

PARIS: A PARALLEL RSA-PRIME INSPECTION TOOL

A Thesis

Presented to

the Faculty of California Polytechnic State University

San Luis Obispo

In Partial Fulfillment

of the Requirements for the Degree

Master of Science in Computer Science

by

Joseph White

June 2013

© 2013

Joseph White

ALL RIGHTS RESERVED

## COMMITTEE MEMBERSHIP

TITLE:                   PARIS: A PARallel RSA-prime InSpection  
                          tool

AUTHOR: Joseph White

DATE SUBMITTED: June 2013

COMMITTEE CHAIR: Chris Lupo, Ph.D.

COMMITTEE MEMBER: Phillip Nico, Ph.D.

COMMITTEE MEMBER: Foaad Khosmood, Ph.D.

## **Abstract**

PARIS: A PARallel RSA-prime InSpection tool

Joseph White

Modern-day computer security relies heavily on cryptography as a means to protect the data that we have become increasingly reliant on. As the Internet becomes more ubiquitous, methods of security must be better than ever. Validation tools can be leveraged to help increase our confidence and accountability for methods we employ to secure our systems.

Security validation, however, can be difficult and time-consuming. As our computational ability increases, calculations that were once considered “hard” due to length of computation, can now be done in minutes. We are constantly increasing the size of our keys and attempting to make computations harder to protect our information. This increase in “cracking” difficulty often has the unfortunate side-effect of making validation equally as difficult.

We can leverage massive-parallelism and the computational power that is granted by today’s commodity hardware such as GPUs to make checks that would otherwise be impossible to perform, attainable. Our work presents a practical tool for validating RSA keys for poor prime numbers: a fundamental problem that has led to significant security holes.

Our implementation using NVIDIA’s CUDA framework offers a 27.5 times speedup over a reference sequential implementation. This level of speedup brings this validation into the realm of runtime reachability.

# Contents

<b>List of Tables</b>	<b>vii</b>
<b>List of Figures</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Related Work . . . . .	2
1.3 Overview of CUDA . . . . .	3
1.4 Algorithm Description . . . . .	4
1.4.1 Binary GCD . . . . .	4
1.4.2 Parallel Functions . . . . .	6
1.4.3 Computational Complexity . . . . .	8
1.4.4 Theoretical Speedup . . . . .	8
1.5 Problem Description . . . . .	9
1.6 Implementation . . . . .	10
1.6.1 Problem Decomposition . . . . .	10
1.6.2 Grid Organization . . . . .	10
1.6.3 Shared Memory . . . . .	12
1.6.4 Occupancy . . . . .	13
1.6.5 Bit-vector . . . . .	14
1.7 Experimental Setup . . . . .	15
1.7.1 Test Machine . . . . .	15
1.7.2 Reference Implementations . . . . .	15
1.7.3 Test Sets . . . . .	16

1.8	Results . . . . .	16
1.9	Conclusion . . . . .	18
1.10	Future Work . . . . .	18
<b>2</b>	<b>Background</b>	<b>21</b>
2.1	Motivation . . . . .	23
2.1.1	Digital Rights Management . . . . .	23
2.1.2	Internet Security and Certificates . . . . .	24
2.1.3	Data Collection . . . . .	26
2.1.4	Data Analysis . . . . .	28
2.1.5	Implications . . . . .	31
2.1.6	Conclusions . . . . .	32
<b>3</b>	<b>Related Works</b>	<b>33</b>
3.1	The Vulnerability . . . . .	33
3.2	The Process . . . . .	34
3.3	Completed Areas . . . . .	35
3.4	Multiple GPUs . . . . .	35
3.5	Large Data . . . . .	37
3.6	Optimization . . . . .	38
3.7	Similar Algorithms . . . . .	39
<b>4</b>	<b>Validation</b>	<b>41</b>
	<b>Bibliography</b>	<b>43</b>

# List of Tables

1.1	Table giving occupancy for various block dimensions . . . . .	13
1.2	Run-times of sequential and CUDA implementations . . . . .	17
2.1	Major CDN Traffic (2010) . . . . .	27
2.2	Alexa top 1M X.509 RSA bit-length . . . . .	30

# List of Figures

1.1	Total percentage of CUDA implementation that is parallel . . . . .	9
1.2	Number of comparisons needed vs. total number of keys in set . . .	11
1.3	Speedup of CUDA Implementation to Sequential C++ . . . . .	17
2.1	Explanation of poor-prime vulnerability . . . . .	22
2.2	X.509 Fields[8] . . . . .	25
2.3	TLS Overview[3] . . . . .	25
2.4	Presence of RSA key or other security in Alexa top 1M sites . . .	29
2.5	Breakdown of Alexa top 1M sites' RSA keys by bit-length . . . . .	30



# Chapter 1

## Introduction

### 1.1 Introduction

RSA is a public key encryption scheme which relies on the difficulty of factoring large numbers. The algorithm is prevalent throughout security, specifically it is used for many web-security applications. An RSA public key is comprised of a modulus  $n$  of specified length (the product of primes  $p$  and  $q$ ), and an exponent  $e$ . The length of  $n$  is given in terms of bits, thus the term “1024-bit RSA key” refers to the number of bits which make up this value. The associated private key uses the same  $n$ , and another value  $d$  such that  $d \cdot e = 1 \bmod \phi(n)$  where  $\phi(n) = (p - 1) \cdot (q - 1)$ [?]. Ideally, given the number of possible primes that may be used to construct a 1024-bit modulus, no random number generators should reuse either prime. Thus, the likelihood of either  $p$  or  $q$  being repeated in a set of keys should be approximately 0. An individual key may be considered secure by itself, but when compared to other keys, might exhibit a weakness which allows each key’s  $d$  to be calculated entirely from public information.

When considering two keys, a weakness exists when the greatest common divisor of both moduli,  $n_1$  and  $n_2$ , is greater than 1. If  $GCD(n_1, n_2) = p$ , then  $p$  must be a shared prime factor of  $n_1$  and  $n_2$ . Thus,  $q_1 = \frac{n_1}{p}$  and  $q_2 = \frac{n_2}{p}$ . Once  $p$  and  $q$  are known,  $d_1$  and  $d_2$  can be directly calculated, yielding both private key pairs.

This weakness is discussed in [?], which showed a significant number of existing RSA keys were susceptible to this exploit. The primary goal of our work was to speedup the most computationally intensive part of their process by implementing the GCD comparisons of RSA 1024-bit keys using NVIDIA’s CUDA platform.

To aid in accomplishing this goal, the work in [7] was expanded and adapted to compare all combinations of keys in a given set. In comparison to their work, larger sections of the overall program were able to be executed in parallel, resulting in further speedup.

## 1.2 Related Work

The work documented in [?] served as inspiration for this work. Here, Lenstra et al. performed a sanity check of a wide array of public RSA keys contained in SSL certificates and SSH host keys. Their discovery that a significant fraction of these keys (roughly 0.2%) were weak led to our desire to parallelize their investigation, in order to make it as efficient as possible with commodity hardware.

The CUDA implementation of the binary GCD algorithm that was built upon (cf. [7]) is an important example of similar work being done. On a fundamental level, our work mirrors theirs as we based the core of our algorithm on their work,

specifically 1024-bit GCDs were calculated in parallel using CUDA. However, we expanded its relevance with modifications in order to automatically divide and parallelize lists of large values to compare.

Another example of work that makes use of the GPU for security applications is solving discrete logarithms as presented in [?]. A set of large-precision operations (768-bit) was necessary for this work, and was thus implemented in CUDA. This was similar to our own starting point due to the currently-limited CUDA support for large-precision numbers. Because these large values have numerous applications in computer security, the work shown here displays another component of computer security where parallelizing work with the large values can be highly advantageous.

Our work is an example of an amalgamation of other related works. It functions as a supplement to the other materials mentioned here, and provides another example of a computer security application that significantly benefits from using parallelization with commodity Single Instruction, Multiple Data (SIMD) multiprocessors. What sets it apart is its use of 1024-bit RSA keys and the method of parallelization implemented.

## 1.3 Overview of CUDA

CUDA is a platform that provides a set of tools along with the ability to write programs that make use of NVIDIA's GPUs (cf. [?]). These massively-parallel hardware devices are capable of processing large amounts of data simultaneously, allowing significant speedups in programs with sections of parallelizable code using the SIMD model. The platform allows for various arrangements of threads to perform work, based on the developer's decomposition of the problem. Our

solution to the problem presented in this paper is discussed in §1.6.2. In general, individual threads are grouped into up-to 3-dimensional blocks to allow sharing of common memory between threads. These blocks can then be organized into a 2-dimensional grid.

The GPU breaks the total number of threads into groups called warps, which consist of 32 threads that will be executed simultaneously on a single *streaming multiprocessor* (SM). The GPU consists of several SMs which are each capable of executing a warp. Blocks are scheduled to SMs until all allocated threads have been executed.

There is also a memory hierarchy on the GPU. There are 3 types of memory that are relevant to this work: global memory is the slowest and largest; shared memory is much faster, but also significantly smaller; and a limited number of registers that each SM has access to. Each thread in a block can access the same section of shared memory.

## 1.4 Algorithm Description

### 1.4.1 Binary GCD

Binary GCD is a well known algorithm for computing the greatest common divisor of two numbers. Instead of relying on costly division operations like Euclid’s algorithm, bit-wise shifts are employed. The implementation presented in this paper follows the outline displayed in Algorithm 1.

---

**Algorithm 1:** Binary GCD algorithm outline

---

**Input:**  $x$  and  $y$ : two integers.

**Output:** The greatest common divisor of  $x$  and  $y$ .

```
1 repeat
2   if  $x$  and  $y$  are both even then
3      $GCD(x, y) = 2 \cdot GCD(\frac{x}{2}, \frac{y}{2})$ ;
4   else if  $x$  is even and  $y$  is odd then
5      $GCD(x, y) = GCD(\frac{x}{2}, y)$ ;
6   else if  $x$  is odd and  $y$  is even then
7      $GCD(x, y) = GCD(\frac{y}{2}, x)$ ;
8   else if  $x$  and  $y$  are both odd then
9     if  $x \geq y$  then
10       $GCD(x, y) = GCD(\frac{x-y}{2}, y)$ ;
11    else
12       $GCD(x, y) = GCD(\frac{y-x}{2}, x)$ ;
13    end
14  end
15 until  $GCD(x, y) = GCD(0, y) = y$ ;
```

---

### 1.4.2 Parallel Functions

To accomplish Algorithm 1 using CUDA, the following three functions had to be parallelized: shift, subtract, and greater-than-or-equal. As outlined in [7], each 1024-bit number is divided across one warp so that each thread has its own 32-bit integer.

The parallel shift function is straightforward: each thread is given an equal-sized piece of the large-precision integer. Then each thread except for Thread 0 grabs a copy of the integer at `threadID - 1`. The variable `threadID` refers to a value between 0 and 31 and corresponds to a thread in a warp. Each thread shifts its value once and uses its copy of the adjacent integer to determine if a bit has shifted between threads. This procedure is outlined in Algorithm 2.

---

**Algorithm 2:** Parallel right shift

---

**Input:**  $x[32]$  is a 1024-bit integer represented as an array of 32 `ints`,

$threadID$  is the 0-31 index of the thread in warp.

```

1 if  $threadID \neq 0$  then
2   |  $temp \leftarrow x[threadID - 1];$ 
   else
3   |  $temp \leftarrow 0;$ 
   end
4  $x \leftarrow x >> 1;$ 
5  $x \leftarrow x \text{ OR } (temp << 31);$ 
```

---

The parallel subtract uses a method called *carry skip* from [7]. First, each thread subtracts its piece and sets the “borrow” flag of `threadID - 1` if an underflow occurred. Next, each thread checks if it was borrowed from and if

so, decrements itself and clears the flag. Then, if another underflow occurs, the borrow flag at `threadID - 1` will be set. This continues until all the borrow flags are cleared. An outline can be found in Algorithm 3.

---

**Algorithm 3:** Parallel subtract using “carry skip”

---

**Input:**  $x$  and  $y$ : two 1024-bit integers,  $threadID$  is the 0-31 index of the thread in warp.

```

1  $x[threadID] \leftarrow x[threadID] - y[threadID];$ 
2 if underflow then
3   | set  $borrow[threadID - 1];$ 
  end
4 repeat
5   | if  $borrow[threadID]$  is set then
6   |    $x[threadID] \leftarrow x[threadID] - 1;$ 
7   |   if underflow then
8   |   | set  $borrow[threadID - 1];$ 
  |   end
9   | clear  $borrow[threadID];$ 
  | end
until all borrow flags are cleared;
```

---

The parallel greater-than-or-equal has each thread check if its integers are equal. If this is the case, then it sets a position variable shared by the warp to the minimum of its `threadID` and the current value stored in the position variable. This is done atomically to ensure the correct value is stored. Finally, all the threads do a greater-than-or-equal comparison with the integers specified by the position variable. This function is outlined in Algorithm 4.

---

**Algorithm 4:** Parallel greater-than-or-equal-to

---

**Input:**  $x$  and  $y$ : two 1024-bit integers,  $threadID$  is the 0-31 index of the thread in warp.

**Output:** *True* if  $x \geq y$ ; else *False*.

```
1 if  $x[threadID] \neq y[threadID]$  then  
2   |    $pos \leftarrow \text{atomicMin}(threadID, pos);$   
   end  
3 return  $x[pos] \geq y[pos]$ 
```

---

### 1.4.3 Computational Complexity

The computational complexity of the binary GCD algorithm has been shown by Stein and Vallée (cf. [?], [?]) to have a worst case complexity of  $\mathcal{O}(n^2)$  where  $n$  is the number of bits in the integer. The worst case is produced when each iteration of the algorithm shifts one of its arguments only once. Since for this application  $n$  is fixed at 1024 bits, the complexity of a single GCD calculation can be considered to be constant time for the worst case.

To compare all the keys together, the amount of GCDs that must be calculated grows at a rate of  $k^2$ , where  $k$  is the number of keys.

### 1.4.4 Theoretical Speedup

Maximum speedup is defined as follows:

$$\text{Max Speedup} = \frac{1}{1 - P} \quad (1.1)$$

where  $P$  is the percentage of the program's execution that can be parallelized. This percentage is a function of the number of keys the program needs to process,



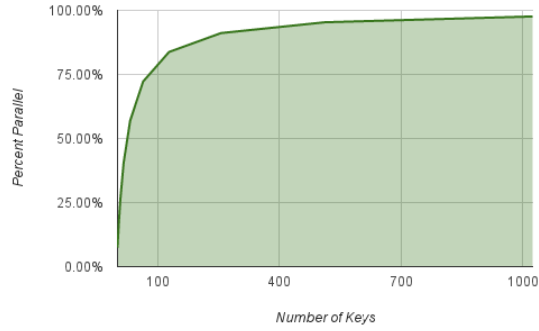


Figure 1.1. Total percentage of CUDA implementation that is parallel and is calculated in Equation 1.2.

$$P = \frac{t \cdot g}{t \cdot g + r \cdot k} \quad (1.2)$$

where

- $t$  = time to process a single GCD
- $g$  = total number of GCD calculations
- $r$  = time to read a single key
- $k$  = total number of keys

Since  $g$  will increase significantly more rapidly than  $k$ ,  $P$  (based on equation 1.2) will approach 1 as  $k$  approaches  $\infty$ . This relationship can be observed in Figure 1.1.

## 1.5 Problem Description

The RSA weakness described above demands that each key in a set be compared with each other key to determine if a GCD greater than 1 exists for any

pair. Given a known set of keys, it is not known before processing the keys which will be likely to have a GCD greater than 1; therefore, there is no way to eliminate comparisons between specific pairs. The natural organization to fulfill this requirement is a comparison matrix of the all keys. Each location in the matrix corresponds to a comparison between two keys.

## 1.6 Implementation

### 1.6.1 Problem Decomposition

Initially, the comparison matrix seems to be an  $n^2$  solution. However, the diagonal of the matrix created consists of unproductive GCD calculations since these entries would compare each key with itself. Furthermore, the matrix is symmetrical over the diagonal. Thus, only the comparisons comprising one of the triangles needs to be performed. Specifically,

$$\text{Total number of GCD compares} = \sum_{i=1}^k i \quad (1.3)$$

This reduction in number of overall compares decreases the work performed significantly, shown in Figure [1.2](#).

### 1.6.2 Grid Organization

One of the most important aspects of any CUDA implementation is the organization of the thread and block array to ensure that the architecture is appropriately used to its full potential. The thread array in this implementation was organized using 3 dimensions. The  $x$  dimension represented the sectioning of a

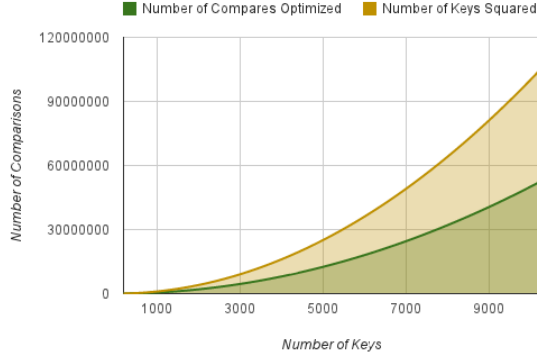


Figure 1.2. Number of comparisons needed vs. total number of keys in set

1024-bit value into individual 32-bit integers, of which there are 32.

$$\frac{1024 \text{ bits per key}}{32 \text{ bits per integer}} = 32 \text{ integers per key}$$

The remaining dimensions,  $y$  and  $z$ , were set to 4, resulting in a block of 512 threads. This design decision was experimentally determined. See §1.6.4 for details about Occupancy optimizations.

$$32 \cdot 4 \times 4 = 512$$

This ensured that each block remained square for algorithmic symmetry and simplicity. The  $y$  and  $z$  dimensions corresponded to how many specific keys within the list of all keys were being compared per block. Thus, two 1024-bit keys were loaded into each 32-thread warp, which was then processed simultaneously as a single comparison. The  $x$  dimension was chosen for two reasons: 1) so one thread in this dimension would represent each of the 32-bit integers inside the key and 2) because there are 32 threads in a single warp. Therefore, this thread-array organization ensured that compares were done using two entire keys (separated into 32 pieces) that were scheduled to the same warp. This eliminated warp divergence since every warp was filled and executed with non-overlapping data.

Blocks were arranged in row-major order based on the key comparisons that they held. The formula for the number of blocks,  $B$ , needed for a vector of keys of size  $k$  can be seen in Equation 1.4.

$$\sum_{i=1}^{\lceil \frac{k}{4} \rceil} i = B \quad (1.4)$$

The limit for a grid in a single dimension is  $2^{16} - 1 = 65535$  and limits the amount of keys that can be processed to 1444. To increase the number of blocks available for computation, a second grid dimension was added. This increased the theoretical maximum number of keys per kernel launch as seen in Equation 1.5.

$$\sum_{i=1}^{\lceil \frac{k}{4} \rceil} i \leq (2^{16} - 1)^2 \quad (1.5)$$

$$k \leq 370716$$

### 1.6.3 Shared Memory

Shared memory was used to load the necessary keys from global memory. Two arrays were created in shared memory, representing the thread-array; both 3-dimensional,  $32 \times 4 \times 4$  and had an integer loaded into each available space. Each array represented which integers would be compared at each location in the matrix. A side effect of this organization was that each key would be repeated 4 times within its integer array. However, this greatly simplified the GCD algorithm so that only a look-up into each array was needed. Since shared memory was not the limiting factor for occupancy, it was not a priority to optimize this aspect of the design and implementation.

Shared memory was also used within the GCD algorithm, specifically in the greater-than-or-equal-to function, and the subtract function. In the greater-than-

Threads per block	128	288	512	800
Occupancy	67%	94%	100%	52%

Table 1.1. Table giving occupancy for various block dimensions

or-equal-to function, a single integer was allocated for each comparison within a block. Within the subtract function, shared memory was utilized to represent the borrow value for each integer.

### 1.6.4 Occupancy

Each SM can be assigned multiple blocks at the same time as long as there are enough free registers and shared memory available. The ratio of active warps to the maximum number of warps supported by a SM is called *occupancy*. On the Fermi architecture, the maximum occupancy is achieved when there are 48 active warps running on a SM at one time. Greater occupancy gives a SM more opportunities to schedule warps in a fashion to hide memory accesses, thus, saturating a SM with many warps decreases performance impact. CUDA Fermi cards have a total of 32768 registers and 49152 bytes of shared memory per SM. The implementation here uses 17 registers and 4762 bytes of shared memory per block and therefore results in a maximum occupancy of 100%.

By using the CUDA occupancy calculator provided by NVIDIA (cf. [?]), a plot can be formed comparing the threads per block with occupancy. To maintain the same block organization outlined above, the block dimensions can be  $2 \times 2$ ,  $3 \times 3$ ,  $4 \times 4$ ,  $5 \times 5$  or 128, 228, 512, 800 threads, respectively. Table 1.1 shows the calculated occupancy for these block sizes. A block size of 512 threads was chosen because it results in the greatest occupancy and thus the best performance.

### 1.6.5 Bit-vector

The initial approach was to allocate a large, multi-dimensional array of integers that would hold the results of the CUDA GCD calculations. This was allocated to the GPU, so each thread could have access as needed; however, since the number of results grew at  $n^2$ , the lack of scalability in this approach was quickly apparent. Additionally, performance decreased due to the large array that was being sent over the PCIe bus. Memory transfers to the GPU are slow, and must be minimized.

After more careful consideration, a new approach was implemented. There would only be a single bit allocated per key-compare to mark whether or not the pair had a GCD greater than 1. In this way, only 2 bytes (16 bits = 1 bit per compare) were necessary per block ( $4 \times 4 = 16$  compares per block), as opposed to the previous  $16 \cdot 32 \cdot 4 = 2048$  bytes. Despite not having access to the answer immediately after returning from the kernel calculation, this approach would be more efficient since there would be a theoretically small number of keys that actually returned with GCDs greater than 1 (i.e. the flag was set). This small set could then be re-processed (GCDs calculated) using a different kernel or using a CPU algorithm. Efficiency would also be increased due to the time saved in memory transfers since there was significantly less memory to transfer before calling the kernel.

## 1.7 Experimental Setup

### 1.7.1 Test Machine

All performance measurements were made on a single machine with an Intel Xeon W3503 dual-core CPU and 4 GB of RAM. This machine has one NVIDIA GeForce GTX 480 GPU with 480 CUDA cores and 1.5 GB of memory. The CUDA driver version present on the machine is 4.2.0, release 302.17, the runtime version is 4.2.9. The CUDA compute capability is version 2.0, and the maximum threads per block is 1024, with each warp having 32 threads.

### 1.7.2 Reference Implementations

In order to check the accuracy of the final implementation, as well as to provide a point of comparison for benchmarking, two reference implementations of this exploit were created. Each was able to use the same format key databases (described in §1.7.3).

The first implementation was written purely in Python using the open source PyCrypto cryptography library. This implementation was able to perform the entire exploit, from finding weak 1024-bit RSA public keys through generating the discovered private keys. This implementation was not used for performance comparison as it was dissimilar to the other two implementations.

A sequential version of the binary GCD algorithm was implemented to serve as a second validation tool for the CUDA implementation. This version sequentially processed the same input as both other implementations and produced output of the same format. Comparison with this implementation ensured that unexpected errors did not result merely from processing the data in parallel.

### 1.7.3 Test Sets

In order to conduct meaningful tests, it was necessary to use an identical data set in all tests. To facilitate this, a tool was written in Python to generate both regular and intentionally weak RSA key pairs using PyCrypto and store them in an SQLite3 database. All keys were generated with a constant  $e$  of 65537, chosen because this was discovered to be a commonly used value (cf. [?]).

The generation process produced a database of RSA key pairs. Intentionally weak keys were evenly distributed.

In order to generate a weak key, this program would generate an initial normal RSA key but save the prime used for  $p$ . For each subsequent bad key,  $p$  would be replaced with this constant, and  $n$  was recalculated. The result was that each weak key would have a GCD greater than 1 when tested against any other weak key.

Using this tool, it was possible to build arbitrarily large test sets with a known number of keys exhibiting the weakness. When these databases were processed using any of the reference implementations, the discovered number of weak keys could be directly compared with the number of keys expected to be found. This allowed both repeatable testing to measure run time, and a method to validate the parallel algorithm was indeed finding GCDs as expected.

## 1.8 Results

The accuracy of the parallel implementation was verified against the sequential implementation by using identical test data sets with known weak keys. Since both implementations found the same set of compromised keys, it was validated



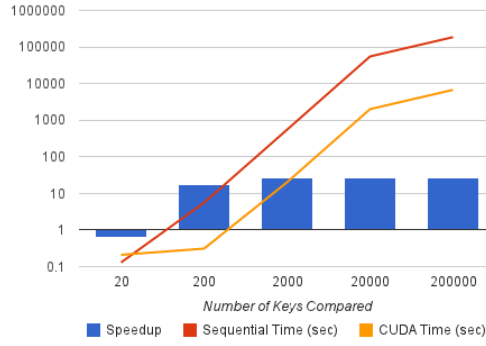


Figure 1.3. Speedup of CUDA Implementation to Sequential C++

Number of Keys	20	200	2000	20000	200000
Sequential Time (sec)	0.13	5.59	550.69	55121.86	185551.91
CUDA Time (sec)	0.21	0.31	20.21	2005.09	6748.23
<b>Speedup</b>	0.6	18.0	27.2	27.5	27.5

Table 1.2. Run-times of sequential and CUDA implementations

that these two implementations were internally consistent. Furthermore, both matched the results of the separate Python reference implementation: supporting the assertion of accurate functionality. The speedup of the CUDA implementation (seen in Figure 1.3) was calculated by comparing its run time with that of the sequential implementation. Compare this with Figure 1.1: this similarity is evidence of the implementation presented here matching with theoretical expectations.

Figure 1.3 shows that speedup increases dramatically with the number of keys until about 2000 comparisons. At this point, the GPU becomes saturated with enough blocks to fully occupy all of the SMs. Speedup remains constant at 27.5 for up to 200000 keys. We have no data beyond this number of keys due to the very long run time of the sequential implementation.

## 1.9 Conclusion

A large speedup resulted directly from writing a CUDA implementation when compared to the sequential implementation. Many more keys are able to be compared in a given amount of time using the CUDA implementation.

A tool was developed to efficiently and completely compare a list of 1024-bit RSA public keys, avoiding repetition and unnecessary work. This tool allows an increased number of keys to be compared compared to prior work, in turn allowing overall execution time to decrease due to the increased parallelism.

The tool described in this paper offers significant advantages over other GCD algorithms in CUDA, and practically applies this for comparison of 1024-bit RSA keys in order to test for a particular weakness. There also exist several areas where the implementation would benefit from further investigation and development including application of the GCDs, and expansion to iterative kernel calls in order to handle even larger sets of keys to compare.

## 1.10 Future Work

The primary limiting factor of this implementation is the amount of memory available on the GPU. Since all key combinations must be computed to expose any potential weakness, the kernel was structured to take a single vector of keys and perform all possible comparisons. In order to process more keys, either a GPU with more memory must be used, or the algorithm must be modified in order to use multiple, iterative kernel launches. The iterative kernel approach would require memory to be separated into two sections that could each be filled with subsections of the large, complete array of keys. All comparisons would then be

performed between the two subsections by calling the kernel that is currently implemented. Upon returning from the kernel, the data in one subsection would be shifted, and all the compares would then be done for those two sets of keys. This process would continue until both vectors have iterated over all keys: specifically, the kernel call would reside in a pair of nested `for` loops. This change would allow the implementation presented here to process an arbitrary number of keys.

To further enhance the above proposed addition, asynchronous memory transfers could also be added to the implementation. When combined with multiple kernel launches, a significant portion of the memory I/O (which is one of the main limiting factors of performance using the GPU) would be able to be masked by simultaneously processing the data currently on the GPU while new data is being copied onto it.

An aspect that was originally intended for this project, but was not implemented was to have the CUDA kernel return the actual GCD of any keys that were found to be “weak” in the sense there existed a GCD greater than 1. Since a large majority of the GCDs found by our implementation are equal to one, memory is wasted if all the results are transferred back to the host. The most memory efficient solution would include dynamically allocating memory for any significant results on the device. This would remove the recalculation step in the current implementation needed to produce private keys.

Recently NVIDIA has released information about a new GPU architecture called Kepler[?]. The Kepler architecture introduces new features that may increase the performance of this implementation.

A feature known as Dynamic Parallelism allows a CUDA kernel to launch new kernels from the GPU. This would allow dynamic allocation of block sizes for different areas of the comparison matrix and remove idle threads from the kernel. Hyper-Q is a new technology that manages multiple CUDA kernels from multiple CPU threads. With the current Fermi architecture, only one CUDA kernel may run on the device at one time. This can lead to under utilization of the GPU hardware. An approach using multiple CPU threads, each running their own CUDA kernel, could greatly increase throughput.

The final missing component of this implementation would be to complete the algorithm and use the GCDs which are being calculated to generate RSA private key pairs. In order to do this, a heterogeneous multi-process approach may be most straightforward. After the parallel run completes, the bit-vector of

# Chapter 2

## Background

RSA is an asymmetric key encryption scheme. Keys come in matched pairs: public keys and private keys. Using the same algorithm, information encrypted with the one key, can be decrypted with the other and vice versa. A party must keep the private key secret, but the public portion of the key can be seen and used by anyone in the world. To generate a pair, an algorithm is performed on two randomly-generated prime numbers whose product is of a certain bit-length (e.g. 1024 bits). These are the primes that our work looks to discover inadequacies with.

The work presented by Lenstra et. al. [\[10\]](#) shows that around 0.2% of RSA keys collected from sources including the SSL Observatory suffer from inadequate random prime generation. They show that when the primes used to build an RSA public-private key pair are not actually random, factoring does not need to be done to break the keys. Instead, only a greater common divisor needs to be calculated between two RSA moduli. When anything greater than 1 is found as a result of this calculation, the keys can be considered broken and offer no protection. This is because once a GCD greater than 1 is found, the private keys

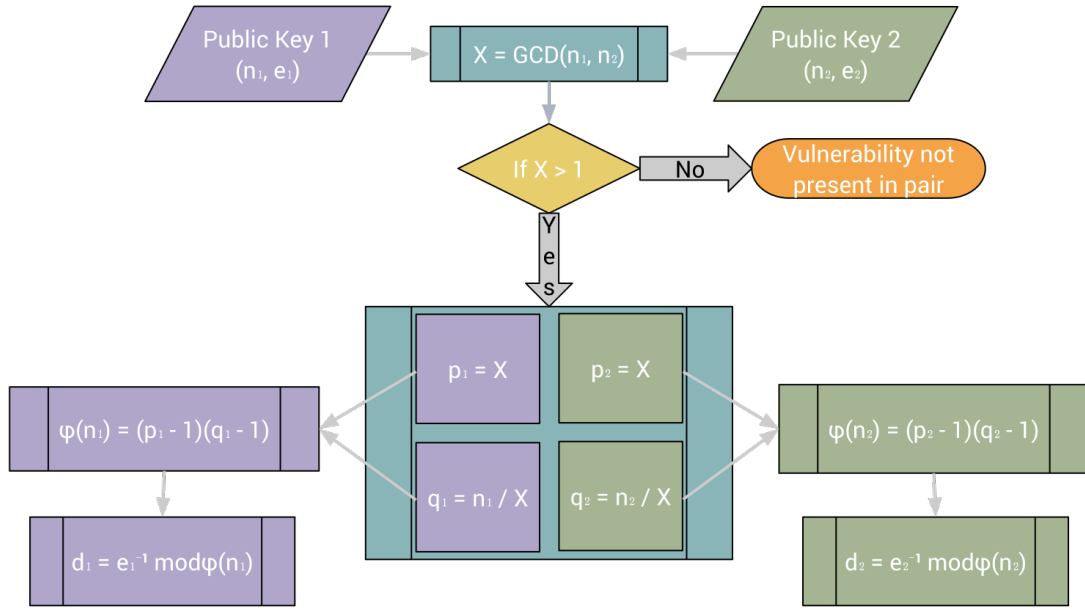


Figure 2.1. Explanation of poor-prime vulnerability

can be generated from the publicly available information. The vulnerability is explained in detail in Figure 2.1.

A tool by Scharfglass et. al. [14] performs a parallelized version of this GCD calculation (using the binary GCD algorithm) between every pair combination in a given set of keys. It structures these calculations so that as many as possible can be done at once. To increase parallelism, and because the keys are so large, they are broken up into 32 chunks, which are also operated on in parallel using the GPU. Through these 2 levels of parallelism, a speedup of 27.5 is achieved. This tool provides a practical tool for checking a given set of keys for the vulnerability we are focusing on.

## 2.1 Motivation

Scharfglass et. al. use NVIDIA’s CUDA framework to significantly speedup a check for the vulnerability presented in [10]. The work in [14] lacks any context or explanation of usage for RSA keys in general, and although their data set is security- and encryption-focused, no discussion is given to the security aspects and the potential impacts of the work. Some research was done in an attempt to add context to the work, and the implications it might have.

Most initial research yielded general statements that did not explain or give appropriate context to RSA (statements including “...all over the Internet...” were accurate, but not informative). More in-depth research yielded two primary use-cases of RSA: Digital Rights Management (DRM) and Secure Socket Layer (SSL) and Transport Layer Security (TLS) Internet security protocols. Other uses included password alternatives (e.g. ssh connections or command line interface tools like `git`) but these are more difficult to collect data for and analyze since they tend to be done on an individual basis, and are managed by individual users.

### 2.1.1 Digital Rights Management

Digital rights management is a protocol used to secure the usage and distribution of various types of media content. It is adopted by content providers and device manufacturers to ensure users don’t misuse or wrongfully share protected, copyrighted content. The RSA Association proposed a protocol that could be applied to various types of media content that secured its usage [1, 2]. These announcements argued that using RSA for DRM offered benefits to the content providers, the device manufacturers, and, arguably, the consumers of the content. It apparently allowed device manufactures and content providers to ensure

proper usage of protected content while allowing users to consume and playback their rightfully-owned media on an array of their own devices.

### 2.1.2 Internet Security and Certificates

Internet security relies on a particular protocol called SSL and its successor TLS. These protocols define how to securely transfer information over the Internet by using encryption and signing mechanisms. One aspect of both the SSL and TLS protocols involve certificates to ensure the expected party is actually being communicated with. These certificates provide information about a server that a user is connected to (e.g. Amazon or Google) and is signed by a certificate authority (e.g. VeriSign). These certificates are also encrypted with a subject's (e.g. Amazon) public key to ensure that the entity is who they say they are. Furthermore, the public key is then used to transfer information back to the subject in a secure way. Figure 2.3 describes the TLS process between a client and server.

The most widely used type of certificate is the X.509 certificate, whose current revision is 3. A breakdown of the major sections of the certificate is shown in Figure 2.2. The “Subject Public Key Information” section of the certificate holds the relevant data (the subject's RSA public key) for our work. The public key portion of the certificates can actually use one of several algorithms defined in the specification. RSA is by far the most popular, accounting for over half of the found certificates (as shown in Figure 2.4). Other options include DSA and Diffie-Hellman. For more detailed information about X.509 certificates and analysis on their infrastructure, see [8].



X509v3 Certificate		
Version	Serial no.	Sig. algo.
Issuer		
Validity	Not Before	Not After
Subject		
Subject Public Key Info		
	Algorithm	Public Key
X509 v3 Extensions		
	CA Flag, EV, CRL, etc.	
Signature		

Figure 2.2. X.509 Fields[8]

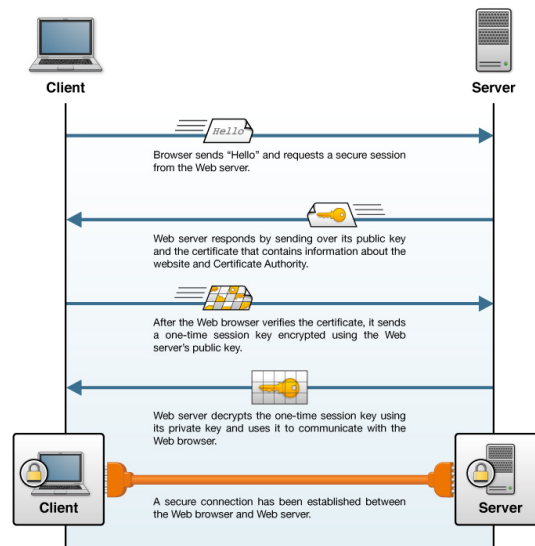


Figure 2.3. TLS Overview[3]

### 2.1.3 Data Collection

The proposed protocol given by the RSA actually modeled the certificate model used by SSL/TLS. For this reason, the SSL certificate model was the primary focus of the data collected. Moreover, the data set of RSA keys in SSL certificates is much larger than DRM implementations that use RSA. Thus, the data set that was collected and then analyzed consisted entirely of X.509 certificates.

Since a primary motivation of this work is the potential impact that vulnerable RSA keys may have, a data set representative of a large portion of traffic was desired. A survey was done over 2007-2010 [9] that gathered some useful data on the Internet as a whole. Relevant to web traffic, this study showed that together, the major content data networks (CDN) (i.e. Google, Facebook, Amazon, etc.) account for nearly 17% of web traffic. Specific breakdowns and conversions can be found in Table 2.1.

Alexa was used to find websites with large amounts of traffic<sup>1</sup>, as it serves as one of the Internet's most prevalent sources for traffic information. In addition to individualized site traffic data, Alexa provides a daily-updated list of the top one million websites ordered by traffic<sup>2</sup>. In regards to the previous point concerning CDNs and traffic breakdowns, it was noted that all web sites hosted by the major CDNs from [9] are present in the top one million list. Additionally, the vast majority of sites from the top one million list were not provided by the CDNs mentioned in [9]. This means that a much larger percentage than 17% is actually represented by the sites in the list, although quantifying precisely how much becomes difficult, and lies beyond the scope of this project.

---

<sup>1</sup><http://www.alexa.com/>

<sup>2</sup><http://s3.amazonaws.com/alexa-static/top-1m.csv.zip>

Table 2.1. Major CDN Traffic (2010)

CDN	Percent of all Internet traffic	Approx. percentage of web traffic
Google	7	12.72
Facebook	0.45	0.818
Amazon CDN (NASA/JPL, PBS,...)	0.55	1
EdgeCast (Yahoo!, Break, Imgur,...)	0.5	0.909
LeaseWeb (Heineken, Starbucks,...)	0.8	1.454

A Python script was used to parse the CSV file provided by Alexa. Each extracted URL was visited on port 443 (corresponding to “https://”) using `openssl`. If a site responded on that port, it provided an X.509 certificate to be parsed and verified. Most of the time, this is done by a user’s web browser which has certificate authorities’ private keys hidden within. However, since the work here was not concerned with the actual web content, the certificate was simply saved. When certificates were found, the “Subject Public Key Algorithm” section was inspected. If this contained some form of RSA, `openssl` was used again to extract the RSA key and save it into a PEM file (a standard format for saving public and private key information in). Additionally, since part of this work’s motivation lies in expanding the work done in [14], the keys were also stored into a SQLite3 database so that later retrieval by the tool was trivial.

After writing an initial implementation of the Python script, it was realized

that the collection process was too slow. Therefore, some of Python’s multi-process capabilities were utilized to speed up the collection of data. Eight processes were spawned to perform the work described above. Seven of the processes read URLs from a deque (double-ended queue) and attempted to obtain an X.509 certificate and RSA key from each. When found, these keys were added to a queue. These data structures not only offered convenient structures for organizing the data, but were also thread-safe and could be used with our updated implementation. An eighth thread read keys from the queue and entered them into the SQLite3 database. This reimplementaion offered a  $15\times$  speedup over the previous, naïve approach.

#### 2.1.4 Data Analysis

After collecting data from each of the top one million Alexa sites, some surface-level analysis was performed in order to gain an overview of the security of these one million sites. This analysis included overall security breakdowns (i.e. whether a site used any kind of security) and bit-length breakdowns of the RSA keys that were found. This data analysis can be found in Figures 2.4 and 2.5. Table 2.2 offer specific breakdown for bit-lengths, and offers details that cannot be seen in the figure.

Interestingly (and slightly worrisome), is that the majority of web sites in this top one million list did not even respond on port 443, implying that they do not allow any option for secure traffic to their site. Nearly 60% of websites in the list fell into this category.

Also of note are the specifics of bit-lengths found. It is a surprising data set for various reasons. First, bit-lengths as small as 384, 512, and 768 were not

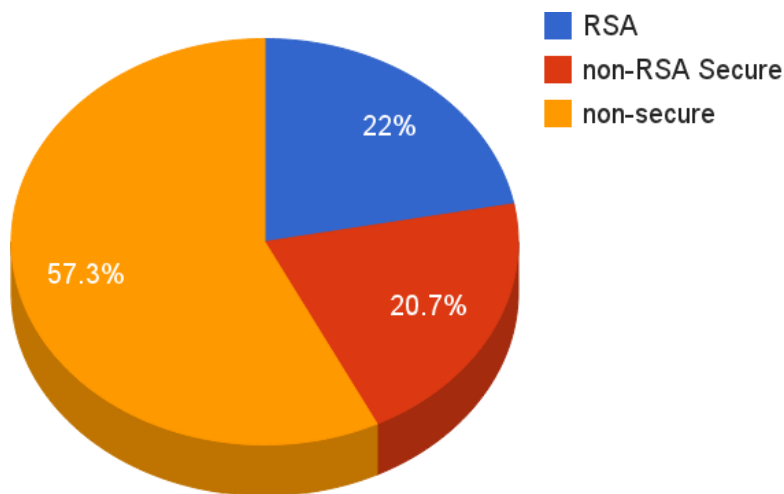


Figure 2.4. Presence of RSA key or other security in Alexa top 1M sites

expected as keys of this size were factored years ago [4], and cannot be expected to offer much security<sup>3</sup>.

On the other hand, the high number of longer bit-lengths was also surprising, but comforting. The fact that the vast majority of keys were 2048-bit was not expected, but means that a large number of high-traffic websites have adequate (for now) encryption. Moreover, there were quite a number of sites with 4096, and, even more surprisingly, 8192-bit.

Further analysis was done using the tool described in [14]. This tool performs a complete check of all RSA keys in a given set for the previously mentioned vulnerability described in [10]. This check is done using NVIDIA’s CUDA framework to parallelize the work as much as possible which allows for a practical security validation tool. At the time of writing, the tool’s current implementation is limited to 1024-bit RSA keys. Due to this limitation, the data set was reduced to 46,736 RSA keys. Due to the  $O(n^2)$  nature of this problem, this check resulted

---

<sup>3</sup>Interesting note: The author attempted to visit the sites 384-bit keys. Google Chrome did not find a web page on port 443 for any of the sites.

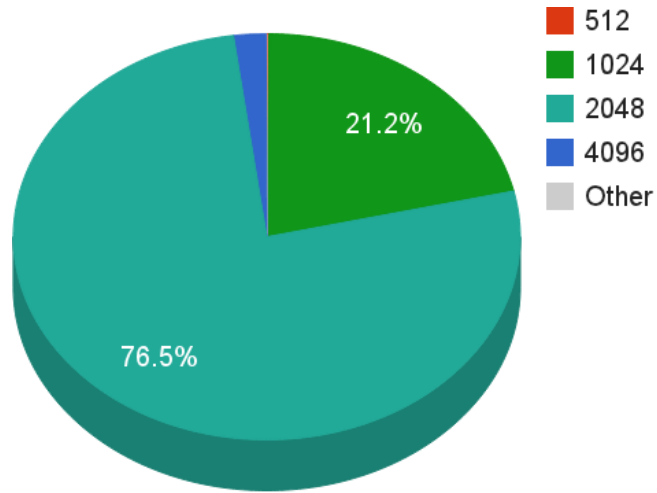


Figure 2.5. Breakdown of Alexa top 1M sites' RSA keys by bit-length

Table 2.2. Alexa top 1M X.509 RSA bit-length

Bits	Count
384	4
512	309
768	4
1024	46736
2048	168391
2432	37
3072	17
4096	4511
8192	14
other	44

in 1,092,150,216 comparisons. The run of the tool using this data set took 138 minutes, 37 seconds or 131,305 comparisons per second. After applying the tool to this set of keys, it was found that no keys were susceptible to the vulnerability (insofar as none of the keys in the set shared primes with a given key). Although no vulnerabilities were detected, this does not mean the set is completely secure from the previously mentioned vulnerability. Since such a limited data set was used, and the security tool relies on large sets to accumulate many primes, the analysis presented can only be considered an initial pass.

### **2.1.5 Implications**

For the same reasons the data set was chosen to be X.509 certificates, presence of the vulnerability mentioned here would have the largest effect on Internet security: specifically in the certificate infrastructure. If servers were unknowingly providing compromised public keys in their certificates, their traffic could offer no security if this became known. This could offer ill-effects since the vulnerability is agnostic of the type of data being encrypted.

Additionally, the password-alternatives that were mentioned briefly before could be particularly impacted. If a system was generating many new RSA keys (e.g. a user is building new accounts for various websites, and tying RSA keys to them) with inadequate primes, potentially all these accounts for this user could end up compromised, despite the sense of stronger security.

Along these lines, a system administrator could potentially suffer from similar issues. If many machines and key-pairs are being setup and generated as part of an initial setup process for a new network of users, poor random number generation could have severe consequences when done at scale. This would likely

create many new, vulnerable keys at once. By being aware that this problem could exist, and that a tool such as [14] exists, an insecure set of users could be avoided or detected much earlier than would have otherwise occurred.

### **2.1.6 Conclusions**

This survey determined that the primary use of RSA keys in practice is the signing of SSL certificates. A data set was gathered that accounts for a significant amount of certificates in the context of web traffic. Some surface-level analysis was done to gain a general understanding of the collected data. A high-performance tool was then applied to a subset of the data to check for a vulnerability within the data subset. No vulnerabilities were found; however, due to limitations in the tool used, as well as limitations in the data set, this cannot be considered a complete analysis and requires further investigation.

Though limited, this survey does successfully add real-world application context to the work done in [14]. Specific implications are able to be drawn from the practicality of the tool, and more tangible examples of what vulnerabilities in RSA keys look like have been realized.



# Chapter 3

## Related Works

### 3.1 The Vulnerability

An RSA vulnerability was discovered and detailed in [10] in early 2012 that forms the basis of the exploit the work presented here builds upon. Namely, that poor random number generation in RSA keys results in insecure encryption using these keys. Specifically, a key's private components can be generated using publicly available information by taking advantage of the work presented in [14].

The exploit explored in [10] makes use of the two prime values that play a fundamental role in an RSA key. These values must be sufficiently random such that repeated values are infeasible and mathematically extremely improbable to encounter due to the scope of the set (this work focuses on 1024-bit keys). However, on some systems and due to some less-than-adequate code, this may not always be the case. When these values are repeated between keys, the greatest-common-divisor algorithm can be applied to find the shared value between keys. This is significant because the GCD calculation is significantly less computa-

tionally expensive than the alternative: factoring the large values. With this new information about each key, their private components are straight-forward to generate since the RSA key-generation process is public and not difficult to implement.

A large sample set was obtained via the OpenSSL Observatory and by crawling the Internet for SSL certificates (which use RSA as an encryption scheme). Since it was thought this would provide an adequate sample set, all these keys were analyzed. The set contained 1024-bit keys, as well as 2048-bit keys. The findings were that 0.2% of these keys were vulnerable to this exploit, and thus offered no security. This included 2048-bit keys as well, despite the conception 2048-bit keys are more secure (this would be true if this vulnerability were avoided).

## 3.2 The Process

As mentioned, the vulnerability is exploited by performing greatest common divisor (GCD) calculations with pairs of keys. The research presented in this paper looks at how to most quickly and efficiently analyze a large database of RSA keys and detect if this vulnerability is present within the set. Traditionally, this would be a tedious process that would likely be infeasible due to time constraints. Since the calculations are independent of one another (between pairs of keys) they can be done in parallel, given the proper system. One such system, and the focus of this work, is NVIDIA's CUDA. This will be detailed later, however, it relates to work done in [7].

This work optimized a specific version of the GCD (binary GCD). This method reduces the algorithm to three repeated operations. This algorithm was

then implemented in CUDA for 1024-bit numbers (something CUDA does not offer native support for). This allows much of each operation to be performed in a parallel fashion, further optimizing the process.

### 3.3 Completed Areas

As a groundwork for the research presented here, the work done in [14] will be used extensively. Here, the ideas present in both [10] and [7] are combined with further optimization and parallelization of the data to achieve a significant speedup of the overall process (when compared to the sequential runtimes). This work makes use of CUDA to parallelize the binary GCD algorithm as presented in [7], but additionally parallelizes each GCD calculation as each is independent of the other (that is, many complete GCD comparisons are performed in parallel). This same approach will be used in the work that follows.

This work, however, does not offer a complete solution and suffers from several limitations. The work in [14] does not actually produce the private components of the RSA keys found to be vulnerable by the algorithm. Furthermore, the set of keys that can be used in the algorithm as presented, is limited by the memory present on the GPU. To offer a complete solution, this must be overcome to allow RSA key-sets of arbitrary size to be used for the validation. The usefulness of such work is severely impacted due to these two limitations.

### 3.4 Multiple GPUs

As an extension to the work in [14], multiple GPU support will be explored. This concept is discussed at length in [13]. Here, several GPU configurations are

investigated to analyze how a distributed system may benefit over a single-GPU design. Work is divided intelligently among the systems in the distributed system to gain more performance increase than would be possible with a single system.

At the time [13] was written, single-system, multiple-GPU configurations were not possible using CUDA over SLI (Scalable Link Interface, NVIDIA’s brand for the interface between multiple GPUs). However, this is possible now, and will be explored at length in our work. This configuration offers benefits including transfer speeds over alternative configurations.

The work done in [15] helps to make a case for the additional speedup that making use of multiple GPUs can offer. This work was done after CUDA had added adequate support for using multiple GPUs simultaneously, and configurations with two-GPU and four-GPU systems were used.

Significant additional speedups were gained, ranging up to three times speedup using four GPUs over the single-GPU system (which translates to 100 times speedup compared to the original, non-GPU benchmark being used). This kind of speedup offers great optimism in how multiple GPUs may offer significant increases to performance of the GCD key algorithm.

Another interesting aspect of the work done in [15] is the varied domain sizes that were used. The data sets were broken down into several different configurations which was each tested on the single, dual, or quad-GPU configuration. This was done to achieve a more accurate answer to how much faster the multi-GPU system could be. The largest domain set ( $1024 \times 32 \times 1024$ ), interestingly, offered the greatest additional speedup between configurations—the quad-GPU configuration was 3 times faster than the single-GPU, compared to 2 times for the  $256 \times 32 \times 256$  domain. This is significant, and implies that exploring different

ways to organize our own data may offer increased advantages when moving to a multi-GPU system.

## 3.5 Large Data

Since one of our primary goals is to remove the limitation present in [14] that prevents larger sets of keys from being analyzed, work investigating how to manage large data sets is relevant. Specifically, large data sets that must be managed by the GPU.

One such work is [16]. In this work, like our own, the data is too large to fit within the extremely-limited memory of the GPU. Thus, the data must be partitioned. In this work, the data was first prepared by organizing it so that it could be effectively partitioned. There was a straightforward way of doing so for this application (K-means).

Additionally, some features of the GPU were considered and used, namely asynchronous memory transfer, a type of multi-tasked streaming. This feature allows memory transfers to be done between the GPU and CPU while the GPU is processing data it already has. This can increase performance significantly when compared to having the GPU move from processing to data transfer (i.e. having a single tasked thread). Optimization of threads was looked into in this work as well, as the authors attempted to find how many of these threads would add the positively impact speedup the most. They found that only 2 threads (i.e. one processing, one transferring) were needed to gain the greatest performance benefit. Adding additional streams after this point did not offer additional benefit.

Decomposition of data is the primary problem when attempting to perform work with a GPU due to the limitations of current hardware. Therefore, much work and thought has been given to optimizing the structure of large data sets. Although a different type of parallelism is looked at in [6], the strategies to data decomposition are relevant to most parallel problem in general, including CUDA. The focus in this survey is matrix multiplication, a standard parallel problem. The most effective way to chunk the data in this work appears to be entire rows/columns of the element matrices. This is logical, and the approach is mimicked (as much as it can be since matrix multiplication is not perfectly synonymous) in [14]. However, a slight modification was found to make this strategy most effective: hold the results in memory until an entire set of iterations completes (all that will fit in memory) and only then transfer the results back to the CPU. Although this does not decrease the amount of memory sent from the GPU to the CPU, it reduces the number of total memory transfers which is often more important. Such fine details must be found and applied to our work here to further increase our performance.

## 3.6 Optimization

As mentioned, data decomposition and organization around the architecture are some of the most important (if not *the* most important) aspects of creating a parallel solution for use with a GPGPU. The work done in [12] gives a detailed overview of how one might optimize work done with CUDA. It discusses the memory model, and how certain techniques such as tiling are used to optimize the performance gains. Although some of these techniques are used in the original implementation found in [14], it does not appear that they were used extensively.

It seems that the research and surveying done in [12] can be applied to our work here and our speedup will benefit from their results.

Unfortunately, this work was done before multi-GPU systems were fully available and accessible using CUDA. The optimization strategies presented in this work will no doubt be useful; however, other strategies also must be employed that considered systems with more than a single GPU (such as those mentioned here from [15]). Leveraging techniques in both realms will offer the greatest amount of performance gain.

### 3.7 Similar Algorithms

A similar algorithm that suffers from the same problems that pairwise comparisons does (each element much be traversed and compared), is search. The work done in [11] explores how search can be optimized for the CUDA framework. Two approaches are given to the bitonic search algorithm that is used: the initial, most straight-forward implementation, and another implementation optimized around the CUDA memory model. The latter offered significant performance increases over the original.

A similar approach will be applied to the work that is used from [14]. Namely, shared memory will be used as efficiently as possible to allow CUDA to perform at its greatest capacity. This is a vital concept in general when optimizing for CUDA; however, the work and challenges from [11] are similar to our own, so the implementation overlaps more than arbitrary examples.

Another important contribution to the field and concepts of distributed systems is the work presented in [5]. MapReduce has become an extremely useful

tool for companies like Google with massive data centers that must distribute the work to perform as reasonable (real-time) levels. MapReduce explores how work might be delegated to multiple workers that exist in a particular distributed system. These concepts apply to the problem of performing the necessary GCD comparisons using multiple GPUs in a single system.



# Chapter 4

## Validation

Accuracy and speedup will both be validated to gain a measure of success.

To validate accuracy, a sequential implementation will be compared against. Much as the work found in [14] performed accuracy validation, a separate, sequential implementation will be created and ran alongside the parallel implementation. Since this work will function as an extension to [14], the validation process will likely be identical or even reused. The results each implementation produces must match for a given set of keys (i.e. both implementations must find the same set of vulnerable keys given identical sets to start with).

An additional limitation realized in [14] was the runtime of the sequential version becoming too long for large sets of keys. Since this is known to be an extremely time-consuming process, it will not always be performed with larger sets of keys. If correct operation of the parallel implementation for small (i.e. sets that can be run sequential in a reasonable amount of time) can be validated using many varying sets of keys of increasing number, it can be assumed the same level of accuracy will continue for larger sets that will not be run sequentially.

Thus, accuracy will be validated for sets no larger than 200,000 keys (as in [14]), and assumed to be accurate (given success of the aforementioned validation) for sets larger than this.

Speedup validation also involves comparing with a sequential implementation; however, this time the runtimes will be compared so an overall speedup may be calculated. Again, much as in [14], speedup will be calculated by taking the ratio of the sequential and parallel runtimes.

As seen in [14], the speedup reached an upper threshold, namely, 27.5. Since this work will expand upon on that, it seems natural to use 27.5 as a baseline. However, there are some differences to consider. First, in [14], sets of keys that could fit entirely in the GPU’s memory were tested and presented. This ensures that the number of memory transfers to the memory card is 1: this significantly reduces memory transfer time, and would help increase the speedup. Since one goal of this work is to allow an arbitrary number of keys to be used in a set, this will likely affect the overall speedup since memory transfers will become smaller and more frequent.

On the other hand, another goal is to add multi-card support which would further increase the total speedup of the final, parallel implementation. Since we will be looking at how fast additional cards make the process (compared to the single-card implementation), the sequential runtime will become less important. It will continue to serve as a baseline to the entire algorithm though.

Considering the modifications done, speedup near 25 seems like a reasonable expectation. Functionality that is expected to decrease the original 27.5, as well as additions that are expected to increase it will both be researched and added. Therefore, a speedup in the same area appears to be a fair expectation of success.

# Bibliography

- [1] RSA security announces new digital rights management solution. [http://www.rsa.com/press\\_release.aspx?id=5159](http://www.rsa.com/press_release.aspx?id=5159). Accessed: 22 February 2013.
- [2] RSA security supports open mobile alliance DRM 2.0 for delivery of secure content. [http://www.rsa.com/press\\_release.aspx?id=3337](http://www.rsa.com/press_release.aspx?id=3337). Accessed: 27 February 2013.
- [3] Transport security layer (TLS) concepts. [http://tech.kaazing.com/documentation/xmpp/3.5/security/c\\_tls.html](http://tech.kaazing.com/documentation/xmpp/3.5/security/c_tls.html). Accessed: 3 March 2013.
- [4] The RSA Factoring Challenge. <http://www.rsa.com/rsalabs/node.aspx?id=2093>, 2007.
- [5] J. Dean and S. Ghemawat. MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [6] C. Fu. Chunking data across multiple C++ AMP kernels. <http://blogs.msdn.com/b/nativeconcurrency/archive/2012/03/05/chunking-data-across-multiple-c-amp-kernels.aspx>, March 2012.
- [7] N. Fujimoto. High-throughput multiple-precision GCD on the CUDA archi-

- ture. In *Signal Processing and Information Technology (ISSPIT), 2009 IEEE International Symposium on*, pages 507–512. IEEE, 2009.
- [8] R. Holz, L. Braun, N. Kammenhuber, and G. Carle. The SSL landscape: a thorough analysis of the x.509 PKI using active and passive measurements. In *Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference*, pages 427–444. ACM, 2011.
- [9] C. Labovitz. Internet Traffic Evolution 2007-2011. In *Global Peering Forum, April*, 2011.
- [10] A. Lenstra, J. Hughes, M. Augier, J. Bos, T. Kleinjung, and C. Wachter. Ron was wrong, Whit is right. *IACR eprint archive*, 64, 2012. [Online; accessed May 2012].
- [11] H. Peters, O. Schulz-Hildebrandt, and N. Luttenberger. Fast in-place, comparison-based sorting with CUDA: a study with bitonic sort. *Concurrency and Computation: Practice and Experience*, 23(7):681–693, 2011.
- [12] S. Ryoo, C. Rodrigues, S. Bagsorkhi, S. Stone, D. Kirk, and W. Hwu. Optimization principles and application performance evaluation of a multithreaded GPU using CUDA. In *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and practice of parallel programming*, pages 73–82. ACM, 2008.
- [13] D. Schaa and D. Kaeli. Exploring the multiple-GPU design space. In *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on*, pages 1–12. IEEE, 2009.
- [14] K. Scharfglass, D. Weng, J. White, and C. Lupo. Breaking weak 1024-bit

- RSA keys with CUDA. In *Parallel and Distributed Computing, Applications and Technologies (PDCAT)*, 2012.
- [15] J. Thibault and I. Senocak. Cuda implementation of a Navier-Stokes solver on multi-GPU desktop platforms for incompressible flows. *Mechanical and Biomedical Engineering Faculty Publications and Presentations*, page 4, 2009.
- [16] R. Wu, B. Zhang, and M. Hsu. Clustering billions of data points using GPUs. In *Proceedings of the combined workshops on UnConventional high performance computing workshop plus memory access workshop*, pages 1–6. ACM, 2009.