# Jitter-based Adaptive True Random Number Generation for FPGAs in the Cloud

Xiang Li, Peter Stanwicks, George Provelengios, Russell Tessier, and Daniel Holcomb Department of Electrical and Computer Engineering University of Massachusetts, Amherst, MA, USA

Abstract—In this paper we present and evaluate a true random number generator (TRNG) design that is compatible with the restrictions imposed by cloud-based FPGA providers such as Amazon Web Services (AWS) EC2 F1. Because cloud FPGA providers disallow the ring oscillator circuits that conventionally generate TRNG entropy, our design is oscillator-free and uses clock jitter as its entropy source. The clock jitter is harvested with a time-to-digital converter (TDC) and a controllable delay line that is continuously tuned to compensate for process, voltage, and temperature variations. After describing the design, we present and validate a stochastic model that conservatively quantifies its worst-case entropy. We deploy and model the design in the cloud on 60 EC2 F1 FPGA instances to ensure sufficient randomness is captured. TRNG entropy is further validated using NIST test suites, and experiments are performed to understand how the TRNG responds to on-die power attacks that disturb the FPGA supply voltage in the vicinity of the TRNG.

#### I. Introduction

Random numbers are fundamental to cryptographic systems and widely used for generating keys, nonces, and initialization vectors. The quality of randomness required in these applications necessitates the use of true-random number generators (TRNGs). TRNGs exploit the inherent physical properties of the system in which they are embedded to generate statistically random and unpredictable numbers. This characteristic makes the outputs of a TRNG unpredictable even to an adversary that knows the current state of the circuit. Physical sources of entropy commonly used in FPGAs by on-chip TRNGs include thermal noise and clock or oscillator jitter. The randomness of the numbers created by TRNGs is typically evaluated using stochastic models and statistical tests [1].

FPGAs are increasingly being used in cloud-based systems for prototyping and acceleration, and to support secure soft processors which require a source of random numbers. Yet to protect their infrastructure from malicious voltage attacks [2], cloud providers such as AWS impose restrictions on the types of circuits that are allowed on their FPGAs. Circuits that deviate from standard digital design flows, including logic-driven clocks and combinational loops as found in ring oscillators (ROs), are detected during bitstream compilation and disallowed from being loaded onto the FPGA [3]. This restriction causes difficulty in creating and characterizing jitter-based TRNG circuits for cloud applications.

In this work, we present a TRNG design and validation procedure that is tailored around the restrictions of cloud-based FPGAs. Our design is able to harvest jitter without creating oscillators or being able to manipulate clocks. The design adjusts to changing environmental conditions and can be characterized

without requiring ground truth delay measurements that are commonly obtained by counting oscillations. We make several specific contributions in this work:

- A TRNG for Virtex UltraScale+ FPGAs used in the cloud is detailed, implemented, and analyzed across numerous AWS EC2 F1 instances. The design, which is based on tunable delay chains and a time-to-digital converter (TDC) that harvests entropy from clock jitter, avoids primitives such as combinational loops that are common in TRNGs but disallowed by AWS and other cloud providers.
- A novel procedure is proposed for computing the minentropy per sample using a stochastic model. The model empirically relates component delays to clock jitter by least-squares fitting. The delays computed by the model during entropy evaluation strongly correlate to the FPGA's own timing report, which supports the validity of our approach.
- The robustness of our TRNG is evaluated by implementing a voltage attack against it on F1 and showing how the TRNG adjusts in response to the attack without compromising its ability to create random numbers. Demonstrating resilience to environmental changes is important for a TRNG that will be used in cloud settings, and is a novel feature of the work.

The remainder of this paper is structured as follows. Section II provides background on previous FPGA TRNG approaches. Section III describes the structure of our TRNG and modeling is discussed in Section IV. Sections V and VI evaluate and discuss the TRNG entropy, resilience, and costs. Section VII concludes the paper.

# II. BACKGROUND AND RELATED WORK

The implementation of TRNGs in FPGAs has been widely studied, although none of the prior approaches address the unique constraints of cloud FPGAs. A large majority of these previous implementations rely on ROs to generate high-frequency signals that exhibit significant jitter. For example, Kohlbrenner and Gaj [4] use two ROs and a sampling circuit to measure jitter and Maiti et al. [5] deploy up to 128 ROs to amplify uncertainty. Some TRNGs augment ROs with delay paths to increase timing sensitivity. Like our approach, Rozic et al. [6] and Yang et al. [7] use carry logic-based delay chains to assist with entropy extraction. These approaches do not include tunable delays to combat environmental factors and an RO is used to excite the delays.

Several non-RO based TRNGs have been built for FPGAs, but they also have limitations that make them inappropriate

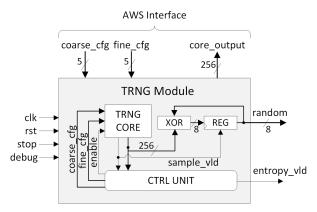


Fig. 1: Structure of TRNG design and interface.

for cloud deployment. Majzoobi et al. [8] use programmable delay lines built from LUTs that can be difficult to characterize on a per-FPGA basis. Deák et al. [9] use the jitter from an on-FPGA phase locked loop (PLL) to create a TRNG using clock settings that are a challenge to replicate in a cloud setting. Perhaps the most similar TRNG approach to ours [10] uses a standard clock input, tunable delay buffers, and a delay path. However, unlike our approach, the delay path is made from a chain of LUTs which have variable logic and routing delays across stages.

#### III. STRUCTURE OF PROPOSED TRNG

Fig. 1 shows the top-level view of our TRNG design. It is a hardware module that serially generates 8-bit random numbers. We instantiate the TRNG module within a hardware testbench for analysis. The design is created for AWS EC2 F1 instances, which contain Xilinx Virtex UltraScale+ VU9P FPGAs. The bitstream is generated using Amazon's Hardware Development Kit (HDK), and then converted to an Amazon FPGA Image (AFI) that is reused for deployment across F1 instances. Amazon provides a runtime tool to interact with the deployed design by reading and writing 32-bit data to or from user-defined registers using AXI4 over PCIe. We make the signals at the top of Fig. 1 accessible to the runtime tool only when the TRNG hardware is set to debug mode. The debug mode allows us to control the TRNG and observe sample values from the TRNG core, which is useful for data collection and analysis in the cloud, but would be insecure if enabled in production.

Internally, the TRNG module gets entropy from the TRNG core, and hashes it into a local entropy pool by XOR operation. The control unit keeps a conservative estimate of the current entropy in the pool by counting the number of valid samples provided to it. Once enough entropy is collected, a signal from the control unit is asserted and the 8-bit random value in the registers can be read out, upon which the entropy count is reset to 0. We now describe in more detail the components of the TRNG core (Fig. 2).

# A. Carry Chain Description

The TDC in our circuit (Fig. 2a) consists of 32 8-bit carry stages, and each output bit from the carry chain is the data input

to a D flip-flop in the same slice. The controller repeatedly generates a single rising edge that propagates into the carry chain with appropriate delay such that it will be propagating through the carry chain when the next rising clock edge occurs. The number of 1-values captured in the 256 flip-flops, i.e. the Hamming weight of the sample, is an indication of how far up the chain the rising edge has propagated by the time of the rising clock edge. The Hamming weight of samples will fluctuate slightly in each trial due to clock jitter, which is our source of randomness.

#### B. Tunable Delay Elements and Feedback Control

Our circuit uses tunable delay elements and feedback to ensure that the rising edge from the delay line is within the TDC chain when the clock arrives. Increasing the propagation delay will cause the rising edge to reach fewer TDC stages and thereby will reduce the Hamming weight of samples. Similarly, decreasing the delay will increase the Hamming weight. In this way the tunable delay circuits are the knob used for adjusting the Hamming weight. The controller uses the delay knob to position the rising edge in the TDC chain during clock arrival, as is required to generate randomness.

Ideally, the Hamming weight of the samples should be centered at around 128, which gives a maximum margin against delay changes in either direction that could detune the circuit. The controller configures the coarse-tuning and fine-tuning stages using the simple feedback scheme of Eq. 1, where (c,f) and (c',f') are the current and next values of the coarse and fine tuning settings, and HW is the Hamming weight of the current sample; note in Eq. 1 that  $f_{mid}$  represents the middle setting for fine tuning, which in our case is 15. Furthermore, as the controller monitors samples it credits entropy to the entropy pool only when the Hamming weight is between 30 and 225 so that samples are not counted as random when the rising edge is approaching either end of the carry chain where jitter may not be captured.

$$(c', f') = \begin{cases} (c+1, f_{mid}) & \text{if } 208 \le HW\\ (c, f+1) & \text{if } 158 \le HW < 208\\ (c, f-1) & \text{if } 48 < HW \le 98\\ (c-1, f_{mid}) & \text{if } HW \le 48 \end{cases}$$
(1)

The coarse and fine tuning stages are implemented as follows. Each coarse tuning stage adds or bypasses a LUT1 primitive that implements a logical buffer, as shown in Fig. 2b. Each fine-tuning stage selects between a shorter and longer pin-to-pin delay of a LUT5, where the enabled path through the LUT is set by the control input, as shown in Fig. 2c. The stages are controlled using thermometer encoding, so that incrementing or decrementing their configuration settings will change only one stage along the delay line, which helps ensure predictable control but has higher area cost than a binary-encoded tunable delay in which each stage has twice the delay of the next. Fig. 3 shows the Hamming weight of samples for all combinations of tuning; note that debug mode is used to generate this plot, as it overrides the feedback of the controller, and allows the samples from the TRNG core to be logged.

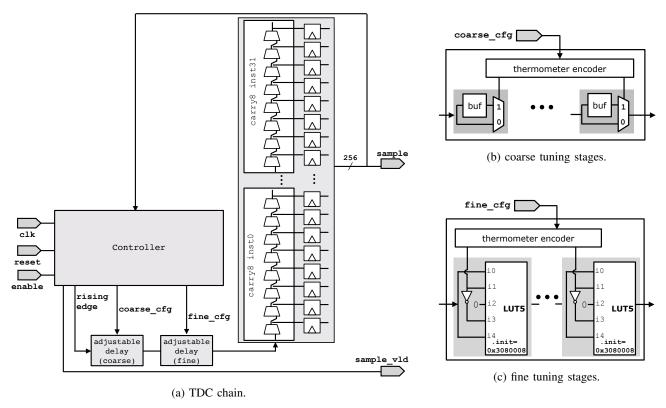


Fig. 2: Components of TRNG core.

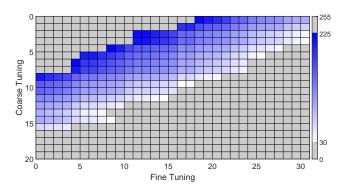


Fig. 3: Heatmap showing the Hamming weight of samples on a single F1 instance for all possible tuning settings. Increasing the coarse or fine tuning reduces the Hamming weight. Sampled values indicative of a poorly-tuned TDC, colored gray, would not cause the entropy count to be incremented.

# C. Post-processing Circuit

The 256-bit samples from the TRNG core are hashed into the 8-bit state of the entropy pool using a simple scheme as shown in Fig. 4. Each update includes a 1-bit circular shift of the 8-bit entropy pool, which ensures that randomness will get distributed through the 8 bits even if always coming from the same position in the 256-bit sample. We have used this particular scheme for simplicity, but it could be replaced with any number of other hash functions for the same effect. A

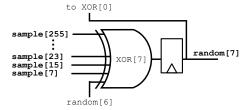


Fig. 4: 256-bit samples are hashed into the 8-bit random signal which is the entropy pool. The hashing uses eight XOR gates, all configured in the same manner as the one that is shown.

counter tracks the number of valid samples produced by the TRNG core, and requires that 80 valid samples are hashed into the entropy pool before its is considered to be random, which assumes 0.1 bits of entropy per sample. Section IV of the paper will show that the actual entropy per sample exceeds this conservative estimate by a factor of more than 2 using both a stochastic model and NIST tests.

#### IV. MODELING OF TRNG

The randomness of the TRNG comes from the samples of the delay line in the TDC. Even if the propagation delay of the delay line does not change across trials, the TDC can produce different samples if it has fine enough time resolution and its sampling clock has sufficient jitter. The time resolution on the TDC is a consequence of the low propagation delay of the stages in the hard carry chains of Xilinx Ultrascale+ devices.

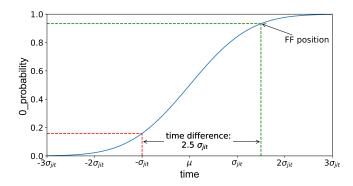


Fig. 5: TDC position in normal distribution CDF for modeling.

A larger clock jitter or finer time resolution will both have the same effect of making the TDC samples more random because both changes will make the jitter relatively larger in comparison to the TDC time resolution. Accordingly, the relevant consideration for modeling a TRNG such as ours is to quantify how the amount of jitter compares to the time resolution of the TDC.

In our stochastic model of the TRNG, we relate jitter to TDC time resolution without relying on conservative timing reports or jitter estimates. We rely in modeling only on the simplifying assumption that jitter is normally distributed, which is in particular consistent with the jitter component caused by thermal noise [11] [12]. We also assume its standard deviation is invariant with respect to the tuning settings of the delay line. From this we calculate the time resolution of each TDC stage in terms of the standard deviation of jitter, which we use as the unit delay in our model. After calculating the time resolution of each stage, we use the same model to calculate a lower bound on min-entropy per sample. The next subsection describes our modeling approach.

# A. Empirical Model Relating TDC Delay to Jitter

In a given trial, the flip-flop associated with each TDC stage will sample a 0 value if its clock arrives before its rising data input from the delay line, and will sample a 1 otherwise. If the delay tuning settings are held constant across samples, the 0-probability (1-probability) of the stage indicates the proportion of trials in which its clock arrives before (after) its rising data input from the delay line.

If the flip-flop of stage i samples 0 with probability 0.159, then 15.9 percent of clock edges arrive there before the rising data input, and 84.1 percent of clock edges arrive after. Under the assumption that jitter is normally distributed, the observation that 15.9 percent of clock edges arrive before the rising data input reveals that rising data input coincides with the clock being  $-1.0*\sigma_{jit}$  away from its mean value, because  $\Phi^{-1}(0.159) = -1.0$ , where  $\Phi^{-1}$  is the inverse CDF of a normal distribution. This scenario is depicted graphically in Fig. 5.

In these same trials, if the flip-flop of stage j samples 0 with probability 0.933, then we similarly conclude that its rising data input coincides with its clock being  $+1.5 * \sigma_{jit}$  away from its mean because  $\Phi^{-1}$  (0.933) = 1.5. If there is no

clock skew between the flip-flops of i and j, then these two findings together indicate that the time difference between the rising data inputs on i and j is equal to  $2.5 * \sigma_{iit}$ . If clock skew is allowed, then we generalize the claim slightly to more formally conclude only the difference in criticality (i.e. timing slack) between the two flops is equal to  $2.5 * \sigma_{iit}$ , although for our purposes it is actually the criticality that matters so we need not worry about skew. The delay or criticality difference between any two stages can therefore be estimated from their 0-probabilities in a set of trials. Because the estimate is noisy when the associated 0-probabilities are close to 0 or close to 1, we apply it only when the 0-probabilities indicate that both stages are within  $\pm 2\sigma_{iit}$  of their means. Note that the delay difference is being calculated in units of  $\sigma_{iit}$  even though the value of  $\sigma_{iit}$  is not known in absolute terms. Using this approach, two flip-flops i and j have arrival times (denoted  $T_i$  and  $T_i$ ) that are related as shown in Eq. 2, where  $\hat{P}_i$  and  $\widehat{P}_i$  are their respective 0-probabilities for a particular tuning

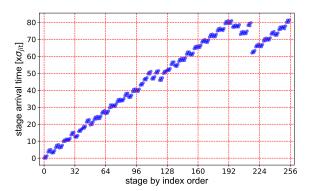
$$T_i - T_j = \left(\Phi^{-1}(\widehat{P}_i) - \Phi^{-1}(\widehat{P}_j)\right)\sigma_{jit} \tag{2}$$

## B. Calculating Stage Arrival Times

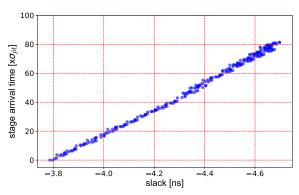
Given that 0-probabilities from each tuning setting will relate the arrival times of some of the TDC stages, and that the same stages can be related to each other by multiple different tunings, we can solve a set of equations to obtain for each stage i the arrival time  $T_i$ . The set of equations is as described in Eq. 3 where T is the n-element column vector of unknown arrival times, A is the m-by-n matrix of coefficients in which all entries are 0 except for a single +1 and single -1 in each row for the two stages that are related, and B is an m-by-1 column vector of the arrival time differences, calculated as shown on the right hand side of Eq. 2. We then find the least-squares solution to AT = B (Eq. 3) which gives stage arrival times in terms of  $\sigma_{iit}$ . Because the formulation deals with differences in arrival times, we adopt the convention that  $T_0$ , the arrival time of the first stage, is 0; other  $T_i$  values therefore represent the arrival time of stage i relative to first stage.

$$AT = B\sigma_{jit} \tag{3}$$

Fig. 6a shows, for one instance of the TRNG, the stage arrival times  $(T_0, \ldots, T_{255})$  obtained by solving Eq. 3. The arrival times across the 256 stages cover a span of approximately 84 times  $\sigma_{iit}$ , which implies  $\sigma_{iit}$  to be around 12ps in absolute terms. While the arrival time generally increases while moving up the TDC chain, note that the trend is not smooth. Although the rising edge does propagate through the carry chains in sequential order, the anomalies in this trend imply that it does not reach the d input of the flip-flops in sequential order. This can occur because the delays between carry chain and flip-flop of each stage are not uniform, and because clock skew at the flip-flop aliases to delay on its data input; both of these artifacts are captured in the timing report so it is instructive to compare the modeled arrival times to the slack from the timing report. Fig 6b shows the same modeled arrival times from Fig. 6a, but now plotted against the reported timing slack for



(a) Modeled arrival time vs stage index.



(b) Modeled arrival time vs slack from timing report.

Fig. 6: TDC characterization.

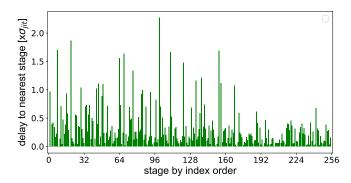


Fig. 7: Precision of each stage in TDC chain from one FPGA instance, obtained from characterization procedure.

the corresponding flip-flop which accounts for all delays and clock skew. There is a high correlation (r=-0.997) between the arrival times from the model and the reported slack. Given that the model was fitted to data measured on the board, the high correlation between the two quantities supports the validity of the model for correctly resolving on-chip delays, and hence also for capturing the difference in criticality between stages. Now that the empirical timing model is validated, we use it as the basis for estimating the worst-case entropy of our TRNG.

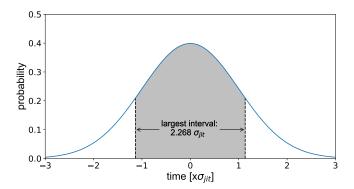


Fig. 8: Largest share of clock arrival times that will cause TDC to sample the same value.

#### C. Stochastic Model for Entropy Estimation

Based on the precision of each stage in Fig. 7, we can estimate a lower bound on min-entropy of the samples. The worst-case min-entropy corresponds to the sampled value that can be produced with highest probability. This would occur when the mean arrival time of the clock coincides with the center of the largest timing gap of any stage, which we denote here as  $\Delta_{max}$ . This is illustrated in Fig. 8 where the shaded region represents all the clock arrival times that would result in the same sample being produced. The probability of producing this sample would then be the probability associated with the shaded region, which we denote as  $P_{max}$ , and calculate from the normal CDF as in Eq. 4. The min-entropy of an outcome with probability  $P_{max}$  is given by Eq. 5. For the specific instance used to generate these results, the largest interval is  $2.268 *\sigma_{jit}$  shown in Fig. 7, which corresponds to a shaded area in Fig. 8 with  $P_{max}$  of 0.743, and hence min-entropy of 0.429 bits.

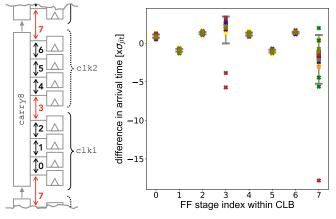
$$P_{max} = \Phi\left(\frac{\Delta_{max}}{2}\right) - \Phi\left(\frac{-\Delta_{max}}{2}\right) \tag{4}$$

$$entropy_{min} = \log_2\left(\frac{1}{P_{max}}\right) \tag{5}$$

#### D. Impact of Routing and Clock Skew on Entropy

The previous subsection explains that worst-case min-entropy is limited by the largest timing gap among all the stages. It is therefore desirable to make all of the timing gaps uniform so that none are unusually large. We now present further results and discussion to explain why routing makes this objective difficult to accomplish in practice.

Fig. 9b shows the difference in arrival time between the FF of each stage and that of the next stage by index. Here, instead of using indices from 0 to 255 to represent the stages across all 32 CLBs as was done in Fig. 6a, we use indices 0-7 for each CLB as annotated in Fig. 9a, and have occurrences of each index from all 32 CLBs. Fig. 9b shows, for each stage index, the 32 differences in arrival time between that stage and the next. The differences in arrival time are predictable for stages 0,1,2 and for stages 4,5,6. In stages 3 and 7, the arrival time difference



(a) CLB stage indices. (b) Differences in arrival times by index.

Fig. 9: Across the 32 CLBs in the TDC, the difference in arrival time between one index and the next is predictable for indices 0,1,2 and 4,5,6. Indices 3 and 7 are each followed by a stage that is on a different clock leaf, and there the difference in arrival time is inconsistent due to clock skew. Error bars extend one standard deviation from the mean.

is inconsistent from CLB to CLB. This inconsistency occurs because Ultrascale+ uses different clock inputs for the upper and lower halves of the CLB, which causes stages 3 and 7 to span two different clock leaf nodes (clk1 and clk2 in Fig. 9a); the skew between the clock leaf nodes aliases to arrival time as discussed in Sec. IV-B. For the TRNG design, one must therefore be careful to avoid large positive clock skew at these points as it can reduce worst-case entropy by causing highly probable outcomes for certain unlucky conditions. Negative skew causes no such problem, as can be observed in Fig. 6, so the one-sided restriction on skew is easy to satisfy in practice.

# V. TRNG QUALITY EVALUATION

In this section, we use three different techniques to test the quality of the random numbers produced by our design. As described in the following three subsections, the results support 1) that our design exceeds the 0.1 bits of min-entropy per trial that was assumed as a security parameter; 2) that our stochastic model gives a reasonable estimate of min-entropy; and 3) the random numbers generated pass tests for statistical randomness.

## A. Stochastic Model Applied Across EC2 F1 Instances

The stochastic model from Section IV-C is our primary strategy for estimating the worst-case min-entropy for each single instance of the TRNG. To test across FPGAs, we load the same bitstream onto 60 different EC2 F1 instances, and on each machine apply our characterization procedure to evaluate the worst-case min-entropy. The distribution of calculated minentropy values (Fig. 10) range from 0.250 to 0.972. These values indicate that across all 60 instances our design exceeds, by at least a margin of  $2.5\times$ , the 0.1 bits of min-entropy per sample that was assumed.

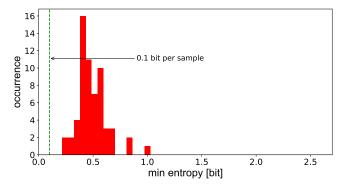


Fig. 10: Entropy from stochastic model on 60 FPGAs.

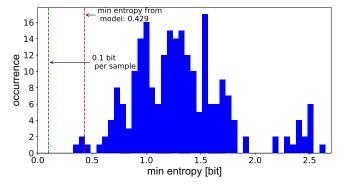


Fig. 11: Entropy per NIST SP800-90B suite for 234 different tunings on one instance.

# B. Stochastic Model vs. NIST Entropy Assessment

Next, to check our stochastic model, we apply the NIST SP800-90B entropy assessment suite [13] to obtain an independently calculated estimate of min-entropy. To generate data for the NIST assessment, we apply on one instance all tuning settings and collect 1,000,000 samples from the TDC with each setting applied [14]. We keep the data from any settings in which the average Hamming weight of samples is in our allowed range of 30 to 225, which corresponds to a total of 234 tuning settings. The NIST assessment is applied separately to each of these 234 datasets, and the distribution of results are shown as a histogram in Fig. 11. The NIST estimate of min-entropy for most of the tuning settings fall above the estimate of 0.429 bits per sample from the stochastic model on this instance. Because the NIST entropy values tend to exceed our estimated worst-case, we gain some confidence that our stochastic model is not overestimating entropy.

## C. End-to-End NIST Statistical Tests

Although the evaluation of the entropy source in the prior subsections is the primary validation for a TRNG, we also apply statistical tests to the post-processed 8-bit values produced by the TRNG as a further validation. The NIST Statistical Test Suite [15], which is widely used with random number generators, applies a collection of statistical tests and for each test reports whether the sequences of bits are consistent with being random. The report shows how often the P-values from

Statistical tests	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	P-value	%
Frequency	10	12	6	12	11	4	11	14	4	16	0.091	98
BlockFrequency	15	17	8	7	7	7	12	9	6	12	0.163	99
CumulativeSums	13	6	10	7	10	15	13	9	9	8	0.596	99
CumulativeSums	11	8	13	8	12	7	15	8	11	7	0.637	98
Runs	9	8	8	8	16	6	15	15	7	8	0.172	98
LongestRun	9	12	8	11	16	6	9	10	11	8	0.657	98
Rank	9	12	13	10	12	5	11	9	13	6	0.637	100
FFT	9	12	12	15	8	5	11	10	12	6	0.494	100
NonOverlap. Template	8	4	7	11	13	9	15	13	10	10	0.401	100
Overlapping Template	8	9	11	8	8	15	12	9	8	12	0.817	98
Universal	8	8	8	12	14	12	14	6	8	10	0.616	99
Approximate Entropy	18	10	12	6	8	14	6	12	6	8	0.109	99
Serial	5	10	8	9	11	12	10	9	16	10	0.616	99
Serial	6	8	10	11	6	13	11	10	12	13	0.740	100
LinearComplexity	13	8	4	14	13	7	15	9	8	9	0.249	100

TABLE I: Results from applying NIST statistical test suite to 100M generated bits shows that the TRNG outputs are evaluated as consistent with being random.

each test fall within uniformly sized bins C1 through C10, and should tend toward being uniformly distributed when enough random data is tested. The test suite is applied to a dataset comprising 100 sequences of 1,000,000 bits from the TRNG and the results are displayed in Tab. I. The final column of the table shows the proportion of sequences that pass the test, indicating that the sequences have statistical properties consistent with being random.

## VI. TRNG PERFORMANCE AND COST EVALUATIONS

Aside from requirement of avoiding circuits such as oscillators that are disallowed in certain clouds, the large capacity of cloud FPGAs implies that the TRNG must also be resilient to any noise, voltage, or temperature fluctuations that are caused by high-powered circuitry around the TRNG.

# A. Resilience to Environmental Fluctuations

We subject the TRNG design to intentional environmental disruptions to check that its feedback is able to adapt appropriately. Specifically, we build a configurable power consumption circuit that is next to the TRNG on the F1 instance. The power waster consists of 32 different levels of power consumption that can be enabled. Each level turns on one instance of a circuit comprising four combinational rounds of the Advanced Encryption Standard (AES) block cipher, with additional feedforward paths added to increase glitching [16]. The power consumption of the circuit is measured as the average power reported by the fpga-describe-local-image command provided in the AWS management tools. Because the reported power updates only once per minute, we perform separate experiments to characterize the consumption of the power wasters instead of measuring their power in real-time when using them to disturb the TRNG. The baseline power consumption of the instance is 8W, and each enabled level of power waster consumes an additional 3W. Turning on the power wasters can disrupt the TRNG by causing heating and voltage droop.

Fig. 12a shows how the TRNG adapts when 56W of power consumption is enabled after the first 1,000 samples are

collected. The blue line shows the Hamming weight of the samples when feedback is disabled, and the orange line shows the Hamming weight when feedback is enabled. The 56W of power draw causes the Hamming weight to have a quick drop, presumably due to voltage droop, followed by a slower decline as the circuit heats up. Both voltage droop and increasing temperature increase propagation delay between the controller and the TDC, which can explain the drop in Hamming weight. The feedback allows the TRNG to compensate for this.

Fig. 12b shows a similar experiment, except now the power wasting circuit is toggled on and off every 1,000 samples, and each time it is switched on an additional five of the 32 power waster instances are enabled, which corresponds to around 15W of additional power consumption. When the feedback is disabled, we can again see by the Hamming weight that the magnitude of power consumption has a direct relation on the delay of the circuit. As before, the controller uses feedback to adapt, and is able to keep the TRNG tuned and operating correctly.

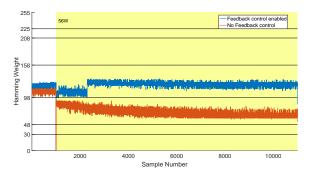
For an end-to-end validation of the TRNG under disturbance from power wasters, we repeat the analysis of Section V-C. As before, the NIST test suite is applied to 100 sequences of 1,000,000 bits, and now 32 power waster instances are running during the data collection. The power wasting circuitry toggles between on and off with every 1,000,000 TRNG bits collected. Similar to Table I, the TRNG again passes the tests, which indicates that these environmental fluctuations are not observed to compromise the TRNG quality.

#### B. Comparison to Prior Work

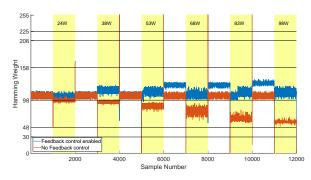
The distinguishing feature of our work is its suitability for, and deployment on, cloud FPGAs. As we've described, this imposes limitations on the types of circuitry that can be used, and increases the importance of the TRNG being robust to environmental changes. Despite these challenges, the costs of our TRNG are found to be reasonable for a large cloud FPGA. Table II compares the throughput and logic utilization of our TRNG to other recently published TRNGs that are implemented on Xilinx FPGAs. Our TRNG design (Fig. 1) consumes 791 LUTs (0.067% of available), 33 CARRY8s, and 559 flip-flops (0.024%) across a total of 184 slices. Among these resources, the controller logic that configures the coarseand fine-tuning consumes 92 LUTs and 34 flip-flops, while the remainder of the resources are consumed by the TRNG core itself. Our design generates random numbers at a rate of 2.43Mbps, which is sufficient for most applications, but could be increased through parallelization if needed.

# VII. CONCLUSION

Cloud FPGAs are commonly used for accelerating computationally expensive cryptographic operations that rely on the generation of random numbers. In this paper, we introduced and evaluated a true random number generator (TRNG) design that is compatible with the design restrictions imposed by cloud-based FPGA providers. The TRNG oscillator-free design that we impose uses a controllable delay and harvests clock jitter as an entropy source using a circuit that is similar to a time-to-digital converter. The effectiveness of the design is supported by



(a) Power consumption starting at sample 1,000.



(b) Power consumption toggling with increasing magnitude.

Fig. 12: Control loop adapts to changes in localized power consumption on the FPGA in order to keep the TRNG tuned.

Work	FPGA type	Throughput	Utilization	Approach
[17]	Spartan 6	100 Mbps	46 slices	self-timed ring
[9]	Spartan 6	14.3 Mbps	67 slices	ring oscillator
[7]	Spartan 6	1.15 Mbps	3 slices	ring oscillator
[18]	Virtex-4	12.5 Mbps	580 slices	RS latches metastability
[19]	Virtex-6	50 Mbps	224 slices	timing nonuniformity
[8]	Virtex-5	2 Mbps	32 slices	metastability
[20]	Spartan 6	3.3 Mbps	27 slices	ring oscillator
This work	Virtex UltraScale+	2.43 Mbps	184 slices	clock jitter

TABLE II: Comparison with related TRNGs implemented on Xilinx FPGAs.

NIST test results and a stochastic model of the entropy source. Furthermore, the design is shown to be able to compensate for voltage droop that may occur during a power attack, and its entropy is not compromised in this scenario. Future work can consider further increases in entropy-per-sample and the impact of advanced clocking features.

## ACKNOWLEDGEMENT

This research was funded in part by NSF grants CNS-1749845 and CNS-1902532. The authors also gratefully acknowledge Gradient Technologies and Amazon for their support and contributions to this work.

#### REFERENCES

- [1] J. Soto, "Statistical testing of random number generators," in *the 22nd National Information Systems Security Conference*, vol. 10, no. 99. NIST Gaithersburg, MD, 1999, p. 12.
- [2] G. Provelengios, D. Holcomb, and R. Tessier, "Characterizing power distribution attacks in multi-user FPGA environments," in 2019 29th International Conference on Field Programmable Logic and Applications (FPL), 2019, pp. 194–201.
- [3] I. Giechaskiel, K. B. Rasmussen, and J. Szefer, "Measuring long wire leakage with ring oscillators in cloud FPGAs," in 2019 29th International Conference on Field Programmable Logic and Applications (FPL). IEEE, 2019, pp. 45–50.
- [4] P. Kohlbrenner and K. Gaj, "An embedded true random number generator for FPGAs," in *Proceedings of the 2004 ACM/SIGDA 12th International Symposium on Field Programmable Gate Arrays.* ACM, 2004, pp. 71–78.
- [5] A. Maiti, R. Nagesh, A. Reddy, and P. Schaumont, "Physical unclonable function and true random number generator: a compact and scalable implementation," in *Proceedings of the 19th ACM Great Lakes Symposium* on VLSI, 2009, pp. 425–428.
- [6] V. Rozic, B. Yang, W. Dehaene, and I. Verbauwhede, "Highly efficient entropy extraction for true random number generators on FPGAs," in 2015 52nd ACM/IEEE Design Automation Conference (DAC). IEEE, 2015, pp. 1–6.
- [7] B. Yang, V. Rožic, M. Grujic, N. Mentens, and I. Verbauwhede, "ES-TRNG: a high-throughput, low-area true random number generator based on edge sampling," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 267–292, 2018.
- [8] M. Majzoobi, F. Koushanfar, and S. Devadas, "FPGA-based true random number generation using circuit metastability with adaptive feedback control," in *International Workshop on Cryptographic Hardware and Embedded Systems*. Springer, 2011, pp. 17–32.
- [9] N. Deák, T. Györfi, K. Márton, L. Vacariu, and O. Cret, "Highly efficient true random number generator in FPGA devices using phase-locked loops," in 2015 20th International Conference on Control Systems and Computer Science. IEEE, 2015, pp. 453–458.
- [10] J.-L. Danger, S. Guilley, and P. Hoogvorst, "High speed true random number generator based on open loop structures in FPGAs," *Microelectronics Journal*, vol. 40, no. 11, pp. 1650–1656, 2009.
- [11] T. J. Yamaguchi, K. Ichiyama, H. X. Hou, and M. Ishida, "A robust method for identifying a deterministic jitter model in a total jitter distribution," in 2009 International Test Conference. IEEE, 2009, pp. 1–10.
- [12] L. Xu, Y. Duan, and D. Chen, "A low cost jitter separation and characterization method," in 2015 IEEE 33rd VLSI Test Symposium (VTS). IEEE, 2015, pp. 1–5.
- [13] M. S. Turan, E. Barker, J. Kelsey, K. A. McKay, M. L. Baish, and M. Boyle, "Recommendation for the entropy sources used for random bit generation," *NIST Special Publication*, vol. 800, no. 90B, 2018.
- [14] C. Celi, "NIST SP800-90B entropy assessment," https://github.com/ usnistgov/SP800-90B\_EntropyAssessment, 2019.
- [15] L. Bassham, Rukhin et al., "A statistical test suite for random and pseudorandom number generators for cryptographic applications," National Institute of Standards and Technology, Tech. Rep., 2010.
- [16] G. Provelengios, D. Holcomb, and R. Tessier, "Power wasting circuits for cloud FPGA attacks," in *International Conference on Field Programmable Logic and Applications (FPL)*, 2020.
- [17] J.-Y. Choe and K.-W. Shin, "A self-timed ring based TRNG with feedback structure for FPGA implementation," in 2020 International Conference on Electronics, Information, and Communication (ICEIC). IEEE, 2020, pp. 1–4.
- [18] H. Hata and S. Ichikawa, "FPGA implementation of ability-based true random number generator," *IEICE Transactions on Information and Systems*, vol. 95, no. 2, pp. 426–436, 2012.
- [19] X. Yang and R. C. Cheung, "A complementary architecture for high-speed true random number generator," in 2014 International Conference on Field-Programmable Technology (FPT). IEEE, 2014, pp. 248–251.
- [20] A. Peetermans, V. Rozic, and I. Verbauwhede, "A highly-portable true random number generator based on coherent sampling," in 2019 29th International Conference on Field Programmable Logic and Applications (FPL). IEEE, 2019, pp. 218–224.