

**МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ ФЕДЕРАЦИИ  
Федеральное государственное автономное образовательное учреждение  
высшего образования Национальный исследовательский Нижегородский  
государственный университет им. Н. И. Лобачевского**

**Институт информационных технологий, математики и механики**

**Кафедра дифференциальных уравнений, математического и численного  
анализа**

Направление подготовки

**02.04.02. Фундаментальная информатика и информационные технологии**

Направленность образовательной программы

**магистерская программа «Компьютерная графика и моделирование живых и  
технических систем»**

**Отчёт**

**Отчет по лабораторной работе**

на тему:

**«Промахи кэша»**

Квалификация (степень)

**магистр**

Форма обучения

**очная**

**Выполнил:**

студент гр. 381706 – 2м

**Бабаев Иван Владимирович**

Нижний Новгород

2018

## Оглавление

Введение.....	3
Кэш промахи .....	3
Профилирование тестовой программы. ....	3
Заключение.....	6
Приложение 1.....	7

## Введение

Кэш микропроцессора — кэш (сверхоперативная память), используемый микропроцессором компьютера для уменьшения среднего времени доступа к компьютерной памяти. Является одним из верхних уровней иерархии памяти. Кэш использует небольшую, очень быструю память (обычно типа SRAM), которая хранит копии часто используемых данных из основной памяти. Если большая часть запросов в память будет обрабатываться кэшем, средняя задержка обращения к памяти будет приближаться к задержкам работы кэша. Данные между кэшем и памятью передаются блоками фиксированного размера, также называемые линиями кэша (англ. cache line) или блоками кэша. Большинство современных микропроцессоров для компьютеров и серверов имеют как минимум три независимых кэша: кэш инструкций для ускорения загрузки машинного кода, кэш данных для ускорения чтения и записи данных и буфер ассоциативной трансляции (TLB) для ускорения трансляции виртуальных (логических) адресов в физические, как для инструкций, так и для данных. Кэш данных часто реализуется в виде многоуровневого кэша (L1, L2, L3).

## Кэш промахи

Cache Miss (промах кэша) случается, когда запрашиваемые данные отсутствуют в кэше и их нужно подгружать из основного источника.

Виды промахов:

- Промах по чтению из кэша инструкций. Обычно дает очень большую задержку, поскольку процессор не может продолжать исполнение программы (по крайней мере, текущего потока исполнения) и вынужден простаивать в ожидании загрузки инструкции из памяти.
- Промах по чтению из кэша данных. Обычно дает меньшую задержку, поскольку инструкции, не зависящие от запрошенных данных, могут продолжать исполняться, пока запрос обрабатывается в основной памяти. После получения данных из памяти можно продолжать исполнение зависимых инструкций.
- Промах по записи в кэш данных. Обычно дает наименьшую задержку, поскольку запись может быть поставлена в очередь и последующие инструкции практически не ограничены в своих возможностях. Процессор может продолжать свою работу, кроме случаев промаха по записи с полностью заполненной очередью.

## Профилирование тестовой программы.

Для демонстрации проблемы кэш-промахов была написана небольшая программа на языке C++ (смотри приложение 1), которая прибавляет ко всем элементам матрицы заданную константу и повторяет заданное число раз (30000).

Профилирование производилось на тестовой машине со следующими характеристиками:

Операционная система: Windows 10

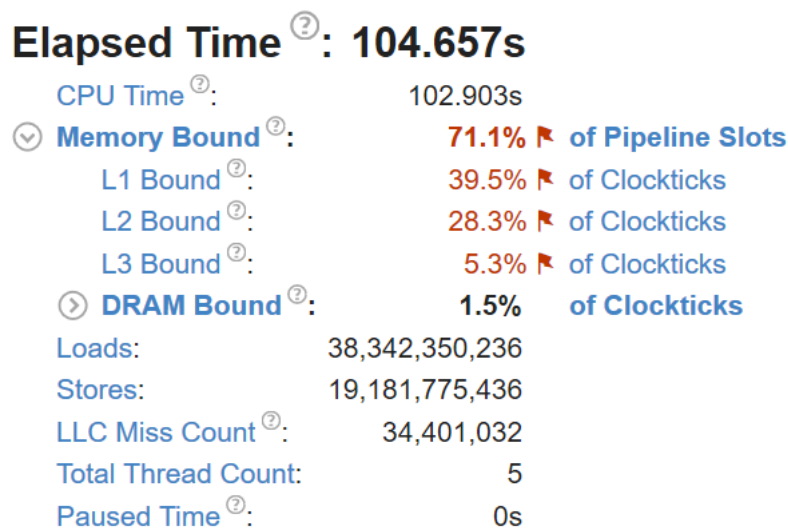
CPU: Intel Core i5-3470, Caches: L1 – 64kb x4, L2- 256Kb x4, L3 – 1024Kb x6.

GPU: NVIDIA GeForce GTX 970, 4058Mb

Оперативная память: 16384Mb DDR3

Профилирование проводилось с помощью инструмента Intel VTune Amplifier Memory Access Tool, который предоставляет набор метрик для выявления проблем, связанных с доступом к памяти.

При запуске анализа работы приложения, Memory Access Tool запускает профилируемое приложение и собирает информацию о его работе. По завершении работы формируется отчет, указывающий критические точки и показывающий общую производительность:



Рассмотрим Memory bound – метрики эффективности операций с памятью:

- **L1 Bound** - процент времени, которое процессор ожидал данные с кэша L1.
- **L2 Bound** - процент времени, которое процессор ожидал данные с кэша L2.
- **L3 Bound** – процент времени, которое процессор ожидал данные с кэша L3.
- **DRAM Bound** – процент времени, которое процессор ожидал данные с основной памяти.

В данном случае наблюдается проблема при работе с кэш-памятью. Memory Access Tool указывает места кода, в которых возникает проблема:

Src	Source	CPU Time: Total	CPU Time: Self	Memory Bound: Total	Memory Bound: Self			
					L1 Bound	L2 Bound	L3 Bound	DRAM Bound
4	<code>#include &lt;memory&gt;</code>							
5								
6	<code>void AddConstToMatrix(int** x, int n, int m, int c)</code>							
7	<code>{</code>							
8	<code>for (int j = 0; j &lt; m; j++)</code>							
9	<code>for (int i = 0; i &lt; n; i++)</code>							
10	<code>x[i][j] = c + x[i][j];</code>							
11	<code>}</code>							
12								
13	<code>void AddConstToMatrixB(int** x, int n, int m, int c)</code>							
14	<code>{</code>							
15	<code>for (int ii = 0; ii &lt; n; ii += 8)</code>							
16	<code>for (int jj = 0; jj &lt; m; jj += 8)</code>							
17	<code>for (int i = ii; i &lt; ii + 8; i++)</code>							
18	<code>for (int j = jj; j &lt; jj + 8; j++)</code>							
19	<code>x[i][j] = x[i][j] + c;</code>							
20	<code>}</code>							
21								
22	<code>int main()</code>							
23	<code>{</code>							
24	<code>int n = 800;</code>							
25	<code>int m = 800;</code>							
26	<code>int **x;</code>							
27	<code>x = new int*[n];</code>							
28	<code>for (int i = 0; i &lt; n; i++) {</code>							
29	<code>x[i] = new int[m];</code>							
30	<code>}</code>							
31								
32	<code>for (int i = 0; i &lt; n; i++) {</code>							
33	<code>for (int j = 0; j &lt; m; j++) {</code>							
34	<code>x[i][j] = 1-i;</code>							
35	<code>}</code>							
36	<code>}</code>							
37								
38	<code>for (int k = 0; k &lt; 30000; k++)</code>							
39	<code>AddConstToMatrix(x, n, m, 2);</code>	102.522s	102.522s	71.3%	39.5%	28.4%	5.3%	1.5%
40								
41	<code>for (int i = 0; i &lt; n; i++) {</code>							
42	<code>delete[] x[i];</code>							
43	<code>}</code>							
44	<code>delete[] x;</code>							
45								
46	<code>return 0;</code>							
47	<code>}</code>							

```

void AddConstToMatrix(int** x, int n, int m, int c)
{
    for (int j = 0; j < m; j++)
        for (int i = 0; i < n; i++)
            x[i][j] = c + x[i][j];
}

```

В данном случае проблема возникает в функции добавления константы ко всем элементам матрицы. А именно, в размещении памяти. Если получать доступ к кэшу в строковом порядке, то будет использоваться вся память кэша. Если идти по столбцам, то кэш закончится прежде, чем память сможет быть повторно используема.



Двумерный массив представляется в памяти в виде:

row,col	0,0	0,1	0,2
	1,0	1,1	1,2
	2,0	2,1	2,2

			0,0	0,1	0,2	1,0	1,1	1,2	2,0	2,1	2,2			
--	--	--	-----	-----	-----	-----	-----	-----	-----	-----	-----	--	--	--

Поэтому, если идти по столбцам, то промежутки между данными будут велики и память кэша быстро закончится. Размер входных данных (размер матрицы):  $800 * 800 * \text{sizeof}(\text{int}) = 2500\text{Kb}$ .

Исправим наш пример:

```
void AddConstToMatrix(int** x, int n, int m, int c)
{
    for (int i = 0; i < n; i++)
        for (int j = 0; j < m; j++)
            x[i][j] = c + x[i][j];
}
```

Запустим Memory Access Tool:

<b>Elapsed Time</b> <sup>?</sup> : 11.614s		
CPU Time <sup>?</sup> :	10.948s	
Memory Bound <sup>?</sup> :	1.4%	of Pipeline Slots
L1 Bound <sup>?</sup> :	0.8%	of Clockticks
L2 Bound <sup>?</sup> :	0.0%	of Clockticks
L3 Bound <sup>?</sup> :	1.1%	of Clockticks
DRAM Bound <sup>?</sup> :	0.0%	of Clockticks
DRAM Bandwidth Bound <sup>?</sup> :	0.0%	of Elapsed Time
Loads:	37,174,315,196	
Stores:	18,594,157,808	
LLC Miss Count <sup>?</sup> :	0	
Total Thread Count:	4	
Paused Time <sup>?</sup> :	0s	

Обращение по строкам дало значительный прирост производительности и сделало процент ожидания данных с L3 достаточно малым, и обнулило ожидание данных с кэша L2.

## Заключение

Тестовая программа показала, что неверное обращение к элементам массива плохо сказывается на чтении/записи данных и влечет потерю производительности из-за долгого чтения. Правильное чтение/запись ускорили работу программы в 9 раз и обнулили промахи кэша L2, практически обнулили промахи кэша L3.

# Литература

- [1] Intel® VTune™ Amplifier, <https://software.intel.com/en-us/vtune-amplifier-help>
- [2] Department of Computer Science at University of Verona, Cache performance presentation, <http://www.di.univr.it/documenti/OccorrenzaIns/matdid/matdid566638.pdf>
- [3] HOW FAST CAN YOU GO – OPTIMIZING MEMORY CACHE PERFORMANCE  
CLAIRE CATES DISTINGUISHED DEVELOPER CLAIRE.CATES@SAS.COM,  
[https://www.cmg.org/wp-content/uploads/2015/10/memory\\_cache.pdf](https://www.cmg.org/wp-content/uploads/2015/10/memory_cache.pdf)
- [4] Why software developers should care about CPU caches, <https://medium.com/software-design/why-software-developers-should-care-about-cpu-caches-8da04355bb8a>

## Приложение 1

Код программы:

```
#include <time.h>
#include <stdio.h>
#include <stdlib.h>
#include <memory>

void AddConstToMatrix(int** x, int n, int m, int c)
{
    for (int i = 0; i < n; i++)
        for (int j = 0; j < m; j++)
            x[i][j] = c + x[i][j];
}

int main()
{
    int n = 800;
    int m = 800;
    int **x;
    x = new int *[n];
    for (int i = 0; i < n; i++) {
        x[i] = new int[m];
    }

    for (int i = 0; i < n; i++) {
        for (int j = 0; j < m; j++) {
            x[i][j] = i;
        }
    }

    for (int k = 0; k < 30000; k++)
        AddConstToMatrix(x, n, m, 2);

    for (int i = 0; i < n; i++) {
        delete[] x[i];
    }
    delete[] x;

    return 0;
}
```