

WROCŁAW UNIVERSITY OF SCIENCE AND TECHNOLOGY  
FACULTY OF FUNDAMENTAL PROBLEMS OF  
TECHNOLOGY

---

FIELD: Computer Science  
SPECIALIZATION: Computer Security

**MASTER OF SCIENCE THESIS**

Timing Attack Resistant Implementation of RSA  
on GPU.

AUTHOR:  
Krzysztof Hamerski

SUPERVISOR:  
dr Maciej Gębala

GRADE:

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Environment Specification . . . . .	2
<b>2</b>	<b>Theory</b>	<b>4</b>
2.1	RSA . . . . .	4
2.2	Side channel attacks . . . . .	5
2.2.1	Timing attacks . . . . .	5
2.3	Parallel programming on CUDA . . . . .	6
2.3.1	Shared Memory . . . . .	6
<b>3</b>	<b>Implementation</b>	<b>7</b>
3.1	Parallel equals . . . . .	12
3.2	Parallel compare . . . . .	12
3.3	Parallel bit length . . . . .	12
3.4	Parallel left shift . . . . .	12
3.5	Parallel right shift . . . . .	12
3.6	Parallel add . . . . .	12
3.7	Parallel subtract . . . . .	12
3.8	Parallel multiply . . . . .	12
3.9	Parallel modulo reduction . . . . .	12
3.10	Parallel multiply modulo . . . . .	12
3.11	Parallel power modulo . . . . .	12
<b>4</b>	<b>Testing</b>	<b>13</b>
4.1	CUDA experiments . . . . .	13
4.1.1	Occupancy . . . . .	13
4.1.2	Theoretical occupancy . . . . .	13
4.1.3	Achieved occupancy . . . . .	15
4.1.4	Occupancy charts . . . . .	15
4.1.5	Instruction statistics . . . . .	17
4.1.6	Branch statistics . . . . .	21
4.1.7	Issue efficiency . . . . .	24
4.1.8	Pipe utilization . . . . .	27
4.1.9	Memory statistics . . . . .	30
4.2	Equals . . . . .	32
	<b>References</b>	<b>32</b>

# Chapter 1

## Introduction

Public key cryptography is the key factor in providing secure communication between two parties. Fast development of distributed system requiring not only security, but also integrity and non-repudiation has pushed cryptography to the limit. Since 1978[6] most commonly used cryptosystem is RSA, which provides asymmetric encryption, as well as generation of digital signatures. The security of RSA is mainly base on the bitwise key length. As computational power of modern CPUs arises, the minimal bit length of RSA key gets significantly bigger to provide sufficient security. At least 4096 bits long keys are considered secure nowadays. This leads to very high workload required to perform encryption/decryption. RSA is mainly based on modular arithmetics and simple computations become infeasible for moder CPUs, when dealing with so large integers.

One of the solution to this problem is to parallelize. GPGPU[9] (for General-Purpose computing on the Graphics Processing Unit) enables the use of GPU for parallel computation other than graphics. GPUs are designed to perform computations in parallel. Since every PC is equipped with some kind of GPU, one can easily exploit its capabilities. NVIDIA has made it even more accessible by creating CUDA (Compute Unified Device Architecture)[7]. It is a parallel computing platform and API, which exposes GPUs true potential.

### 1.1 Environment Specification

This project was developed in Microsoft Visual Studio 2015[5] with NSight plugin and CUDA API.[7] Source code is written mainly in C++ and inline assembly - PTX.[8] Table 1.1 more precisely illustrates GPU specification on which the program and all tests are run.

Table 1.1: Device specification

Device name:	GeForce GTX 960
CUDA Driver Version:	9.0
CUDA Runtime Version:	8.0
CUDA Capability version number:	5.2
Total amount of global memory:	4096 MBytes (4294967296 bytes)
GPU Max Clock rate:	1304 MHz (1.30 GHz)
Total amount of shared memory per block:	49152 bytes
Warp size:	32

Table below presents full platform and system information.

Table 1.2: PC specification

Operating System:	Windows 10 Education 64-bit
Motherboard:	Gigabyte Technology Co., Ltd. P55A-UD4
CPU:	Intel(R) Core(TM) i5 CPU 750 @ 2.67GHz (4 CPUs), 2.7GHz
RAM Memory: :	4096MB

# Chapter 2

## Theory

This chapter basically presents minimal amount of theory required to fully understand the concepts, problems and solutions which were applied and implemented within this project.

### 2.1 RSA

The RSA algorithm was created by Ronald Rivest, Adi Shamir, and Leonard Adleman in 1970s. The security of the algorithm is based on the problem of integer factorization.

- Key generation[4]

- Choose two large and distinct prime numbers  $p$  and  $q$ .
- Compute the modulus  $n$

$$n = pq$$

- Compute

$$\lambda(n) = lcm(\lambda(p), \lambda(q)) = lcm(p-1, q-1)$$

where  $\lambda$  is the Carmichael Totient Function[2]. This value must be kept in private.

- Choose an integer  $e$  such that:

$$1 < e < \lambda(n)$$

and

$$gcd(e, \lambda(n)) = 1$$

This means  $e$  and  $\lambda(n)$  are coprime.[10]

- Choose value  $d$  for decryption and solve for  $d$ :

$$de \equiv 1(mod \lambda(n))$$

$e$  is used for encryption. Usually this is done the other way around.  $e$  is chosen first so it has short bit length and small Hamming weight,[3] which provides for much faster encryption. In most cases  $e = 3, 5$  or  $7$ . High security is provided if larger number is used. Most common is Fermat's four:  $e = 2^{16} + 1 = 65537 = 0x10001$

– Public key is then:

$$(e, n)$$

and the private key:

$$(d, n)$$

- Encryption

For message  $m$  the ciphertext  $c$  is in relation:

$$c \equiv m^e \bmod n$$

- Decryption

$$m \equiv c^d \bmod n \equiv (m^e)^d \bmod n \equiv m \bmod n$$

## 2.2 Side channel attacks

Side-channel attacks are very powerful attacks against cryptographic implementations. They take the advantage of implementation flaws, gathering informations about almost everything that can be measured during execution time. If no care is taken, side-channel attacks can be used to compromise virtually any mathematically secure system.

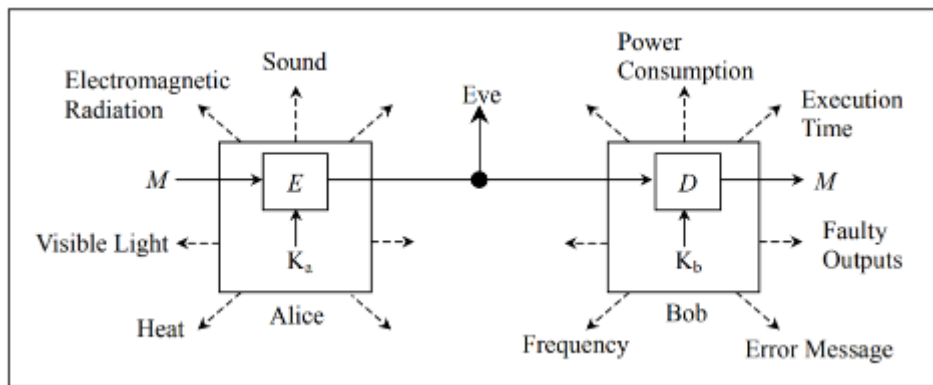


Figure 2.1: An example of cryptosystem including side-channel information leakage

Side-channel attacks aim to retrieve secret data from a cryptographic system by observing factors outside the normal computation. Power consumption, execution time and even electromagnetic radiation, sound, light and heat can leak some information about data being computed. Running statistical analysis on this partial information makes even the most sophisticated systems vulnerable to breakage.

### 2.2.1 Timing attacks

Implementations of cryptographic algorithms often perform computations in non-constant time, due to performance optimizations. If such operations involve secret parameters, these timing variations can leak some information and, provided enough knowledge of the implementation is at hand, a careful statistical analysis could even lead to the total recovery of these secret parameters.[11]

## 2.3 Parallel programming on CUDA

### 2.3.1 Shared Memory

"Threads within a block can cooperate by sharing data through shared memory and by synchronizing their execution to coordinate memory accesses. Because it is on-chip, shared memory has much higher bandwidth and much lower latency than local or global memory. For that shared memory is equivalent to a user-managed cache: The application explicitly allocates and accesses it. A typical programming pattern is to stage data coming from device memory into shared memory, process the data in shared memory while sharing the shared view of the data across the threads of a block, and write the results back to device memory.

To achieve high bandwidth, shared memory is divided into equally-sized memory modules, called banks, which can be accessed simultaneously. Any memory read or write request made of  $n$  addresses that fall in  $n$  distinct memory banks can therefore be serviced simultaneously, yielding an overall bandwidth that is  $n$  times as high as the bandwidth of a single module. However, if two addresses of a memory request fall in the same memory bank, there is a bank conflict and the access has to be serialized. The hardware splits a memory request with bank conflicts into as many separate conflict-free requests as necessary, decreasing throughput by a factor equal to the number of separate memory requests. If the number of separate memory requests is  $n$ , the initial memory request is said to cause  $n$ -way bank conflicts. To get maximum performance, it is therefore important to minimize bank conflicts.

Shared memory has 32 banks that are organized such that successive 32bit words map to successive banks. A shared memory request for a warp does not generate a bank conflict between two threads that access any address within the same 32bit word (even though the two addresses fall in the same bank).

For devices of compute capability 3.x, shared memory has 32 banks with two addressing modes that can be configured using `cudaDeviceSetSharedMemConfig()`. Each bank has a bandwidth of 64 bits per clock cycle. In 64bit mode, successive 64bit words map to successive banks. A shared memory request for a warp does not generate a bank conflict between two threads that access any sub-word within the same 64bit word (even though the addresses of the two sub-words fall in the same bank). In 32bit mode (default), successive 32bit words map to successive banks. A shared memory request for a warp does not generate a bank conflict between two threads that access any sub-word within the same 32bit word or within two 32bit words whose indices  $i$  and  $j$  are in the same 64word aligned segment (i.e., a segment whose first index is a multiple of 64) and such that  $j=i+32$  (even though the addresses of the two sub-words fall in the same bank)."[1]

# Chapter 3

## Implementation

The code was written in C++ language. In order to minimize data movement between host and the device, most of logics and computation are executed fully on GPU. Implementation of RSA encryption requires only one function - modular exponentiation of a multi precision integer. Implemented Big Integer class provides much more functions, to properly handle data and validate results. Code listing below shows the header of class BigInteger.

```
class BigInteger
{
//fields
public:

// 4096 bits
static const int ARRAY_SIZE = 128;

private:
// Magnitude array in little endian order.
// Most-significant int is mag[length-1].
// Least-significant int is mag[0].
// Allocated on the device.
unsigned int* deviceMagnitude;

// same array allocated on the host
// provides faster access if nothing was changed
unsigned int* hostMagnitude;

// flag indicating if hostMagnitude matches deviceMagnitude
bool upToDate;

// Device wrapper instance different for every integer
// to provide parallel execution
DeviceWrapper* deviceWrapper;

// methods
public:
BigInteger();
BigInteger(const BigInteger& x);
BigInteger(unsigned int value);
~BigInteger();
const unsigned int& operator[](int index);

// factory
static BigInteger* fromHexString(const char* string);
```



```

static BigInteger* createRandom(int bitLength);

// setters, getters
void set(const BigInteger& x);
unsigned int* getDeviceMagnitude(void) const;

// arithmetics
void add(const BigInteger& x);
void subtract(const BigInteger& x);
void multiply(const BigInteger& x);
void square(void);
void mod(const BigInteger& modulus);
void multiplyMod(const BigInteger& x, const BigInteger& modulus);
void squareMod(const BigInteger& modulus);
void powerMod(BigInteger& exponent, const BigInteger& modulus);

// logics
void shiftLeft(int bits);
void shiftRight(int bits);

// extras
bool equals(const BigInteger& value) const;
int compare(const BigInteger& value) const;
int getBitwiseLengthDifference(const BigInteger& value) const;
int getBitwiseLength(void) const;
int getLSB(void) const;
bool testBit(int bit);
void synchronize(void);
char* toHexString(void);
void print(const char* title);

//timer
void startTimer(void);
unsigned long long stopTimer(void);

// async calls
// must call synchronize to read from
void modAsync(const BigInteger& modulus);
void multiplyModAsync(const BigInteger& x, const BigInteger& modulus);
void squareModAsync(const BigInteger& modulus);

private:

void setMagnitude(const unsigned int* magnitude);
void clear(void);
void updateDeviceMagnitiude(void);
void updateHostMagnitiude(void);
static unsigned int random32(void);

/*
Parses hex string to unsigned int type.
Accepts both upper and lower case, no "0x" at the beginning.
E.g.: 314Da43F
*/
static unsigned int parseUnsignedInt(const char* hexString);
};

```

Listing 3.1: BigInteger.h

Big Integer class handles mathematical logic, validates input / output, and provides comfortable and readable interface.

The intermediary between Big Integer and GPU is another class - Device Wrapper. It handles GPU kernel launches, synchronization and data movement. Class definition is listed below.

```

class DeviceWrapper
{
private:
    // main stream for kernel launches
    cudaStream_t mainStream;

    // lauch config
    dim3 block_1, block_2, block_4;
    dim3 thread_warp, thread_2_warp, thread_4_warp;

    // 4 ints to help store results
    int* deviceWords;

    // auxiliary arrays
    unsigned int* device4arrays;
    unsigned int* device128arrays;
    unsigned int* deviceArray;

    unsigned long long* deviceStartTime;
    unsigned long long* deviceStopTime;

public:
    DeviceWrapper();
    ~DeviceWrapper();

    // sync
    unsigned int* init(int size) const;
    unsigned int* init(int size, const unsigned int* initial) const;
    void updateDevice(unsigned int* device_array, const unsigned int*
        host_array, int size) const;
    void updateHost(unsigned int* host_array, const unsigned int*
        device_array, int size) const;
    void free(unsigned int* device_x) const;

    // extras
    void clearParallel(unsigned int* device_x) const;
    void cloneParallel(unsigned int* device_x, const unsigned int* device_y)
        const;
    int compareParallel(const unsigned int* device_x, const unsigned int*
        device_y) const;
    bool equalsParallel(const unsigned int* device_x, const unsigned int*
        device_y) const;
    int getLSB(const unsigned int* device_x) const;
    int getBitLength(const unsigned int* device_x) const;
    void synchronize(void);

    // measure time
    void startClock(void);
    unsigned long long stopClock(void);

```

```

// logics
void shiftLeftParallel(unsigned int* device_x, int bits) const;
void shiftRightParallel(unsigned int* device_x, int bits) const;

// arithmetics
void addParallel(unsigned int* device_x, const unsigned int* device_y)
    const;
void subtractParallel(unsigned int* device_x, const unsigned int*
    device_y) const;
void multiplyParallel(unsigned int* device_x, const unsigned int*
    device_y) const;
void squareParallel(unsigned int* device_x) const;
void squareParallelAsync(unsigned int* device_x) const;
void modParallel(unsigned int* device_x, unsigned int* device_m) const;
void modParallelAsync(unsigned int* device_x, unsigned int* device_m)
    const;
void multiplyModParallel(unsigned int* device_x, const unsigned int*
    device_y, const unsigned int* device_m) const;
void multiplyModParallelAsync(unsigned int* device_x, const unsigned int*
    * device_y, const unsigned int* device_m) const;
void squareModParallel(unsigned int* device_x, const unsigned int*
    device_m) const;
void squareModParallelAsync(unsigned int* device_x, const unsigned int*
    device_m) const;

private:
void inline addParallelWithOverflow(unsigned int* device_x, const
    unsigned int* device_y, int blocks) const;
};

```

Listing 3.2: DeviceWrapper.h

Project also contains Test class to validate computations, measure times and simulate encryption. RSA class contains single "encrypt" function, which encrypts provided value. Full class diagram is presented on figure 3.1.

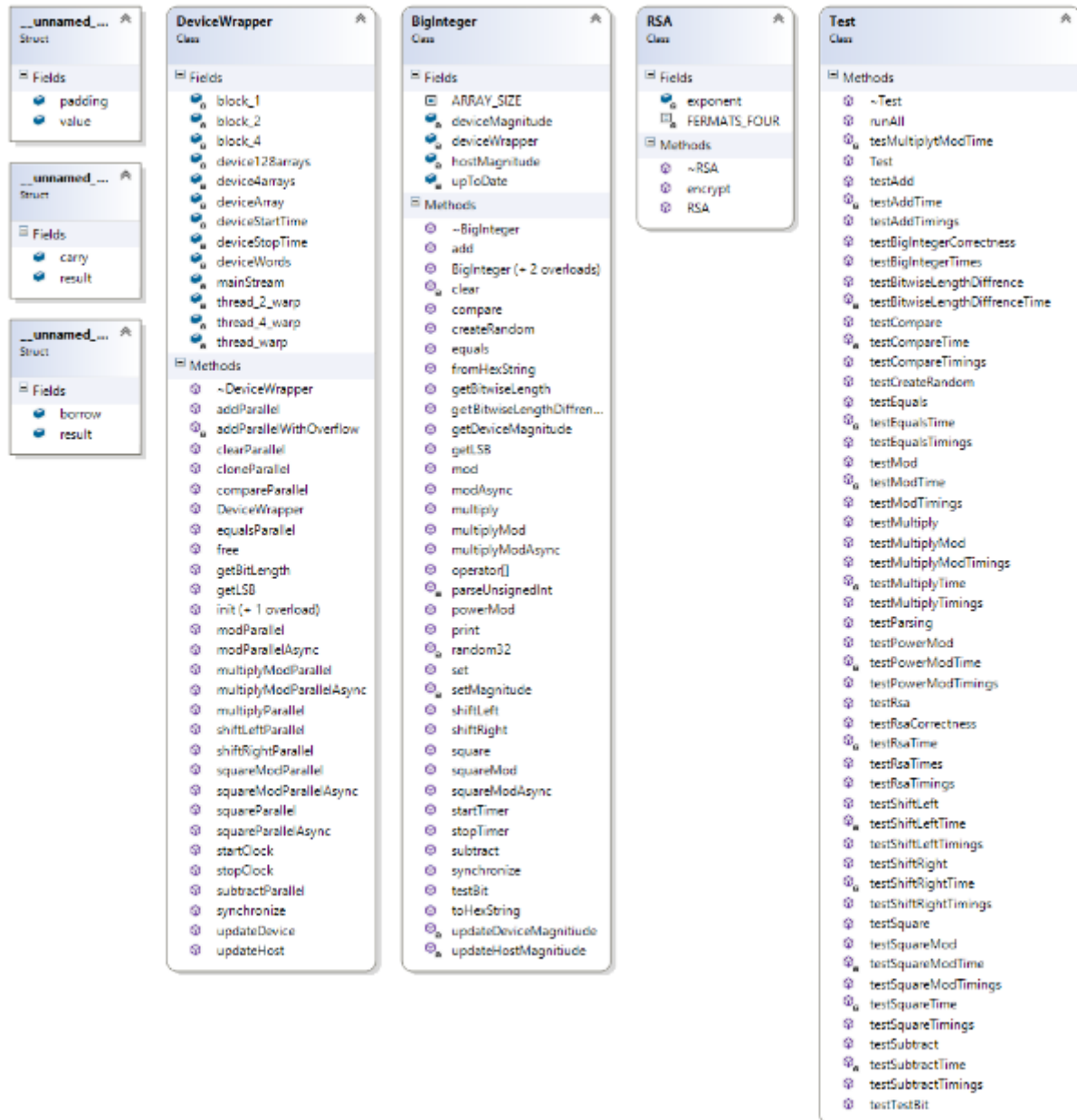


Figure 3.1: Class Diagram

- 
- 3.1 Parallel equals
  - 3.2 Parallel compare
  - 3.3 Parallel bit length
  - 3.4 Parallel left shift
  - 3.5 Parallel right shift
  - 3.6 Parallel add
  - 3.7 Parallel subtract
  - 3.8 Parallel multiply
  - 3.9 Parallel modulo reduction
  - 3.10 Parallel multiply modulo
  - 3.11 Parallel power modulo

# Chapter 4

## Testing

This chapter presents the behavior of implemented algorithms based on the inputs' lengths. All of the functions are designed to handle data up to 4096 bits long. Each test is repeated 50-100 times (depends on the complexity) to provide rather average results.

Other test were executed using CUDA Profiler Activity. "This tool gathers detailed performance information, in addition to timing and launch configuration details. A CUDA Profiler activity consists of a kernel filter and a set of profiler experiments. Profile experiments are directed analysis tests targeted at collecting in-depth performance information for an isolated instance of a kernel launch."[1]

### 4.1 CUDA experiments

#### 4.1.1 Occupancy

"A warp (32 threads) is considered active from the time its threads begin executing to the time when all threads in the warp have exited from the kernel. There is a maximum number of warps which can be concurrently active on a Streaming Multiprocessor (SM). Occupancy is defined as the ratio of active warps on an SM to the maximum number of active warps supported by the SM. Occupancy varies over time as warps begin and end, and can be different for each SM.

Low occupancy results in poor instruction issue efficiency, because there are not enough eligible warps to hide latency between dependent instructions. When occupancy is at a sufficient level to hide latency, increasing it further may degrade performance due to the reduction in resources per thread. An early step of kernel performance analysis should be to check occupancy and observe the effects on kernel execution time when running at different occupancy levels.

#### 4.1.2 Theoretical occupancy

There is an upper limit for active warps, and thus also for occupancy, derivable from the launch configuration, compile options for the kernel, and device capabilities. Each block of a kernel launch gets distributed to one of the SMs for execution. A block is considered active from the time its warps begin executing to the time when all warps in the block have exited from the kernel. The number of blocks which can execute concurrently on an SM is limited by the factors listed below. The upper limit for active warps is the product of the upper limit for active blocks and the number of warps per block. Thus, the upper

limit for active warps can be raised by increasing the number of warps per block (defined by block dimensions), or by changing the factors limiting how many blocks can fit on an SM to allow more active blocks. The limiting factors are:

- Warps per SM

The SM has a maximum number of warps that can be active at once. Since occupancy is the ratio of active warps to maximum supported active warps, occupancy is 100% if the number of active warps equals the maximum. If this factor is limiting active blocks, occupancy cannot be increased. For example, on a GPU that supports 64 active warps per SM, 8 active blocks with 256 threads per block (8 warps per block) results in 64 active warps, and 100% theoretical occupancy. Similarly, 16 active blocks with 128 threads per block (4 warps per block) would also result in 64 active warps, and 100% theoretical occupancy.

- Blocks per SM

The SM has a maximum number of blocks that can be active at once. If occupancy is below 100% and this factor is limiting active blocks, it means each block does not contain enough warps to reach 100% occupancy when the device's active block limit is reached. Occupancy can be increased by increasing block size. For example, on a GPU that supports 16 active blocks and 64 active warps per SM, blocks with 32 threads (1 warp per block) result in at most 16 active warps (25% theoretical occupancy), because only 16 blocks can be active, and each block has only one warp. On this GPU, increasing block size to 4 warps per block makes it possible to achieve 100% theoretical occupancy.

- Registers per SM

The SM has a set of registers shared by all active threads. If this factor is limiting active blocks, it means the number of registers per thread allocated by the compiler can be reduced to increase occupancy. Kernel execution time and average eligible warps should be monitored carefully when adjusting registers per thread to control occupancy. The performance gain from improved latency hiding due to increased occupancy may be outweighed by the performance loss of having fewer registers per thread, and spilling to local memory more often. The best-performing balance of occupancy and registers per thread can be found experimentally by tracing the kernel compiled with different numbers of registers per thread.

- Shared memory per SM

The SM has a fixed amount of shared memory shared by all active threads. If this factor is limiting active blocks, it means the shared memory needed per thread can be reduced to increase occupancy. Shared memory per thread is the sum of static shared memory, the total size needed for all `'__shared__'` variables, and dynamic shared memory, the amount of shared memory specified as a parameter to the kernel launch. For some CUDA devices, the amount of shared memory per SM is configurable, trading between shared memory size and L1 cache size. If such a GPU is configured to use more L1 cache and shared memory is the limiting factor for occupancy, then occupancy can also be increased by choosing to use less L1 cache and more shared memory.

### 4.1.3 Achieved occupancy

Theoretical occupancy shows the upper bound active warps on an SM, but the true number of active warps varies over the duration of the kernel, as warps begin and end. An SM contain one or more warp schedulers. Each warp scheduler attempts to issue instructions from a warp on each clock cycle. To sufficiently hide latencies between dependent instructions, each scheduler must have at least one warp eligible to issue an instruction every clock cycle. Maintaining as many active warps as possible (a high occupancy) throughout the execution of the kernel helps to avoid situations where all warps are stalled and no instructions are issued. Achieved occupancy is measured on each warp scheduler using hardware performance counters to count the number of active warps on that scheduler every clock cycle. These counts are then summed across all warp schedulers on each SM and divided by the clock cycles the SM is active to find the average active warps per SM. Dividing by the SM's maximum supported number of active warps gives the achieved occupancy per SM averaged over the duration of the kernel. Averaging across all SMs gives the overall achieved occupancy.

### 4.1.4 Occupancy charts

- Varying Block Size

Shows how varying the block size while holding other parameters constant would affect the theoretical occupancy. The circled point shows the current number of threads per block and the current upper limit of active warps. If the chart's line goes higher than the circle, changing the block size could increase occupancy without changing the other factors.

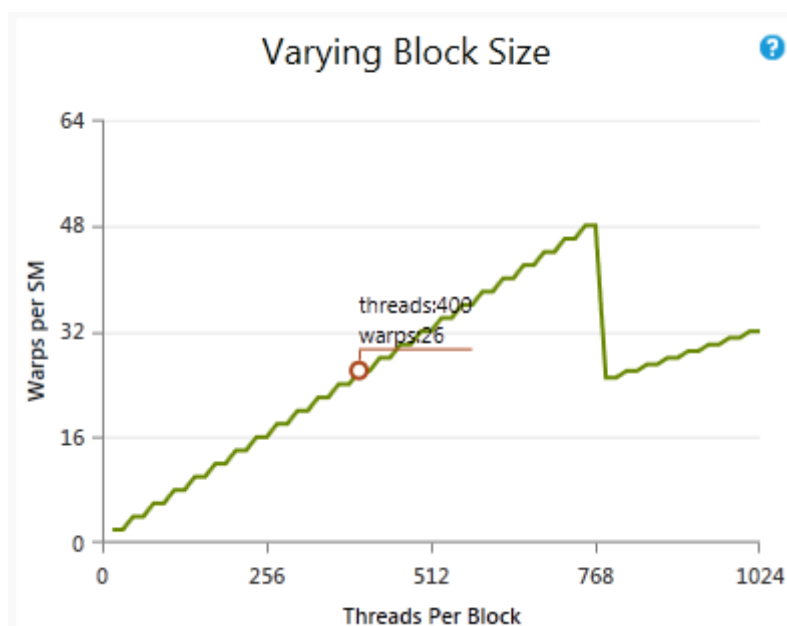


Figure 4.1: Varying block size chart example

- Varying Register Count

Shows how varying the register count while holding other parameters constant would affect the theoretical occupancy. The circled point shows the current number of registers per thread and the current upper limit of active warps. If the chart's



line goes higher than the circle, changing the number of registers per thread could increase occupancy without changing the other factors.

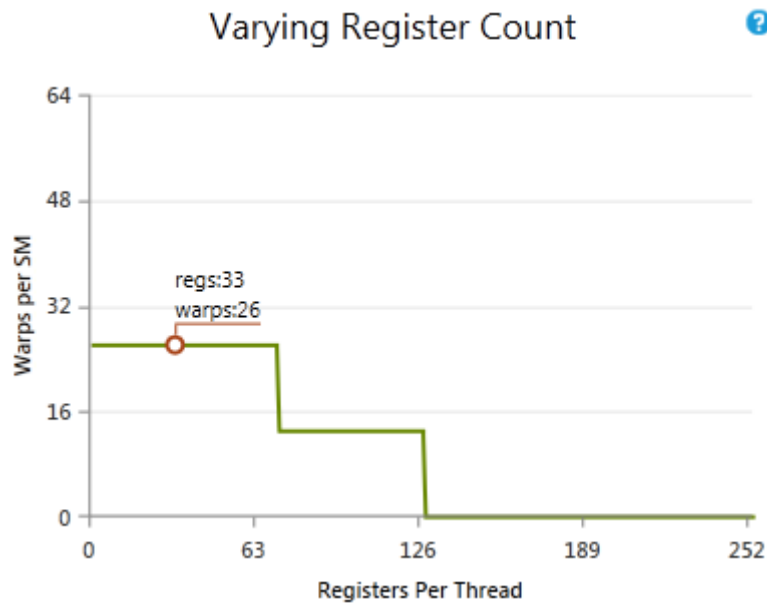


Figure 4.2: Varying register count chart example

- Varying Shared Memory Usage

Shows how varying the shared memory usage while holding other parameters constant would affect the theoretical occupancy. The circled point shows the current amount of shared memory per block and the current upper limit of active warps. If the chart's line goes higher than the circle, changing the amount of shared memory per block could increase occupancy without changing the other factors.

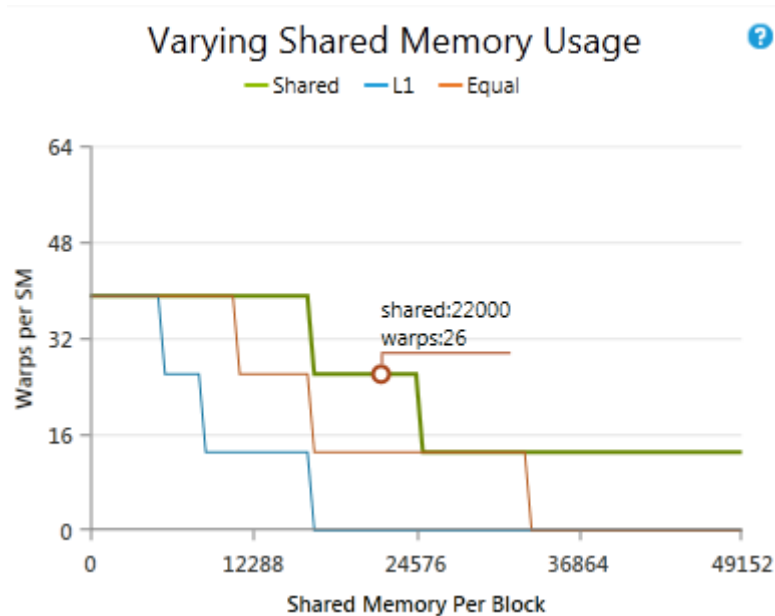


Figure 4.3: Varying shared memory usage chart example

- Achieved Occupancy Per SM

The achieved occupancy for each SM. The values reported are the average across all warp schedulers for the duration of the kernel execution. The line across all bars is the average, which is the number reported as Achieved Occupancy in the other tables.

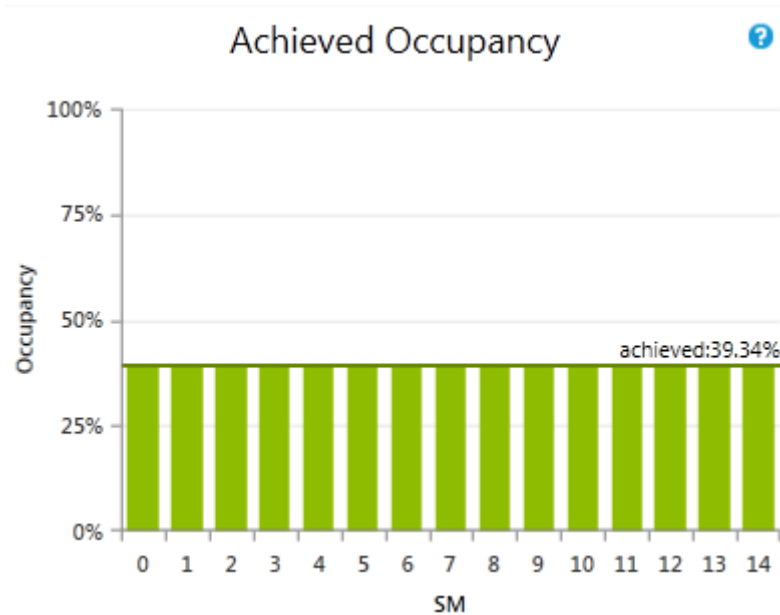


Figure 4.4: Achieved occupancy per SM chart example

#### 4.1.5 Instruction statistics

The Instruction statistics experiment provides a first level triage for understanding the overall utilization of the target device when executing the kernel.

The global work distribution engine schedules thread blocks to Streaming Multiprocessors. A multiprocessor is considered to be active if it has at least one warp assigned to it. The total number of cycles a SM was active for the duration of the kernel execution is defined as Active Cycles.

Each multiprocessor exposes multiple warp schedulers that are able to execute at least one instruction per cycle. At every instruction issue time, each warp scheduler selects one warp that is able to make forward process from its assigned list of warps. For this selected warp the scheduler then issues either the next single instruction or the next two instructions.

A warp scheduler might need to issue an instruction multiple times to actually complete the execution for all 32 threads of a warp. The two primary reasons for this difference between Instructions Issued and Instructions Executed are: First, address divergence and bank conflicts on memory operations. Second, assembly instructions that can only be issued for a half-warp per cycle and thus need to be issued twice. Double floating-point instructions are the prime example for such instructions. As each executed instruction needs to be at least issued once, the following statement holds true in all cases:

$$\text{Instructions Issued} \geq \text{Instructions Executed}$$

Issuing an instruction multiple times is also referred to as Instruction Replay. Each replay iteration takes away the ability to make forward progress by issuing new instructions

on that warp scheduler. Also the compute resources required to process the instruction are consumed for every instruction replay. In short, the more instruction replay iterations are required the higher is the performance impact on the kernel execution.

- Instructions Per Clock (IPC)

A z-ordered column graph showing the achieved instructions throughputs per SM for both, issued instructions and executed instructions. The theoretical maximum peak IPC is a device limit and defined by the compute capabilities of the target device. The y-axis is scaled to this peak value.

- Issued IPC

The average number of issued instructions per cycle accounting for every iteration of instruction replays. Optimal if as close as possible to the Executed IPC. Some assembly instructions require to be multi-issued. hence the instruction mix affects the definition of the optimal target for this metric.

- Executed IPC

The average number of executed instructions per cycle. Higher numbers indicate more efficient usage of the available all resources. As each warp scheduler of a multiprocessor can execute instructions independently, a target goal of executing one instruction per cycle means executing on average with an IPC equal to the number of warp schedulers per SM. The maximum achievable target IPC for a kernel is dependent on the mixture of instructions executed.

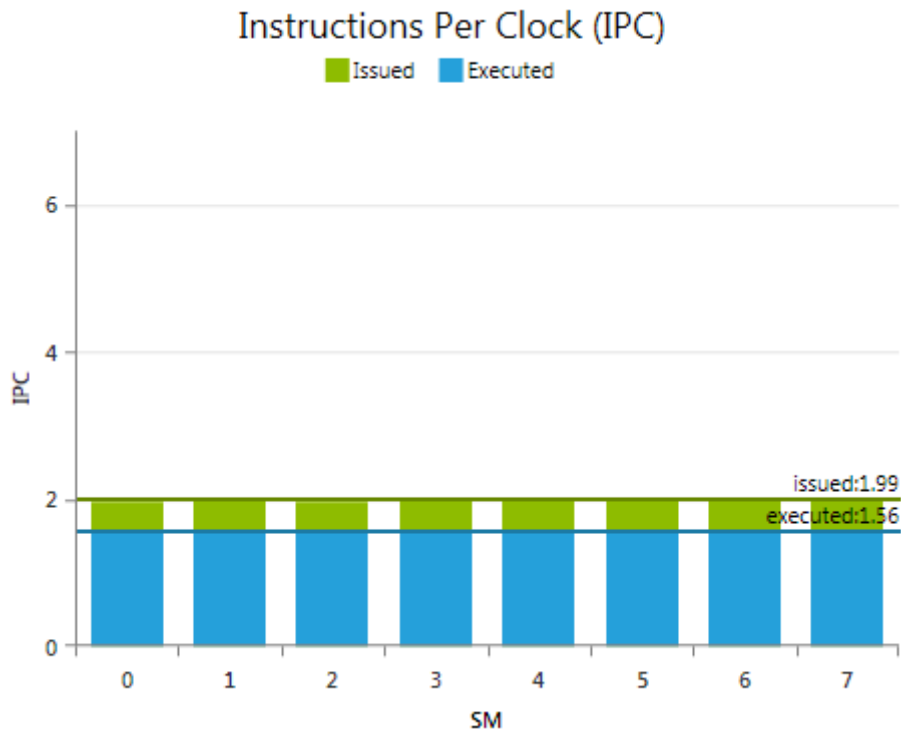


Figure 4.5: Instructions per clocks chart example

- SM Activity Shows the percentage of time each multiprocessor was active during the duration of the kernel launch. A multiprocessor is considered to be active if at least one warp is currently assigned for execution. An SM can be inactive - even

though the kernel grid is not yet completed - due to high workload imbalances. Such uneven balancing between the SMs can be caused by a few factors: Different execution times for the kernel blocks, variations between the number of scheduled blocks per SM, or a combination of the two.

The observable result of a load imbalance are highly different activity values across the multiprocessors; simply caused by the fact that some SMs are still busy executing work, while others SMs already completed their share of work and stay idle as no more work items are left to be scheduled. This is typically referred to as a "tail effect". Small kernel grids with a low number of blocks are more likely to be affected by a tail effects.

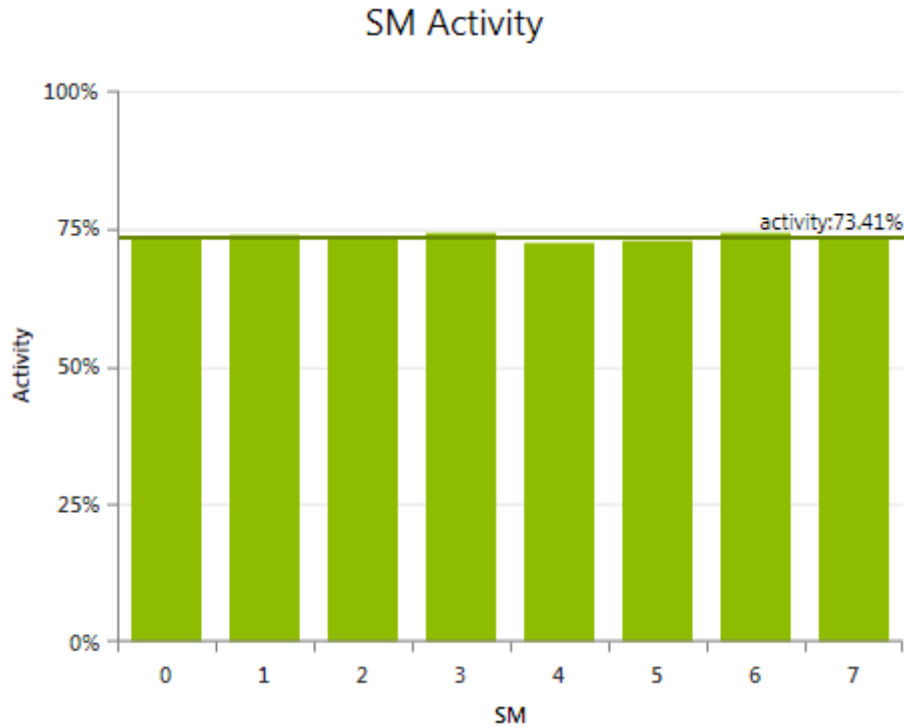


Figure 4.6: SM activity chart example

- Instructions Per Warp (IPW) Shows the average executed instructions per warp for each multiprocessor. High variations in the IPW metric across the SMs indicate non-uniform workloads for the blocks of the kernel grid. While such imbalance does not necessarily have to result in low performance, IPW is very useful to understand the cause of variations in SM Activity.

The most common code pattern to cause high variations in IPW is conditionally executed code blocks where the conditional expression is dependent on the block index. Examples include: special pre-processing or post-processing operations executed for a single block only, or costly detection and handling of edge conditions that are only triggered for some subset of the grid.

- Warps Launched Shows the total number of warps launched per multiprocessor for the executed kernel grid. Large differences in the number of warps launched per SM are most commonly the result of providing an insufficient amount of parallelism with the kernel grid. More specifically, the number of kernel blocks is too low to make good use of all available compute resources. A high variation in the number

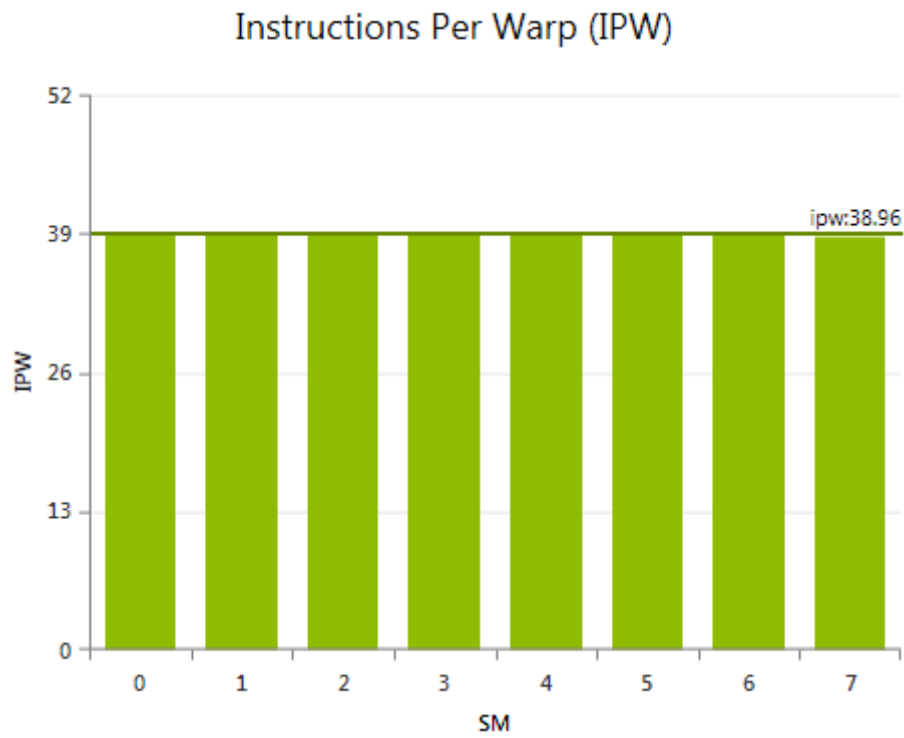


Figure 4.7: Instructions per warp chart example

of warps launched is only a concern if the SM Activity is low on one or more SMs. In this case, partitioning workload that either result in less variance in execution duration per warp or in the execution of more thread blocks. Both cases will help the work distributor in dispatching the given work more evenly.

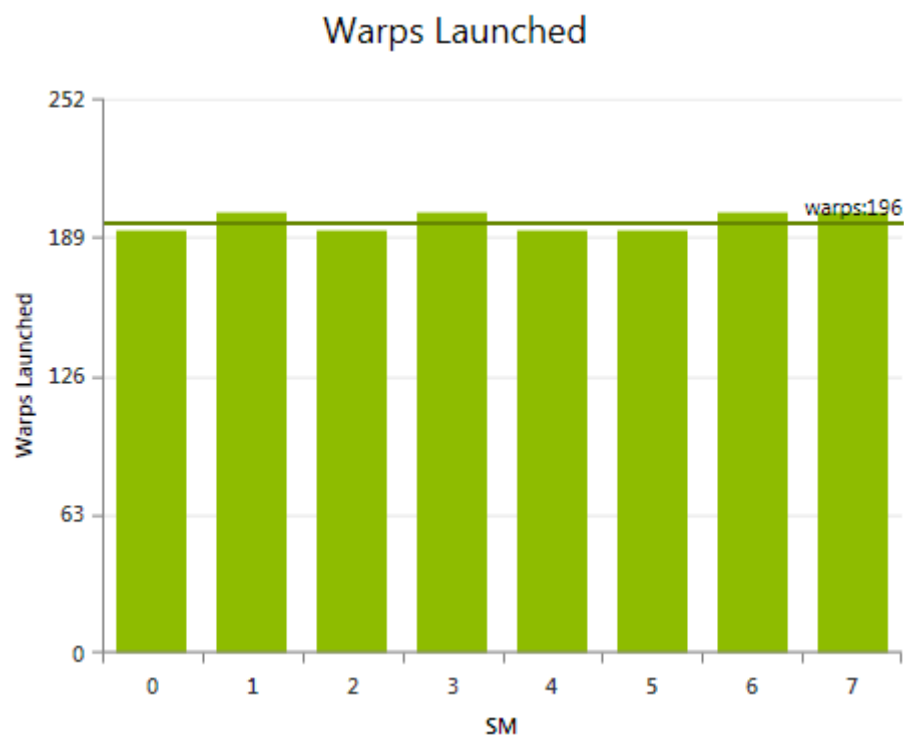


Figure 4.8: Warps launched chart example

### 4.1.6 Branch statistics

Flow control can have a serious impact on the efficiency of executing a kernel. Especially if a lot of flow control decisions are divergent, forcing the threads of a warp into very different control paths through the kernel code. The Branch Statistics experiment helps answering the question of how often flow control instructions were executed, how many of them were uniform versus divergent, and how much the flow control impacted the overall kernel execution performance. A flow control instruction is considered to be divergent if it forces the threads of a warp to execute different execution paths. If this happens, the different execution paths must be serialized, since all of the threads of a warp share a program counter; this increases the total number of instructions executed for this warp. When all the different execution paths have completed, the threads converge back to the same execution path. Conditional expressions that evaluate to a uniform decision across all threads of a warp, do only execute the single, selected code path - consequently causing a lot less overhead. The ratio of executed uniform flow control decisions over all executed conditionals is defined as Branch Efficiency.

The actual performance impact caused by divergent flow control is proportional to the combination of how many different code paths need to be evaluated and how expensive the serialized code segments are. One way of capturing this is to track for all executed instructions how many of the threads in a warp were actually participating in the execution, i.e. how many threads were not predicated off. This is typically referred to as Control Flow Efficiency. By definition this is independent of the number of flow control decisions made, but rather states an upper limit of the utilization of the available compute resources due to flow control.

- Efficiency

The efficiency chart shows the two primary metrics for evaluating the impact of flow control.

- Branch Efficiency

States the ratio of uniform control flow decisions over all executed branch instructions. Shown per-SM (the bars) and averaged over all SMs (the Branch line). Higher values are better, as warps more often take a uniform code path. A value lower than 100% is a necessary, but not sufficient indicator for a negative impact on the kernel execution performance, since the metric does not have any knowledge about the size of the code regions enclosed by the conditionals. For example, one divergent flow control decision out of ten executed branches may be negligible if it encloses very few lines of code only; but it may have a huge impact if it forced the warp to execute many different code paths with thousands of instructions.

- Control Flow Efficiency

Defined as the ratio of active threads that are not predicated off over the maximum number of threads per warp for each executed instruction. Gets lower with fewer threads per warp being active per instruction; therefore serving as a metric for the efficiency in using the available processing units. Lower control flow efficiency can be caused by: Launching warps with less than 32 threads active. Terminating some threads in a warp earlier than others. Or executing instructions with only a subset of the threads enabled.

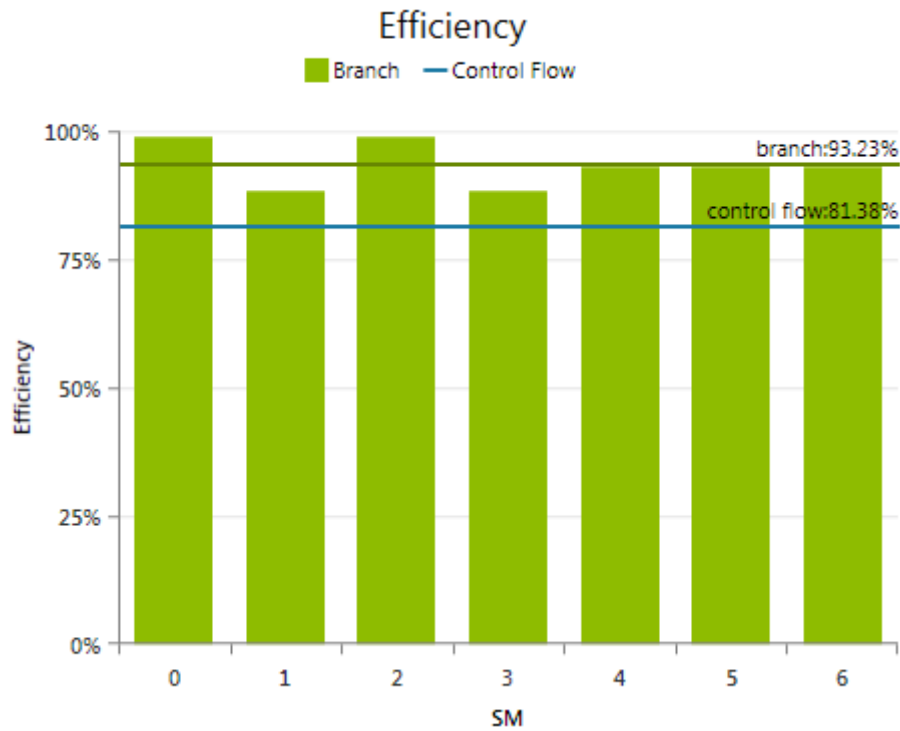


Figure 4.9: Efficiency chart example

- Branches per warp

Shows the average count of executed branch instructions per warp per SM grouped by the outcome of the evaluation of the conditional statement. Useful to investigate the total amount of flow control instructions executed for the warps of the kernel grid.

- Not Taken / Taken

Average number of executed branch instructions with a uniform control flow decision per warp; that is all active threads of a warp either take or not take the branch

- Diverged

Average number of executed branch instruction per warp for which the conditional resulted in different outcomes across the threads of the warp. All code paths with at least one participating thread get executed sequentially. Lower numbers are better.

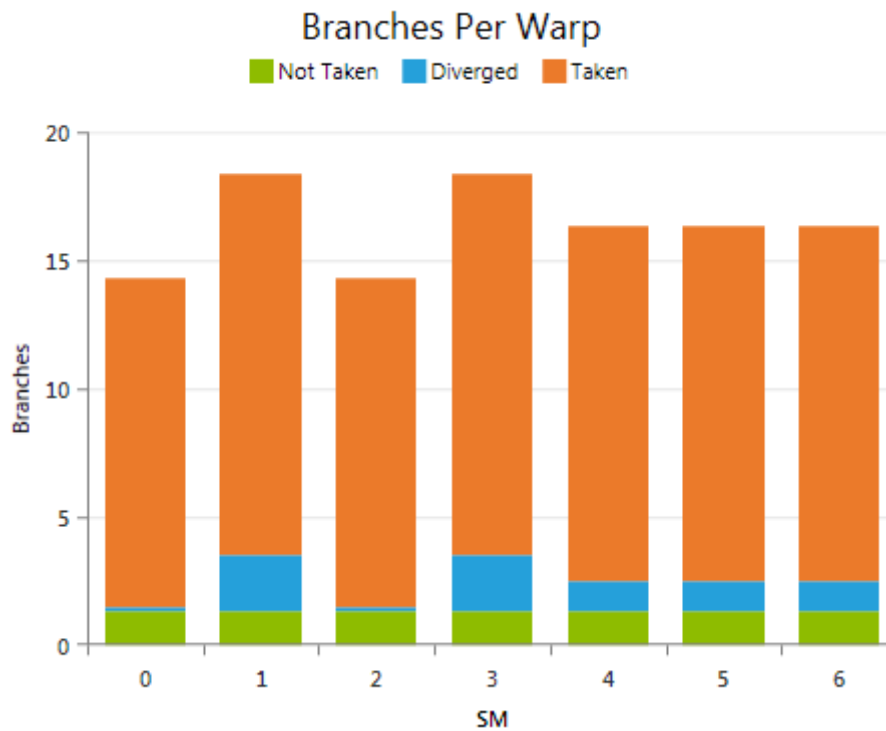


Figure 4.10: Branches per warp chart example

- Branches condition

Shows the distribution of executed branches that were uniform versus divergent aggregated across all warps of the kernel grid.

- Not Taken / Taken

Total number of executed branch instructions with a uniform control flow decision; that is all active threads of a warp either take or not take the branch.

- Diverged

Total number of executed branch instruction for which the conditional resulted in different outcomes across the threads of the warp. All code paths with at least one participating thread get executed sequentially. Lower numbers are better.



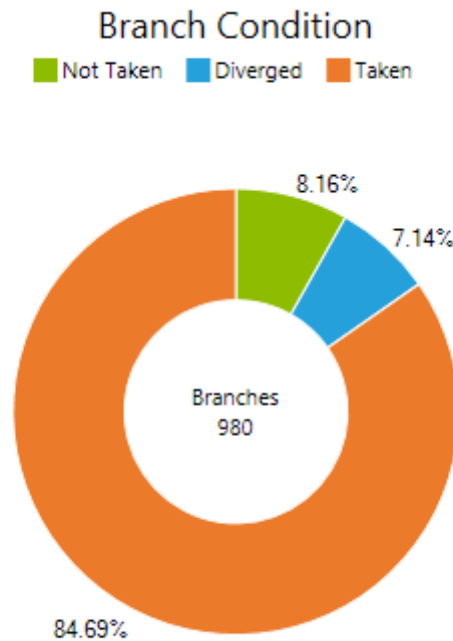


Figure 4.11: Branch condition chart example

#### 4.1.7 Issue efficiency

The issue efficiency experiment provides information about the device's ability to issue the instructions. The key takeaway is the answer to the question if the device was able to issue instructions every cycle. Not being able to do so inevitably lowers the potential peak performance of the kernel.

- Warps Per SM

The metrics are reported as average values across the complete kernel execution for each individual SM of the target device. The y-Axis is scaled to the device limit.

- Active Warps

A warp is active from the time it is scheduled on a multiprocessor until it completes the last instruction. Each warp scheduler maintains its own list of assigned active warps. This assignment of warps to the schedulers is done once at the time a warp becomes active and is valid for the lifetime of the warp.

- Eligible Warps

An active warp is considered eligible if it is able to issue the next instruction. Each warp scheduler will select the next warp to issue an instruction from the pool of eligible warps. Warps that are not eligible will report an Issue Stall Reason. The target is to have at least one eligible warp per scheduler per cycle.

- Theoretical Occupancy

The theoretical occupancy acts as upper limit to active warps and consequently also eligible warps per SM. It is defined by the execution configuration of the kernel launch.

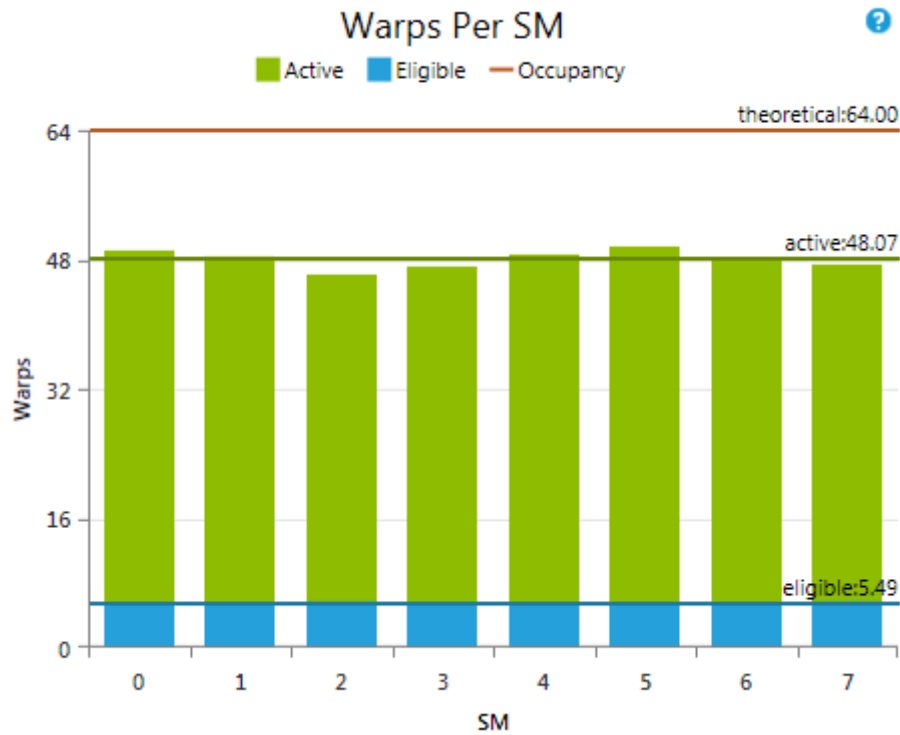


Figure 4.12: Warps per SM chart example

- Warp Issue Efficiency

The metrics are reported as average values across the complete kernel execution for each individual SM of the target device. The y-Axis is scaled to the device limit.

- Active Warps

On every clock cycle, a warp scheduler tries to issue an instruction from one of its warps. When a warp issues an instruction, it takes at least a few cycles before it becomes eligible to issue again, so many warps should be active on a warp scheduler to ensure it can issue an instruction from some warp on every cycle. In this experiment, the profiler counts whether an instruction was issued or not for each clock cycle on each warp scheduler. The Warp Issue Efficiency chart shows the average across all warp schedulers over the duration of kernel execution.

- No Eligible

The number of cycles that a warp scheduler had no eligible warps to select from and therefore did not issue an instruction. The lower the percentage of cycles with no eligible warp the more efficient the code runs on the target device.

- One or More Eligible

The number of cycles that a warp scheduler had at least one eligible warps to select from. This metric is equal to total number of cycles an instruction was issued summed across all warp schedulers. Better if the value is higher with a target of getting close to 100%.

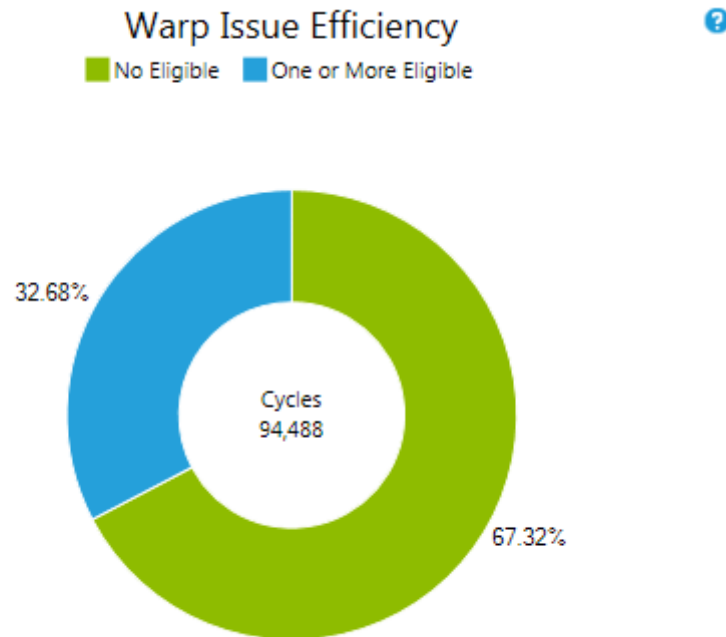


Figure 4.13: Warp Issue Efficiency chart example

- Issue Stall Reasons

The issue stall reasons capture why an active warp is not eligible. On devices of compute capability 3.0 and higher, every stalled warp increments its most critical stall reason by one on every cycle. The sum of the stall reasons, hence increment per multiprocessor per cycle, by a value between zero (if all warps are eligible) and the number of active warps (if all warps are stalled). The update of the stall reason counters occurs for all stalled warps independent of being able to issue an instruction that cycle or not.

- Pipeline Busy

The compute resources required by the instruction are not yet available.

- Texture

The texture subsystem is fully utilized, or has too many outstanding requests

- Constant

A constant load is blocked due to a miss in the constants cache.

- Instruction Fetch

The next assembly instruction has not yet been fetched.

- Memory Throttle

A large number of pending memory operations prevent further forward progress. These can be reduced by combining several memory transactions into one.

- Memory Dependency

A load/store cannot be made because the required resources are not available or are fully utilized, or too many requests of a given type are outstanding. Memory dependency stalls can potentially be reduced by optimizing memory alignment and access patterns.

- Synchronization

The warp is blocked at a `_syncthreads()` call.

- Execution Dependency

An input required by the instruction is not yet available. Execution dependency stalls can potentially be reduced by increasing instruction-level parallelism.

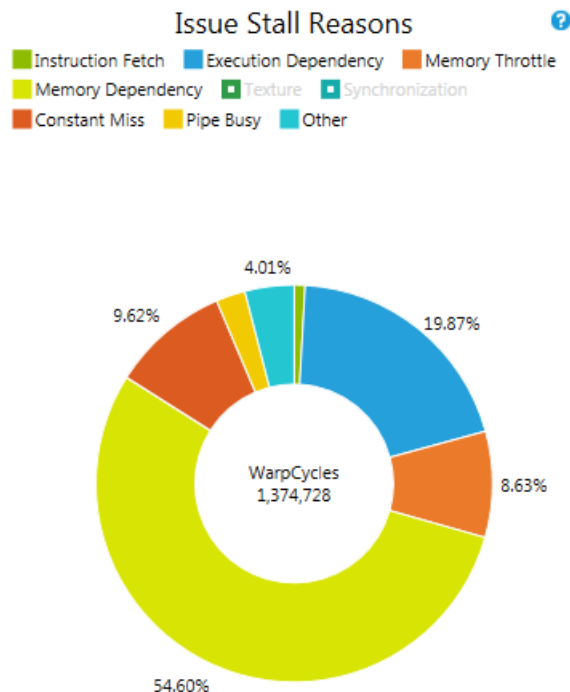


Figure 4.14: Issue Stall Reason chart example

### 4.1.8 Pipe utilization

Each Streaming Multiprocessor (SM) of a CUDA device features numerous hardware units that are specialized in performing specific task. At the chip level those units provide execution pipelines to which the warp schedulers dispatch instructions to. For example, texture units provide the ability to execute texture fetches and perform texture filtering. Load/Store units fetch and save data to memory. Understanding the utilization of those pipelines and knowing how close they are to the peak performance of the target device are key information for analyzing the efficiency of executing a kernel; and also allows to identify performance bottlenecks caused by oversubscribing to a certain type of pipeline.

Pipeline Utilization metrics report the observed utilization for each pipeline at run-time. High pipeline utilization states that the corresponding compute resources were used heavily and kept busy often during the execution of the kernel. Low values indicate that the pipeline is not frequently used and resources were idle. The results for individual pipelines are independent of each other; summing up two or more pipeline utilization percentages does not result in a meaningful value. As the pipeline metrics are reported as an average over the duration of the kernel launch, a low value does not necessarily rule out that the pipeline was a bottleneck at some point in time during the kernel execution.

- Pipe Utilization Chart

Shows the average utilization of the four major logical pipelines of the SMs during the execution of the kernel. Useful for investigating if a pipeline is oversubscribed and therefore is limiting the kernel's performance. Also helpful to estimate if adding more work will scale well or if a pipeline limit will be hit. In this context adding more work may refer to adding more arithmetic workload (for example by increasing the accuracy of some calculations), increasing the number of memory operations (including introducing register spilling), or increasing the number of active warps per SM with the goal of improving instruction latency hiding.

– Load / Store

Covers all issued instructions that trigger a request to the memory system of the target device - excluding texture operations. Accounts for load and store operations to global, local, shared memory as well as any atomic operation. Also includes register spills. Devices of compute capability 3.5 and higher support loading global memory through the read-only data cache (LDG); those operations do not contribute to the load/store group, but are accounted for in the texture pipeline utilization instead.

– Texture

Covers all issued instructions that perform a texture fetch and, for devices of compute capability 3.5 and higher, global memory loads via the read-only data cache.

– Control Flow

Covers all issued instructions that can have an effect on the control flow, such as branch instructions (BRA,BRX), jump instructions (JMP,JMX), function calls (CAL,JCAL), loop control instructions (BRK,CONT), return instructions (RET), program termination (EXIT), and barrier synchronization (BAR).

– Arithmetic

Covers all issued floating point instructions, integer instructions, conversion operations, and movement instructions.

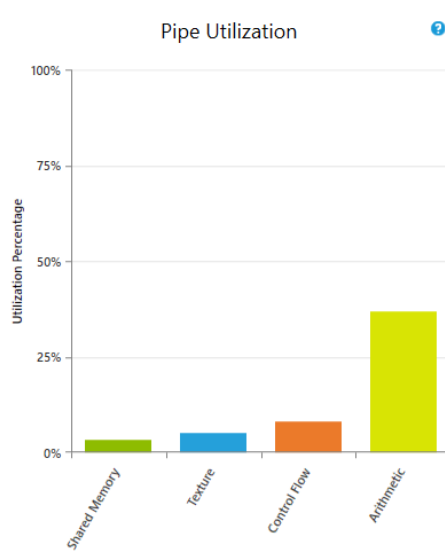


Figure 4.15: Pipe Utilization chart example

- Arithmetic Workload

Provides the distribution of estimated costs for numerous classes of arithmetic instructions. The cost model is based on the issue count weighted by the reciprocal of the corresponding instruction throughput.

- FP16

Estimated workload for all 16-bit floating-point add (HADD), multiply (HMUL), multiply-add (HFMA) instructions.

- FP32

Estimated workload for all 32-bit floating-point add (FADD), multiply (FMUL), multiply-add (FMAD) instructions.

- FP64

Estimated workload for all 64-bit floating-point add (DADD), multiply (DMUL), multiply-add (DMAD) instructions.

- FP32 (Special)

Estimated workload for all 32-bit floating-point reciprocal (RCP), reciprocal square root (RSQ), base-2 logarithm (LG2), base 2 exponential (EX2), sine (SIN), cosine (COS) instructions.

- I32 (Add)

Estimated workload for all 32-bit integer add (IADD), extended-precision add, subtract, extended-precision subtract, minimum (IMNMX), maximum instructions.

- I32 (Mul)

Estimated workload for all 32-bit integer multiply (IMUL), multiply-add (IMAD), extended-precision multiply-add, sum of absolute difference (ISAD), population count (POPC), count of leading zeros, most significant non-sign bit (FLO).

- I32 (Shift)

Estimated workload for all 32-bit integer shift left (SHL), shift right (SHR), funnel shift (SHF) instructions.

- Cmp/Min/Max

Estimated workload for all comparison operations

- I32 (Bitfield/Rev)

Estimated workload for all 32-bit integer bit reverse, bit field extract (BFE), and bit field insert (BFI) instructions

- I32 (Bitwise Logic)

Estimated workload for all logical operations (LOP).

- Warp Shuffle

Estimated workload for all warp shuffle (SHFL) instructions.

- Video SIMD

Estimated workload for all video vector instructions

- Conv (From I8/I16 to I32)

Estimated workload for all type conversions from 8-bit and 16-bit integer to 32-bit types (subset of I2I).

- Conv (To/From FP64)

Estimated workload for all type conversions from and to 64-bit types (subset of I2F, F2I, and F2F).

- Conv (All Other)

Estimated workload for all all other type conversions (remaining subset of I2I, I2F, F2I, and F2F).

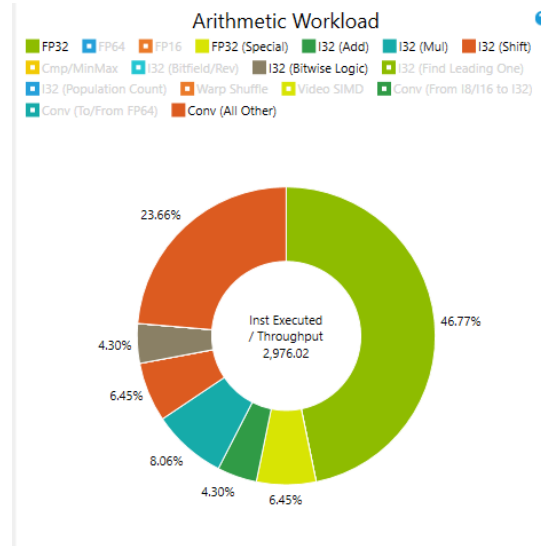


Figure 4.16: Arithmetic Workload chart example

### 4.1.9 Memory statistics

- Overview Chart

Shows a summary view of the memory hierarchy of the CUDA programming model. Key metrics are reported for the areas that were covered by memory experiments during the data collection. The nodes in the diagram depict either a logical memory space (global, local, shared, ...) or an actual hardware unit on the chip (caches, shared memory, device memory). For the various caches the reported percentage number states the cache hit rate; that is the ratio of requests that could be served with data locally available to the cache over all requests made. Requests that hit data in the cache are served much faster than requests that miss the cache; missed data needs to be fetched from another layer of the memory hierarchy.

Links between the nodes in the diagram depict the data paths between the SMs to the memory spaces into the memory system.

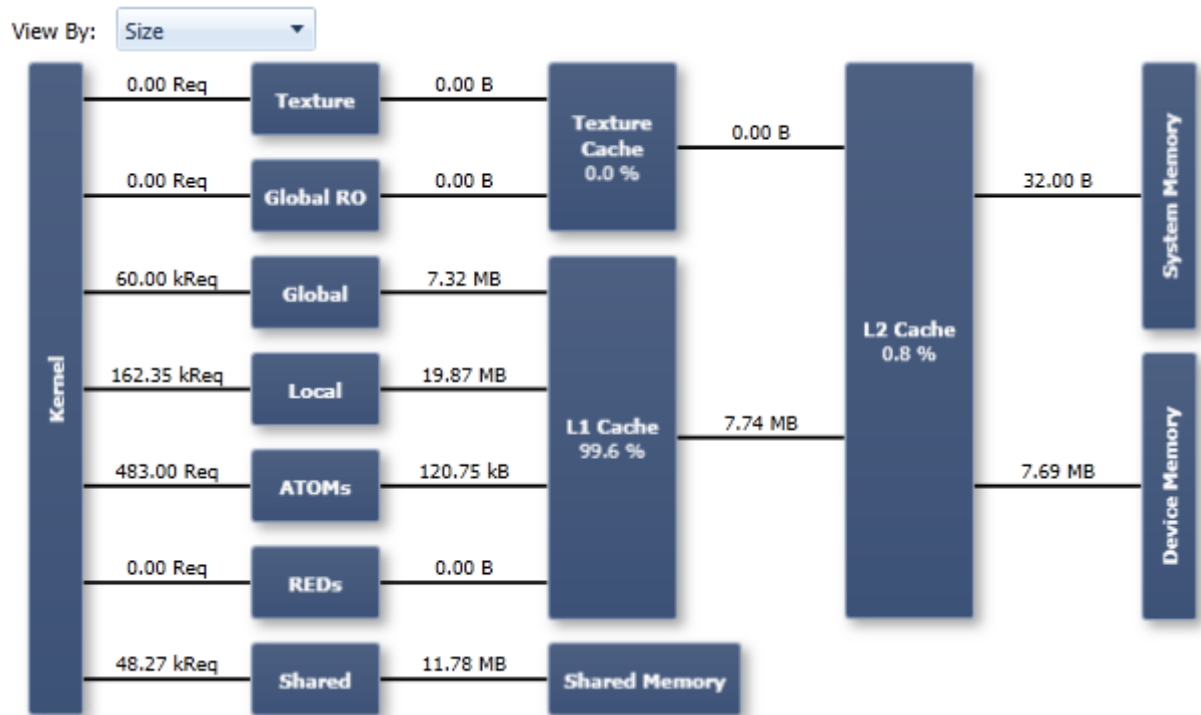


Figure 4.17: Memory Statistics chart example

- Shared memory Chart

The number of Load/Store Requests equals the amount of shared memory instructions executed. When a warp executes an instruction that accesses shared memory, it resolves the bank conflicts. Each bank conflict forces a new memory transaction. The more transactions are necessary, the more unused words are transferred in addition to the words accessed by the threads, reducing the instruction throughput accordingly. Each memory transaction also requires the assembly instruction to be issued again; causing instruction replays if more than one transaction is required to fulfill the request of a warp.

The Transactions Per Request chart shows the average number of shared memory transactions required per executed shared memory instruction, separately for load and store operations. Lower numbers are better; the target for a single shared memory operation on a full warp is 1 transaction for both, a 4byte access and an 8byte access, and 2 transactions for a 16byte access."



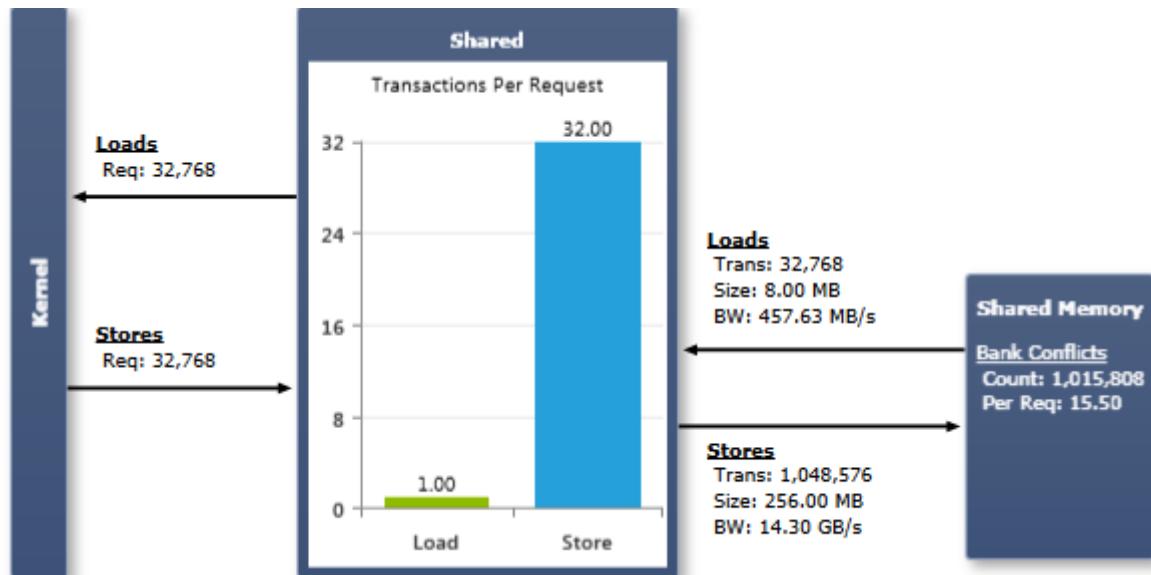
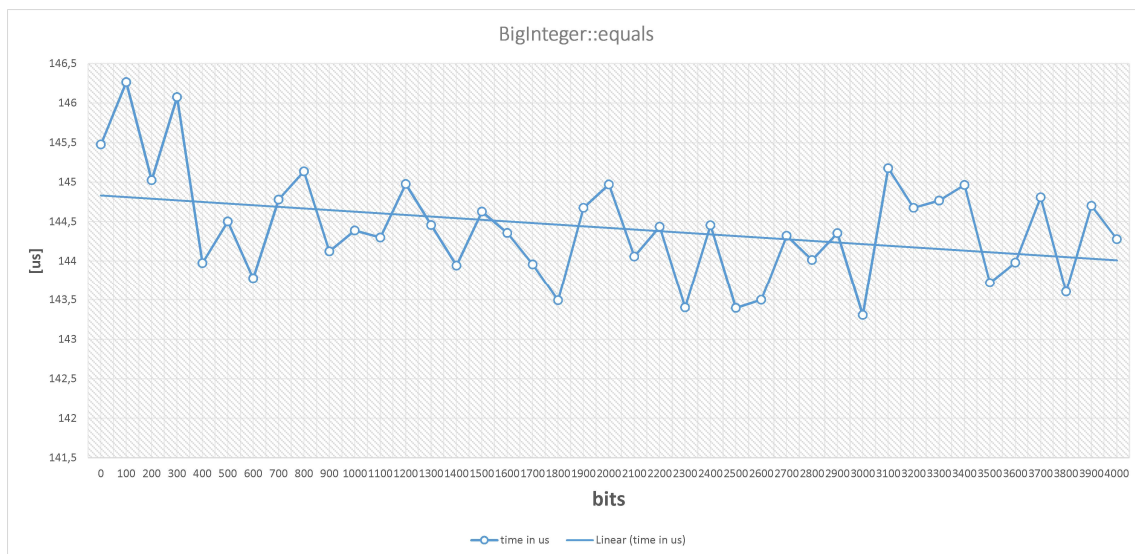


Figure 4.18: Shared Memory Statistics chart example

## 4.2 Equals

Figure 4.19 presents execution times of the function Equals in  $\mu\text{s}$  based on input's bitwise length. The trend line shows the average time actually drops as the input's length rises. This is probably caused by caching techniques implemented on the device. It is safe to state the function runs in constant time, and is resistant to timing attacks.

Figure 4.19: Execution time of function Equals in  $\mu\text{s}$  based on input's bitwise length

# Bibliography

- [1] N. CORPORATION. Nvidia nsight visual studio edition 5.3 user guide. [http://docs.nvidia.com/nsight-visual-studio-edition/5.3/Nsight\\_Visual\\_Studio\\_Edition\\_User\\_Guide.htm#Nsight\\_Visual\\_Studio\\_Edition\\_User\\_Guide.htm%3FTocPath%3D\\_\\_\\_\\_\\_1](http://docs.nvidia.com/nsight-visual-studio-edition/5.3/Nsight_Visual_Studio_Edition_User_Guide.htm#Nsight_Visual_Studio_Edition_User_Guide.htm%3FTocPath%3D_____1). Available: 2017-08-28.
- [2] Y. Ge. A note on the carmichael function. <http://yimin-ge.com/doc/carmichael.pdf>. Available: 2017-08-28.
- [3] E. O. Iskra Nunez. Generalized hamming weights for linear codes. <http://www.uprh.edu/~simu/Reports2001/NOU.pdf>, 2001. Available: 2017-08-28.
- [4] R. Laboratories. Rsaes-oeap encryption scheme algorithm specification and supporting documentation. [http://www.inf.pucrs.br/~calazans/graduate/TPVLSI\\_I/RSA-oeap\\_spec.pdf](http://www.inf.pucrs.br/~calazans/graduate/TPVLSI_I/RSA-oeap_spec.pdf). Available: 2017-08-28.
- [5] Microsoft. Visual studio. <https://www.visualstudio.com>. Available: 2017-08-28.
- [6] E. Milanov. The rsa algorithm. [https://sites.math.washington.edu/~morrow/336\\_09/papers/Yevgeny.pdf](https://sites.math.washington.edu/~morrow/336_09/papers/Yevgeny.pdf), June 2009. Available: 2017-08-27.
- [7] NVIDIA. Cuda zone. <https://developer.nvidia.com/cuda-zone>. Available: 2017-08-27.
- [8] NVIDIA. Using inline ptx assembly in cuda. [https://www.cs.cmu.edu/afs/cs/academic/class/15668-s11/www/cuda-doc/ptx\\_isa\\_2.2.pdf](https://www.cs.cmu.edu/afs/cs/academic/class/15668-s11/www/cuda-doc/ptx_isa_2.2.pdf). Available: 2017-08-28.
- [9] J. D. Owens, D. Luebke, N. Govindaraju, M. Harris, J. Krüger, A. Lefohn, T. J. Purcell. A survey of general-purpose computation on graphics hardware. *Computer Graphics Forum*, 26(1):80–113, 2007.
- [10] W. Stein. Elementary number theory: Primes, congruences, and secrets. <http://wstein.org/ent/ent.pdf>, 2017. Available: 2017-08-28.
- [11] D. F. YongBin Zhou. Side-channel attacks: Ten years after its publication and the impacts on cryptographic module security testing. <https://eprint.iacr.org/2005/388.pdf>. Available: 2017-08-28.

# List of Figures

2.1	An example of cryptosystem including side-channel information leakage . . .	5
3.1	Class Diagram . . . . .	11
4.1	Varying block size chart example . . . . .	15
4.2	Varying register count chart example . . . . .	16
4.3	Varying shared memory usage chart example . . . . .	16
4.4	Achieved occupancy per SM chart example . . . . .	17
4.5	Instructions per clocks chart example . . . . .	18
4.6	SM activity chart example . . . . .	19
4.7	Instructions per warp chart example . . . . .	20
4.8	Warps launched chart example . . . . .	20
4.9	Efficiency chart example . . . . .	22
4.10	Branches per warp chart example . . . . .	23
4.11	Branch condition chart example . . . . .	24
4.12	Warps per SM chart example . . . . .	25
4.13	Warp Issue Efficiency chart example . . . . .	26
4.14	Issue Stall Reason chart example . . . . .	27
4.15	Pipe Utilization chart example . . . . .	28
4.16	Arithmetic Workload chart example . . . . .	30
4.17	Memory Statistics chart example . . . . .	31
4.18	Shared Memory Statistics chart example . . . . .	32
4.19	Execution time of function Equals in $\mu$ m based on input's bitwise length . .	32

# List of Tables

1.1	Device specification . . . . .	2
1.2	PC specification . . . . .	3