acessos a memória “row-wise”.

A primeira versão serão, com seria de esperar, apresentará problemas de performance relacionados com o acesso à memória. É importante também realçar que as aplicações são ambas compiladas com a versão de compilador GCC 4.9.0., com as seguintes flags de compilação -O2 -ggdb -g”. Denote que a flag de otimização -O3, que activa um conjunto de otimizações de código² está desativada, muito pela necessidade de manter o código (apesar de pouco eficiente) legível e com capacidade de ser compreendido pelo programador.

2. Caracterização do Hardware do ambiente de testes

Tão importante com especificar as versões de código e as propriedades é especificar os ambientes de teste nos quais pretendemos realizar as benchmarks.

Através da análise do hardware disponível no Search6³, uma das nossas plataformas de teste, foi seleccionado o nó do tipo compute-431, sendo a disponibilidade global do mesmo e o correcto funcionamento do perf os principais factores. Na tabela 1 encontram-se especificadas as principais características dos sistemas em teste.

Tabela 1: Características de Hardware do nó 431

Sistema	compute-431
# CPUs	2
CPU	Intel® Xeon® X5650
Arquitectura de Processador	Nehalem
# Cores por CPU	6
# Threads por CPU	12
Freq. Clock	2.66 GHz
Cache L1	192KB (32KB por Core)
Cache L2	1536KB (256KB por Core)
Cache L3	12288KB (partilhada)
Ext. Inst. Set	SSE4.2
#Memory Channels	3
Memória Ram Disponível	48GB
Peak Memory BW Fab. CPU	32 GB/s

3. Determinação do tamanho dos datasets

Talvez dos pontos mais importantes presente na tabela de caracterização de hardware seja o tamanho da Cache L3, dado que ambos os algoritmos dependem massivamente de leitura/escrita na memória principal. Assim sendo, se pretendemos realçar as penalizações resultantes de diferentes tipos de acesso aos dados devemos realizar dois tipos distintos de testes:

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

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

3 matrizes

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

1. <http://sandsoftwaresound.net/perf/perf-tutorial-hot-spots/>

2. -finline-functions, -funswitch-loops, -fpredictive-commoning, -fgcse-after-reload, -ftree-vectorize, -fvect-cost-model, -ftree-partial-pre and -fipa-cp-clone options

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.

Dorante, 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 commando:

```
[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:

```
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
```

```
Performance counter stats for './naive':  
  
201.696023      cpu-clock (msec)  
  
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
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ←  
stat -e cpu-clock,faults ./naive  
  
Performance counter stats for './naive':  
  
196.395754      cpu-clock (msec)  
844            faults  
  
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.

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ←  
record -e cpu-clock,faults ./naive  
[ perf record: Woken up 1 times to write data ]  
[ perf record: Captured and wrote 0.046 MB perf. data (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 prof.data (ficheiro default no qual é gravada a informação):

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

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ←  
report --stdio  
# To display the perf.data header info, please use --header/--header-only options.  
#  
# Samples: 808 of event 'cpu-clock'  
# Event count (approx.): 808  
#  
# Overhead Command Shared Object Symbol  
# .....  
#  
94.68% naive naive [.] multiply_matrices  
2.10% naive naive [.] initialize_matrices  
1.36% naive libc-2.12.so [.] __random
```

```

0.87% naive libc-2.12.so      [.] ←
__random_r
0.25% naive [kernel.kallsyms] [k] ←
__mem_cgroup_commit_charge
0.25% naive naive           [.] ←
rand@plt
0.12% naive [kernel.kallsyms] [k] ←
__call_rcu
0.12% naive [kernel.kallsyms] [k] ←
__mem_cgroup_uncharge_common
0.12% naive libc-2.12.so      [.] ←
malloc
0.12% naive libc-2.12.so      [.] rand

# Samples: 22 of event 'faults'
# Event count (approx.): 1157
#
# Overhead Command Shared Object Symbol
# .....
#
#
63.61% naive naive           [.] ←
initialize_matrices
29.56% naive libc-2.12.so      [.] ←
__init_cpu_features
5.62% naive ld-2.12.so         [.] ←
dl_main
0.52% naive ld-2.12.so         [.] ←
_dl_setup_hash
0.26% naive ld-2.12.so         [.] ←
_dl_start
0.17% naive [kernel.kallsyms] [k] ←
__clear_user
0.17% naive libc-2.12.so      [.] ←
__strchr_sse2
0.09% naive ld-2.12.so         [.] ←
_start

```

Existem outras alternativas a esse mesmo comando sendo esta:

```
[a57816@compute-431-l 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:

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

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

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

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

```

[a57816@search6 ESC_FLAME_GRAPH]$ perf report ←
stdio --sort comm,dso
[kernel.kallsyms] with build id 09←
b53ce5dfe09c7blad9473e1356d7d922512e27 not ←
found, continuing without symbols
# To display the perf.data header info, please ←
use --header/--header-only options.
#
# Samples: 808 of event 'cpu-clock'
# Event count (approx.): 808
#

```

```

# Overhead Command Shared Object
# .....
#
97.03% naive naive
2.48% naive libc-2.12.so
0.50% naive [kernel.kallsyms]

# Samples: 22 of event 'faults'
# Event count (approx.): 1157
#
# Overhead Command Shared Object
# .....
#
63.61% naive naive
29.73% naive libc-2.12.so
6.48% naive ld-2.12.so
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
```

```

[a57816@search6 ESC_FLAME_GRAPH]$ perf report ←
stdio --sort comm,dso --header
[kernel.kallsyms] with build id 09←
b53ce5dfe09c7blad9473e1356d7d922512e27 not ←
found, continuing without symbols
# =====
# captured on: Mon May 23 23:00:56 2016
# hostname : compute-431-l.local
# os release : 2.6.32-279.14.1.el6.x86_64
# perf version : 4.0.0
# arch : x86_64
# nrcpus online : 24
# nrcpus avail : 24
# cpudesc : Intel(R) Xeon(R) CPU X5650 @ 2.67GHz
# cpuid : GenuineIntel,6,44,2
# total memory : 12319324 kB
# cmdline : /share/jade/SOFT/perf/perf record -e←
cpu-clock,faults ./naive
# event : name = cpu-clock, type = 1, config = 0←
x0, config1 = 0x0, config2 = 0x0, excl_usr ←
= 0, excl_kern = 0, excl_host = 0, ←
excl_guest = 1, precise_ip = 0, attr_mmap2 ←
= 0, attr
# event : name = faults, type = 1, config = 0x2,←
config1 = 0x0, config2 = 0x0, excl_usr = ←
0, excl_kern = 0, excl_host = 0, excl_guest←
= 1, precise_ip = 0, attr_mmap2 = 0, ←
attr_mm
# HEADER_CPU_TOPOLOGY info available, use -I to ←
display
# HEADER_NUMA_TOPOLOGY info available, use -I to←
display
# pmu mappings: cpu = 4, tracepoint = 2, ←
software = 1

```

```
# =====
#
# Samples: 808 of event 'cpu-clock'
# Event count (approx.): 808
#
# Overhead Command Shared Object
# .....
#
# 97.03% naive naive
# 2.48% naive libc-2.12.so
# 0.50% naive [kernel.kallsyms]
#
# Samples: 22 of event 'faults'
# Event count (approx.): 1157
#
# Overhead Command Shared Object
# .....
#
# 63.61% naive naive
# 29.73% naive libc-2.12.so
# 6.48% naive ld-2.12.so
# 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
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ←
report --stdio --dsos=naive,libc-2.12.so
# To display the perf.data header info, please ←
use --header/--header-only options.
#
# Samples: 808 of event 'cpu-clock'
# Event count (approx.): 808
#
# Overhead Command Shared Object Symbol
# .....
#
# 94.68% naive naive [...] ←
# multiply_matrices
# 2.10% naive naive [...] ←
# initialize_matrices
# 1.36% naive libc-2.12.so [...] __random
# 0.87% naive libc-2.12.so [...] ←
# __random_r
# 0.25% naive naive [...] rand@plt
# 0.12% naive libc-2.12.so [...] malloc
# 0.12% naive libc-2.12.so [...] rand
#
# Samples: 22 of event 'faults'
# Event count (approx.): 1157
#
# Overhead Command Shared Object Symbol
# .....
#
# 63.61% naive naive [...] ←
# initialize_matrices
# 29.56% naive libc-2.12.so [...] ←
# __init_cpu_features
```

```
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 $512 \times 512 \times (4 \text{ bytes}) = 1024 \text{ KB}$, multiplicado por 3 (3 matrizes), representa um total de 3072 KB, que representa um total de páginas $\frac{3072 \text{ KB}}{4 \text{ KB por página}} = 768$. Ora, 63.61% do total de **page faults** contabilizadas representa $0.6361 \times 1157 = 735,9677 = 736 \text{ páginas}$, como espectado.

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:

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ←
annotate --stdio --dsos=naive --symbol=←
multiply_matrices
Percent | Source code & Disassembly of naive for←
cpu-clock
:
:
:
: Disassembly of section .text:
:
: 0000000000400810 <multiply_matrices>:
: multiply_matrices():
:
: }
: }
:
: void multiply_matrices()
: {
0.00 : 400810: pxor %xmm2,%xmm2
0.00 : 400814: mov $0x7e97c0,%edi
0.00 : 400819: mov %rdi,%r8
```

```

0.00 : 40081c: xor %esi,%esi
0.00 : 40081e: sub $0x7e97c0,%r8
0.00 : 400825: nopl (%rax)
0.13 : 400828: lea 0x6f5580(%rsi)←
,%rax
0.00 : 40082f: lea 0x7e97c0(%rsi)←
,%rcx
0.00 : 400836: mov %rdi,%rdx
0.00 : 400839: movaps %xmm2,%xmm1
0.00 : 40083c: nopl 0x0(%rax)
:
: for (i = 0 ; i < MSIZE ; i++)←
: {
: for (j = 0 ; j < ←
MSIZE ; j++) {
: float sum = ←
0.0 ;
: for (k = 0 ; ←
k < MSIZE ; k++) {
: sum =←
sum + (matrix_a[i][k] * matrix_b[←
k][j]) ;
24.97 : 400840: movss (%rdx),%xmm0
0.00 : 400844: add $0x7d0,%rax
0.00 : 40084a: mulss -0x7d0(%rax),%←
xmm0
18.56 : 400852: add $0x4,%rdx
: int i, j, k ;
:
: for (i = 0 ; i < MSIZE ; i++)←
: {
: for (j = 0 ; j < ←
MSIZE ; j++) {
: float sum = ←
0.0 ;
: for (k = 0 ; ←
k < MSIZE ; k++) {
22.35 : 400856: cmp %rcx,%rax
: sum =←
sum + (matrix_a[i][k] * matrix_b[←
k][j]) ;
0.00 : 400859: addss %xmm0,%xmm1
: int i, j, k ;
:
: for (i = 0 ; i < MSIZE ; i++)←
: {
: for (j = 0 ; j < ←
MSIZE ; j++) {
: float sum = ←
0.0 ;
: for (k = 0 ; ←
k < MSIZE ; k++) {
33.86 : 40085d: jne 400840 <←
multiply_matrices+0x30>
: sum =←
sum + (matrix_a[i][k] * matrix_b[←
k][j]) ;
: }
: matrix_r[i][j]←
] = sum ;
0.00 : 40085f: movss %xmm1,0x601340←
(%r8,%rsi,1)
0.13 : 400869: add $0x4,%rsi
: void multiply_matrices()
: {
: int i, j, k ;
:
: for (i = 0 ; i < MSIZE ; i++)←
: {
: for (j = 0 ; j < ←
MSIZE ; j++) {
0.00 : 40086d: cmp $0x7d0,%rsi
0.00 : 400874: jne 400828 <←
multiply_matrices+0x18>
0.00 : 400876: add $0x7d0,%rdi
:
: void multiply_matrices()

```

```

: {
: int i, j, k ;
:
: for (i = 0 ; i < MSIZE ; i++)←
: {
0.00 : 40087d: cmp $0x8dda00,%rdi
0.00 : 400884: jne 400819 <←
multiply_matrices+0x9>
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
```

```

Percent | Source code & Disassembly of naive ←
for cpu-clock
:
:
: Disassembly of section .text:
:
: 0000000000400810 <multiply_matrices>:
: multiply_matrices():
0.00 : 400810: pxor %xmm2,%xmm2
0.00 : 400814: mov $0x7e97c0,%edi
0.00 : 400819: mov %rdi,%r8
0.00 : 40081c: xor %esi,%esi
0.00 : 40081e: sub $0x7e97c0,%r8
0.00 : 400825: nopl (%rax)
0.13 : 400828: lea 0x6f5580(%rsi)←
,%rax
0.00 : 40082f: lea 0x7e97c0(%rsi)←
,%rcx
0.00 : 400836: mov %rdi,%rdx
0.00 : 400839: movaps %xmm2,%xmm1
0.00 : 40083c: nopl 0x0(%rax)
24.97 : 400840: movss (%rdx),%xmm0
0.00 : 400844: add $0x7d0,%rax
0.00 : 40084a: mulss -0x7d0(%rax),%←
xmm0
18.56 : 400852: add $0x4,%rdx
22.35 : 400856: cmp %rcx,%rax
0.00 : 400859: addss %xmm0,%xmm1
33.86 : 40085d: jne 400840 <←
multiply_matrices+0x30>
0.00 : 40085f: movss %xmm1,0x601340←
(%r8,%rsi,1)
0.13 : 400869: add $0x4,%rsi
0.00 : 40086d: cmp $0x7d0,%rsi
0.00 : 400874: jne 400828 <←
multiply_matrices+0x18>
0.00 : 400876: add $0x7d0,%rdi
0.00 : 40087d: cmp $0x8dda00,%rdi
0.00 : 400884: jne 400819 <←
multiply_matrices+0x9>
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]) ;
```

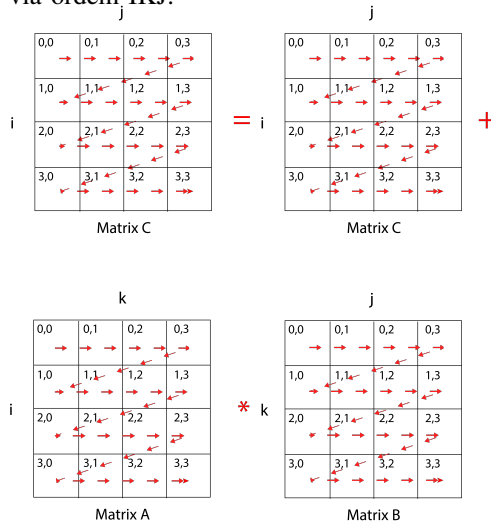


```

for (j = 0 ; j < MSIZE ; j++) {
    matrix_r[i][j] = matrix_r[i][j] ←
        +
        (matrix_a[i][k] * matrix_b[k][j]) ;
}
}
}

```

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:

8 Tabela 2: Performance events (naive vs. interchange) para o
9 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
dTLB-store-misses	446	9
iTLB-load-misses	0	0
branch-loads	129163483	129898962
branch-load-misses	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 LLCA-CHE (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 com 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 doravante 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.

Analiseemos 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)	10.43	65.61

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.

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

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

# EVENT NAME	NAIVE	INTERCHANGE
Elapsed time		
instructions	38376000000	41776400000
cycles	78390700000	12384800000
cache-references	4744200000	14400000
cache-misses	4008300000	11200000
LLC-loads	4815000000	14800000
LLC-load-misses	4073600000	14400000
dTLB-load-misses	1100000	100000
branches	3923900000	3847500000
branch-misses	2200000	1900000

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

Tabela 5: Sampling mode: large_naive vs. large_interchange para o nó compute-431

# EVENT NAME	NAIVE	INTERCHANGE
Elapsed time		
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

7. Conclusão

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ← 1
evlist -F
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ← 1
evlist -F
cpu-clock: sample_freq=8000 2
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ← 1
list | wc -l
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ← 1
list | wc -l
784 2
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ← 1
stat -e cpu-cycles,instructions ./naive
```

```
[a57816@compute-431-l ESC_FLAME_GRAPH]$ perf ← 1
stat -e cpu-cycles,instructions ./naive
Performance counter stats for './naive':
2
3
4
5
6
7
8
516402889      cpu-cycles
908042685      instructions ←
#      1.76  insns per ←
cycle
0.200239282 seconds time elapsed
```



```
[a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf stat -e cpu-cycles,instructions,cache-references,cache-misses,branch-instructions,branch-misses,bus-cycles,L1-dcache-loads,L1-dcache-load-misses,L1-dcache-stores,L1-dcache-store-misses,LLC-loads,LLC-load-misses,LLC-stores,LLC-store-misses,dTLB-load-misses,dTLB-store-misses,iTLB-load-misses,branch-loads,branch-load-misses ./naive
```

```
[a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf stat -e cpu-cycles,instructions,cache-references,cache-misses,branch-instructions,branch-misses,bus-cycles,L1-dcache-loads,L1-dcache-load-misses,L1-dcache-stores,L1-dcache-store-misses,LLC-loads,LLC-load-misses,LLC-stores,LLC-store-misses,dTLB-load-misses,dTLB-store-misses,iTLB-load-misses,branch-loads,branch-load-misses ./naive
```

Performance counter stats for './naive':

```
535187277      cpu-cycles      [20.25%]
1044692763     instructions      # 1.95  insns per cycle
8196140        cache-references
[25.37%]
36522          cache-misses      # 0.446 % of all cache refs
126101720      branch-instructions
[27.67%]
258384         branch-misses      # 0.20% of all branches
0              bus-cycles
246027409      L1-dcache-loads
[22.50%]
56436199       L1-dcache-load-misses # 22.94% of all L1-dcache hits
9973628        L1-dcache-stores
[22.27%]
322982         L1-dcache-store-misses
[22.16%]
7391770        LLC-loads
[22.05%]
2671           LLC-load-misses # 0.04% of all LLC-cache hits
218407         LLC-stores
[10.92%]
18512          LLC-store-misses
[10.86%]
2239           dTLB-load-misses
[16.21%]
446            dTLB-store-misses
[21.51%]
0              iTLB-load-misses
[21.41%]
129163483      branch-loads
```

```
[21.22%]
5688441        branch-load-misses
[20.73%]
0.210071197 seconds time elapsed
```

```
[a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf stat -e cpu-cycles,instructions,cache-references,cache-misses,branch-instructions,branch-misses,bus-cycles,L1-dcache-loads,L1-dcache-load-misses,L1-dcache-stores,L1-dcache-store-misses,LLC-loads,LLC-load-misses,LLC-stores,LLC-store-misses,dTLB-load-misses,dTLB-store-misses,iTLB-load-misses,branch-loads,branch-load-misses ./interchange
```

```
[a57816@compute-431-1 ESC_FLAME_GRAPH]$ perf stat -e cpu-cycles,instructions,cache-references,cache-misses,branch-instructions,branch-misses,bus-cycles,L1-dcache-loads,L1-dcache-load-misses,L1-dcache-stores,L1-dcache-store-misses,LLC-loads,LLC-load-misses,LLC-stores,LLC-store-misses,dTLB-load-misses,dTLB-store-misses,iTLB-load-misses,branch-loads,branch-load-misses ./interchange
```

Performance counter stats for './interchange':

```
399561216      cpu-cycles      [21.57%]
1152237507     instructions      # 2.88  insns per cycle
429971         cache-references
[27.74%]
43034          cache-misses      # 10.009 % of all cache refs
132065934      branch-instructions
[28.66%]
249858         branch-misses      # 0.19% of all branches
0              bus-cycles
253077242      L1-dcache-loads
[21.87%]
7577858        L1-dcache-load-misses # 2.99% of all L1-dcache hits
128034804      L1-dcache-stores
[20.60%]
106020         L1-dcache-store-misses
[20.38%]
262810         LLC-loads
[22.23%]
1001           LLC-load-misses # 0.38% of all LLC-cache hits
69369          LLC-stores
[10.96%]
0              LLC-store-misses
[10.88%]
950            dTLB-load-misses
[16.21%]
```

9	dTLB-store-misses ↵	21	
	[21.47%]		↵
0	iTLB-load-misses ↵	22	
	[21.33%]		↵
129898962	branch-loads ↵	23	↵
	[21.19%]		
5560030	branch-load-misses ↵	24	↵
	[21.05%]		
0.159788849	seconds time elapsed	25	
		26	