

Fast Random Integer Generation in an Interval

DANIEL LEMIRE, Université du Québec (TELUQ), Canada

In simulations, probabilistic algorithms and statistical tests, we often generate random integers in an interval (e.g., $[0, s)$). For example, random integers in an interval are essential to the Fisher-Yates random shuffle. Consequently, popular languages like Java, Python, C++, Swift and Go include ranged random integer generation functions as part of their runtime libraries.

Pseudo-random values are usually generated in words of a fixed number of bits (e.g., 32 bits, 64 bits) using algorithms such as a linear congruential generator. We need functions to convert such random words to random integers in an interval $([0, s))$ without introducing statistical biases. The standard functions in programming languages such as Java involve integer divisions. Unfortunately, division instructions are relatively expensive. We review an unbiased function to generate ranged integers from a source of random words that avoids integer divisions with high probability. To establish the practical usefulness of the approach, we show that this algorithm can multiply the speed of unbiased random shuffling on x64 processors. Our proposed approach has been adopted by the Go language for its implementation of the shuffle function.

CCS Concepts: • **Theory of computation** → **Pseudorandomness and derandomization**; • **Software and its engineering** → **Software performance**;

Additional Key Words and Phrases: Random number generation, Rejection method, Randomized algorithms

1 INTRODUCTION

There are many efficient techniques to generate high-quality pseudo-random numbers such as Mersenne Twister [28], Xorshift [27, 32], linear congruential generators [6, 9, 20, 21, 24] and so forth [22, 26]. Many pseudo-random number generators produce 32-bit or 64-bit words that can be interpreted as integers in $[0, 2^{32})$ and $[0, 2^{64})$ respectively: the produced values are practically indistinguishable from truly random numbers in $[0, 2^{32})$ or $[0, 2^{64})$. In particular, no single value is more likely than any other.

However, we often need random integers selected uniformly from an interval $[0, s)$, and this interval may change dynamically. It is useful for selecting an element at random in an array containing s elements, but there are less trivial uses. For example, the Fisher-Yates random shuffle described by Knuth [8, 16] (see Algorithm 1) requires one random integer in an interval for each value in an array to be shuffled. Ideally, we would want these values to be generated without bias so that all integer values in $[0, s)$ are equally likely. Only then are all permutations equally likely. A related algorithm is *reservoir sampling* [38] (see Algorithm 2) which randomly selects a subset of values from a possibly very large array, even when the size of the array is initially unknown.

We use random permutations as part of simulation algorithms [1, 5, 7, 14, 29, 31]. The performance of randomized permutation algorithms is important. Various non-parametric tests in statistics and machine learning repeatedly permute randomly the original data. In some contexts, Hinrichs et al. found that the computational burden due to random permutations can be prohibitive [15]. Unsurprisingly, much work has been done on parallelizing permutations [13, 18, 34, 35, 40] for greater speed.

One might think that going from fixed-bit pseudo-random numbers (e.g., 32-bit integers) to pseudo-random numbers in an interval is a minor, inexpensive operation. However, this may not be true, at least when the interval is a changing parameter and we desire a uniformly-distributed result.

ALGORITHM 1: Fisher-Yates random shuffle: it shuffles an array of size n so that $n!$ possible permutations are equiprobable.

Require: array A made of n elements indexed from 0 to $n - 1$

```
1: for  $i = n - 1, \dots, 1$  do
2:    $j \leftarrow$  random integer in  $[0, i]$ 
3:   exchange  $A[i]$  and  $A[j]$ 
4: end for
```

ALGORITHM 2: Reservoir sampling: returns an array R containing k distinct elements picked randomly from an array A of size n so that all $\binom{n}{k}$ possible samples are equiprobable.

Require: array A made of n elements indexed from 0 to $n - 1$

Require: integer k ($0 < k \leq n$)

```
1:  $R \leftarrow$  array of size  $k$ 
2: for  $i = 0, \dots, k - 1$  do
3:    $R[i] \leftarrow A[i]$ 
4: end for
5: for  $i = k, \dots, n - 1$  do
6:    $j \leftarrow$  random integer in  $[0, i]$ 
7:   if  $j < k$  then
8:      $R[j] \leftarrow A[i]$ 
9:   end if
10: end for
11: return  $R$ 
```

- Let us consider a common but biased approach to the problem of converting numbers from a large interval $[0, 2^n)$ to numbers in a subinterval $[0, s)$ ($s \leq 2^n$): the modulo reduction $x \rightarrow x \bmod s$. On x64 processors, this could be implemented through the division (`div`) instruction when both s and x are parameters. When applied to 32-bit registers, this instruction has a latency of 26 cycles [10]. With 64-bit registers, the latency ranges from 35 to 88 cycles, with longer running times for small values of s .
- Another biased but common approach consists in using a fixed-point floating-point representation consisting of the following step:
 - we convert the random word to a floating-point number in the interval $[0, 1)$,
 - we convert the integer s into a floating-point number,
 - we multiply the two resulting floating-point numbers,
 - and we convert the floating-point result to an integer in $[0, s)$.

When using the typical floating-point standard (IEEE 754), we can at best represent all values in $[0, 2^{24})$ divided by 2^{24} using a 32-bit floating-point number. Thus we do not get the full 32-bit range: we cannot generate all numbers in $[0, s)$ if $s > 2^{24}$. To do so, we must use double precision floating-point numbers, and then we can represent all values in $[0, 2^{53})$ divided by 2^{53} . Moreover converting between floating-point values and integers is not without cost: the corresponding instructions on x64 processors (e.g., `cvtsi2ss`, `cvtss2si`) have at least six cycles of latency on Skylake processors [10].

While generating a whole new 64-bit pseudo-random number can take as little as a handful of cycles [33], transforming it into an integer in an interval ($[0, s)$ for $s \in [0, 2^{64})$) without bias can take an order of magnitude longer when using division operations.

There is a fast technique that avoids division and does not require floating-point numbers. Indeed, given an integer x in the interval $[0, 2^L)$, we have that the integer $(x \times s) \div 2^L$ is in $[0, s)$ for any integer $s \in [0, 2^L]$. If the integer x is picked at random in $[0, 2^L)$, then the result $(x \times s) \div 2^L$ is a random integer in $[0, s)$ [30]. The division by a power of two ($\div 2^L$) can be implemented by a bit shift instruction, which is inexpensive. A multiplication followed by a shift is much more economical on current processors than a division, as it can be completed in only a handful of cycles. It introduces a bias, but we can correct for it efficiently using the rejection method (see § 4). This multiply-and-shift approach is similar in spirit to the multiplication by a floating-point number in the unit interval $([0, 1))$ in the sense that $(sx) \div 2^L$ can be intuitively compared with $s \times x/2^L$ where $x/2^L$ is a random number in $[0, 1)$.

Though the idea that we can avoid divisions when generating numbers in an interval is not novel, we find that many standard libraries (Java, Go, ...) use an approach that incurs at least one integer division per function call. We believe that a better default would be an algorithm that avoids division instructions with high probability. We show experimentally that such an algorithm can provide superior performance.

2 MATHEMATICAL NOTATION

We let $\lfloor x \rfloor$ be the largest integer smaller than or equal to x , we let $\lceil x \rceil$ be the smallest integer greater than or equal to x . We let $x \div y$ be the integer division of x by y , defined as $\lfloor x/y \rfloor$. We define the remainder of the division of x by y as $x \bmod y$: $x \bmod y \equiv x - (x \div y)y$.

We are interested in the integers in an interval $[0, 2^L)$ where, typically, $L = 32$ or $L = 64$. We refer to these integers as L -bit integers.

When we consider the values of $x \bmod s$ as x goes from 0 to 2^L , we get the 2^L values

$$\underbrace{0, 1, \dots, s-1}_{s \text{ values}}, \underbrace{0, 1, \dots, s-1}_{s \text{ values}}, \underbrace{0, 1, \dots, s-1}_{s \text{ values}}, \underbrace{0, 1, \dots, (2^L \bmod s) - 1}_{2^L \bmod s \text{ values}}.$$

$(2^L \div s)s$ values

We have the following lemma by inspection.

LEMMA 2.1. *Given integers $a, b, s > 0$, there are exactly $(b - a) \div s$ multiples of s in $[a, b)$ whenever s divides $b - a$. More generally, there are exactly $(b - a) \div s$ integers in $[a, b)$ having a given remainder with s whenever s divides $b - a$.*

A geometric distribution with success probability $p \in [0, 1]$ is a discrete distribution taking value $k \in \{1, 2, \dots\}$ with probability $(1 - p)^{k-1}p$. The mean of a geometric distribution is $1/p$.

3 EXISTING UNBIASED TECHNIQUES FOUND IN COMMON SOFTWARE LIBRARIES

Assume that we have a source of uniformly-distributed L -bit random numbers, i.e., integers in $[0, 2^L)$. From such a source of random numbers, we want to produce a uniformly-distributed random integer y in $[0, s)$ for some integer $s \in [1, 2^L]$. That is all integers from the interval are equally likely: $P(y = z) = 1/s$ for any integer $z \in [0, s)$. We then say that the result is *unbiased*.

If s divides 2^L , i.e., it is a power of two, then we can divide the random integer x from $[0, 2^L)$ by $2^L/s = 2^L \div s$. However, we are interested in the general case where s may be any integer value in $[0, 2^L)$.

We can achieve an unbiased result by the rejection method [39]. For example, we could generate random integers $x \in [0, 2^L)$ until x is in $[0, s)$, rejecting all other cases.¹ Rejecting so many values

¹For some applications where streams of random numbers need to be synchronized, the rejection method is not applicable [2, 4, 12, 19, 23].

ALGORITHM 3: The OpenBSD algorithm.

Require: source of uniformly-distributed random integers in $[0, 2^L)$

Require: target interval $[0, s)$ with $s \in [0, 2^L)$

```
1:  $t \leftarrow (2^L - s) \bmod s$   $\{(2^L - s) \bmod s = 2^L \bmod s\}$ 
2:  $x \leftarrow$  random integer in  $[0, 2^L)$ 
3: while  $x < t$  do {Application of the rejection method}
4:    $x \leftarrow$  random integer in  $[0, 2^L)$ 
5: end while
6: return  $x \bmod s$ 
```

is wasteful. Popular software libraries use more efficient algorithms. We provide code samples in Appendix A.

3.1 The OpenBSD Algorithm

The C standard library in OpenBSD and macOS have an `arc4random_uniform` function to generate unbiased random integers in an interval $[0, s)$. See Algorithm 3. The Go language (e.g., version 1.9) has adopted the same algorithm for its `Int63n` and `Int31n` functions, with minor implementation differences [37]. The GNU C++ standard library (e.g., version 7.2) also relies on the same algorithm [11].

The interval $[2^L \bmod s, 2^L)$ has size $2^L - (2^L \bmod s)$ which is divisible by s . From Lemma 2.1, if we generate random integers from integer values in $[2^L \bmod s, 2^L)$ as remainders of a division by s ($x \bmod s$), then each of the integers occur for $2^L \div s$ integers $x \in [0, 2^L)$. To produce integers in $[2^L \bmod s, 2^L)$, we use the rejection method: we generate integers in $x \in [0, 2^L)$ but reject the result whenever $x < 2^L \bmod s$. If we have a source of unbiased random integers in $[0, 2^L)$, then the result of Algorithm 3 is an unbiased random integer in $[0, s)$.

The number of random integers consumed by this algorithm follows a geometric distribution, with a success probability $p = 1 - (2^L \bmod s)/2^L$. On average, we need $1/p$ random words. This average is less than two, irrespective of the value of s . The algorithm always requires the computation of two remainders.

There is a possible trivial variation on the algorithm where instead of rejecting the integer from $[0, 2^L)$ when it is part of the first $2^L \bmod s$ values (in $[0, 2^L \bmod s)$), we reject the integer from $[0, 2^L)$ when it is part of the last $2^L \bmod s$ values (in $[2^L - (2^L \bmod s), 2^L)$).

3.2 The Java Approach

It is unfortunate that Algorithm 3 always requires the computation of two remainders, especially because we anticipate such computations to have high latency. The first remainder is used to determine whether a rejection is necessary ($x < (2^L - s) \bmod s$), and the second remainder is used to generate the value in $[0, s)$ as $x \bmod s$.

The Java language, in its `Random` class, uses an approach that often requires a single remainder (e.g., as of OpenJDK 9). Suppose we pick a number $x \in [0, 2^L)$ and we compute its remainder $x \bmod s$. Having both x and $x \bmod s$, we can determine whether x is allowable (is in $[0, 2^L - (2^L \bmod s))$) without using another division. When it is the case, we can then return $x \bmod s$ without any additional computation. See Algorithm 4.

The number of random words and remainders used by the Java algorithm follows a geometric distribution, with a success probability $p = 1 - (2^L \bmod s)/2^L$. Thus, on average, we need $1/p$ random words and remainders. Thus when s is small ($s \ll 2^L$), we need ≈ 1 random words and remainders. This compares favorably to the OpenBSD algorithm that always requires the computation of two

ALGORITHM 4: The Java algorithm.

Require: source of uniformly-distributed random integers in $[0, 2^L)$

Require: target interval $[0, s)$ with $s \in [0, 2^L)$

1: $x \leftarrow$ random integer in $[0, 2^L)$

2: $r \leftarrow x \bmod s$

$$\{x - r > 2^L - s \Leftrightarrow x \in [0, 2^L - (2^L \bmod s))\}$$

3: **while** $x - r > 2^L - s$ **do** {Application of the rejection method}

4: $x \leftarrow$ random integer in $[0, 2^L)$

5: $r \leftarrow x \bmod s$

6: **end while**

7: **return** r

remainders. However, the maximal number of divisions required by the OpenBSD algorithm is two, whereas the Java approach could require infinitely many divisions in the worst case.

4 AVOIDING DIVISION

Though arbitrary integer divisions are relatively expensive on common processors, bit shifts are less expensive, often requiring just one cycle. When working with unsigned integers, a bit shift is equivalent to a division by a power of two. Thus we can compute $x \div 2^k$ quickly for any power of two 2^k . Similarly, we can compute the remainder of the division by a power of two as a simple bitwise logical AND: $x \bmod 2^k = x \text{ AND } (2^k - 1)$.

Moreover, common general-purpose processors (e.g., x64 or ARM processors) can efficiently compute the full result of a multiplication. That is, when multiplying two 32-bit integers, we can get the full 64-bit result and, in particular, the most significant 32 bits. Similarly, we can multiply two 64-bit integers and get the full 128-bit result or just the most significant 64 bits when using a 64-bit processor. Most modern programming languages (C, C++, Java, Swift, Go...) have native support for 64-bit integers. Thus it is efficient to get the most significant 32 bits of the product of two 32-bit integers. To get the most significant 64 bits of the product of two 64-bit integers in C and C++, we can use either intrinsics (e.g., `__umulh` when using the Visual Studio compiler) or the `__uint128_t` extension supported by the GNU GCC and LLVM's clang compilers. The Swift language has the `multipliedFullWidth` function that works with both 32-bit and 64-bit integers. It gets compiled to efficient binary code.

Given an integer $x \in [0, 2^L)$, we have that $(x \times s) \div 2^L \in [0, s)$. By multiplying by s , we take integer values in the range $[0, 2^L)$ and map them to multiples of s in $[0, s \times 2^L)$. By dividing by 2^L , we map all multiples of s in $[0, 2^L)$ to 0, all multiples of s in $[2^L, 2 \times 2^L)$ to one, and so forth. The $(i + 1)^{\text{th}}$ interval is $[i \times 2^L, (i + 1) \times 2^L)$. By Lemma 2.1, there are exactly $\lfloor 2^L/s \rfloor$ multiples of s in intervals $[i \times 2^L + (2^L \bmod s), (i + 1) \times 2^L)$ since s divides the size of the interval $(2^L - (2^L \bmod s))$. Thus if we reject the multiples of s that appear in $[i \times 2^L, i \times 2^L + (2^L \bmod s))$, we get that all intervals have exactly $\lfloor 2^L/s \rfloor$ multiples of s . We can formalize this result as the next lemma.

LEMMA 4.1. *Given any integer $s \in [0, 2^L)$, we have that for any integer $y \in [0, s)$, there are exactly $\lfloor 2^L/s \rfloor$ values $x \in [0, 2^L)$ such that $(x \times s) \div 2^L = y$ and $(x \times s) \bmod 2^L \geq 2^L \bmod s$.*

Algorithm 5 is a direct application of Lemma 4.1. It generates unbiased random integers in $[0, s)$ for any integer $s \in (0, 2^L)$.

This algorithm has the same probabilistic distribution of random words as the previously presented algorithms, requiring the same number of L -bit uniformly-distributed random words on average. However, the number of divisions (or remainders) is more advantageous. Excluding divisions by a power of two, we have a probability $\frac{2^L - s}{2^L}$ of using no division at all. Otherwise, if the

ALGORITHM 5: An efficient algorithm to generate unbiased random integers in an interval.

Require: source of uniformly-distributed random integers in $[0, 2^L)$
Require: target interval $[0, s)$ with $s \in [0, 2^L)$

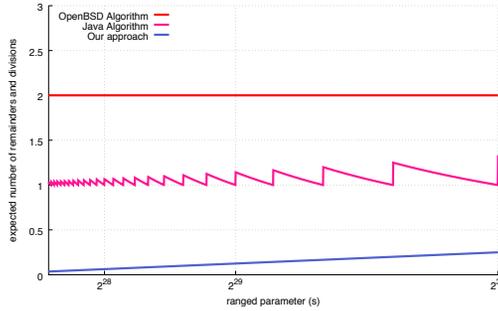
```

1:  $x \leftarrow$  random integer in  $[0, 2^L)$ 
2:  $m \leftarrow x \times s$ 
3:  $l \leftarrow m \bmod 2^L$ 
4: if  $l < s$  then {Application of the rejection method}
5:    $t \leftarrow (2^L - s) \bmod s$   $\{2^L \bmod s = (2^L - s) \bmod s\}$ 
6:   while  $l < t$  do
7:      $x \leftarrow$  random integer in  $[0, 2^L)$ 
8:      $m \leftarrow x \times s$ 
9:      $l \leftarrow m \bmod 2^L$ 
10:  end while
11: end if
12: return  $m \div 2^L$ 

```

 Table 1. Number of remainder computations (excluding those by known powers of two) in the production an unbiased random integer in $[0, s)$.

	expected number of remainders per integer in $[0, s)$	expected number of remainders when the interval is small ($s \ll 2^L$)	maximal number of remainders per integer in $[0, s)$
OpenBSD (Algorithm 3)	2	2	2
Java (Algorithm 4)	$\frac{2^L}{2^L - (2^L \bmod s)}$	1	∞
Our approach (Algorithm 5)	$\frac{s}{2^L}$	0	1


 Fig. 1. Expected number of remainder computations (excluding those by known powers of two) in the production an unbiased random integer in $[0, s)$ for 32-bit integers.

initial value of $l = (x \times s) \bmod 2^L$ is less than s , then we incur automatically the cost of a single division (to compute t), but no further division (not counting divisions by a power of two). Table 1 and Fig. 1 compare the three algorithms in terms of the number of remainder computations needed.

5 EXPERIMENTS

We implemented our software in C++ on a Linux server with an Intel (Skylake) i7-6700 processor running at 3.4 GHz. This processor has 32 kB of L1 data cache, 256 kB of L2 cache per core with 8 MB of L3 cache, and 32 GB of RAM (DDR4 2133, double-channel). We use the GNU GCC 5.4 compilers with the “-O3 -march=native” flags. To ensure reproducibility, we make our software freely available.² Though we implement and benchmark a Java-like approach (Algorithm 4), all our experiments are conducted using C++ code.

For our experiments, we use a convenient and fast linear congruential generator with the recurrence formula $X_{n+1} = c \times X_n \bmod 2^{128}$ where $c = 15750249268501108917$ to update the 128-bit state of the generator ($X_i \in [0, 2^{128})$ for $i = 0, 1, 2, \dots$), returning $X_{n+1} \div 2^{64}$ as a 64-bit random word [21]. We start from a 128-bit seed X_0 . This well-established generator passes difficult statistical tests such as Big Crush [25]. It is well suited to x64 processors because they have fast 64-bit multipliers.

We benchmark the time required, per element, to randomly shuffle arrays of integers having different sizes. We can consider array indexes to be either 32-bit or 64-bit values. When working with 64-bit indexes, we require 64-bit integer divisions which are slower than 32-bit integer divisions on x64 processors. We always use the same 64-bit random-number generator, but in the 32-bit case, we only use the least significant 32 bits of each random word. For reference, in both the 32-bit and 64-bit figures, we include the results obtained with the shuffle functions of the standard library (`std::shuffle`), implemented using our 64-bit random-number generator. For small arrays, the `std::shuffle` has performance similar to our implementation of the OpenBSD algorithm, but it becomes slightly slower when shuffling larger arrays.

We present our experimental results in Fig. 3. We report the wall-clock time averaged over at least 5 shuffles. When the arrays fit in the cache, we expect them to remain in the cache. The time required is normalized by the number of elements in the array. As long as the arrays fit in the CPU cache, the array size does not affect the performance. As the arrays grow larger, the latency of memory access becomes a factor and the performance decreases.

- In the 32-bit case, the approach with few divisions can be almost twice as fast as the Java-like approach which itself can be at least 50% faster than the OpenBSD-like approach.
- When shuffling with 64-bit indexes as opposed to 32-bit indexes, our implementations of the OpenBSD-like and Java-like algorithms become significantly slower (up to three times) due to the higher cost of the 64-bit division. Thus our approach can be more than three times faster than the Java version in the 64-bit case.

The relative speed differences between the different algorithms become less significant when the arrays grow larger. In Fig. 2, we present the ratio of the OpenBSD-like approach with our approach. We see that the relative benefit of our approach diminishes when the array size increases. In the 32-bit case, for very large arrays, our approach is merely 50% faster whereas it is nearly three times faster for small arrays. Using Linux `perf`, we estimated the number of cache misses to shuffle an array containing 100 million integers and found that the OpenBSD approach generates about 50% more cache misses than our approach.

To help the processor prefetch memory and reduce the number of cache misses, we can compute the random integers in small blocks, and then shuffle while reading the precomputed integers (see Algorithm 6). The resulting buffered algorithm is equivalent to the conventional Fisher-Yates random shuffle, and it involves the computation of the same number of random indexes, but it differs on how the memory accesses are scheduled. In Fig. 4, we see that the OpenBSD-like approach

²<https://github.com/lemire/FastShuffleExperiments>

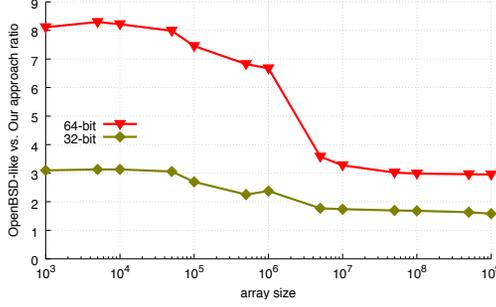
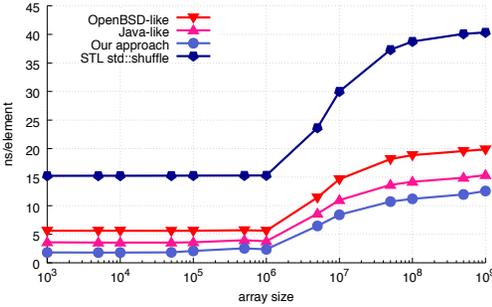
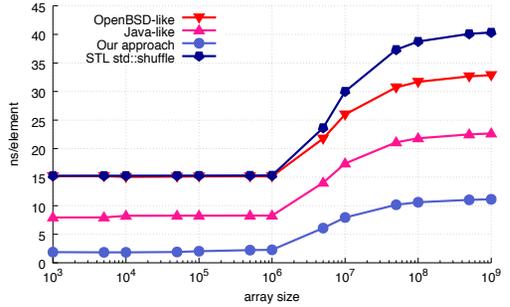


Fig. 2. Ratio of the timings of the OpenBSD-like approach and of our approach.

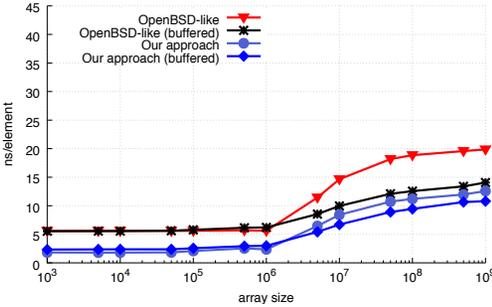


(a) 32-bit indexes

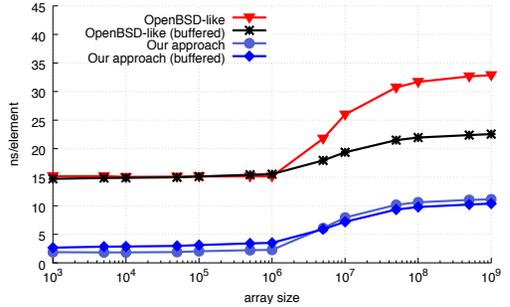


(b) 64-bit indexes

Fig. 3. Wall-clock time in nanoseconds per element to shuffle arrays of random integers.



(a) 32-bit indexes



(b) 64-bit indexes

Fig. 4. Wall-clock time in nanoseconds per element to shuffle arrays of random integers using either regular shuffles or buffered shuffles with a buffer of size 256 (see Algorithm 6).

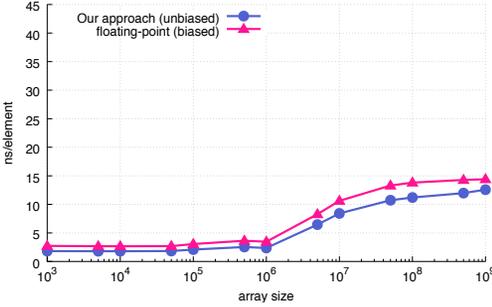
benefits from the buffering when shuffling large arrays. A significant fraction of the running time of the regular OpenBSD-like implementation is due to caching issues. When applied to our approach, the benefits of the buffering are small, and for small to medium arrays, the buffering is slightly harmful.

ALGORITHM 6: Buffered version of the Fisher-Yates random shuffle.

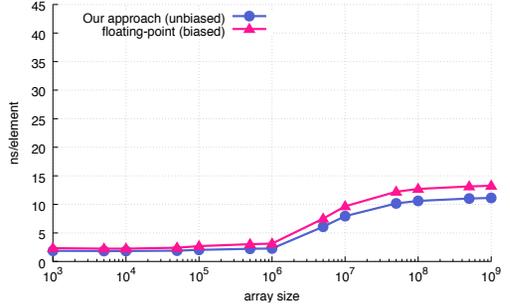
Require: array A made of n elements indexed from 0 to $n - 1$

Require: $B \leftarrow$ a small positive constant (the buffer size)

```
1:  $i \leftarrow n - 1$ 
2:  $Z \leftarrow$  some array with capacity  $B$  (the buffer)
3: while  $i \geq B$  do
4:   for  $k \in \{i, i - 1, \dots, i - B + 1\}$  do
5:      $Z_k \leftarrow$  random integer in  $[0, k]$ 
6:   end for
7:   for  $k \in \{i, i - 1, \dots, i - B + 1\}$  do
8:     exchange  $A[k]$  and  $A[Z_k]$ 
9:   end for
10:   $i \leftarrow i - B$ 
11: end while
12: while  $i > 0$  do
13:   $j \leftarrow$  random integer in  $[0, i]$ 
14:  exchange  $A[i]$  and  $A[j]$ 
15:   $i \leftarrow i - 1$ 
16: end while
```



(a) 32-bit indexes



(b) 64-bit indexes

Fig. 5. Wall-clock time in nanoseconds per element to shuffle arrays of random integers using either Algorithm 5 or an approach based on floating-point multiplications.

We also compare Algorithm 5 to the floating-point approach we described in § 1 where we represent the random number as a floating-point value in $[0, 1)$ which we multiply by s to get a number in $[0, s)$. As illustrated in Fig. 5, the floating-point approach is slightly slower (10%–30%) whether we use 32-bit or 64-bit indexes. Moreover, it introduces a small bias and it is limited to $s \leq 2^{24}$ in the 32-bit case and to $s \leq 2^{53}$ in the 64-bit case.

6 CONCLUSION

We find that the algorithm often used in OpenBSD and macOS (through the `arc4random_uniform` function) requires two divisions per random number generation. It is also the slowest in our tests. The Java approach that often requires only one division, can be faster. We believe that it should be preferred to the to the OpenBSD algorithm.

As we have demonstrated, we can use nearly no division at all and at most one in the worst case. Avoiding divisions can multiply the speed of unbiased random shuffling functions on x64 processors.

For its new random shuffle function, the Go programming language adopted our proposed approach [36]. The Go authors justify this decision by the fact that it results in 30% better speed compared with the application of the OpenBSD approach (Algorithm 3).

Our results are only relevant in the context where the generation of random integers is fast compared with the latency of a division operation. In the case where the generation of random bits is likely the bottleneck, other approaches would be preferable [3]. Moreover, our approach may be not applicable to specialized processors such as Graphics Processing Units (GPUs) that lack support for the computation of the full multiplication [17].

A CODE SAMPLES

```
// returns value in [0,s)
// random64 is a function returning random 64-bit words
uint64_t openbsd(uint64_t s, uint64_t (*random64)(void)) {
    uint64_t t = (-s) % s;
    uint64_t x;
    do {
        x = random64();
    } while (x < t);
    return x % s;
}

uint64_t java(uint64_t s, uint64_t (*random64)(void)) {
    uint64_t x = random64();
    uint64_t r = x % s;
    while (x - r > UINT64_MAX - s + 1) {
        x = random64();
        r = x % s;
    }
    return r;
}

uint64_t nearlydivisionless(uint64_t s, uint64_t (*random64)(void)) {
    uint64_t x = random64();
    __uint128_t m = (__uint128_t) x * (__uint128_t) s;
    uint64_t l = (uint64_t) m;
    if (l < s) {
        uint64_t t = -s % s;
        while (l < t) {
            x = random64();
            m = (__uint128_t) x * (__uint128_t) s;
            l = (uint64_t) m;
        }
    }
    return m >> 64;
}
```

ACKNOWLEDGMENTS

The work is supported by the Natural Sciences and Engineering Research Council of Canada under grant number RGPIN-2017-03910. The author would like to thank R. Startin and J. Epler for independently reproducing the experimental results and providing valuable feedback.

REFERENCES

- [1] Michael Amrein and Hans R. Künsch. 2011. A Variant of Importance Splitting for Rare Event Estimation: Fixed Number of Successes. *ACM Trans. Model. Comput. Simul.* 21, 2, Article 13 (Feb. 2011), 20 pages. <https://doi.org/10.1145/1899396.1899401>
- [2] Søren Asmussen and Peter W Glynn. 2007. *Stochastic simulation: algorithms and analysis*. Springer Science & Business Media, New York, NY, USA.
- [3] Axel Bacher, Olivier Bodini, Hsien-Kuei Hwang, and Tsung-Hsi Tsai. 2017. Generating Random Permutations by Coin Tossing: Classical Algorithms, New Analysis, and Modern Implementation. *ACM Trans. Algorithms* 13, 2, Article 24 (Feb. 2017), 43 pages. <https://doi.org/10.1145/3009909>
- [4] Paul Bratley, Bennet L Fox, and Linus E Schrage. 1987. *A guide to simulation* (2 ed.). Springer Science & Business Media, New York, NY, USA.
- [5] James M. Calvin and Marvin K. Nakayama. 1998. Using Permutations in Regenerative Simulations to Reduce Variance. *ACM Trans. Model. Comput. Simul.* 8, 2 (April 1998), 153–193. <https://doi.org/10.1145/280265.280273>
- [6] A De Matteis and Simonetta Pagnutti. 1988. Parallelization of random number generators and long-range correlations. *Numer. Math.* 53, 5 (1988), 595–608.
- [7] Luc Devroye. 1997. Random Variate Generation for Multivariate Unimodal Densities. *ACM Trans. Model. Comput. Simul.* 7, 4 (Oct. 1997), 447–477. <https://doi.org/10.1145/268403.268413>
- [8] Richard Durstenfeld. 1964. Algorithm 235: Random Permutation. *Commun. ACM* 7, 7 (July 1964), 420–. <https://doi.org/10.1145/364520.364540>
- [9] George Fishman. 2013. *Monte Carlo: concepts, algorithms, and applications* (corrected edition ed.). Springer Science & Business Media, Berlin.
- [10] Agner Fog. 2016. *Instruction tables: Lists of instruction latencies, throughputs and micro-operation breakdowns for Intel, AMD and VIA CPUs*. Technical Report. Copenhagen University College of Engineering, Copenhagen, Denmark. http://www.agner.org/optimize/instruction_tables.pdf [last checked October 2017].
- [11] Free Software Foundation. 2017. The GNU C++ Library API Reference. <https://gcc.gnu.org/onlinedocs/libstdc++/api.html> [last checked October 2017]. (2017).
- [12] Michael C. Fu and Shane G. Henderson. 2017. History of seeking better solutions, AKA simulation optimization. In *Proceedings of the 2017 Winter Simulation Conference*. IEEE Press, New York, NY, USA, 131–157.
- [13] Jens Gustedt. 2008. *Engineering Parallel In-Place Random Generation of Integer Permutations*. Springer Berlin Heidelberg, Berlin, Heidelberg, 129–141. https://doi.org/10.1007/978-3-540-68552-4_10
- [14] Alejandro S. Hernandez, Thomas W. Lucas, and Matthew Carlyle. 2012. Constructing Nearly Orthogonal Latin Hypercubes for Any Nonsaturated Run-variable Combination. *ACM Trans. Model. Comput. Simul.* 22, 4, Article 20 (Nov. 2012), 17 pages. <https://doi.org/10.1145/2379810.2379813>
- [15] Chris Hinrichs, Vamsi K. Ithapu, Qinyuan Sun, Sterling C. Johnson, and Vikas Singh. 2013. Speeding Up Permutation Testing in Neuroimaging. In *Proceedings of the 26th International Conference on Neural Information Processing Systems (NIPS'13)*. Curran Associates Inc., USA, 890–898.
- [16] Donald E. Knuth. 1969. *Seminumerical Algorithms*. The Art of Computer Programming, Vol. 2. Addison-Wesley, Boston.
- [17] W. B. Langdon. 2009. A Fast High Quality Pseudo Random Number Generator for nVidia CUDA. In *Proceedings of the 11th Annual Conference Companion on Genetic and Evolutionary Computation Conference: Late Breaking Papers (GECCO '09)*. ACM, New York, NY, USA, 2511–2514. <https://doi.org/10.1145/1570256.1570353>
- [18] Daniel Langr, Pavel Tvrdík, Tomáš Dytrych, and Jerry P. Draayer. 2014. Algorithm 947: Paraperm—Parallel Generation of Random Permutations with MPI. *ACM Trans. Math. Softw.* 41, 1, Article 5 (Oct. 2014), 26 pages. <https://doi.org/10.1145/2669372>
- [19] Averill M. Law. 2007. *Simulation Modeling and Analysis* (fourth edition ed.). McGraw-Hill, New York, NY, USA.
- [20] Pierre L'Ecuyer. 1990. Random Numbers for Simulation. *Commun. ACM* 33, 10 (Oct. 1990), 85–97. <https://doi.org/10.1145/84537.84555>
- [21] Pierre L'Ecuyer. 1999. Tables of linear congruential generators of different sizes and good lattice structure. *Mathematics of Computation of the American Mathematical Society* 68, 225 (1999), 249–260.
- [22] Pierre L'Ecuyer. 2015. History of uniform random number generation. In *Proceedings of the 2017 Winter Simulation Conference*. IEEE Press, New York, NY, USA, 1–17.

- [23] Pierre L'Ecuyer. 2015. Random number generation with multiple streams for sequential and parallel computing. In *Proceedings of the 2015 Winter Simulation Conference*. IEEE Press, New York, NY, USA, 31–44.
- [24] Pierre L'Ecuyer, François Blouin, and Raymond Couture. 1993. A Search for Good Multiple Recursive Random Number Generators. *ACM Trans. Model. Comput. Simul.* 3, 2 (April 1993), 87–98. <https://doi.org/10.1145/169702.169698>
- [25] Pierre L'Ecuyer and Richard Simard. 2007. TestU01: A C Library for Empirical Testing of Random Number Generators. *ACM Trans. Math. Softw.* 33, 4, Article 22 (Aug. 2007), 40 pages. <https://doi.org/10.1145/1268776.1268777>
- [26] Pierre L'Ecuyer. 2012. Random number generation. In *Handbook of Computational Statistics*. Springer, Berlin, 35–71.
- [27] George Marsaglia. 2003. Xorshift RNGs. *Journal of Statistical Software* 8, 14 (2003), 1–6.
- [28] Makoto Matsumoto and Takuji Nishimura. 1998. Mersenne Twister: A 623-dimensionally Equidistributed Uniform Pseudo-random Number Generator. *ACM Trans. Model. Comput. Simul.* 8, 1 (Jan. 1998), 3–30. <https://doi.org/10.1145/272991.272995>
- [29] Takayuki Osogami. 2009. Finding Probably Best Systems Quickly via Simulations. *ACM Trans. Model. Comput. Simul.* 19, 3, Article 12 (Aug. 2009), 19 pages. <https://doi.org/10.1145/1540530.1540533>
- [30] Mark Overton. 2011. Fast, High-Quality, Parallel Random-Number Generators: Comparing Implementations. <http://www.drdoobs.com/tools/fast-high-quality-parallel-random-number/231000484> [last checked October 2017]. (2011).
- [31] Art B. Owen. 1998. Latin Supercube Sampling for Very High-dimensional Simulations. *ACM Trans. Model. Comput. Simul.* 8, 1 (Jan. 1998), 71–102. <https://doi.org/10.1145/272991.273010>
- [32] François Panneton and Pierre L'Ecuyer. 2005. On the Xorshift Random Number Generators. *ACM Trans. Model. Comput. Simul.* 15, 4 (Oct. 2005), 346–361. <https://doi.org/10.1145/1113316.1113319>
- [33] Mutsuo Saito and Makoto Matsumoto. 2008. *SIMD-Oriented Fast Mersenne Twister: a 128-bit Pseudorandom Number Generator*. Springer Berlin Heidelberg, Berlin, Heidelberg, 607–622. https://doi.org/10.1007/978-3-540-74496-2_36
- [34] Peter Sanders. 1998. Random permutations on distributed, external and hierarchical memory. *Inform. Process. Lett.* 67, 6 (1998), 305–309. [https://doi.org/10.1016/S0020-0190\(98\)00127-6](https://doi.org/10.1016/S0020-0190(98)00127-6)
- [35] Ivo D. Shterev, Sin-Ho Jung, Stephen L. George, and Kouros Owzar. 2010. permGPU: Using graphics processing units in RNA microarray association studies. *BMC Bioinformatics* 11, 1 (16 Jun 2010), 329. <https://doi.org/10.1186/1471-2105-11-329>
- [36] Josh Bleecher Snyder. 2017. math/rand: add Shuffle. <https://go-review.golang.org/c/go/+51891> [last checked October 2017]. (2017).
- [37] The Go authors 2017. Package rand implements pseudo-random number generators. <https://github.com/golang/go/blob/master/src/math/rand/rand.go> [last checked October 2017]. (2017).
- [38] Jeffrey S. Vitter. 1985. Random Sampling with a Reservoir. *ACM Trans. Math. Softw.* 11, 1 (March 1985), 37–57. <https://doi.org/10.1145/3147.3165>
- [39] John Von Neumann. 1951. Various techniques used in connection with random digits. *National Bureau of Standards Series* 12 (1951), 36–38.
- [40] Michael Waechter, Kay Hamacher, Franziska Hoffgaard, Sven Widmer, and Michael Goesele. 2012. *Is Your Permutation Algorithm Unbiased for $n \neq 2^m$?* Springer Berlin Heidelberg, Berlin, Heidelberg, 297–306. https://doi.org/10.1007/978-3-642-31464-3_30