

COMP 360

FINAL PROJECT

RSA Moduli Pitfalls

Authors:

Matt Adelman

Evan Carmi

Adam Forbes

December 14, 2012

CONTENTS

1 History:	3
2 Original Ideas:	3
3 Lit Review:	3
4 Algorithm Implementation:	4
5 Data Collection:	5
6 Experiments:	6
7 Analysis of Data:	6
8 Problems:	8
9 Conclusion:	8
10 Further Research:	8
A Appendix of code:	9

Abstract

The algorithm RSA is one of the most widely implemented algorithms in the world. The security of RSA implementations has long been a question of study. In this paper we aim to examine a vulnerability in public keys of RSA found by Lenstra et al. which states that RSA offers only a 99.8% security because of poorly generated prime numbers. We collect data on the top 1 million websites and attempt to reanalyze this claim.

Keywords: RSA, security, prime numbers, Euclidean algorithm, factoring, data collection

1. HISTORY:

RSA is a public key encryption algorithm developed by Ron Rivest, Adi Shamir, and Leonard Adleman at MIT in 1977. Rivest and Shamir were both computer scientists that were working on an “unbreakable” public key encryption method. Rivest and Shamir worked on many different codes, and would pose them as a challenge to Adleman. Forty two of these codes were presented, and Adleman broke them all. Finally on attempt number 43, they created what is now known as the RSA scheme. Incidentally, English Professor Clifford Cocks developed the exact same encryption system in 1973, but it was classified as top-secret, so it was not released until 1997. The RSA algorithm was released for public domain in 1997.

The algorithm operates by using two distinct, large prime numbers to generate public and private keys. Anyone can use the public key to encrypt a message, but only someone with the private key can decrypt it. Breaking the RSA requires factorization of the prime numbers used to create the keys. Since all known prime factorization algorithms are exponential in time, if the prime numbers are large enough RSA is secure.

2. ORIGINAL IDEAS:

When we first started out with this project, our ideas were fairly straightforward. We were going to implement the RSA algorithm in JavaScript, put it on a website, and address some of the common implementation issues

that cause security flaws. A few that we were going to address were timing attacks, public exponent flaws and bad prime generation. This is where things got interesting. We were given the paper by Lenstra et al “Ron Was Wrong, What Is Right.” While more on this paper will be addressed in the next section, the experiments performed by Lenstra exposed how big of an issue bad prime generation was. Our goal then shifted away from rigorous implementation, although we still did create a sufficient implementation, to recreating the experiments from the paper and attempting to find some bad moduli of our own.

3. LIT REVIEW:

As stated earlier, “Ron was Wrong, What is Right” was not the initial focus of our project. However, upon further reading and the study of Lenstra et al. proved to be instrumental in the direction of our project. The paper was published on February 12, 2012 and created a splash in the field of cryptography. Its most influential discovery was that two out of every thousand RSA moduli offer no security at all. Leading to the conclusion that RSA offers a disconcertingly inadequate 99.8% security. This was discovered by utilizing the Euclidean Algorithm and running it against their huge dataset of certificates in order to find moduli with shared primes.

Their findings were shocking to us. We were inspired to recreate their experiments on the state of the web nine months later in the hopes of finding one of two possible conclu-

sions. Either we find broken moduli or we find none. The former would imply a lack of real-world implementation despite the influence of their paper. The latter would indicate an improvement of online security since the data collection of Lenstra et al. However, in order to begin our recreation we needed to first review their work which will be summarized below.

They gathered their data from a collection of public-key databases. Of these the two most significant databases being the SSL Observatory (a project run by the Electronic Frontier Foundation) and one from Massachusetts Institute of Technology. Their method of collection differed greatly from ours as they “did not engage in web crawling, extensive ssh-session monitoring, or other data activities that may be perceived as intrusive or aggressive”, a philosophy we did not share. Their research led them to spend a long time discussing the differences between the public keys they collected and the security differences between them. They did extensive clustering of moduli to describe them as a graph. However, many of their findings are only tangentially related to the work we have done and therefore will be left out of this section.

The most pertinent information we are concerned with lies in the size of their final dataset and the number of compromised RSA moduli agnostic of the types of certificates they came from. They collected 11666704 public keys of which 6386984 had distinct RSA moduli with the rest being split evenly between ElGamal and DSA. Of these 6386984 moduli they found 12934 different ones which could be factored by 14901 distinct primes and were therefore compromised.

4. ALGORITHM IMPLEMENTATION:

For our implementation, we chose JavaScript as our environment because of its ease for web development. We really wanted to end

up with something that we could show off and display in a public setting. The code to the implementation of our RSA algorithm can be found in Appendix A. For our implementation, we relied heavily on the BigInt library found here: <http://leemon.com/crypto/BigInt.js>. While there were other BigInt libraries, out there but this is the one that had all of the functions we needed, and was implemented in an easy format for us to deal with. The main functions we used from this library were the following: `str2bigInt` which takes a string and a base that represents the big integer and converts it for us. `inverseMod` which takes a big integer and a modulus as a big integer, and returns the big integer that is the inverse if one exists, otherwise null. `mult` which has the obvious function of multiplying two big integers. `greater` which returns a boolean depending on whether or not the first argument is greater than the second. `modPow(a,b,c)` which takes three big integers and returns $a^b \bmod c$. `addInt` which takes a big integer and an integer and adds the integer to the big integer. This was used in the computation of $\phi(n)$.

The functions we created are as follows: `generate_key(p:String, q:String, e:String) = [[n:BigInt, e:BigInt], [n:BigInt, d:BigInt]]` where p and q are our distinct primes and e is the public exponent all represented as strings. We then return n which is the modulus, with $n = pq$. Together with e the public exponent and d the private exponent. These together make up our keys. `crypt(key:[BigInt, BigInt], message:BigInt) = message':BigInt` where the *key* is a modulus with an exponent, and the *message* is something to be encrypted or decrypted depending on what *key* with which we call `crypt`.

The next step for our implementation was getting everything to play nice with our web implementation (found here: [4](http://carmi.</p>
</div>
<div data-bbox=)

github.com/rsa/code/). Our plan was to create some helper functions that would access the data from our form, check for errors on input and even randomly generate pseudo-random primes based on a bit length. The code for these helper functions can be found in Appendix A. Some of the big ones are the `gen` function which takes the information given in the form, tests for errors such as bad public exponents, or if the public exponent is not relatively prime to $\phi(n)$. We also have functions to encrypt and decrypt messages given that we have public and private keys entered into our form.

To ensure that future changes do not introduce bugs into our already working codebase we use tests. For our JavaScript implementation of RSA we used QUnit tests, a testing framework in JavaScript used by jQuery among others. This allows our tests to be run easily in the browser along side our implementation. If future changes break a function, then the associated test will fail alerting us to the portion of code that is broken. This included minor tests for generation, encryption, decryption, and even a few number theory functions which were in our original implementation of the RSA algorithm that only worked on integers. This is actually when we realized we needed to switch to arbitrary precision integers, because the exponents get very large, very quickly.

5. DATA COLLECTION:

Any analysis of RSA certificates requires a set of data with which to test. For our project we needed a large set of data so as to best mimic the data that Lenstra used in their analysis of public keys. Lenstra stated in their paper that they used the Electronic Frontier Foundation's SSL Observatory - <https://www.eff.org/observatory>. We initially looked at the SSL Observatory as a possibly source of data for our project. However, unfortunately the

SSL Observatory has not been updated in a few years. Furthermore, it includes no complete guide on how to use its source, and is simply a collection of python scripts that interact with a MySQL database in distinct ways. Thus, we were unable to successfully use the SSL Observatory source code as a means to gather RSA public keys. However, there exists data dumps from the SSL Observatory from 2010. We gathered this data, available through bittorrent, and used it as a starting point.

In order to see the effectiveness and change since the Lenstra paper was published in early 2012 we needed to scrap for fresh data. (The paper briefly mentions that the latest data they looked at was from February 2012.) The SSL Observatory scans all IPv4 space for servers that respond to https requests. We did not have the computational time/power to do this, so we decided on the Alexa top 1 million domains as a source of hosts. Then for each one of these domains we attempted to connect via HTTPS and collect the public key used for the connection. The specific method of collecting these moduli was the This was done by `curl`'ing the domain at port 443 with a HTTP request. If we received a response, we used `openssl s_client` to get the decoded certificate and `openssl x509` to decode the certificate into a easily passable format. However, we noticed that less than 10% of domains responded to our requests. We added some further error handling for these situations. We followed redirects and tried a few common subdomains to improve our rate of responses to over 12%. Out of the 1000000 domains we received 125723 X.509 certificates each containing a moduli. This was the source of our data for our experiments.

In order to deal with our large directory of certificates we needed combine them into a single csv that for the sake of efficiency discards all the impertinent data. We wrote a parsing program in Java that reads in every file and parses it into a combined csv file. The

biggest challenge with dealing with the certificates was handling the exceptions that inevitably arise when dealing with large numbers of only mostly standardized data.

Once we had all of the certificates in a csv file, the next step was to parse out the information we needed. Since the modulus was in a different column for each certificate file, therefore it was in a different column for each line in the csv file. However, there was the unique character pattern “)Modulus(” on each line of the file. Therefore we could match for that using the substring method of the String class in Java. Once we had that, we stripped the modulus of all extraneous characters such as white space and colons. Then since all of the other information we needed in the csv file was before the modulus, we just used a `split(“,”)` on the line to split around commas and easily obtained the domain name, and a unique identifier for each modulus.

Finally we needed code to create a random sample of moduli to more swiftly run an analysis on. We also wanted to make sure that these moduli were unique. We did this by looping through and adding the moduli we chose to a set to make sure we never chose two that were the same. The parsing and random sample getting code is seen in Appendix A.

6. EXPERIMENTS:

The main idea of our experiments is as follows: if we have two moduli m, n we can find if these two moduli share a prime p simply by taking the gcd of the two moduli. This is very fast and is linear in the number of bits. To do around 100000 gcd operations on the machines we were using took around 5 seconds each. Therefore our code, seen in Appendix A, does the following: It takes a file that contains the list of moduli, along with other identifying information. It also takes a smaller file which only contains a small section of the moduli.

This was done so the processing could be broken up onto multiple machines. Once the two files are read into memory, the program stores the moduli in an array, and a unique identifier corresponding to the moduli in another array. We then loop through the data comparing every modulus in the smaller file with every modulus in the bigger file.

We used Java for our experiments due to speed, and the ease with which we can allocate memory to the virtual machine. The default JVM heap space is not large enough to read in the files we were working with (some of which were over 100MB) so we allocated extra memory with the java `-Xmx2g` option which allocates an extra 2GB of data to the virtual machine. We probably did not need this much more allocated, but it was used just as a precaution.

The algorithm was designed to overlook all of moduli that are the same, because while this is interesting information, it gave us no additional data. If the algorithm came across two moduli n, m such that $n \neq m$ and $\gcd(n, m) \neq 1$ then we know we have a prime in common and have factored the modulus, rendering it insecure. Once it finds one, it prints out the gcd, the unique identifier for each modulus and the result of calling `isProbablePrime(80)` which returns true if there is less than a 2^{-80} chance that the gcd is composite. We then can use this data to try and factor the modulus.

7. ANALYSIS OF DATA:

Using the results of Lenstra we analyzed the significance of our findings. Here is a summary of both of our analyses:

Lenstra et al.

- 11666704 public keys
- 6386984 distinct RSA moduli
- 12934 distinct broken moduli
- 14901 distinct primes

Us

- 105984 distinct RSA moduli
- ~60000 moduli analyzed
- 0 broken moduli
- 0 primes

These results are reassuring. They show that out of the 105984 distinct RSA moduli observed by our scraping we could not find any that were compromised. Unfortunately, due to a lack of time and resources (we ran analysis on six computers for approximately two weeks) we were only able to verify around 60000 of the moduli against the set of 105984. Our analysis is ongoing and therefore 60000 is constantly increasing.

The set of 60000 were spread over six clusters arranged alphabetically on their url. We will assume that these were chosen randomly as the original list of the top 1 million domain names was outputted alphabetically and not by ranking. In addition to this prefixes to domains and subdomains further splits up the set of certificates collected.

Our first attempt at calculating the statistical significance of our results was to ask the following question: if there is a .002 chance of finding a broken moduli in the set of Lenstra et al. how many moduli should we expect to find?

- Let n be the size of our sample set of distinct moduli
- Let p be the probability of finding a broken moduli
- Let m be the expected number of broken moduli
- $np = m$
- $105984 \cdot .002 = 212$

We would expect to find 212, and we found 0. Is this significant? Yes it is, because the chance of finding no broken moduli is $(.998)^{60000} = 6.7995 \cdot 10^{-53}$.

However, there are many problems with this analysis. The largest being that we were only able to scrape a dataset with a size that's 1.66% of that analyzed by Lenstra et al. This

wouldn't be a large problem if you could identify the broken moduli by themselves, but by the very nature of our analysis you need to run the Euclidean Algorithm on every single moduli against every other moduli, hence the unsurpassable amount of time and computing power required to factor the moduli. This means that the chances of finding a broken moduli scale very badly with a smaller sample size.

Lenstra et al. found a large range of compromised moduli clusters. However, for this analysis let's assume the worst case: every modulus has a pair that shares a prime. In the case of the research done by Lenstra et al. we have 14901 distinct primes and therefore a total of 29802 pairs. This is no the actual case because their research shows clusters of moduli that share the same prime. If we were to include the statistical probabilities of every group size the model would become too cumbersome and therefore analyzing the worst case scenario provides a good estimate.

Due to a lack of probability experience we have written a program that simulates our statistical problem and determines the expected number of pairs. The program runs as follows:

- Populate an array of size 6386984 with 0s
- Loop 14901 times. Each time randomly assigning two unassigned elements to be a and negative a . These will act as the pairs of moduli that share a prime.
- Create another array of size 105984 and populate it with a random selection from the first array.
- Count the number of pairs in the new array

This program was run 100000 times and found that, on average, there are 8.188 pairs in the second array. Meaning that we should expect to find at 8 pairs in our dataset - a fact we did not. Now that we know we should expect 8 pairs in our dataset we need to calculate the

probability of our findings. This code is seen in Appendix A.

First we need to calculate the probability any given moduli should be broken. $p = 8/105984 = 7.73 \cdot 10^{-5}$ But if we only completed analyzing 60000 moduli, how many should we find to be broken? $\#broken = (7.73 \cdot 10^{-5}) \cdot 60000 = 4.638$ Let's round this up to a probable 5 broken moduli.

Therefore we can say with confidence that 8 out of every 100000 should be broken in our smaller set and we should expect to find a little under 4 broken moduli in our analysis. Notice that this chance of pairs is based on the research of Lenstra et al. and yet gives us a substantially worse chance of finding broken moduli. We have a 0.0000773 chance of finding a compromised moduli out of 105984 whereas their research found a 0.002 chance of finding a compromised moduli out of 6386984 moduli.

8. PROBLEMS:

Computation Power - The largest problem we encountered was computational power. Simply put, we did not have the computing power to analyze all the data we gathered. There primary function we wanted to do was run a gcd across every moduli we encountered. And we couldn't do this in time. There were a few methods we took to try and solve this. We split our data into multiple sections and ran it across a few different machines. Nonetheless, we were unable to process enough data fast enough. We simply needed a compute cluster, or more time. At some points we had data running on five different machines.

SSL Observatory - We attempted to contact the EFF SSL Observatory maintainer for help getting the Observatory running, and to ask if there was a newer data dump than the 2010 data that was on bittorrent. We never received a response. It's likely that the not-for-profit EFF didn't have anyone working on the

project who wanted to help us.

Same Moduli - Something we did not expect but ran into quite frequently were the existence of many moduli that were the same. For instance out of the 125723 moduli that we collected, just over 104000 were unique. Therefore around 18% of the moduli we collected were not unique. In addition, around 10% of that was all from the same domain *.googleusercontent.com. The implication of this is that we could not run as many comparisons as we would have liked, because we had to throw out a lot of same data. The fact that Google uses the same moduli for so many of its domains really speaks to the security of this algorithm

9. CONCLUSION:

We did not find any moduli that did not offer security. What are the chances of this happening? $p = (1 - 0.0000773)^{60000} = 0.00968$ There is a 1% chance of not finding any moduli in our set, and therefore we can say with 99% confidence that our findings show an improvement in RSA security because we were unable to find any broken moduli.

One possible explanation for this can be found in the bias of our domain names. We scraped from the top 1 million most trafficked websites online, and therefore one can hypothesize that there is a lot of overlap between this list and the list of most maintained websites online. The highest probability of finding a bad modulus may lie in looking at the top million least visited websites.

10. FURTHER RESEARCH:

The first priority would be to collect a larger dataset of moduli closer to that of Lenstra et al. as it would greatly improve the statistical significance of our research. With the data we collected, including the older 2010 data dump from the SSL Observatory, we have a large

set of moduli. We discussed the possibility of creating a service that would check any new moduli against the entire every seen set of moduli to ensure there are never duplicates. This would, in effect, allow us verify that certificates are safe prior to them being used by

websites. In fact, we hope that some certificate authorities are already doing this to ensure that within their own network of certificates there are no duplicates. The fact that gcd is so fast makes a database, and a test like this incredibly feasible.

A. APPENDIX OF CODE:

Listing 1: A ruby script that interfaces with unix command line utilities to collect X.509 certificates given a list of hosts.

```

1 # A ruby shell script to scrape for https certificates and decode them.
2 # @author - Evan Carmi
3
4 # a csv file with a number and host on each line. Ex: 1,google.com.
5 # Often useful: http://s3.amazonaws.com/alexa-static/top-1m.csv.zip
6
7 HOST_FILE = 'unseen.txt'
8 PORT = 443
9
10 # Directory for downloaded certificates
11 CERT_DIR = 'certs'
12
13 # Directory for decoded certificates
14 DECODED_DIR = 'decoded'
15
16 'mkdir -p #{CERT_DIR}/'
17 'mkdir -p #{DECODED_DIR}/'
18
19 def cmd_runner(cmd)
20   out = '#{cmd}'
21   unless $? .success?
22     puts "\n\n\nERROR\nFailed while running #{cmd} with error: #{out}\n\n"
23   end
24   $? .success?
25 end
26
27 File.open(HOST_FILE).each_line do |line|
28   domain = line.strip
29
30   # first check if server responds to https request
31   cmd = "curl -sI https://#{domain} --connect-timeout 2 -m 2"
32   worked = cmd_runner(cmd)
33   unless worked
34     domain.prepend("www.")
35     cmd = "curl -sI https://#{domain} --connect-timeout 2 -m 2"
36   end
37   next unless cmd_runner(cmd)
38
39   # download certificate
40   cmd1 = "echo -n | openssl s_client -connect #{domain}:#{PORT} | sed -ne '/--BEGIN CERTIFICATE-/,/--END CERTIFICATE-/p' > #{CERT_DIR}/#{domain}.cert"
41   next unless cmd_runner(cmd1)
42
43   # unpack certificate
44   cmd2 = "openssl x509 -in #{CERT_DIR}/#{domain}.cert -text -noout > #{DECODED_DIR}/#{domain}.cert.decoded"
45   next unless cmd_runner(cmd2)
46
47   cmd3 = "echo '#{domain}' >> downloaded_hosts.txt"
48   next unless cmd_runner(cmd3)
49
50 end

```

Listing 2: A Java class that takes the csv file of certificates and extracts all the useful information.

```

1 /*
2  * Matt Adelman
3  * RSA Final Project
4  */
5
6 import java.math.BigInteger;
7 import java.util.Scanner;
8 import java.io.FileNotFoundException;
9 import java.io.FileReader;

```

```

10
11 public class FormatNewData {
12
13     public static void main(String[] args) {
14
15         String[] lines = new String[125723];
16
17         try {
18             // Getting the file
19             String fileName = "formatted-certs-all.csv";
20             Scanner file = new Scanner(new FileReader(fileName));
21             // Looping through the file
22             for (int i = 0; i < lines.length; i++) {
23                 String line = file.nextLine();
24                 String add = "";
25                 int start = line.indexOf("Modulus");
26                 String rest = line.substring(start);
27                 start = rest.indexOf(":");
28                 rest = rest.substring(start);
29                 String modulus = (rest.split(",")[0]).replaceAll(":", "");
30                 String[] splits = line.split(",");
31                 add = add + splits[0] + ",";
32                 add = add + splits[2] + ",";
33                 add = add + modulus + "\n";
34                 lines[i] = add;
35             }
36
37             for (int i = 0; i < lines.length; i++) {
38                 System.out.print(lines[i]);
39             }
40
41         }
42
43
44
45         // Catching a FileNotFoundException exception
46         catch (FileNotFoundException fnfe) {
47             System.out.println("File was not found");
48             System.exit(0);
49         }
50
51     }
52
53 }

```

Listing 3: A Java class to get the random sample of moduli.

```

1 /*
2  * Matt Adelman
3  * RSA Final Project
4  */
5
6 import java.util.Scanner;
7 import java.io.FileNotFoundException;
8 import java.io.FileReader;
9 import java.util.HashSet;
10 import java.util.Random;
11
12 public class GetRandomSample {
13
14
15     public static void main(String[] args) {
16
17         int count = 0;
18         int samples = 10000;
19
20         String[] allLines = new String[125723];
21
22         String[] sample = new String[10000];
23         String[] modArr = new String[allLines.length];
24         Random rand = new Random();
25         HashSet mods = new HashSet(20000);
26
27         try {
28             // Getting the file
29             String fileName = "new-form-data.csv";
30             Scanner file = new Scanner(new FileReader(fileName));
31             // Looping through the file
32             for (int i = 0; i < allLines.length; i++) {
33                 String line = file.nextLine();
34                 allLines[i] = line + "\n";
35                 modArr[i] = line.split(",")[2];
36             }
37         }
38
39         catch (FileNotFoundException fnfe) {

```

```

40         System.out.println("File was not found");
41         System.exit(0);
42     }
43
44     while (count < samples) {
45         int index = rand.nextInt(allLines.length);
46         if (modArr[index].length() > 10 && mods.add(modArr[index])) {
47             sample[count] = allLines[index];
48             count++;
49         }
50     }
51
52     for (int i = 0; i < sample.length; i++) {
53         System.out.print(sample[i]);
54     }
55
56 }
57
58 }

```

Listing 4: A Java class that reads in two files containing our moduli and other information and prints out a prime if we find one between two moduli.

```

1  /*
2   * Matt Adelman
3   * RSA Final Project
4   */
5
6  import java.math.BigInteger;
7  import java.util.Scanner;
8  import java.io.FileNotFoundException;
9  import java.lang.NumberFormatException;
10 import java.lang.StringIndexOutOfBoundsException;
11 import java.io.FileReader;
12
13 public class RandomSampleCheck {
14
15     public static void main(String[] args) {
16
17         int count1 = 0;
18         int count2 = 0;
19         // Public moduli
20         BigInteger[] keys1 = new BigInteger[2000];
21         BigInteger[] keys2 = new BigInteger[125723];
22         // associated data
23         String[] data1 = new String[keys1.length];
24         String[] data2 = new String[keys2.length];
25         BigInteger big;
26
27         try {
28             // Getting the file
29             // Change each time
30             String fileName1 = "xaa";
31             Scanner file1 = new Scanner(new FileReader(fileName1));
32             String fileName2 = "new-form-data.csv";
33             Scanner file2 = new Scanner(new FileReader(fileName2));
34             // Looping through the file
35             for (int i = 0; i < keys1.length; i++) {
36                 // Getting the line of the file with pertinent info
37                 String line = file1.nextLine();
38                 // Break the line, because it's a CSV
39                 //
40
41                 String[] lineScan = line.split(",");
42                 data1[i] = lineScan[1]; //id number
43                 String modulus = lineScan[2];
44                 big = new BigInteger(modulus, 16);
45                 keys1[count1] = big;
46                 count1++;
47             }
48             for (int i = 0; i < keys2.length; i++) {
49                 // Getting the line of the file with pertinent info
50                 String line = file2.nextLine();
51                 // Break the line, because it's a CSV
52                 //
53
54                 String[] lineScan = line.split(",");
55                 data2[i] = lineScan[1]; //id number
56                 String modulus = lineScan[2];
57                 big = new BigInteger(modulus, 16);
58                 keys2[count2] = big;
59                 count2++;
60             }
61             System.out.println("Done!");
62             // Timing

```

```

63     Long startTime = new Long(System.currentTimeMillis());
64     // Pairwise GCD of moduli
65     for (int i = 0; i < count1; i++) {
66         if ((i % 1000) == 0) {
67             System.out.println("Working on i = " + i);
68         }
69         for (int j = 0; j < count2; j++) {
70             BigInteger gcd = keys1[i].gcd(keys2[j]);
71             BigInteger one = BigInteger.ONE;
72             // If the Moduli are equal we get no useful info
73             if (!(one.equals(gcd)) && !(keys1[i].equals(keys2[j]))) {
74                 System.out.println(data1[i] + ", " + data2[j] +
75                     ", " + gcd + gcd.isProbablePrime(80));
76             }
77         }
78     }
79
80     Long endTime = new Long(System.currentTimeMillis());
81
82     //System.out.println(endTime.intValue() - startTime.intValue());
83     //System.out.println(count);
84     file1.close();
85     file2.close();
86 }
87
88 // Catching a FileNotFoundException exception
89 catch (FileNotFoundException fnfe) {
90     System.out.println("File was not found");
91     System.exit(0);
92 }
93 catch (NumberFormatException nfe) {
94     System.out.println(count1 + 1);
95     big = BigInteger.ONE;
96 }
97 catch (StringIndexOutOfBoundsException sioobe) {
98     System.out.println(count1 + 1);
99     System.out.println("out of bounds");
100 }
101
102
103
104 }
105
106
107 }

```

Listing 5: A Java class that simulates the probable number of pairs of bad moduli we should have found.

```

1  /*
2  Adam Forbes
3  RSA Project
4  */
5
6
7  import java.math.*;
8  import java.util.Arrays;
9
10 public class statTest {
11     public static void main(String[] args) {
12         //The following program is written to figure out the chances of getting
13         //no pairs in our smaller set of moduli collected.
14         //Lenstra et al. used 6 386 984 moduli and found 14901 distinct primes
15         int numTests = 100000;
16         int lenstraSetSize = 6386984;
17         int numberOfPairs = 14901*2;
18         int ourSetSize = 105984;
19         int[] numPairs = new int[numTests];
20         populateNegOne(numPairs);
21
22         double sum = 0;
23         for (int i = 0; i < numTests; i++) {
24             int[] lenstraSet = new int[lenstraSetSize];
25             populateZero(lenstraSet);
26             populatePairs(lenstraSet, numberOfPairs);
27
28             int[] ourSet = new int[ourSetSize];
29             populateSubset(lenstraSet, ourSet);
30
31             numPairs[i] = numPairs(ourSet);
32
33             sum += numPairs[i];
34
35             if (i % 100 == 0) {
36                 System.out.println("Test number: " + i + " of " + numTests);
37                 System.out.println("Average number of pairs: " + sum/i);
38             }
39         }
40     }
41
42     public static void populatePairs(int[] arr, int numPairs) {
43         //We will use the following notation:
44         //0: an unbroken moduli
45         //k: a broken moduli p1
46         //-k: the pair of p1
47
48         //we don't want to use 0 here because our default slot is 0!
49         for (int i = 1; i <= numPairs; i++) {
50             int p1 = (int)(Math.random() * arr.length);
51             int p2 = (int)(Math.random() * arr.length);
52
53             while (arr[p1] != 0) {
54                 p1 = (int)(Math.random() * arr.length);
55             }
56             while (arr[p2] != 0 && p2 != p1) {
57                 p2 = (int)(Math.random() * arr.length);
58             }
59             //Now we have chosen two empty spots the pair
60
61             arr[p1] = i;
62             arr[p2] = -i;
63         }
64     }
65
66     public static void populateSubset(int[] arr, int[] subArr) {
67         for (int i = 0; i < subArr.length; i++) {
68             //Because the array is populated at random we can assume
69             //that by choosing the first subArr.length number of
70             //moduli we are getting a random sample
71             subArr[i] = arr[i];
72         }
73     }
74
75     public static void populateZero(int[] arr) {
76         for (int i = 0; i < arr.length; i++){
77             arr[i] = 0;
78         }
79     }
80
81     public static void populateNegOne(int[] arr) {
82         for (int i = 0; i < arr.length; i++){
83             arr[i] = -1;
84         }
85     }
86 }

```

```
87 public static int numPairs(int[] arr) {
88     //this method will check for the number of pairs in an array
89     int[] sortedArr = Arrays.copyOf(arr, arr.length);
90     Arrays.sort(sortedArr);
91     int numPairs = 0;
92
93     //in order to check for pairs we need only to iterate from the start
94     //of the sorted array to the beginning of the 0s, each time checking
95     //if there exists a match to the negative number in the positive side
96     //A custom built sorting algorithm would probably be the most efficient
97     //method to do this ... but time constraints have always been a problem
98     int i = 0;
99     while (sortedArr[i] < 0) {
100         //if the search is < 0 then there is no pair!!
101         if (Arrays.binarySearch(sortedArr, sortedArr[i]*-1) > 0) {
102             numPairs++;
103         }
104         i++;
105     }
106
107     return numPairs;
108 }
109
110 public static void printArr(int[] arr) {
111     for (int i = 0; i < arr.length; i++) {
112         if (i % 20 != 0) {
113             System.out.print(" " + arr[i]);
114         }
115         else {
116             System.out.println();
117         }
118     }
119     System.out.println();
120 }
121 }
```