

**Note: Fully integrated 3.2 Gbps quantum random number generator with real-time extraction**

Xiao-Guang Zhang, You-Qi Nie, Hongyi Zhou, Hao Liang, Xiongfeng Ma, Jun Zhang<sup>\*</sup>, and Jian-Wei Pan

Citation: *Rev. Sci. Instrum.* **87**, 076102 (2016); doi: 10.1063/1.4958663

View online: <http://dx.doi.org/10.1063/1.4958663>

View Table of Contents: <http://aip.scitation.org/toc/rsi/87/7>

Published by the [American Institute of Physics](#)

---

---

## Note: Fully integrated 3.2 Gbps quantum random number generator with real-time extraction

Xiao-Guang Zhang,<sup>1,2</sup> You-Qi Nie,<sup>1,2</sup> Hongyi Zhou,<sup>3</sup> Hao Liang,<sup>1,2</sup> Xiongfeng Ma,<sup>3</sup>  
 Jun Zhang,<sup>1,2,a)</sup> and Jian-Wei Pan<sup>1,2</sup>

<sup>1</sup>*Hefei National Laboratory for Physical Sciences at the Microscale and Department of Modern Physics, University of Science and Technology of China, Hefei, Anhui 230026, China*

<sup>2</sup>*CAS Center for Excellence and Synergetic Innovation Center in Quantum Information and Quantum Physics, University of Science and Technology of China, Hefei, Anhui 230026, China*

<sup>3</sup>*Center for Quantum Information, Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing 100084, China*

(Received 24 March 2016; accepted 29 June 2016; published online 11 July 2016)

We present a real-time and fully integrated quantum random number generator (QRNG) by measuring laser phase fluctuations. The QRNG scheme based on laser phase fluctuations is featured for its capability of generating ultra-high-speed random numbers. However, the speed bottleneck of a practical QRNG lies on the limited speed of randomness extraction. To close the gap between the fast randomness generation and the slow post-processing, we propose a pipeline extraction algorithm based on Toeplitz matrix hashing and implement it in a high-speed field-programmable gate array. Further, all the QRNG components are integrated into a module, including a compact and actively stabilized interferometer, high-speed data acquisition, and real-time data post-processing and transmission. The final generation rate of the QRNG module with real-time extraction can reach 3.2 Gbps. *Published by AIP Publishing.* [<http://dx.doi.org/10.1063/1.4958663>]

Random numbers are required in many applications such as numerical simulations, cryptography, and even lotteries. Quantum random number generators (QRNGs),<sup>1</sup> exploiting the basic principles of quantum physics, can produce true random numbers which are unpredictable, irreproducible, and unbiased. So far, various QRNG schemes have been proposed and experimentally demonstrated.<sup>2–22</sup> These QRNG schemes can be simply sorted into three categories. The first one is the beam splitter (BS) scheme by measuring the path selection when a single photon passes through a beam splitter.<sup>2–4</sup> In such scheme, one bit at most is generated per photon detection and QRNG speed is limited by the count rates of single-photon detectors. The second one is time measurement scheme by measuring and digitizing photon arrival times,<sup>5–8,11,18</sup> and QRNG speed in this scheme reaches roughly 100 Mbps.<sup>7,11,18</sup> The third one is quantum fluctuation scheme by measuring vacuum states<sup>10,12–14,16</sup> or measuring laser phase fluctuations,<sup>9,17,19,22</sup> in which classical photodetectors are used instead of single-photon detectors and thus the generation rate can be greatly increased up to Gbps.

In the scheme of laser phase fluctuations,<sup>9,17,22</sup> the randomness originates from laser spontaneous emission. Given a laser operated around its threshold level, the contribution ratio between spontaneous and stimulated emissions can be pretty high so that quantum noise dominates the phase fluctuations. Further, phase fluctuations can be converted into intensity fluctuations using an interferometer, which can be measured by a fast photodetector. The output signals of the photodetector are digitized to generate raw

random data. To generate final random numbers the min-entropy of the raw random data is evaluated and the bias is removed with randomness extraction. Owing to the fast photodetector and the high-speed digitizer as well, the final random bit rate can be extremely high. For instance, a record QRNG speed of 68 Gbps has been recently reported based on the scheme of laser phase fluctuations.<sup>22</sup>

However, we remark that such high-speed QRNGs have their limitations for practical use.<sup>23</sup> In the previous experiments,<sup>9,17,22</sup> the photodetector outputs were digitized by oscilloscopes, and the sampled raw random numbers were temporarily stored in the limited memories of the oscilloscopes, which were then post-processed offline. In addition, compact and integrated QRNG modules are highly required for practical applications.

To close the gap between experimental demonstration and practical use, here we report a real-time, fully integrated and standalone instrument of QRNG with a generation rate of 3.2 Gbps based on the scheme of laser phase fluctuations. The random bitstream is transmitted via a small form-factor pluggable (SFP) optical transceiver. This real-time bit rate is higher than the fastest commercially available QRNG product with over 20 times.

The design diagram of the QRNG module is described in Fig. 1(a). The key components of the QRNG module include a stable interferometer that is assembled inside a compact box and two printed circuit boards (PCBs). One PCB includes a laser diode driver with temperature control and a circuit for phase stabilization, whilst the other PCB is designed for raw data acquisition, real-time post-processing, and final random data transmission. All the components are integrated into a metal box with a size of 304 mm × 250 mm × 78 mm as shown in Fig. 1(b).

<sup>a)</sup>Electronic mail: zhangjun@ustc.edu.cn

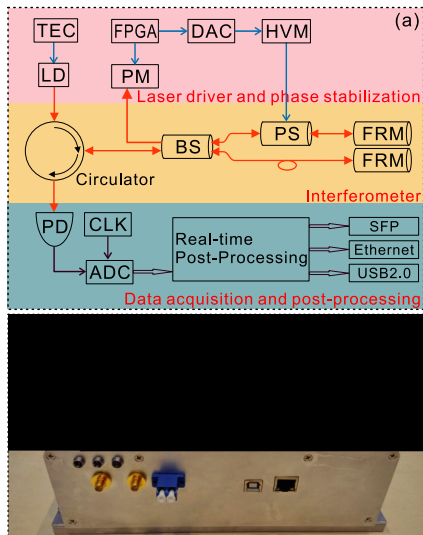


FIG. 1. Design diagram (a) and photo (b) of the QRNG module.

A 1550 nm laser diode (LD) is driven by constant current that is slightly above its threshold, and a thermoelectric cooler (TEC) is employed to stabilize the temperature of LD. The emitted continuous wavelength photons enter an unbalanced interferometer via an optical circulator. One port of the circulator is connected with a 50/50 beam splitter (BS), whose output ports are connected with two Faraday rotator mirrors (FRMs) to construct a polarization-insensitive Michelson interferometer. A phase shifter (PS) is inserted in one arm of the interferometer, and the time difference between the two arms is around 0.8 ns which is much smaller than the coherence time of LD. One output port of the interferometer is detected by a 9.5 GHz InGaAs photodetector (PD) after passing through the circulator, whilst the other output port is monitored by a power meter (PM). The PM is implemented using another PD to measure the optical power and an analog-to-digital converter (ADC) to digitize the measured signal. A field-programmable gate array (FPGA) reads out the PM data, and sends feedback data to a digital-to-analog converter (DAC) to form a voltage signal after the computation by a proportional-integral-derivative (PID) algorithm. The voltage signal regulates the output of a high-voltage module (HVM), which results in automatic adjustments to the PS. Due to real-time and fast responses of the PID algorithm, the interferometer can be highly stable, which guarantees continuous operations of the QRNG module.

To obtain raw random data, the voltage output of the PD is amplified and then digitized by an 8-bit ADC (TI ADC083000) with a clock (CLK) of 1 GSa/s, and the sampled data are then fed into a high-speed FPGA (Xilinx Virtex-6). A pipeline post-processing algorithm based on Toeplitz hashing matrix is implemented in this FPGA for real-time randomness extraction. The extracted random numbers are transmitted in real-time via the interface of SFP with 3.2 Gbps. For the applications requiring lower bit rates, either the Gbps Ethernet port or the universal serial bus (USB) 2.0 port can be used.

To quantify the randomness of the raw random data min-entropy evaluation is applied. The min-entropy is defined

as  $H_{min}(X) = -\log_2(\max_{x \in \{0,1\}^N} P_r[X = x])$ , which can be exploited to quantify the extraction ratio between the raw random bits and the final random bits given a probability distribution of  $\{0,1\}^N$ .<sup>17,24</sup> Detailed modeling and analysis of min-entropy for the scheme of laser phase fluctuations can be found in the literatures.<sup>9,17,22,23</sup> The key parameter for the evaluation is the ratio of the quantum phase fluctuations to classical noise ( $\gamma$ ). For our QRNG module, we follow the same experimental approach to measure  $\gamma$  by tuning the laser power<sup>17,22</sup> and repeat this measurement many times. The value of optimal  $\gamma$  is stable, and a typical measured value is 6.87. The variance of quantum phase fluctuations can be calculated<sup>9,17,22,23</sup> by  $\sigma_q^2 = \frac{\gamma}{\gamma+1} \langle V(t)^2 \rangle$ , where  $\sigma_q$  is the standard deviation of the raw random data, and  $V(t)$  is the voltage output amplitude of PD. In our QRNG module, the measured intensity variance is around 8311  $mV^2$ , therefore,  $\sigma_q$  is 85.2 mV. Then, one can calculate that the maximum probability in the whole distribution is 0.011 since the quantum signal follows a Gaussian distribution.<sup>24</sup>  $H_{min}(X)$  is thus calculated as 6.5 bits per sample or 0.8 bits per raw bit, which means that 6.5 random bits can be generated from each sample.

Since the capability limit of real-time post-processing in FPGA is 5 Gbps, 3 bits are abandoned for each sample including the least significant bit and the left two bits. Considering the worst case there are still 3.5 bits that can be extracted from each sample, which corresponds to a new min-entropy of 0.7 bits per raw bit.

After the min-entropy evaluation, a Toeplitz hashing extractor is applied to distill the raw random data, and the rigorous discussion about Toeplitz hashing extractor can be found in Ref. 24. Given a binary Toeplitz matrix with a size of  $m \times n$ ,  $m$  final random bits are extracted by multiplying the matrix and  $n$  raw bits. The Toeplitz matrix can be constructed by a sequence of  $m + n - 1$  random bits called matrix building seeds, due to its characteristic that all the elements of each descending diagonal from left to right are the same.

For our implementation of Toeplitz hashing extraction, we choose  $m = 1024$  and  $n = 1520$  so that the extraction ratio is  $m/n = 0.67$ . According to the leftover hash lemma,<sup>25</sup>  $m = nH_{min}(X) - 2 \log_2(\frac{1}{\epsilon})$ , one can calculate the information theoretic security bound  $\epsilon = 2^{-20}$ , that is, the statistical distance between the extracted random sequence and the uniform sequence is bounded by  $\epsilon = 2^{-20}$ . The extraction efficiency is  $m/[nH_{min}(X)] = 0.96$ .

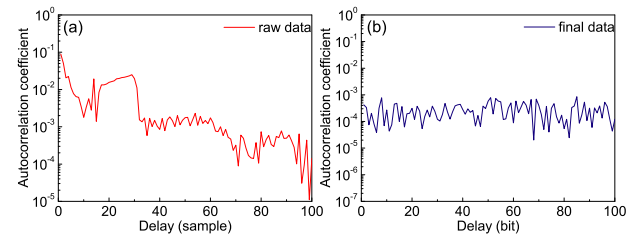
In the QRNG module, the Toeplitz hashing randomness extraction is implemented in FPGA. By taking advantage of concurrent computation capability in FPGA, the real-time post-processing speed can be extremely improved compared with the software implementation in computer.<sup>22</sup> Due to the resource limit in FPGA, it is impossible to directly compute such a large matrix rapidly. Therefore, we propose a concurrent pipeline algorithm to achieve the real-time extraction. In particular, in order to implement real-time randomness extraction, we design three modules including matrix building, submatrix multiplication and vector accumulation in FPGA, as shown in Fig. 2. These modules are operated in a pipeline mode with a synchronized clock of 62.5 MHz. In the matrix building module, the  $m + n - 1$  ( $=2543$ ) bit seeds are used to construct



FIG. 2. FPGA implementation of Toeplitz hashing randomness extraction.

the complete Toeplitz matrix. These seeds can be refreshed periodically. The whole Toeplitz matrix is divided into  $n/k$  submatrices. Considering the resource consumption in FPGA, the value of  $k$  is selected to be 80. In each clock period, a submatrix with a size of  $m \times k$  and  $k$  raw bits is multiplied to output a temporary column vector in the submatrix multiplication module. After  $n/k$  clocks, all the calculations for  $n$  raw bits are finished. The  $n/k$  temporary column vectors are then added in the vector accumulation module and  $m$  final random bits are generated. We note that binary multiplications and binary additions can be realized by bitwise AND and bitwise XOR operations, respectively. With such a configuration, real-time Toeplitz hashing randomness extraction is achieved in the FPGA, and the generation rate limit of final random bits can reach  $5 \text{ Gbps} \times 1024/1520 \sim 3.36 \text{ Gbps}$ . The final generation rate in FPGA can be tuned to match the interfaces for the real-time transmission. With a SFP optical transceiver, the real-time rate of the final random numbers reaches 3.2 Gbps. Apart from the SFP, optional interfaces including Gbps Ethernet port and USB 2.0 port are also designed for the applications requiring lower random bit rates, whose average speeds are tested as 968.7 Mbps and 259.5 Mbps, respectively.

To test the randomness of the final data, we perform an autocorrelation comparison between the raw data and the extracted random data, as shown in Fig. 3. The relatively large autocorrelation existing in the raw data is mainly due to the short sampling interval.<sup>23</sup> For each sample, the interferometer output can be regarded as an interference from two certain points in the same beam with a time difference of  $\tau$ , where  $\tau$  stands for the time delay between the two arms of the interferometer. The two points form a wave interval. When the sampling rate is high enough, one interval overlaps with adjacent intervals, which results in a large autocorrelation. Fig. 3 shows that after post-processing this autocorrelation can be significantly reduced. For the randomness quantification, the theory given in Ref. 23 can be applied, which aims at a high sampling rate situation. Finally, the standard NIST statistical

FIG. 3. Autocorrelation analysis of  $10^7$  raw samples (a) and  $10^7$  extracted random bits (b). The autocorrelation existing in the raw data is completely eliminated by the Toeplitz hashing extraction.

tests are applied to test the randomness of final data. Typically, 10 final random data files with each file size of 1 Gbits are tested, and all the files can well pass the test items. We note that given the items that produce multiple outcomes the  $p$ -values are processed by a Kolmogorov-Smirnov uniformity test and the proportions are averaged.

In summary, we have developed a real-time and fully integrated QRNG module based on the scheme of laser phase fluctuations. We propose and implement a pipeline post-processing algorithm based on Toeplitz hashing randomness extraction in FPGA, which can extremely increase the real-time generation rate of final random numbers up to 3.2 Gbps.

This work has been financially supported by the National Basic Research Program of China Grant No. 2013CB336800, the National Natural Science Foundation of China Grant No. 61275121, and the Chinese Academy of Sciences. X.-G. Zhang and Y.-Q. Nie contributed equally to this work.

<sup>1</sup>X. Ma *et al.*, e-print [arXiv:1510.08957](https://arxiv.org/abs/1510.08957) (2015).<sup>2</sup>J. Rarity *et al.*, *J. Mod. Opt.* **41**, 2435 (1994).<sup>3</sup>A. Stefanov *et al.*, *J. Mod. Opt.* **47**, 595 (2000).<sup>4</sup>T. Jennewein *et al.*, *Rev. Sci. Instrum.* **71**, 1675 (2000).<sup>5</sup>H.-Q. Ma *et al.*, *Appl. Optics* **44**, 7760 (2005).<sup>6</sup>J. Dynes *et al.*, *Appl. Phys. Lett.* **93**, 031109 (2008).<sup>7</sup>M. A. Wayne and P. G. Kwiat, *Opt. Express* **18**, 9351 (2010).<sup>8</sup>M. Fürst *et al.*, *Opt. Express* **18**, 13029 (2010).<sup>9</sup>B. Qi *et al.*, *Opt. Lett.* **35**, 312 (2010).<sup>10</sup>C. Gabriel *et al.*, *Nat. Photonics* **4**, 711 (2010).<sup>11</sup>M. Wahl *et al.*, *Appl. Phys. Lett.* **98**, 171105 (2011).<sup>12</sup>T. Symul *et al.*, *Appl. Phys. Lett.* **98**, 231103 (2011).<sup>13</sup>M. Jofre *et al.*, *Opt. Express* **19**, 20665 (2011).<sup>14</sup>P. J. Bustard *et al.*, *Opt. Express* **19**, 25173 (2011).<sup>15</sup>Y. Jian *et al.*, *Rev. Sci. Instrum.* **82**, 073109 (2011).<sup>16</sup>A. Marandi *et al.*, *Opt. Express* **20**, 19322 (2012).<sup>17</sup>F. Xu *et al.*, *Opt. Express* **20**, 12366 (2012).<sup>18</sup>Y.-Q. Nie *et al.*, *Appl. Phys. Lett.* **104**, 051110 (2014).<sup>19</sup>Z. L. Yuan *et al.*, *Appl. Phys. Lett.* **104**, 261112 (2014).<sup>20</sup>B. Sanguinetti *et al.*, *Phys. Rev. X* **4**, 031056 (2014).<sup>21</sup>Q. Yan *et al.*, *Rev. Sci. Instrum.* **85**, 103116 (2014).<sup>22</sup>Y.-Q. Nie *et al.*, *Rev. Sci. Instrum.* **86**, 063105 (2015).<sup>23</sup>H. Zhou *et al.*, *Phys. Rev. A* **91**, 062316 (2015).<sup>24</sup>X. Ma *et al.*, *Phys. Rev. A* **87**, 062327 (2013).<sup>25</sup>R. Impagliazzo *et al.*, in *Proceedings of the Twenty-First Annual ACM Symposium on Theory of Computing* (ACM, 1989).