

# Generating safe primes

Joachim von zur Gathen and Igor E. Shparlinski

Communicated by Spyros Magliveras

*Dedicated to Professor Tran Van Trung on the occasion of his 65th birthday*

**Abstract.** Safe primes and safe RSA moduli are used in several cryptographic schemes. The most common notion is that of a prime  $p$ , where  $(p - 1)/2$  is also prime. The latter is then a Sophie Germain prime. Under appropriate heuristics, they exist in abundance and can be generated efficiently. But the modern methods of analytic number theory have – so far – not even allowed to prove that there are infinitely many of them. Also for other notions of safe primes, there is no algorithm in the literature that is unconditionally proven to terminate, let alone to be efficient. This paper considers a different notion of safe primes and moduli. They can be generated in polynomial time, without any unproven assumptions, and are good enough for the cryptographic applications that we are aware of.

**Keywords.** Safe prime, Sophie Germain prime, Hofheinz–Kiltz–Shoup cryptosystem.

**2010 Mathematics Subject Classification.** 11N25, 94A60.

## 1 Introduction

There are various notions of *safe primes*  $p$  and *safe RSA moduli*  $N = pq$ , where  $p$  and  $q$  are safe primes, in the cryptographic literature. Two standard conditions on an integer (greater than 1) are  $y$ -smooth (with all prime factors at most  $y$ ) and  $y$ -rough (with all prime factors greater than  $y$ ). The condition “ $y$ -rough” implies “not  $y$ -smooth”, but not conversely. A *Sophie Germain prime*  $\ell$  is such that  $p = 2\ell + 1$  is also prime, and we call its “elder sister”  $p$  a *Marie Germain prime*<sup>1</sup>. More generally, we call  $p = 2\ell + 1$  a *Marie Germain<sub>i</sub> prime* if  $\ell$  is the product of exactly  $i$  distinct primes. Our safe primes will be in the set  $\text{MG}_{\leq 2} = \text{MG}_1 \cup \text{MG}_2$  of primes  $p$  with  $(p - 1)/2$  having either one or two prime factors. The set  $\text{MG}_1$  appears in several cryptographic protocols, but it is still unknown whether it is

---

The first author’s work was supported by the B-IT Foundation and the Land Nordrhein-Westfalen. The second author’s work was supported by the Australian Research Council grant DP110100628 and Singapore National Research Foundation grant CRP2-2007-03.

<sup>1</sup> This alludes to Marie-Madeline Germain, the elder sister of the mathematician Marie-Sophie Germain.

infinite – although there is no reason to doubt this. It seems unlikely that this long-standing open problem will be resolved in the near future.

The purpose of this paper is to show that the safe primes in  $MG_{\leq 2}$

- exist in abundance,
- can be effectively generated,
- can be successfully used instead of  $MG_1$  primes in many protocols.

In order to place this in context, we assemble four properties of primes  $p$  from the literature in the following table, where  $\ell = (p - 1)/2$ .

Marie Germain	$\ell$ prime
Marie Germain $_{\leq 2}$	$\ell$ is prime or product of two primes
rough	$\ell$ is a product of large primes
not smooth	$\ell$ has a large prime divisor

We have not quantified the four notions, and at this point have the implications

$$\text{Marie Germain} \implies \text{Marie Germain}_{\leq 2} \implies \text{rough} \implies \text{not smooth}.$$

We now add more detail. The most stringent notion asks for  $\ell$  to be a *Sophie Germain prime*. This is also the most common one in the cryptographic literature. It is put forth by Menezes, van Oorschot and Vanstone [25, Section 4.6.1] and by Galbraith [11], and used in Shoup [30] and Hofheinz, Kiltz and Shoup [18] (and even in the Wikipedia entry on “safe primes”). Furthermore, Naccache [26] also uses this notion, but assumes incorrect heuristics for the probability of random  $p$  and  $(p - 1)/2$  to be prime. The work of Damgård and Koprowski [7] first appeals to this notion, but says, correctly, on [7, p. 153] that “we do not even know if there are infinitely many safe primes”. Indeed, it is conjectured that there are about  $cx/\ln^2 x$  Sophie Germain primes up to  $x$ , for some explicit constant  $c$ ; see Conjecture 6.1 below. If the conjecture is true, then one can efficiently generate such primes. This works quite well in practice. However, the currently available methods of analytic number theory have not even allowed to show that there are infinitely many Sophie Germain primes. Thus no cryptosystem assuming an infinite supply of them can be proven unconditionally to work.

As usual, we assume a security parameter  $n$ . Then “large” usually means values which are exponential in  $n$ , for example primes with about  $n/2$  bits, and moduli with about  $n$  bits. “Small” means polynomial in  $n$ , that is, with  $O(\log n)$  bits.

$MG_2$  primes have not been used in cryptography, but we show in this paper that they provide a useful concept: it can be efficiently sampled, and is good enough for the cryptographic applications that we are aware of.

In the third condition,  $y$ -rough is used with a small  $y$ , in the sense above. The signature scheme of Gennaro, Halevi and Rabin [13] works with this notion, as does the distributed moduli generation of Damgård and Koprowski [7] and of Fouque and Stern [10]; they assume in the proof of their Theorem 2 that primes are uniformly distributed. The roughness condition is also employed in [7, 22, 27], where  $y$  is the number of players in a certain multiparty computation, and [28, Section 3.5] says that “safe primes are unfortunately less dense than unrestricted primes”, but no lower bound is attempted in this experimental paper. The work of Joye and Paillier [22] gives a heuristic improvement for this task by sieving modulo the product of small primes greater than 2 and making sure that  $(p - 1)/2$  has no small factor by  $p$  being a nonsquare modulo each small prime. Clearly the second notion implies this one, even with a value for  $y$  that may be as “large” as is allowed in the second one.

The fourth notion requires that  $k$  have at least one “large” prime factor. This used to be required for RSA moduli in some scenarios, in order to resist certain factorization methods. It follows from the third notion when the smoothness and roughness parameters are chosen identically, but is now considered obsolete. A result of Baker and Harman [1] implies that the set of primes  $p$  for which  $p - 1$  is not  $p^{0.677}$ -smooth is of positive relative density in the set of all primes.

Now take some  $p \in \text{MG}_{\leq 2}$  with only large prime factors in  $(p - 1)/2$ , and a different prime  $q$  with the same property. Then the modulus  $N = pq$  enjoys the following properties:

- $N$  is a Blum integer,
- $\varphi(N)/4$  has only large prime factors,
- the square of a random element in the residue ring  $\mathbb{Z}_N$  is a generator of the subgroup of squares in the unit group  $\mathbb{Z}_N^\times$  with probability exponentially close to 1.

For the Marie Germain $_{\leq 2}$  primes, the topic of this paper, a quantified version of these properties is proven in Theorem 5.2 and Corollary 5.3. Algorithmically, the crux is to show that a certain rejection sampling process producing such integers works in polynomial time, that is, the primes we generate form an inverse polynomial fraction of all integers up to some bound. Our main technical tool is a result of Heath-Brown [16].

Our interest in this theme has been spawned by the cryptosystem of Hofheinz, Kiltz and Shoup [18], whose breaking (in the sense of IND-CCA2 security) is equivalent to factoring the modulus  $N$ . This reduction is free of any unproved hypotheses. They present two versions. One of them uses the Goldreich–Levin predicate. The other one, simpler and more natural, takes the modulus as  $N =$

$pq$ , where  $p$  and  $q$  are two distinct Marie Germain primes. As noted above, this works well heuristically, but no proof for it is in sight. We show here that modern analytic number theory is still powerful enough to close this gap in the argument of Hofheinz, Kiltz and Shoup [18].

In the following, the letters  $\ell$ ,  $p$ , and  $q$  with or without subscripts denote prime numbers, and for real  $x < y$ , we use  $[x..y]$  and  $[x..y]_{\mathbb{R}}$  to denote the sets of integers and reals between  $x$  and  $y$ , respectively. We also use  $(x..y)$  and  $(x..y)_{\mathbb{R}}$  in similar meanings. All our random choices are made uniformly from finite sets, unless explicitly stated otherwise. We use the big-Oh notation and its relatives as sets, so that

$$O(g) = \{f : \exists c |f(x)| \leq c \cdot g(x) \text{ for sufficiently large } x\},$$

where  $f$  and  $g$  are real functions, and pointwise operators, so that, for example,  $g + O(h) = \{g + f : f \in O(h)\}$ .

Section 2 through 5 build up the required algorithmic machinery in several steps. Just after Algorithm 3.1, we explain why a naive sampling method fails and something like the approach of Section 2 is needed; see the end of Section 7 for a simpler but less efficient version. Sections 6 and 7 present extensions and variations of our method, heuristics for the number of Marie Germain $_{\leq 2}$  primes, and, assuming these heuristics, more precise information about the runtime of our algorithms.

## 2 Sampling the harmonic distribution

Our goal in this section is to sample the *harmonic distribution* which gives to an integer  $k$  in some finite interval a probability proportional to  $1/k$ . This yields an approximately uniform sampling of integer points under a hyperbola in Algorithm 2.3, which in turn leads to the random unbalanced moduli in Algorithm 3.1.

We first recall the *inversion method* for the continuous version  $\mathcal{D}^*$  of this distribution; see [23, Section 3.4.1] and [9, Section III.2.2 B]. We have two positive real numbers  $A < B$ , set

$$a^* = \frac{1}{\ln(B/A)},$$

and take the continuous distribution  $\mathcal{D}^*$  on  $(A..B)_{\mathbb{R}}$  with (cumulative) density function

$$F^*(x) = \begin{cases} 0 & \text{if } x < A, \\ a^* \ln(x/A) & \text{if } A \leq x \leq B, \\ 1 & \text{if } x > B. \end{cases}$$

Thus, choosing  $x$  according to  $\mathcal{D}^*$  is equivalent to  $\text{prob}(x \leq y) = F^*(y)$  for all  $y \in (A..B]_{\mathbb{R}}$ . On  $(A..B]_{\mathbb{R}}$ ,  $F^*$  takes values between 0 and 1, increases strictly monotonically, and its functional inverse is

$$G^*(x) = Ae^{x/a^*}.$$

In particular,  $G^*(x) \in (A..B]$  for  $x \in (0..1]_{\mathbb{R}}$ .

The inversion method takes a sample  $u$  from the uniform real distribution on  $(0..1]_{\mathbb{R}}$  and sets  $x = G^*(u)$ . Then for any  $y \in (A..B]_{\mathbb{R}}$ , we have

$$\text{prob}(x \leq y) = \text{prob}(G^*(u) \leq y) = \text{prob}(u \leq F^*(y)) = F^*(y).$$

Thus  $x$  samples  $\mathcal{D}^*$ .

In several steps, we now transform this method into a discrete algorithm that approximately samples the harmonic distribution. In a first step, we take positive real numbers  $A < B$ , the harmonic number

$$H_n = \sum_{1 \leq k \leq n} 1/k$$

for an integer  $n$ , set

$$a = \frac{1}{H_{\lfloor B \rfloor} - H_{\lfloor A \rfloor}} \tag{2.1}$$

and consider the discrete harmonic distribution  $\mathcal{D}$  with

$$\text{prob}_{\mathcal{D}}(k) = \frac{a}{k}$$

on the integers  $k \in (A..B]$ . We use the Euler–Mascheroni constant  $\gamma \approx 0.57721$  and the bounds

$$0 < H_n - (\ln n + \gamma) < \frac{1}{2n}; \tag{2.2}$$

see [14] for sharper estimates which imply (2.2).

The value  $H_A$  has a standard definition also for non-integral  $A$ . In order to avoid confusion, we write  $H(A)$  for  $H_{\lfloor A \rfloor}$  in the following. For our rounded arguments, we use  $|\ln(1+z) - z| \leq z^2$  if  $|z| \leq 1/2$ , and for  $A \geq 4$

$$\begin{aligned} |H(A) - (\ln A + \gamma)| &= |H_{\lfloor A \rfloor} - \ln A + \gamma| \\ &\leq |H_{\lfloor A \rfloor} - (\ln \lfloor A \rfloor + \gamma)| + \left| \ln \frac{\lfloor A \rfloor}{A} \right| \\ &\leq \frac{1}{2\lfloor A \rfloor} + \left| \ln \left( 1 - \frac{A - \lfloor A \rfloor}{A} \right) \right| \leq \frac{1}{2\lfloor A \rfloor} + \frac{1}{A^2} < \frac{1}{A}. \end{aligned}$$

For  $4 \leq A < B$ , we have

$$\begin{aligned} |H(B) - H(A) - \ln(B/A)| &= |H(B) - (\ln B + \gamma) - (H(A) - (\ln A + \gamma))| \\ &< B^{-1} + A^{-1} < 2A^{-1}. \end{aligned}$$

If also

$$\ln(B/A) \geq A^{-1} + (1 + A^{-2})^{1/2}, \quad (2.3)$$

then  $(\ln(B/A) - 2A^{-1}) \cdot \ln(B/A) \geq 1$  and

$$\begin{aligned} |a - a^*| &= \left| \frac{1}{H(B) - H(A)} - \frac{1}{\ln(B/A)} \right| \\ &= \frac{|\ln(B/A) - (H(B) - H(A))|}{|H(B) - H(A)| \cdot \ln(B/A)} \\ &< \frac{2A^{-1}}{|\ln(B/A) - 2A^{-1}| \cdot \ln(B/A)} \leq 2A^{-1}. \end{aligned} \quad (2.4)$$

One can check that the densities of  $\mathcal{D}^*$  on integers and  $\mathcal{D}$  agree closely, but we do not need this here.

We now consider the inversion method, where  $u \in (0..1]_{\mathbb{R}}$  is chosen uniformly at random and the integer

$$t^*(u) = \lfloor G^*(u) \rfloor = \lfloor Ae^{u/a^*} \rfloor$$

is produced. Then  $A \leq t^*(u) \leq B$ . One can check that any  $k \in [A..B]$  is returned with probability  $a^* \ln(1 + 1/k) \sim a^*/k \cdot (1 + O(A^{-1}))$ .

In the next step, we replace  $u$  by a discrete approximation and again round  $G^*(u)$  down. So we have a (large) positive integer  $M$ , choose an integer  $v \in (0..M]$  uniformly at random, so that  $v/M$  is an approximation of  $u$  as above, and produce

$$t(v) = \lfloor G^*(v/M) \rfloor = \lfloor Ae^{v/a^*M} \rfloor. \quad (2.5)$$

For  $k \in (A..B]$ , we let

$$V = \{v \in (0..M] : t(v) = k\} = t^{-1}(k),$$

so that  $b_1(k) = \#V/M$  is the probability with which  $k = t(v)$ . We now claim that  $b_1(k)$  is close to  $a^*/k$ , and provide an error estimate for the approximation quality. For real  $x < y$ , the number of integers in  $(x..y]$  satisfies

$$|\#\{x..y\} - (y - x)| \leq 1. \quad (2.6)$$

Using  $F^*(k+1) - F^*(k) = a^* \ln(1 + k^{-1})$ , we have

$$\begin{aligned} v \in V &\iff k \leq G^*\left(\frac{v}{M}\right) < k+1 \\ &\iff F^*(k) \leq \frac{v}{M} < F^*(k+1) \\ &\iff MF^*(k) \leq v < MF^*(k+1). \end{aligned}$$

Therefore

$$|\#V - a^* M \ln(1 + k^{-1})| \leq 1.$$

For  $k \geq 2$  we find

$$\begin{aligned} \left|b_1(k) - \frac{a^*}{k}\right| &= \left|\frac{\#V}{M} - \frac{a^*}{k}\right| \\ &\leq \left|\frac{\#V}{M} - a^* \ln\left(1 + \frac{1}{k}\right)\right| + \left|a^* \ln\left(1 + \frac{1}{k}\right) - \frac{a^*}{k}\right| \\ &\leq \frac{1}{M} + \frac{a^*}{k^2} = \frac{a^*}{k} \left(\frac{1}{k} + \frac{k}{a^* M}\right) \leq \frac{a^*}{k} \left(\frac{1}{A} + \frac{B}{a^* M}\right). \end{aligned} \quad (2.7)$$

None of the methods sketched so far can be literally implemented on a computer, since we cannot compute a real number like  $e^{v/aM}$  exactly. So we now consider floating-point computations with real numbers using  $m_0$  bits of precision. We take some  $n$  so that all quantities in the algorithm are absolutely at most  $2^n$ . We assume  $m_0 > n$  and set  $m = m_0 - n$ . Then if a real value  $r$  is to be computed, the algorithm computes some  $\tilde{r}$  with  $|r - \tilde{r}| < 2^{-m} = \varepsilon$ , which we call an  $m$ -bit approximation (which in standard parlance is an  $m_0$ -bit approximation). We assume some standard representation where  $\lfloor \tilde{r} \rfloor$  can be computed exactly, by truncating after the decimal (or binary) point. This leads to the following algorithm.

**Algorithm 2.1** (Sampling the harmonic distribution).

INPUT: Positive real numbers  $A$  and  $B$  with  $4 \leq A < B$ , and positive integers  $M$  and  $m$ .

OUTPUT: An integer in  $[A..B]$ .

- (1) Choose an integer  $v \in [1..M]$  uniformly at random.
- (2) Calculate an  $m$ -bit approximation  $\tilde{s}$  to  $Ae^{v \ln(B/A)/M}$ .
- (3) Return  $\lfloor \tilde{s} \rfloor$ .

**Theorem 2.2.** *We assume that  $B < M < 2^n$  and*

$$m \geq \log_2 \left( \frac{M}{A \ln(B/A)} \right), \quad (2.8)$$

*and we calculate numerically with precision  $m + n$ .*

(i) *Let  $k \in (A..B]$  and*

$$\delta_1 = \frac{a^* + 2}{a^* A} + \frac{3B}{a^* M} + \frac{1}{2^{m-2}}, \quad (2.9)$$

*and assume that (2.3) holds. Then the probability  $b_2(k)$  that  $k$  is returned by the algorithm satisfies*

$$\left| b_2(k) - \frac{a}{k} \right| \leq \frac{a^* \delta_1}{k}.$$

(ii) *Any output of Algorithm 2.1 is in  $[A..B]$  and the algorithm uses time polynomial in  $m + n$ .*

*Proof.* For  $1 \leq v \leq M$ , we write  $s(v) = A(B/A)^{v/M} = Ae^{v/a^*M} = G^*(v/c)$ ,  $t(v) = \lfloor s(v) \rfloor$  as in (2.5), and  $\tilde{s}(v)$  for the approximation to  $s(v)$  calculated in step (2). Since rounding down is exact,  $\tilde{t}(v) = \lfloor \tilde{s}(v) \rfloor$  is returned in step (3).

(ii): By assumption, we have  $A/a^*M \geq \varepsilon = 2^{-m}$ . For any  $v \geq 1$ , we have

$$\tilde{s}(v) \geq s(v) - \varepsilon \geq s(1) - \varepsilon = Ae^{1/a^*M} - \varepsilon \geq A \left( 1 + \frac{1}{a^*M} \right) - \varepsilon \geq A$$

and  $\tilde{t}(v) \geq A$ . For any  $v \leq M$ , we have

$$\tilde{s}(v) \leq s(v) + \varepsilon \leq s(M) + \varepsilon = Ae^{\ln(B/A)} + \varepsilon = B + \varepsilon$$

and  $\tilde{t}(v) \leq B$ . It is well known how to compute numerically the required approximations to  $\ln(B/A)$ ,  $e^{v \ln(B/A)/M} = (B/A)^{v/M}$ , and  $s(v)$  in time polynomial in  $m + n$ ; see [5, Sections 4.2.5, 4.4] for some details.

(i): We take some  $k \in (A..B]$  and want to show that  $b_2(k)$  is close to  $a/k$ . We define the five real intervals

$$\begin{aligned} S &= [k..k+1)_{\mathbb{R}}, \\ S_{0,+} &= [k..k+\varepsilon)_{\mathbb{R}}, \\ S_{0,-} &= [k-\varepsilon..k)_{\mathbb{R}}, \\ S_{1,+} &= [k+1..k+1+\varepsilon)_{\mathbb{R}}, \\ S_{1,-} &= [k+1-\varepsilon..k+1)_{\mathbb{R}}. \end{aligned}$$



We take  $t^{-1}(k) = \{v \in [1..M]: t(v) = k\}$ , and similarly for  $\tilde{t}$ . Then

$$\begin{aligned} t^{-1}(k) &= s^{-1}(S), \\ s^{-1}(S \setminus (S_{0,+} \cup S_{1,-})) &\subseteq \tilde{t}^{-1}(k) \subseteq s^{-1}(S \cup S_{0,-} \cup S_{1,+}). \end{aligned} \quad (2.10)$$

We start by considering the case  $s(v) \in S_{0,+}$ . Then  $s(v) = k + \gamma_0$  with  $0 \leq \gamma_0 < \varepsilon$ . Setting  $\gamma_1 = \gamma_0/k$ , we have  $0 \leq \gamma_1 < \varepsilon/k < 1/2$  and  $G^*(v/M) = s(v) = k(1 + \gamma_1)$ . Furthermore,

$$\frac{v}{M} = F^* \circ G^* \left( \frac{v}{M} \right) = F^*(k(1 + \gamma_1)) = F^*(k) + a^* \ln(1 + \gamma_1)$$

and

$$0 \leq \ln(1 + \gamma_1) \leq 2\gamma_1 < \frac{2\varepsilon}{k}.$$

Therefore

$$v = MF^*(k) + a^* M \ln(1 + \gamma_1) \in \left[ MF^*(k) .. MF^*(k) + \frac{2\varepsilon a^* M}{k} \right)$$

which in turn implies that

$$\#s^{-1}(S_{0,+}) \leq \frac{2\varepsilon a^* M}{k} + 1.$$

One finds the same bound for  $\#s^{-1}(S_{1,-})$ ,  $\#s^{-1}(S_{0,-})$ , and  $\#s^{-1}(S_{1,+})$ . It now follows that

$$\#s^{-1}(S_{0,+} \cup S_{1,-}), \#s^{-1}(S_{0,-} \cup S_{1,+}) \leq \frac{4\varepsilon a^* M}{k} + 2.$$

Let  $\mathcal{A}$  be the algorithm described by (2.5), which is just the exact version of Algorithm 2.1, with  $s(v)$  in step (2) calculated exactly and output  $t(v)$ . Now  $\mathcal{A}$  works well and returns  $k$  with probability  $b_1(k)$  close to  $a^*/k$ , namely satisfying (2.7). By (2.10), this happens if and only if  $s(v) \in S$ , so that

$$\left| \frac{\#s^{-1}(S)}{M} - \frac{a^*}{k} \right| = \left| b_1(k) - \frac{a^*}{k} \right| \leq \frac{a^*}{k} \left( \frac{1}{A} + \frac{B}{a^* M} \right),$$

Moreover,

$$\begin{aligned} \#s^{-1}(S) - \left( \frac{4\varepsilon a^* M}{k} + 2 \right) &\leq \#s^{-1}(S \setminus (S_{0,+} \cup S_{1,-})) \leq \#\tilde{t}^{-1}(k) \\ &\leq s^{-1}(S \cup S_{0,-} \cup S_{1,+}) \\ &\leq \#s^{-1}(S) + \frac{4\varepsilon a^* M}{k} + 2. \end{aligned}$$

Using (2.4), it follows that

$$\begin{aligned}
 \left| b_2(k) - \frac{a}{k} \right| &\leq \left| \frac{\#\tilde{t}^{-1}(k)}{M} - \frac{a^*}{k} \right| + \left| \frac{a^* - a}{k} \right| \\
 &\leq \left| \frac{\#\tilde{t}^{-1}(k) - \#s^{-1}(S)}{M} \right| + \left| \frac{\#s^{-1}(S)}{M} - \frac{a^*}{k} \right| + \frac{2}{Ak} \\
 &\leq \frac{4\varepsilon a^*}{k} + \frac{2}{M} + \frac{a^*}{k} \left( \frac{1}{A} + \frac{B}{a^*M} \right) + \frac{2}{Ak} \\
 &\leq \frac{a^*}{k} \left( 4\varepsilon + \frac{2B}{a^*M} + \frac{1}{A} + \frac{B}{a^*M} + \frac{2}{a^*A} \right) = \frac{a^* \delta_1}{k}. \quad \square
 \end{aligned}$$

Using the penultimate rather than the last bound in (2.7), we can replace the summand  $\sigma = 1/A + 3B/a^*M$  by  $\sigma_k = 1/k + 3k/a^*M$ . As a function of a real variable  $k$  on  $(A..B]_{\mathbb{R}}$ ,  $\sigma_k$  is convex, and  $\sigma$  can be replaced by  $\max\{\sigma_A, \sigma_B\}$ . This equals  $\sigma_B$  if and only if  $a^*M \geq 3AB$ .

We now present a method to generate almost uniformly random pairs of positive integers under a hyperbola  $xy = D$ .

**Algorithm 2.3** (Random integers under a hyperbola).

INPUT: Positive real numbers  $8 \leq A < B < C < D < 2^n$ , and integers  $M$  and  $m$  greater than 1.

OUTPUT: A pair  $(k_1, k_2)$  of integers with  $k_1 \in [A..B]$  and  $k_1 k_2 \in [C..D]$ .

- (1) Call Algorithm 2.1 with inputs  $A, B, M$ , and  $m$ , and output  $k_1$ , calculating numerically with precision  $m + n$ .
- (2) Choose a uniformly random integer  $k_2 \in [C/k_1..D/k_1]$ .
- (3) Return  $(k_1, k_2)$ .

For any pair  $(k_1, k_2)$  of integers, we let  $b_3(k_1, k_2)$  be the probability with which it is returned by the algorithm, and write

$$K = \{(k_1, k_2) : k_1 \in [A..B], k_1 k_2 \in [C..D]\}. \quad (2.11)$$

For  $A = C = 1$  and  $B = D$ ,  $\#K = \sum_{k \leq D} \tau(k)$  equals  $2D$  times the average value of Dirichlet's divisor function  $\tau$  on  $[1..D]$ . In our application,  $B \approx D^{1/2}$  is much smaller than  $D$ .

**Theorem 2.4.** *We assume that (2.3) and (2.8) hold,  $a$  is as in (2.1),  $\delta_1$  as in (2.9), and set*

$$\delta_2 = \frac{(a+1)B + 2\delta_1(D - C + aB)}{D - C - B}.$$

- (i) Any output of Algorithm 2.3 satisfies the output specification and for any  $(k_1, k_2) \in K$  we have

$$\left| b_3(k_1, k_2) - \frac{1}{\#K} \right| < \frac{\delta_2}{\#K}.$$

- (ii) The algorithm uses time polynomial in  $m$ ,  $n$ , and  $\ln M$ .

*Proof.* For  $k_1 \in [A..B]$ , the number  $d(k_1)$  of choices in step (2) satisfies by (2.6)

$$\left| d(k_1) - \frac{D-C}{k_1} \right| \leq 1.$$

It follows that

$$d(k_1) \geq \frac{D-C}{k_1} - 1 \geq \frac{D-C-B}{k_1}.$$

Hence

$$\begin{aligned} \left| \frac{1}{d(k_1)} - \frac{k_1}{D-C} \right| &= \left| \frac{D-C-k_1d(k_1)}{d(k_1)(D-C)} \right| \\ &\leq \frac{k_1}{d(k_1)(D-C)} \leq \frac{k_1^2}{(D-C)(D-C-B)}. \end{aligned}$$

Furthermore,

$$\begin{aligned} \left| \#K - \frac{D-C}{a} \right| &= \left| \sum_{A < k_1 \leq B} d(k_1) - (D-C)(H(B) - H(A)) \right| \\ &= \left| \sum_{A < k_1 \leq B} \left( d(k_1) - \frac{D-C}{k_1} \right) \right| \leq \sum_{A < k_1 \leq B} 1 < B. \end{aligned} \quad (2.12)$$

Since  $b_3(k_1, k_2) = b_2(k_1)/d(k_1)$ , we have

$$\begin{aligned} \#K \cdot b_3(k_1, k_2) - 1 &= b_2(k_1) \cdot \frac{\#K}{d(k_1)} - 1 \\ &= \frac{a}{k_1} \cdot \frac{D-C}{a} \cdot \frac{k_1}{D-C} + \frac{(b_2(k_1) - a/k_1)\#K}{d(k_1)} + \frac{a(\#K - (D-C)/a)}{k_1d(k_1)} \\ &\quad + \frac{(D-C)(d(k_1)^{-1} - k_1/(D-C))}{k_1} - 1 \\ &= \frac{(k_1b(k_1) - a)\#K}{k_1d(k_1)} + \frac{a\#K - (D-C)}{k_1d(k_1)} + \frac{(D-C) - k_1d(k_1)}{k_1d(k_1)}. \end{aligned}$$

Therefore,

$$\begin{aligned} & k_1 d(k_1) \cdot |\#K \cdot b_3(k_1, k_2) - 1| \\ & \leq a^* \delta_1 \#K + aB + (D - C)d(k_1) \cdot \frac{k_1^2}{(D - C)(D - C + B)}. \end{aligned}$$

Finally, we derive

$$\begin{aligned} \#K \cdot \left| b_3(k_1, k_2) - \frac{1}{\#K} \right| & \leq \frac{a^* \delta_1 \#K + aB + k_1^2 d(k_1)/(D - C + B)}{k_1 d(k_1)} \\ & = \frac{aB + a^* \delta_1 \#K}{k_1 d(k_1)} + \frac{k_1}{D - C + B} \\ & \leq \frac{aB + (a + 2A^{-1})\delta_1((D - C)/a + B)}{k_1 \cdot (D - C - B)/k_1} \\ & \quad + \frac{B}{D - C + B} \\ & < \frac{(a + 1)B + 2\delta_1(D - C + aB)}{D - C - B} = \delta_2. \end{aligned}$$

The time bound claimed in (ii) follows from Theorem 2.2.  $\square$

### 3 Uniform unbalanced moduli

As a next step, we present an efficient algorithm to generate (almost) uniformly random unbalanced RSA moduli, where one prime factor is allowed to be considerably smaller than the other one. The natural rejection sampling process cannot be proved to be efficient in this situation, as explained after stating the algorithm.

This is used to generate safe primes in Section 4, which in turn leads to the safe moduli of Section 5, the ultimate goal of this paper.

**Algorithm 3.1** (Random unbalanced moduli).

INPUT: Positive real numbers  $8 \leq A < B < C < D$ .

OUTPUT: Primes  $\ell_1 \neq \ell_2$  so that  $A \leq \ell_1 \leq B$  and  $C \leq \ell_1 \ell_2 \leq D$ .

(1) Compute  $M = \lceil 6AB \rceil$  and  $m = \lceil \log_2 6B \rceil$ .

(2) Repeat steps (2a) and (2b) until primes  $\ell_1$  and  $\ell_2$  are found.

a. Call Algorithm 2.3 with the input values as above and output  $(k_1, k_2)$ .

b. If  $k_1$  and  $k_2$  are distinct primes, set  $\ell_1 = k_1$  and  $\ell_2 = k_2$ .

(3) Return  $(\ell_1, \ell_2)$ .

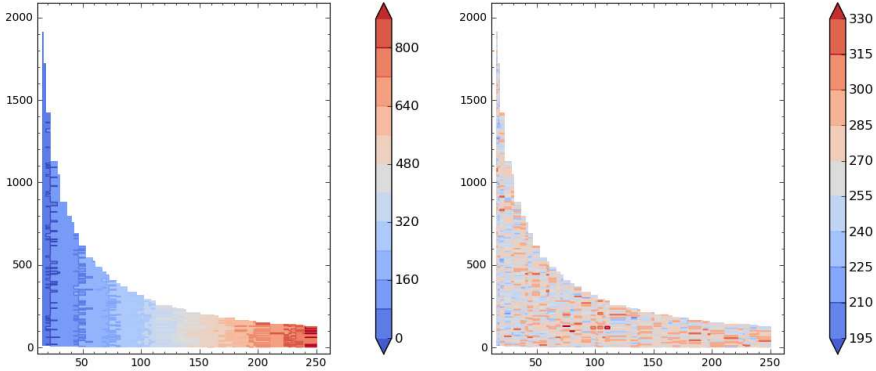


Figure 1. Distribution of 1 000 000 prime pairs  $(\ell_0, \ell_1)$  obtained from different sampling algorithms for  $x = 2^{16}$ .

For simplicity, the algorithm does not return the unbalanced modulus  $N = \ell_1 \ell_2$  explicitly.

As an alternative one might consider the naive approach of generating uniformly randomly first  $\ell_1 \in [A..B]$  and then  $\ell_2 \in [C/\ell_1..D/\ell_1]$ .

This naive approach fails for the following reason. Assume that  $C = o(D)$ ,  $A$  and  $B$  are of different orders of magnitude, smaller than that of  $D$ , and let  $\pi(x)$  denote the number of primes up to  $x$ . Then there are about  $\pi(D/\ell_1) \sim D/\ell_1 \ln(D/\ell_1) \sim D/\ell_1 \ln D$  possibilities for  $\ell_2$ , and a particular  $\ell_2$  is chosen with probability proportional to  $\ell_1$ . The same holds for any particular  $(\ell_1, \ell_2)$ . This value is exponentially smaller for small values of  $\ell_1$  (close to  $A$ ) than for large ones (close to  $B$ ). Thus the distribution on  $(\ell_1, \ell_2)$  is highly nonuniform, as illustrated by the left heat diagram in Figure 1 which is taken from the work of Ziegler and Zollmann [35]. In Algorithm 4.1 below we need (approximately) uniform  $(\ell_1, \ell_2)$ , and this is actually delivered by Algorithm 3.1 as shown on the right of Figure 1.

In the usual RSA key generation, one typically uses primes from dyadic intervals, that is with  $B = 2A$ , and then this nonuniformity is not much of a problem.

We recall the set  $K$  of integer pairs from (2.11), and  $\delta_2$  from Theorem 2.4. For our values of  $A, B, C$ , and  $D$ , we set

$$X = \{(\ell_1, \ell_2) : \ell_1 \text{ and } \ell_2 \text{ distinct primes,} \\ \ell_1 \in [A..B], \ell_1 \ell_2 \in [C..D]\} \subseteq K, \quad (3.1)$$

$$a_1 = \frac{1}{\ln \ln B - \ln \ln A} = \frac{1}{\ln(\log_A(B))}. \quad (3.2)$$

We have the following bounds on  $a_1$  and the size of  $X$ .

**Lemma 3.2.** *Let  $8 \leq A < B < C < D$  be real numbers with*

$$2a_1 \leq \ln^2 A, \quad (3.3)$$

$$82A \ln A \leq 41B \leq D, \quad (3.4)$$

$$2^{34} \leq D. \quad (3.5)$$

Then

$$\left| a_1^{-1} - \sum_{A \leq \ell \leq B} \frac{1}{\ell} \right| \leq \frac{1}{\ln^2 A} \quad (3.6)$$

and

$$\frac{D}{5a_1 \ln D} \leq \frac{D}{5a_1 \ln(D/A)} \leq \#X \leq \frac{2D}{a_1 \ln(D/B)}. \quad (3.7)$$

*Proof.* Mertens' theorem, see [29, Theorem 5], implies that

$$\ln \ln B + c - \frac{1}{2 \ln^2 B} \leq \sum_{\ell \leq B} \frac{1}{\ell} \leq \ln \ln B + c + \frac{1}{2 \ln^2 B}$$

for some constant  $c$ , from which (3.6) follows. An explicit formula and the approximation 0.26150 for  $c$  are given in [32, § I.1.6]. There are at most  $D^{1/2}$  pairs  $(k_1, k_2)$  with  $k_1 = k_2$  and  $k_1^2 \leq D$ . We first bound from below the size of  $X$  as

$$\begin{aligned} \#X &\geq \sum_{A \leq \ell_1 \leq B} (\pi(D/\ell_1) - \pi(C/\ell_1)) - D^{1/2} \\ &\geq \sum_{A \leq \ell_1 \leq B} (\pi(D/\ell_1) - \pi(D/2\ell_1)) - D^{1/2}. \end{aligned} \quad (3.8)$$

The asymptotic value of  $\#X$  is given in Theorem 6.2 below. Now since  $D/2\ell_1 \geq D/2B \geq 20.5$  by (3.4), we see that [29, Corollary 3] implies that

$$\pi(D/\ell_1) - \pi(D/2\ell_1) \geq \frac{3D}{5\ell_1 \ln(D/\ell_1)} > \frac{D}{2\ell_1 \ln(D/A)}.$$

It follows that

$$\begin{aligned} D^{1/2} + \#X &\geq \frac{D}{2 \ln(D/A)} \sum_{A \leq \ell_1 \leq B} \frac{1}{\ell_1} \geq \frac{D}{2 \ln(D/A)} \left( \frac{1}{a_1} - \frac{1}{\ln^2 A} \right) \\ &\geq \frac{D}{2 \ln(D/A)} \left( \frac{1}{a_1} - \frac{1}{2a_1} \right) = \frac{D}{4a_1 \ln(D/A)}. \end{aligned}$$

The assumptions (3.3) and (3.5) yield

$$20a_1 \ln(D/A) \leq 10 \ln^2 A \ln D < 10 \ln^3 D \leq D^{1/2}$$

and

$$\frac{D}{5a_1 \ln D} \leq \frac{D}{5a_1 \ln(D/A)} \leq \frac{D}{4a_1 \ln(D/A)} - D^{1/2} \leq \#X.$$

Furthermore, (3.3) and [29, Corollary 2] imply that

$$\begin{aligned} \#X &\leq \sum_{A \leq \ell_1 \leq B} \pi\left(\frac{D}{\ell_1}\right) \leq \sum_{A \leq \ell_1 \leq B} \frac{5D}{4\ell_1 \ln(D/\ell_1)} \leq \frac{5D}{4 \ln(D/B)} \sum_{A \leq \ell_1 \leq B} \frac{1}{\ell_1} \\ &\leq \frac{5D}{4 \ln(D/B)} \left(a_1^{-1} + \frac{1}{\ln^2 A}\right) < \frac{2D}{a_1 \ln(D/B)} \end{aligned}$$

which concludes the proof.  $\square$

**Theorem 3.3.** *Let  $8 \leq A < B < C < D \leq 2^n$  be real numbers, assume that (3.3), (3.4), and (3.5) hold, and also*

$$42 \ln B \leq A, \quad AB + C \leq D. \quad (3.9)$$

Then the following hold.

- (i) *We have  $\delta_2 \leq (14 \ln B)/A \leq 1/3$  and Algorithm 3.1 returns an element of  $X$ , and for any pair  $(\ell_1, \ell_2) \in X$ , the probability  $b_4(\ell_1, \ell_2)$  with which it is returned satisfies*

$$\left| b_4(\ell_1, \ell_2) - \frac{1}{\#X} \right| \leq \frac{3\delta_2}{\#X}.$$

- (ii) *The algorithm performs an expected number in  $O(a_1 n^2)$  of repetitions of steps (2a) and (2b), and has expected runtime polynomial in  $a_1 n$ .*

*Proof.* We begin with numerical computations that verify some assumptions in previous results. For starters, the condition (2.3) follows from

$$\ln(B/A) \geq \ln 4 > 8^{-1} + (1 + 8^{-2})^{1/2} \geq A^{-1} + (1 + A^{-2})^{1/2}.$$

Next, we find an upper bound on  $\delta_1$ :

$$\begin{aligned} \frac{a^* + 2}{a^* A} &= \frac{1 + 2 \ln(B/A)}{A} < \frac{2 \ln B}{A}, \\ \frac{1}{2^{m-2}} &< \frac{1}{A} < \frac{\ln B}{2A}, \\ \frac{3B}{a^* M} &< \frac{\ln B}{2A}, \\ \delta_1 &< \frac{3 \ln B}{A}. \end{aligned}$$

For the upper bound on  $\delta_2$ , we have from (2.4) that

$$a \leq a^* + 2A^{-1} = (\ln(B/A))^{-1} + 2A^{-1} \leq 1/\ln 4 + 1/4 < 1,$$

so that  $(a + 2)B < 3B \leq D - C$  and  $(D - C + aB)/(D - C - B) \leq 2$ . Furthermore,

$$AB/(2 \ln B) + B \leq AB/4 + B < AB \leq D - C$$

and

$$(a + 1)B < 2B \leq (2 \ln B)(D - C - B)/A,$$

which, together with (3.9), implies the claimed bounds on  $\delta_2$ .

We also have

$$\frac{M}{a \ln(B/A)} < \frac{7AB}{a \ln(B/A)} < \frac{7B}{\ln 4} < 6B,$$

so that (2.8) is satisfied.

For the analysis of the algorithm we have, by the equation (2.12) and with  $a$  as in (2.1), that

$$\#K \leq \frac{D - C + aB}{a}.$$

Let  $(k_1, k_2) \in K$ . Since the inequality (2.3) holds, Theorem 2.4 says that the probability  $b_3(k_1, k_2)$  that  $(k_1, k_2)$  is returned by Algorithm 2.3 satisfies

$$\left| b_3(k_1, k_2) - \frac{1}{\#K} \right| \leq \frac{\delta_2}{\#K}.$$

Thus some element of  $X$  is returned in one execution of steps (2a) and (2b) with probability at least

$$\frac{(1 - \delta_2)\#X}{\#K} \geq \frac{(1 - \delta_2)aD}{5a_1(D - C + aB) \ln D} \geq \frac{(1 - \delta_2)a}{10a_1 \ln D}.$$

Since  $a^{-1} \in O(\ln B)$ , the expected number of executions until success is at most

$$\frac{10a_1 \ln D}{(1 - \delta_2)a} \in O(a_1 \ln D \cdot \ln B) \subseteq O(a_1 n^2).$$

Now let  $r = \#X/\#K$  and denote as  $t$  the random variable counting the number of repetitions until success. For any  $k = (k_1, k_2) \in K$ , the probability for output  $k$  satisfies

$$(1 - \delta_2)r \leq \text{prob}(k \in X) \leq (1 + \delta_2)r,$$

$$1 - (1 + \delta_2)r \leq \text{prob}(k \notin X) \leq 1 - (1 - \delta_2)r.$$



Since the random choices in different repetitions are independent, we have for any  $i \geq 0$

$$\text{prob}(t > i) = (\text{prob}(k \notin X))^i \geq (1 - (1 + \delta_2)r)^i.$$

Thus for any  $(\ell_1, \ell_2) \in X$ , we have

$$\text{prob}(t = i + 1 \text{ and } (\ell_1, \ell_2) \text{ is found}) \geq (1 - (1 + \delta_2)r)^i \frac{1 - \delta_2}{\#K},$$

and

$$b_4(\ell_1, \ell_2) \geq \frac{1 - \delta_2}{\#K} \sum_{i \geq 0} (1 - (1 + \delta_2)r)^i = \frac{1 - \delta_2}{(1 + \delta_2)\#X} \geq \frac{1 - 2\delta_2}{\#X}.$$

A similar calculation shows that

$$b_4(\ell_1, \ell_2) \leq \frac{1 + \delta_2}{(1 - \delta_2)\#X} \leq \frac{1 + 3\delta_2}{\#X}.$$

Primality of an integer can be tested in random polynomial time, see [6, § 3.4], and the claimed cost bound follows.  $\square$

## 4 Safe primes

We use the following notation, for any  $\alpha$  with  $0 \leq \alpha < 1/2$ .

$P$  = set of primes,

$\text{MG}_1 = \{p \in P : (p - 1)/2 \text{ prime}\},$

$\text{MG}_{2,\alpha} = \{p \in P : (p - 1)/2 = \ell_1 \ell_2 \text{ with } \ell_1 < \ell_2 \text{ primes and } (\ell_1 \ell_2)^\alpha \leq \ell_1\}.$

Furthermore,  $P(x)$ ,  $\text{MG}_1(x)$ , and  $\text{MG}_{2,\alpha}(x)$  are the corresponding subsets of those  $p$  with  $p \leq x$ , and  $\pi(x)$ ,  $\pi_1(x)$ , and  $\pi_{2,\alpha}(x)$  denote the respective cardinalities. An integer  $\ell$  is a Sophie Germain prime if and only if  $2\ell + 1 \in \text{MG}_1$  is a Marie Germain prime, as defined in the Introduction, and  $\text{MG}_2 = \text{MG}_{2,0}$ . The last condition in our definition of  $\text{MG}_{2,\alpha}$  is equivalent to  $\ell_2 \leq \ell_1^{1/\alpha-1}$  when  $\alpha \neq 0$ .

**Algorithm 4.1** (Generating a safe prime).

INPUT: Positive bounds  $x$  and  $\alpha < 1/2$ .

OUTPUT: A prime  $p$  with  $x/\ln^2 x < p \leq x$ .

(1) Compute  $y_0 = (x - \ln^2 x)/2 \ln^2 x$  and  $y_1 = (x - 1)/2$ .

- (2) Repeat steps (3)–(6) until a prime is returned.
- (3) Choose a random prime  $\ell \in [y_0 \dots y_1]$ .
- (4) If  $2\ell + 1$  is prime then return  $p = 2\ell + 1$ .
- (5) Using Algorithm 3.1 with inputs  $(A, B, C, D) = (x^\alpha, x^{1/2}, y_0, y_1)$ , choose an approximately uniformly random sample from the pairs  $(\ell_1, \ell_2)$  of primes with  $\ell_1 \in [x^\alpha \dots x^{1/2}]$  and  $\ell_1\ell_2 \in [y_0 \dots y_1]$ .
- (6) If  $2\ell_1\ell_2 + 1$  is prime then return  $p = 2\ell_1\ell_2 + 1$ .

The interesting question is how many repetitions we expect to perform.

- Theorem 4.2.** (i) *Any  $p$  returned by the algorithm satisfies the output specification, and  $p \equiv 3 \pmod{4}$ . For any output  $p$ ,  $(p - 1)/2$  is squarefree with at most two prime divisors, and each of them is at least  $x^\alpha$ .*
- (ii) *Let  $0.25 \leq \alpha < 0.276$  and let  $n$  be sufficiently large. For an input  $x \in [2^{n-1} \dots 2^n]_{\mathbb{R}}$ , the expected number of repetitions made by Algorithm 4.1 until an output is returned is  $O(n)$ , and the expected runtime of the algorithm is polynomial in  $n$ .*

*Proof.* An  $\ell$  leading to an output in step (4) satisfies  $\ell \geq y_0 > x^\alpha$ . Suppose that  $(\ell_1, \ell_2)$  is chosen in step (5) and that  $p = 2\ell_1\ell_2 + 1$  is returned in step (6). Then  $x^\alpha \leq \ell_1 \leq x^{1/2}$  and  $\ell_2 \geq y_0/\ell_1 \geq (x - \ln^2 x)/2x^{1/2} \ln^2 x > x^\alpha$ . Also  $\ell_1 \neq \ell_2$ , since 3 divides  $2\ell^2 + 1$  for every prime  $\ell \neq 3$ . Thus any output has the stated properties. The primality tests can be performed in polynomial time. The assumptions of Theorem 3.3 are easily checked, and thus also one execution of step (5) uses polynomial time. It is sufficient to show the claim about the expected number of repetitions.

Our main tool is a result of Heath-Brown [16, Lemma 1] taken with  $k = 1$ ,  $K = 2$ , any  $u \in \{3, 7, 11, 15\}$ , and  $v = 16$ , which implies that

$$\pi_1(x) + \pi_{2,0.276}(x) \geq \frac{cx}{\ln^2 x} \quad (4.1)$$

for some constant  $c > 0$ . No explicit value for  $c$  is known.

By (4.1), at least one (possibly both) of the alternatives

$$\pi_1(x) \geq \frac{cx}{2 \ln^2 x} \quad (4.2)$$

and

$$\pi_{2,0.276}(x) \geq \frac{cx}{2 \ln^2 x} \quad (4.3)$$

holds. In this and the next section, we make no attempt to optimize constants, and even use the trivial form

$$\frac{x}{2 \ln x} \leq \pi(x) = \#P(x) \leq \frac{2x}{\ln x}$$

of the Prime Number Theorem.

We consider the “shifted multiplication” map  $\mu$  with  $\mu((\ell_1, \ell_2)) = 2\ell_1\ell_2 + 1$ , and

$$\begin{aligned} X_1 &= \left\{ 2\ell + 1 : \frac{x}{\ln^2 x} < 2\ell + 1 \leq x, \ell \text{ prime} \right\}, \\ Y_1 &= X_1 \cap P, \\ X_2 &= \left\{ (\ell_1, \ell_2) : \ell_1 \in [x^\alpha \dots x^{1/2}] \text{ and } \ell_2 \in \left[ \frac{y_0}{\ell_1} \dots \frac{y_1}{\ell_1} \right] \text{ primes} \right\}, \\ Y_2 &= \mu(X_2) \cap P. \end{aligned} \tag{4.4}$$

By the Prime Number Theorem, we have  $\#X_1 \leq \pi(x/2) < 2x/\ln x$ . Since  $2y_0 + 1 = x/\ln^2 x$ , we have  $\text{MG}_1(x) \subseteq Y_1 \cup P(x/\ln^2 x)$ , the integers  $2\ell + 1$  tested in step (4) are uniformly distributed in  $X_1$ , and  $2\ell + 1$  is returned if it is in  $Y_1$ . Since  $n$  is sufficiently large, we may assume  $n > 1 + 18/c$  and  $x > e^{12/c}$ , so that for

$$a_2 = \frac{1}{c - 6/\ln x}$$

we have  $0 < a_2 \leq 2/c$ . If (4.2) holds, then with

$$\begin{aligned} \#Y_1 &\geq \pi_1(x) - \pi\left(\frac{x}{\ln^2 x}\right) \geq \frac{cx}{2 \ln^2 x} - \frac{2x}{\ln^2 x \cdot \ln(x/\ln^2 x)} \\ &\geq \frac{cx}{2 \ln^2 x} - \frac{3x}{\ln^3 x} = \frac{x}{2a_2 \ln^2 x} \geq \frac{cx}{4 \ln^2 x}. \end{aligned} \tag{4.5}$$

Thus the expected number of repetitions of steps (3) and (4) is at most

$$\frac{\#X_1}{\#Y_1} < \frac{2x/\ln x}{cx/4 \ln^2 x} = \frac{8 \ln x}{c} \in O(n).$$

The claim follows in this case.

Now we assume that (4.2) does not hold, so that (4.3) does. We have

$$1.4 < a_1 = \frac{1}{-\ln 2\alpha} < 1.7.$$

Now  $X_2$  is the set  $X$  defined in (3.1) for our input parameters  $(A, B, C, D) = (x^\alpha, x^{1/2}, y_0, y_1)$ . By Lemma 3.2, we have

$$\#X_2 \leq \frac{2x}{a_1 \ln((x-1)/2x^{1/2})} \leq \frac{2x}{a_1(0.5 \ln x - 1)} \leq \frac{4x}{\ln x}.$$

An asymptotic formula for  $\#X_2$  is provided in Theorem 6.2, but at this stage, the above upper bound suffices.

Now let  $\beta = 0.276$  and  $p = 2\ell_1\ell_2 + 1 \in \text{MG}_{2,\beta}(x)$  with  $\ell_1 < \ell_2$ . Then either  $\ell_1\ell_2 < y_0$  or  $(y_0 \leq \ell_1\ell_2 \leq y_1$  and  $\ell_1 < \ell_2 \leq \ell_1^{1/\beta-1})$ . In the latter case, we have  $\ell_1^2 < \ell_1\ell_2 \leq y_1$  and  $\ell_1 \leq y_1^{1/2} < x^{1/2}$ . Furthermore,  $\alpha < \beta$  and  $x^\alpha < y_0^\beta \leq (\ell_1\ell_2)^\beta \leq \ell_1$ , so that  $(\ell_1, \ell_2) \in X_2$ . It follows that  $p \in Y_2$ , and hence  $\text{MG}_{2,\beta}(x) \subseteq Y_2 \cup P(2y_0 + 1)$  and, similar to (4.5),

$$\#Y_2 \geq \pi_{2,\beta}(x) - \pi(2y_0 + 1) \geq \frac{cx}{2 \ln^2 x} - \pi\left(\frac{x}{\ln^2 x}\right) \geq \frac{cx}{4 \ln^2 x}.$$

Since  $Y_2 \subseteq \mu(X_2)$ , we have  $\mu(\mu^{-1}(Y_2)) = Y_2$  and  $\#\mu^{-1}(Y_2) \geq \#Y_2$ . Success occurs in step (6) if  $(\ell_1, \ell_2) \in \mu^{-1}(Y_2)$ , and any  $(\ell_1, \ell_2)$  is chosen in step (5) with probability at least  $(1 - 3\delta_2)/\#X_2$  by Theorem 3.3, whose assumptions are readily checked. Then the probability of success is at least

$$\frac{1 - 3\delta_2}{\#X_2} \cdot \#\mu^{-1}(Y_2) \geq \frac{(1 - 3\delta_2)\#Y_2}{\#X_2}.$$

It follows that the expected number of repetitions of steps (5) and (6) until success is at most

$$\frac{\#X_2}{(1 - 3\delta_2)\#Y_2} \leq \frac{4x/\ln x}{(1 - 3\delta_2) \cdot cx/4 \ln^2 x} \in O(n),$$

which concludes the proof.  $\square$

More detailed calculations show that any  $n \geq \max\{18/c, 155\}$  is sufficiently large for the conclusions to hold, for  $\alpha = 1/4$ .

On the other hand, for an asymptotic result we can replace the lower bound in the output specification by  $p \geq a(y) \cdot y/\ln y$  for any function  $a \in o(1)$ .

An output  $p$  from step (4) is uniformly random in the set  $Y_1$  from (4.4), and by Theorem 3.3 (i), an output from step (6) is approximately random in  $Y_2$ . After Conjecture 6.3 below, we address the question of how to make  $p$  in either case uniformly random in  $Y_1 \cup Y_2$ .

## 5 Safe moduli

We generate a safe modulus  $N = pq$  from two executions of Algorithm 4.1.

**Algorithm 5.1** (Generating a safe modulus).

INPUT: Positive bounds  $y$  and  $\alpha$  with  $1/4 \leq \alpha < 0.276$ .

OUTPUT: A modulus  $N$  with  $y/\ln^2 y \leq N \leq y$ .

- (1) Repeat steps (2)–(4) until a modulus is returned.
- (2) Call Algorithm 4.1 with inputs  $x = y^{1/2}$  and  $\alpha$ , and output  $p$ .
- (3) Call Algorithm 4.1 with inputs  $x = y/p$  and  $\alpha$ , and output  $q$ .
- (4) If  $\gcd((p-1)/2, (q-1)/2) = 1$ , then return  $N = pq$ .

A *Blum integer* is a product of two primes, both congruent to 3 modulo 4. These have been introduced by Blum, Blum and Shub [4].

**Theorem 5.2.** (i) *Any output  $N$  satisfies the output specification and is a Blum integer. Furthermore,  $\varphi(N)/4$  is squarefree with at most four prime factors, and each of these is at least  $y^{\alpha/2}$ .*

- (ii) *For  $n$  sufficiently large and  $y \in [2^{n-1} \dots 2^n]_{\mathbb{R}}$ , the expected number of repetitions in Algorithm 5.1 until an output  $N$  is returned is at most  $1 + y^{-1/8} \ln^2 y < 2$ , and the expected runtime of the algorithm is polynomial in  $n$ .*

*Proof.* (i) The claims follow from Theorem 4.2, using the fact that  $y/p \geq y^{1/2}$ .

(ii) We write  $\ell_0$  or  $(\ell_1, \ell_2)$  for the choice that leads to an output  $p$  in step (2), depending on whether step (4) or step (6) of the call to Algorithm 4.1 is successful, and similarly  $r_0$  or  $(r_1, r_2)$  for step (3). We denote as  $R_0$ ,  $R_1$ , and  $R_2$  the sets of possible values for  $r_0$ ,  $r_1$ , and  $r_2$ , respectively. Since  $\ell_1 \neq \ell_2$  and  $r_1 \neq r_2$ , the gcd condition in step (4) is violated only if one of  $r_0, r_1, r_2$  equals one of  $\ell_0, \ell_1, \ell_2$ . We write  $b_5$  for the probability of this event, conditioned on some output  $p$  of step (2), and distinguish two cases.

In the first case, the second call to Algorithm 4.1 returns from its step (4). Then  $r_0$  has to avoid at most two values, so that  $b_5 \leq 2/\#R_0$ , since  $r_0$  is chosen uniformly in  $R_0$ .

In the second case, Algorithm 4.1 returns from its step (6). Then Theorem 3.3 implies that  $(r_1, r_2)$  assumes any specific value with probability at most  $(1 + 3\delta_2)/\#X < 2/\#X$  with  $X$  from (3.1). There are at most two values that both

$r_1$  and  $r_2$  have to avoid. This makes for a total of at most  $2\#R_1 + 2\#R_2$  pairs to be avoided by  $\#R_1$  choices and

$$b_5 \leq \frac{4\#R_1 + 4\#R_2}{\#X}. \quad (5.1)$$

We now prove upper bounds on the three  $\#R_i$ . By Theorem 4.2, the output  $p$  of step (2) satisfies

$$\frac{4y^{1/2}}{\ln^2 y} = \frac{y^{1/2}}{\ln^2 y^{1/2}} \leq p \leq y^{1/2}.$$

In step (2) of Algorithm 4.1, when called in step (3) of Algorithm 5.1, we have  $y_0 \leq y_1/2$  and

$$\begin{aligned} y_1 &= \frac{y/p - 1}{2}, \\ \#R_0 &= \pi(y_1) - \pi(y_0) \geq \pi(y_1) - \pi(y_1/2) \\ &\geq \frac{3 \cdot y_1/2}{5 \cdot \ln(y_1/2)} \geq \frac{y^{1/2}}{10 \ln(y^{1/2}/6)} > \frac{y^{1/2}}{5 \ln y}, \\ b_5 &< \frac{10 \ln y}{y^{1/2}} < \frac{\ln y}{2y^{1/8}}. \end{aligned}$$

For the second case, we have in step (4) of the same call to Algorithm 4.1

$$\begin{aligned} (y/p)^\alpha &\leq r_1 \leq (y/p)^{1/2}, \\ \#R_1 &= \pi((y/p)^{1/2}) - \pi((y/p)^{\alpha/2}) < \pi\left(\frac{y^{1/4} \ln y}{2}\right) < 4y^{1/4}, \\ \frac{y/p - 1}{2r_1} &< \frac{y}{2pr_1} \leq \frac{y^{1-\alpha}}{2p^{1-\alpha}} \leq \frac{y^{(1-\alpha)/2} (\ln y)^{2-2\alpha}}{5} \leq \frac{y^{3/8} \ln^{3/2} y}{5}. \end{aligned}$$

Finally

$$\begin{aligned} \#R_2 &= \pi\left(\frac{y/p - 1}{2r_1}\right) - \pi\left(\frac{y/p - \ln^2(y/p)}{2r_1 \ln^2(y/p)}\right) \\ &< \pi\left(\frac{y^{3/8} \ln^{3/2} y}{5}\right) \leq 2y^{3/8} \ln^{1/2} y. \end{aligned}$$

From (3.7) with the parameters  $B = (y/p)^{1/2} = A^{1/2\alpha}$  and  $D = y/p \geq y^{1/2}$  we obtain

$$\#X \geq \frac{y^{1/2} \cdot (-\ln(2\alpha))}{5 \cdot \ln y^{1/2}} > \frac{y^{1/2}}{2 \ln y}.$$

It follows that in the second case we have in (5.1)

$$b_5 \leq \frac{4 \cdot (4y^{1/4} + 2y^{3/8} \ln^{1/2} y) \cdot 2 \ln y}{y^{1/2}} \leq \frac{\ln^2 y}{2y^{1/8}}.$$

The gcd condition holds in any case with probability at least  $1 - b_5$ , and the number of iterations is at most  $(1 - b_5)^{-1} \leq 1 + 2b_5 \leq 1 + y^{-1/8} \ln^2 y$ .  $\square$

Section 7 presents more details on the runtime of Algorithms 4.1 and 5.1.

**Corollary 5.3.** *Let  $N$  be an output of Algorithm 5.1. For uniformly random  $a \in \mathbb{Z}_N^\times$ ,  $a^2$  generates the group of squares in  $\mathbb{Z}_N^\times$  with probability at least  $1 - 4N^{-1/8}$ .*

*Proof.* We have

$$\begin{aligned} \mathbb{Z}_p^\times &\cong \mathbb{Z}_2 \times \mathbb{Z}_{(p-1)/2}, \\ \mathbb{Z}_N^\times &\cong \mathbb{Z}_2^2 \times \mathbb{Z}_{(p-1)/2} \times \mathbb{Z}_{(q-1)/2} \cong \mathbb{Z}_2^2 \times \mathbb{Z}_{\varphi(N)/4}, \end{aligned}$$

by the condition in step (4).

The set  $\square_N$  of squares in  $\mathbb{Z}_N^\times$  is isomorphic to  $\mathbb{Z}_{\varphi(N)/4}$  and hence cyclic with  $\varphi(\varphi(N)/4)$  generators. The squaring map from  $\mathbb{Z}_N^\times$  to  $\square_N$  is 4-to-1, and for a uniformly random  $a \in \mathbb{Z}_N^\times$  we have

$$\text{prob}\{a \in \mathbb{Z}_N^\times : a^2 \text{ generates } \square_N\} = \frac{\varphi(\varphi(N)/4)}{\varphi(N)/4}. \quad (5.2)$$

We first consider the case where both  $p$  and  $q$  come from step (6) of the respective call to Algorithm 4.1. We can then write  $\varphi(N)/4 = \ell_1 \ell_2 r_1 r_2$ , with four distinct primes  $\ell_1, \ell_2, r_1, r_2$ , by Theorem 4.2 and the condition in step (4) of Algorithm 5.1. These primes are all at least  $y^{\alpha/2}$ , and

$$\begin{aligned} \frac{\varphi(\varphi(N)/4)}{\varphi(N)/4} &= \left(1 - \frac{1}{\ell_1}\right) \left(1 - \frac{1}{\ell_2}\right) \left(1 - \frac{1}{r_1}\right) \left(1 - \frac{1}{r_2}\right) \\ &\geq 1 - \left(\frac{1}{\ell_1} + \frac{1}{\ell_2} + \frac{1}{r_1} + \frac{1}{r_2}\right) \\ &\geq 1 - \frac{4}{y^{\alpha/2}} \geq 1 - \frac{4}{N^{\alpha/2}} \geq 1 - \frac{4}{N^{1/8}}. \end{aligned}$$

The last estimate also holds for the other possibilities for the factors of  $p - 1$  and  $q - 1$ . Together with (5.2), this concludes the proof.  $\square$

As  $y$  grows, the lower bound of Corollary 5.3 comes exponentially close to 1.

**Corollary 5.4.** *Let  $n$  be a sufficiently large integer. Then we can generate in expected time polynomial in  $n$  an RSA modulus  $N = pq$  so that*

- $2^n/n^2 \leq N \leq 2^n$ ,
- $\varphi(N)/4 = (p-1)(q-1)/4$  is squarefree with at most four prime factors,
- each such prime factor is at least  $2^{n/8}$ ,
- for uniformly random  $a \in \mathbb{Z}_N^\times$ ,  $a^2$  generates the group of squares in  $\mathbb{Z}_N^\times$  with probability at least  $1 - 4 \cdot 2^{-n/8}$ .

It follows that the moduli presented here can be used in the encryption scheme of Hofheinz, Kiltz and Shoup [18]. If  $g$  generates  $\square_N$ , then  $|g|$  generates the group  $\mathbb{QR}_N^+$  of signed quadratic residues, in their notation. The factoring assumption then refers to the moduli generated by Algorithm 5.1.

For various notions of RSA integers, their number is estimated in [8, 24, 34].

## 6 Heuristic estimates and extended range

It is interesting to compare the number of Marie Germain primes to that of the Marie Germain<sub>2</sub> primes generated in steps (5) and (6) of Algorithm 4.1. For Sophie Germain primes, we take a prime  $\ell \leq x$ , of which there are about  $x/\ln x$  many. If  $2\ell + 1$  is also prime, then  $\ell$  is a Sophie Germain prime. Since the density of primes up to  $2x$  is about  $1/\ln(2x) \sim 1/\ln x$ , one might naively expect there to be about  $x/\ln^2 x$  of them. Although this gives the right order of magnitude, the asymptotics is false as it ignores so-called “local” (or divisibility) conditions. The argument of Bateman and Horn [3] in this special case suggests the following.

We take a prime  $q \neq 2, \ell$ . Then  $2\ell + 1 \not\equiv 1 \pmod q$ , while general primes are allowed to be  $1 \pmod q$ . This consideration also applies to twin primes and leads to the standard heuristics on the number of Sophie Germain primes, namely that there are about  $2C_2x/\ln^2 x$  of them up to  $x$ , where

$$C_2 = \prod_{p \geq 3} \left(1 - \frac{1}{(p-1)^2}\right) \approx 0.66016$$

is the *twin prime constant*. In our situation,  $x$  denotes an upper bound on  $2\ell + 1$ , so that we consider the Sophie Germain primes  $\ell \leq x/2$ . We thus rephrase the heuristics as follows for our purposes.

We use  $\sim$  to denote the asymptotic equality of two quantities, so that  $f \sim g$  means that  $|f(x) - g(x)| \in o(g(x))$  as  $x \rightarrow \infty$ . This relation is symmetric in  $f$  and  $g$ .



**Conjecture 6.1** (Sophie Germain prime conjecture). We have

$$\pi_1(x) \sim \frac{2C_2}{\ln x} \cdot \pi(x/2) \sim \frac{C_2 x}{\ln^2 x}.$$

The local behavior of  $2\ell_1\ell_2 + 1$ , as  $(\ell_1, \ell_2)$  ranges over the set  $X_2$  as defined in (4.4), is the same as that of  $2\ell + 1$  in this range. The following conjecture is the natural analog of the previous one:

$$\pi_{2,\alpha}(x) \sim \frac{2C_2}{\ln x} \cdot \#\mu(X_2), \tag{6.1}$$

where the map  $\mu$  is as in the proof of Theorem 4.2. It now remains to estimate the magnitude of  $\#\mu(X_2)$ . This can be done unconditionally.

**Theorem 6.2.** *We have*

$$\#\mu(X_2) \sim \#X_2 \sim \frac{\ln(\alpha^{-1} - 1)}{2} \cdot \frac{x}{\ln x}.$$

*Proof.* We start with  $\#X_2$  and set

$$r_1 = \sum_{x^\alpha \leq \ell \leq x^{1/2}} \pi\left(\frac{x}{2\ell}\right), \quad r_2 = \sum_{x^\alpha \leq \ell \leq x^{1/2}} \pi\left(\frac{x}{2\ell \ln^2 x}\right).$$

Then the asymptotic version of (3.8), ignoring the term  $D^{1/2}$  of small order, says that

$$\#X_2 \sim r_1 - r_2.$$

For  $r_1$ , the Prime Number Theorem and partial summation imply

$$\begin{aligned} r_1 &\sim \frac{x}{2} \int_{x^\alpha}^{x^{1/2}} \frac{1}{\lambda \ln(x/2\lambda) \ln \lambda} d\lambda \sim \frac{x}{2} \int_{x^\alpha}^{x^{1/2}} \frac{1}{\lambda \ln(x/\lambda) \ln \lambda} d\lambda \\ &= \frac{x}{2} \int_{x^\alpha}^{x^{1/2}} \frac{1}{(\ln x - \ln \lambda) \ln \lambda} d \ln \lambda = \frac{x}{2} \int_{\alpha u}^{u/2} \frac{1}{(u-v)v} dv, \end{aligned}$$

where  $u = \ln x$ . Therefore,

$$r_1 \sim \frac{x}{2u} \int_{\alpha u}^{u/2} \left(\frac{1}{u-v} + \frac{1}{v}\right) dv = \frac{x}{2u} (\ln(1-\alpha) - \ln \alpha) = \frac{\ln(\alpha^{-1} - 1)}{2} \cdot \frac{x}{\ln x}.$$

For  $r_2$ , we recall from (3.2) and (3.6) that

$$\sum_{x^\alpha \leq \ell \leq x^{1/2}} \frac{1}{\ell} \sim \ln \frac{1}{2\alpha}.$$

Thus

$$\begin{aligned} r_2 &\leq \sum_{x^\alpha \leq \ell \leq x^{1/2}} \frac{x}{2\ell \ln^2 x \cdot \ln(x/2\ell \ln^2 x)} \\ &\leq \frac{x}{2\ln^2 x} \sum_{x^\alpha \leq \ell \leq x^{1/2}} \frac{1}{\ell \ln(x^{1/2}/2\ell \ln^2 x)} \leq \frac{x \ln(1/2\alpha)}{\ln^3 x} \in o(r_1). \end{aligned}$$

It follows that

$$\#X_2 \sim r_1. \tag{6.2}$$

Now it remains to determine the size of  $\mu(X_2)$ . The map  $\mu : X_2 \rightarrow \mu(X_2)$  is sometimes 1-to-1 and sometimes 2-to-1. The latter happens if and only if  $(\ell_1, \ell_2), (\ell_2, \ell_1) \in X_2$ . For every such pair  $(\ell_1, \ell_2)$ , we have  $\ell_1, \ell_2 \leq x^{1/2}$ , and thus there are at most  $\pi(x^{1/2})^2 \in O(x/\ln^2 x) \subseteq o(r_1)$  of them. Now (6.2) implies that  $\#\mu(X_2) \sim \#X_2 \sim r_1$ .  $\square$

We have  $\ln((1/4)^{-1} - 1) = \ln 3 \approx 1.098$ . The methods of Loebenger and Nüsken [24] provide upper and lower bounds on  $\#\mu(X_2)$  of the form  $\text{const} \cdot x / \ln x$ . From the heuristic argument (6.1), we derive the following.

**Conjecture 6.3** (MG<sub>2</sub> prime conjecture). For  $1/4 \leq \alpha < 1/2$ , we have

$$\pi_{2,\alpha}(x) \sim \frac{C_2 \ln(\alpha^{-1} - 1) x}{\ln^2 x}.$$

Under the two conjectures, there are roughly 10% more MG<sub>2</sub>- than MG<sub>1</sub>-primes for  $\alpha = 1/4$ .

We now extend the range of applicability of our method. Lemma 3.2 and thus Theorem 4.2 depend on a version of Mertens' theorem over certain intervals, that is, on good estimates of the sum

$$M(A, B) = \sum_{\ell \in (A \dots B]} \frac{1}{\ell}.$$

In our present application,  $B$  is substantially larger than  $A$ , so that classical results allow us to handle this sum. However, for future extensions and possible ramifications of our ideas and results, it might be useful to study this sum in a range of  $A$  and  $B$  which is as wide as possible. For two quantities  $x$  and  $y$ , we write  $x \asymp y$  if for some positive constants  $c_1 \leq c_2$  we have  $c_1 y \leq x \leq c_2 y$ .

**Theorem 6.4.** For real  $A$  and  $B$  with  $B \geq A + A^{0.525} \geq 3$  we have

$$M(A, B) \asymp \ln \frac{\ln B}{\ln A}.$$

*Proof.* We assume that  $A$  is large enough and let  $\Delta = B - A$ . A result of Baker, Harman and Pintz [2] on primes in short interval (see also [15, Theorem 7.2]) implies that

$$\pi(A + \Delta) - \pi(A) \asymp \Delta / \ln A \tag{6.3}$$

for  $\Delta \geq A^{0.525}$ . Only the lower bound on the left-hand side in (6.3) requires this restriction; the upper bound holds for any  $\Delta$  by the celebrated Brun–Titchmarsh theorem, see [21, Theorem 6.6].

As a first case, we consider  $A + A^{0.525} \leq B < A + A / \ln A$ . Then

$$\frac{1}{B}(\pi(B) - \pi(A)) \leq M(A, B) \leq \frac{1}{A}(\pi(B) - \pi(A)),$$

and (6.3) implies

$$M(A, B) \asymp \frac{\Delta}{A \ln B}. \tag{6.4}$$

Since  $\Delta = o(A)$  and  $\ln(1 + z) \sim z$  as  $z \rightarrow 0$ , we also have

$$\frac{\ln B}{\ln A} = 1 + \frac{\ln(B/A)}{\ln A} = 1 + \frac{\ln(1 + \Delta/A)}{\ln A} \sim 1 + \frac{\Delta}{A \ln A}$$

and

$$\ln \frac{\ln A}{\ln B} \sim \frac{\Delta}{A \ln A} \sim \frac{\Delta}{A \ln B}.$$

Together with (6.4), the result follows in this case.

We now come to the case where  $A + A / \ln A \leq B$ . Then

$$\frac{\ln B}{\ln A} \geq \frac{\ln A + \ln(1 + 1/\ln A)}{\ln A} \sim \frac{\ln A + 1/\ln A}{\ln A} = 1 + \frac{1}{\ln^2 A}$$

and  $1/\ln^2 A \in O(\ln(\ln B / \ln A))$ . A slight modification of the argument of Vinogradov [33], coupled with [32, Theorem 8, § I.1], implies that

$$\sum_{p \leq x} \frac{1}{p} \in \ln \ln x + \gamma + O(\exp(-c_0(\ln x)^{3/5}))$$

for some absolute constant  $c_0 > 0$ . Furthermore, we have

$$\exp(-c_0(\ln A)^{3/5}) \in o\left(\frac{1}{\ln^2 A}\right) \subseteq o\left(\ln \frac{\ln B}{\ln A}\right).$$

Thus we find

$$M(A, B) \in \ln \frac{\ln B}{\ln A} + O(\exp(-c_0(\ln A)^{3/5})).$$

Hence

$$M(A, B) \sim \ln \frac{\ln B}{\ln A},$$

which concludes the proof also in this case. □

Until the result of Baker, Harman and Pintz [2] is improved, there is clearly no chance to extend the above range of  $A$  and  $B$ .

A famous result of Huxley [19], see also [17], shows that for  $\Delta \geq A^{7/12}$ , we can replace  $\asymp$  in (6.3) by  $\sim$ . That is, in the  $\asymp$  notation, both “constants”  $c_1$  and  $c_2$  can be chosen in  $1 + o(1)$ . Arguing as above, we find that in this range, Theorem 6.4 also holds with  $\sim$ . We have  $7/12 = 0.58333\dots$

## 7 Cost estimates

We now analyze the cost of Algorithm 4.1 more closely. We denote by  $M(n)$  a number of bit operations with which arithmetic (addition, multiplication, division with remainder) can be performed on  $n$ -bit integers. Thus we have  $M(n) \in O(n^2)$  with classical and  $M(n) \in O^\sim(n)$  with fast arithmetic, where the  $O^\sim$  notation hides logarithmic factors. We also denote as  $T(n)$  the cost of testing an  $n$ -bit integer for primality. There are several choices.

- Deterministic primality test:  $T(n) \in O^\sim(n^6)$  (see [6, § 4.5]).
- Probabilistic primality test:  $T(n) \in O^\sim(n^4)$  (see [6, § 4.5]).
- Probabilistic compositeness test:  $T(n) \in O(tnM(n))$  for some  $t$ , the number of iterations of a single test, see [12, Theorem 18.6] or [6, Algorithm 3.4.7].

For practical purposes, one will use the last type of test in the algorithm. After Algorithm 5.1 has produced an output, one can test the pseudoprimes involved ( $p$ ,  $q$ , and the factors of  $p - 1$  and  $q - 1$ ) by a probabilistic primality test. This cost is within the time bound of the algorithm.

We now examine separately  $MG_1$  prime generation, that is, Algorithm 4.1 with steps (5) and (6) removed, and  $MG_2$  prime generation, that is, Algorithm 4.1 with steps (3) and (4) removed.

In  $MG_1$  prime generation, we choose a uniformly random  $k \leq x$ , test it for primality, and on success, test  $2k + 1$  for primality. The probability finding a prime  $p = 2\ell + 1$  with  $\ell$  prime is approximately  $\pi_1(x)/x$  and the expected cost of producing such a  $p$  is  $T(n) \cdot x/\pi_1(x)$ .

In  $MG_2$  prime generation, we choose uniformly random  $(k_1, k_2) \in K$ , test both for primality, and on success, test  $2k_1k_2 + 1$  for primality. The probability finding a prime  $p$  is approximately  $\pi_2(x)/\#K$  and the expected cost of producing such a  $p$  is  $T(n) \cdot \#K/\pi_2(x)$ . For  $a$  from (2.1), we have  $a \approx 1/\ln(x^{1/2}/x^\alpha) = 2/(1-2\alpha) \ln x$ . Using (2.12), we find  $\#K \approx x/a \approx ((1-2\alpha)x \ln x)/2$ .

We know from (4.1) that at least one of  $\pi_1(x)$  and  $\pi_1(x)$  is bounded from below by  $cx/2 \ln^2 x$ . Algorithm 4.1 runs both  $MG_1$  and  $MG_2$  generation in tandem and

bit-length	128	256	512	1024	2048
MG <sub>1</sub>	0.2 s	1.1 s	12.6 s	298.1 s	5798.8 s
MG <sub>2</sub>	20.8 s	164.7 s	2582.5 s	31892.8 s	558600.2 s
MG <sub>2</sub> (fast)	0.4 s	1.8 s	11.7 s	144.9 s	2147.5 s

Table 1. Time needed for finding a safe prime, depending on the bit-length of  $x$ . (Average over at least 100 findings, except for MG<sub>2</sub> of bit-lengths  $\geq 1024$ ; hardware: single core Intel Xeon, 3.00 GHz.)

halts whenever one of them succeeds. Its expected cost is therefore the minimum of the two costs. We have shown the following.

**Theorem 7.1.** *Algorithms 4.1 and 5.1 for the generation of safe primes and safe moduli with nearly  $n$  bits, respectively, take an expected number of  $O(n^3\mathbb{T}(n))$  bit operations.*

With probabilistic compositeness tests, this comes to  $O(n^4M(n))$  operations, that is,  $O(n^6)$  operations with classical and  $O^\sim(n^5)$  with fast arithmetic.

If Conjectures 6.1 and 6.3 hold, then the cost comes to  $O(n^3M(n))$  for MG<sub>1</sub> prime generation and  $O(n^4M(n))$  for MG<sub>2</sub>. We might run MG<sub>1</sub> prime generation  $n$  times for each execution of MG<sub>2</sub> prime generation. Then we are in the best of two worlds:

- The algorithm provably terminates.
- If the Sophie Germain Conjecture holds, then it does not take much more time than pure MG<sub>1</sub> prime generation.

The advantage is that pure MG<sub>1</sub> prime generation is not proven to terminate.

In a “fast” variant of MG<sub>2</sub> generation, instead of rejecting  $(k_1, k_2)$  if one of them is composite, we first test  $k_1$  and on success keep generating values for  $k_2$  until we find a prime. This works well in practice as shown in Table 1 which is taken from [35].

The primes generated have  $n$  or slightly fewer bits. If exactly  $n$  bits are required, one can reject the smaller ones. Under the conjectures, this also works in polynomial time.

For a simplified version of our method, we recall that Algorithms 2.3 and 3.1 rely on an asymptotically uniform sampling of points under a hyperbola given by Algorithm 2.1. We now exhibit a slower but simpler “dyadic” algorithm, using  $y_0$  and  $y_1$  from Algorithm 4.1.

Let  $n = \lceil \log_2 x \rceil$  and  $j_0 = \lfloor \alpha \log_2 x \rfloor$ . We know that either (4.2) or (4.3) holds. In the latter case, for some  $j \in [j_0 \dots n/2]$  there are at least

$$\frac{\pi_{2,\alpha}(x)}{n/2 - j_0} \geq \frac{cx}{n \ln^2 x}$$

primes  $p = 2\ell_1\ell_2 + 1 \in \text{MG}_{2,\alpha}(x)$  with primes  $\ell_1$  and  $\ell_2$  such that  $\ell_1 \in [2^j, 2^{j+1}]$ . A simple version of Algorithm 4.1 repeats on input  $x$  and  $\alpha$  the following steps until success:

- choose  $\ell \in [y_0 \dots y_1]$  uniformly at random and test  $\ell$  and  $2\ell + 1$  for primality,
- for  $j \in [j_0 \dots n/2 - 1]$ , choose  $\ell_1 \in [2^j \dots 2^{j+1}]$  and  $\ell_2 \in [y_0/\ell_1 \dots y_1/\ell_1]$  uniformly at random and test  $\ell_1, \ell_2$ , and  $2\ell_1\ell_2 + 1$  for primality.

This dyadic method has an expected cost of  $O(n^4 M(n))$  bit operations. If one performs about  $n$  iterations of the first step for each execution of the second one, this bound still holds and the algorithm provably terminates. But if Conjecture 6.1 is true, this version uses about as much time as  $\text{MG}_1$  prime generation. Its advantage is its greater simplicity, when compared to Algorithm 4.1.

## 8 Comments and open questions

From the point of view of algorithmic applications, it would be nice to have a version of the result of Heath-Brown [16, Lemma 1] with an effective (or, even better, an explicitly computed) constant  $c$  for which (4.1) holds. This, however, may be a nontrivial task and may only work for values of  $\beta$  smaller than 0.276.

The proofs in this paper rely on fairly deep results in analytic number theory. But the resulting algorithm is quite simple. In any efficient prime generation method, one will need a (probabilistic) primality test. This is sufficient for Sophie Germain prime generation. For  $\text{MG}_2$  prime generation, one only needs in addition a variable-precision numerical package to approximate  $A(B/A)^{v/M}$  in step (2) of Algorithm 2.1 with sufficient accuracy. Table 1 gives some experimental results.

Besides making our main results and algorithms stronger, Algorithm 2.1 may have more applications. For example, one can consider various approximate counting problems with positive integer points  $(m, n)$  in the hyperbolic domain  $mn \leq x$ . The exact determination of the total number of such points is treated by Tao, Croot and Helfgott [31, Theorem 2.1 and Section 2.1]. This number can also be approximated, with the currently best known error bound  $x^{131/416+o(1)}$ , see [20]. However, these methods do not apply to counting integer points  $(m, n)$  under a hyperbola if some additional restrictions are imposed on  $m$  and  $n$  that might be expressed as congruence conditions or properties of  $b$ -ary expansions (to some

fixed base  $b \geq 2$ ) or a combination of both. For such questions, Algorithm 2.1 may lead to effective probabilistic estimation algorithms.

**Acknowledgments.** The authors thank Richard Brent and Paul Zimmermann for hints about numerical approximations, Konstantin Ziegler and Johannes Zollmann for useful discussions and permission to use some figures and tables they produced, and Konstantin Ziegler for suggesting to allude to Marie Germain.

## Bibliography

- [1] R. C. Baker and G. Harman, Shifted primes without large prime factors, *Acta Arith.* **83** (1998), 331–361.
- [2] R. C. Baker, G. Harman and J. Pintz, The difference between consecutive primes. II, *Proc. Lond. Math. Soc.* **83** (2001), 532–562.
- [3] P. T. Bateman and R. A. Horn, A heuristic asymptotic formula concerning the distribution of prime numbers, *Math. Comp.* **16** (1962), 363–367.
- [4] L. Blum, M. Blum and M. Shub, A simple unpredictable pseudo-random number generator, *SIAM J. Comp.* **15** (1986), 364–383.
- [5] R. P. Brent and P. Zimmermann, *Modern Computer Arithmetic*, Cambridge Monogr. Appl. Comput. Math. 18, Cambridge University Press, 2011.
- [6] R. Crandall and C. Pomerance, *Prime Numbers: A Computational Perspective*, Springer, New York, 2005.
- [7] I. Damgård and M. Kopolowski, Practical threshold RSA signatures without a trusted dealer, in: *Advances in Cryptology (EUROCRYPT 2001)*, Lecture Notes in Comp. Sci. 2045, Springer, Berlin (2001), 152–165.
- [8] A. Decker and P. Moree, Counting RSA-integers, *Results Math.* **52** (2008), 35–39.
- [9] L. Devroye, *Non-Uniform Random Variate Generation*, Springer, New York, 1986.
- [10] P.-A. Fouque and J. Stern, Fully distributed threshold RSA under standard assumptions, in: *Advances in Cryptology (ASIACRYPT 2001)*, Lecture Notes in Comp. Sci. 2248, Springer, Berlin (2001), 310–330.
- [11] S. Galbraith, *Mathematics of Public Key Cryptography*, Cambridge University Press, New York, 2012.
- [12] J. von zur Gathen and J. Gerhard, *Modern Computer Algebra*, Cambridge University Press, Cambridge, 2013.
- [13] R. Gennaro, S. Halevi and T. Rabin, Secure hash-and-sign signatures without the random oracle, in: *Advances in Cryptology (EUROCRYPT '99)*, Lecture Notes in Comp. Sci. 1592, Springer, Berlin (1999), 123–139.

- [14] B.-N. Guo and F. Qi, Sharp bounds for harmonic numbers, *Applied Math. Comp.* **218** (2011), 991–995.
- [15] G. Harman, *Prime-Detecting Sieves*, Princeton University Press, Princeton, 2007.
- [16] D. R. Heath-Brown, Artin’s conjecture for primitive roots, *Quart. J. Math.* **37** (1986), 27–38.
- [17] D. R. Heath-Brown, The number of primes in a short interval, *J. Reine Angew. Math.* **389** (1988), 22–63.
- [18] D. Hofheinz, E. Kiltz and V. Shoup, Practical chosen ciphertext secure encryption from factoring, *J. Cryptology* **26** (2013), 375–441.
- [19] M. N. Huxley, On the difference between consecutive primes, *Invent. Math.* **15** (1972), 164–170.
- [20] M. N. Huxley, Exponential sums and lattice points III, *Proc. London Math. Soc.* **87** (2003) 591–609.
- [21] H. Iwaniec and E. Kowalski, *Analytic Number Theory*, American Mathematical Society, Providence, 2004.
- [22] M. Joye and P. Paillier, Fast generation of prime numbers on portable devices: An update, in: *Cryptographic Hardware and Embedded Systems (CHES 2006)*, Lecture Notes in Comp. Sci. 4249, Springer, Berlin (2006), 160–173.
- [23] D. E. Knuth, *The Art of Computer Programming, vol. 2, Seminumerical Algorithms*, Addison Wesley, Reading, 1981.
- [24] D. Loebenberg and M. Nüsken, Notions for RSA integers, *Intern. J. Applied Crypto.* (2013), 1753-053.
- [25] A. J. Menezes, P. C. van Oorschot and S. A. Vanstone, *Handbook of Applied Cryptography*, CRC Press, 1997.
- [26] D. Naccache, Double-speed safe primes generation, preprint (2003), <http://eprint.iacr.org/2003/175>.
- [27] T. Nishide and K. Sakurai, Distributed Paillier cryptosystem without trusted dealer, in: *Information Security Applications (WISA 2010)*, Lecture Notes in Comp. Sci. 6513, Springer, Berlin (2011), 44–60.
- [28] E. Ong and J. Kubiawicz, Optimizing Robustness while generating shared secret safe primes, in: *Public Key Cryptography (PKC 2005)*, Lecture Notes in Comp. Sci. 3386, Springer, Berlin (2005), 120–137.
- [29] J. B. Rosser and L. Schoenfeld, Approximate formulas for some functions of prime numbers, *Illinois J. Math.* **6** (1962), 64–94.
- [30] V. Shoup, Practical threshold signatures, in: *Advances in Cryptology (EUROCRYPT 2000)*, Lecture Notes in Comp. Sci. 1807, Springer, Berlin (2000), 207–220.



- 
- [31] T. Tao, E. Croot and H. Helfgott, Deterministic methods to find primes, *Math. Comp.* **81** (2012), 1233–1246.
  - [32] G. Tenenbaum, *Introduction to Analytic and Probabilistic Number Theory*, Cambridge University Press, 1995.
  - [33] A. I. Vinogradov, On the remainder in Mertens' formula (in Russian), *Dokl. Akad. Nauk SSSR* **148**, 262–263.
  - [34] B. de Weger, Estimates for RSA moduli counting functions, preprint (2008).
  - [35] K. Ziegler and J. Zollmann, Fast and uniform generation of safe RSA moduli — Extended Abstract. *WEWoRC 2013 — Book of Abstracts* (2013), 15–19. Karlsruher Institut für Technologie (KIT), Karlsruhe.

Received 8 March, 2013; accepted 5 August, 2013.

#### **Author information**

Joachim von zur Gathen, B-IT, Universität Bonn, 53113 Bonn, Germany.

E-mail: [gathen@bit.uni-bonn.de](mailto:gathen@bit.uni-bonn.de)

Igor E. Shparlinski, Department of Pure Mathematics, University of New South Wales, NSW 2109, Australia.

E-mail: [igor.shparlinski@unsw.edu.au](mailto:igor.shparlinski@unsw.edu.au)