

# *Q*-Class Authentication System for Double Arbiter PUF

Risa YASHIRO<sup>†a)</sup>, Student Member, Takeshi SUGAWARA<sup>†</sup>, Mitsugu IWAMOTO<sup>†</sup>,  
and Kazuo SAKIYAMA<sup>†</sup>, Members

**SUMMARY** Physically Unclonable Function (PUF) is a cryptographic primitive that is based on physical property of each entity or Integrated Circuit (IC) chip. It is expected that PUF be used in security applications such as ID generation and authentication. Some responses from PUF are unreliable, and they are usually discarded. In this paper, we propose a new PUF-based authentication system that exploits information of unreliable responses. In the proposed method, each response is categorized into multiple classes by its unreliability evaluated by feeding the same challenges several times. This authentication system is named *Q*-class authentication, where *Q* is the number of classes. We perform experiments assuming a challenge-response authentication system with a certain threshold of errors. Considering 4-class separation for 4-1 Double Arbiter PUF, it is figured out that the advantage of a legitimate prover against a clone is improved from 24% to 36% in terms of success rate. In other words, it is possible to improve the tolerance of machine-learning attack by using unreliable information that was previously regarded disadvantageous to authentication systems.

**key words:** *Physically Unclonable Function, authentication, machine-learning attack*

## 1. Introduction

Physically Unclonable Function (PUF) has been studied since its concept was proposed in [1], [2]. PUF's output is based on physical information that is generated in manufacturing process. Due to the nature of unclonability, PUF has been increasing attention as a new cryptographic primitive. PUFs are classified into two types; strong PUF and weak PUF. Especially strong PUF, of which challenge space is relatively large, is of great interest as it could be a key component in authentication systems.

Arbiter PUF (APUF) is one of the strong PUFs. It exploits the propagation time difference between two signal paths in a circuit [3]. These two paths are determined by sequential selector pairs that are controlled by challenge bits. The propagation time difference is evaluated using a latch or flip-flop called arbiter. As a result, one-bit output is generated as response. It is experimentally verified that such propagation time difference is unique between Integrated Circuit (IC) chips because of unavoidable manufacturing variation.

An obstacle for APUF is machine-learning (ML) attack [4], [5]. In this attack, an attacker collects challenge-

response pairs (CRPs), and makes a clone of the legitimate PUF. To address this problem, *n*-XOR APUF, in which multiple APUF outputs are XORed, was proposed to improve the attack tolerance in [6].

However, it is reported that even XOR APUF can be cloned by using unreliable responses [7], [8]. The unreliability of responses becomes valid information when modeling APUF. An *n*-XOR APUF's response is generated by APUF's responses, because it is tightly related to the propagation time difference [8]. Moreover, if a response of APUF is unstable, *n*-XOR APUF's response also becomes unstable. In [7], the author applied the divide-and-conquer ML attack to *n*-XOR APUF. More specifically, Covariance Matrix Adaptation Evolution Strategy (CMA-ES) is used to find a better model for each APUF one by one. It is worth mentioning that non-deterministic feature of CMA-ES improves the accuracy of APUF model.

Meanwhile yet another APUF variant called double APUF (DAPUF) was proposed in order to improve uniqueness [9]. DAPUF combines multiple APUFs similarly to *n*-XOR APUF but in a different way. DAPUF is also evaluated under the ML attack. In [10], the authors reported that an attack on DAPUF using Deep Learning (DL) was unsuccessful. However, they also reported that DAPUF has a lot of unstable responses, which makes it unsuitable for conventional authentication system. In summary, DAPUF is advantageous in terms of resistance against ML attack, but is more unstable compared to APUF and *n*-XOR APUF.

We address the above-mentioned disadvantage of DAPUF by proposing a new APUF-based authentication system. That is inspired by a conventional technique for weak PUF [11] in which the unreliable response is used as the third value. Classification to three classes is not necessarily appropriate for strong PUF because there is threat of ML attack. In our proposed method, the number of classes is increased. More specifically, unreliability of response is estimated by feeding the same challenges several times. Then the response is categorized into multiple classes depending on the estimated unreliability. This authentication system is named *Q*-class authentication, where *Q* is the number of classes. The results show that the proposed authentication system using DAPUF increases the tolerance against DL, and strengthens the prover's advantage compared with the conventional 2-class authentication system.

Manuscript received March 22, 2017.

Manuscript revised July 4, 2017.

<sup>†</sup>The authors are with the Department of Informatics, The University of Electro-Communications, Chofu-shi, 182-8585 Japan.

a) E-mail: yashiro@uec.ac.jp

DOI: 10.1587/transfun.E101.A.129

**Table 1** Notation.

Notation	Explanation
$T$	The number of trials
$L$	Challenge-bit length in one trial of authentication
$N$	The number of devices
$b_{t,l}$	The $l$ -th response bit of the $t$ -th trial (or device)
$k$	The number of responses used in authentication system
$U$	Uniqueness value
$P$	Prediction rate
$S$	Steadiness value
$C$	Correctness value
$R$	Randomness value
$M$	Secure-Operation Margin
$Q$	Total class number in authentication
$m_R$	The number to execute PUF in registration phase
$m_V$	The number to execute PUF in verification phase
$r_s$	The number of 1s in the raw responses when the PUF executed $m_R$ ( $m_V$ ) times
$\tilde{c}$	The challenge
$c_l$	The $l$ -th challenge bit of the challenge of length $i$
$\mathcal{D}$	Maximum advantage of legitimate prover against clone
$A^P$	Prover's success rate
$A^C$	Clone's success rate

## 2. Related Work

### 2.1 Quantitative Metrics for PUF

Table 1 shows the notation in this paper. Correctness<sup>†</sup>, uniqueness, and randomness are defined in [12]. Correctness indicates response stability when the same challenge is provided to a certain PUF. Let  $T$  be the number of trials,  $L$  be the number of total challenges, i.e., challenge-bit length in one trial of the authentication experiment, and  $N$  be the number of devices. Correctness is defined with a response bit as

$$C = 1 - \frac{2}{TL} \sum_{t=2}^T \sum_{l=1}^L (b_{1,l} \oplus b_{t,l}). \quad (1)$$

The ideal number of correctness is 1, and if the number of correctness is 0, the PUF returns a random response.

Uniqueness indicates extent of responses difference when inputting the same challenge to each PUF. Uniqueness is expressed as

$$U = \frac{4}{NL} \sum_{t=2}^N \sum_{l=1}^L (b_{1,l} \oplus b_{t,l}). \quad (2)$$

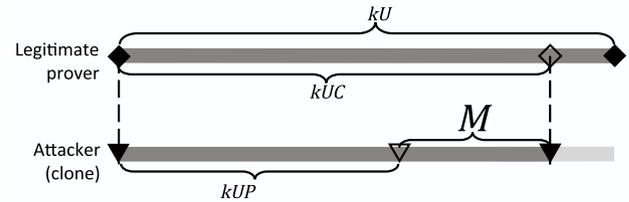
When uniqueness is 0, all the PUFs return the same response.

Finally, randomness is claiming the ratio of 0 and 1 in response bits. Randomness is expressed as the min-entropy of binary probability distribution ( $p, 1 - p$ ) is,

$$R = -\log_2 \max\{p, 1 - p\}, \quad (3)$$

where

<sup>†</sup>It is named reliability in [13].

**Fig. 1** Visual description of Secure-Operation Margin.

$$p = -\frac{1}{TL} \sum_{t=1}^T \sum_{l=1}^L b_{t,l}. \quad (4)$$

#### 2.1.1 Secure-Operation Margin

In [10], the authors proposed a quantitative metric called Secure-Operation Margin (SOM) for APUF-based authentication. This metric indicates the advantage of legitimate prover over clone in the number of bits. SOM is defined as

$$M' = kU(1 - S - P), \quad (5)$$

where  $U$ ,  $S$ , and  $P$  are the uniqueness, the steadiness, and the prediction rate. The prediction rate indicates the accuracy of clones. The steadiness means the stability of responses when the same challenge is repeatedly fed to a PUF. The steadiness is defined as

$$S = \frac{1}{TL} \sum_{t=2}^T \sum_{l=1}^L (b_{1,l} \oplus b_{t,l}). \quad (6)$$

In this paper, we use a slightly modified version of SOM given by

$$M = kU(C - P). \quad (7)$$

In the previous definition in Eq. (5), it is assumed that at least 50% of CRPs are stable. However, the assumption is not satisfied in some practical cases. The obstacle can be avoided in the new definition using  $C$ . SOM expresses an advantage rate of prover in  $k$ -bit authentication system as shown in Fig. 1.  $kU$  represents the number of PUF-specific unique bits found in  $k$  responses. Within  $kU$ , prover stably regenerates  $kUC$  bits, meanwhile clone successfully predicts  $kUP$  bits. Thus,  $M$ , which is given as the difference between  $kUC$  and  $kUP$ , is an advantage of prover over clone.

### 2.2 Double Arbiter PUF

In  $n$ -XOR APUF, an arbiter competes two signals from  $n$ -stage selector pairs. Then,  $n$ -bit responses from the arbiters are XORed to generate one-bit response. As an example, block diagram of 3-XOR APUF with 64 stages is shown in Fig. 2.

In  $n$ -1 DAPUF, propagation time of two signals that propagate through the wires that are the same in terms of physical layout are compared. A one-bit response is generated by XORing the output signals from the arbiters. Figure 2

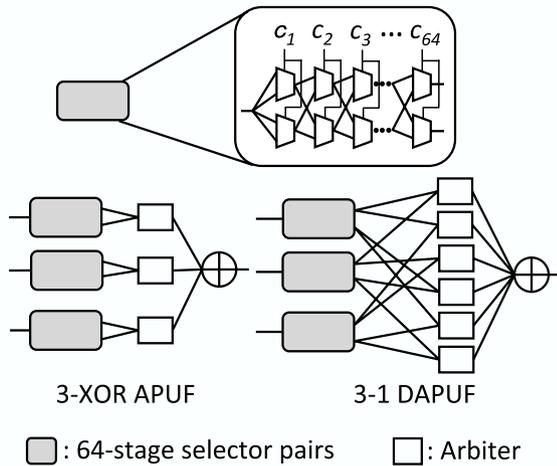


Fig. 2 Block diagrams for 3-XOR APUF and 3-1 DAPUF.

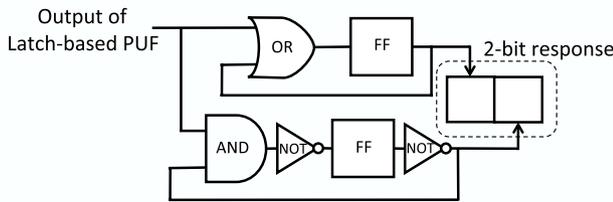


Fig. 3 Schematic view of detection circuit for latch-based PUF.

shows a block diagram of 3-1 DAPUF with 64 stages. DAPUF can have more arbiters compared to  $n$ -XOR APUF, meanwhile the hardware cost of  $n$ -1 DAPUF is almost the same as that of  $n$ -XOR APUF. It should be noted that DAPUF has more unreliable responses than APUF [15]. The unreliability of DAPUF could be a significant problem for the conventional authentication system.

## 2.3 Multivalued Responses of PUF

### 2.3.1 Latch-Based PUF

The authors of [11] proposed a latch-based PUF that outputs multivalued responses. The latch-based PUF has ternary outputs, namely 00, 11, and 10, determined by invoking a raw PUF multiple times. If the responses from the raw PUF are constantly zero or one, then 00 or 11 is used as a response. In the remaining case in which the raw PUF behaves randomly, then the PUF responds with 10. The translation can be achieved using a simple circuit shown in Fig. 3.

The authors showed that the latch-based PUF using ternary is advantageous to conventional ones in which random responses had been either corrected using ECC or simply ignored. It is noted that the latch-based PUF is a weak PUF used for key (ID) generation. However, its application to strong PUF in which ML attack is concerned was remained open.

### 2.3.2 RG-DTM PUF

In [14], response generation according to the delay time measurement (RG-DTM) PUF that outputs multiple-valued responses is proposed. In RG-DTM PUF, delay difference between two paths in APUF is measured more precisely using an array of capacitors. Therefore, multiple-valued responses can be assigned to the amount of delay differences.

The authors claim that uniqueness can be improved by the proposed method. The authors also mention that resistance against ML attack can be improved by using multiple classes. However, the number of classes should be determined before fabrication because delay difference (cf. unreliability) is used for classification.

## 3. Proposed Authentication System

### 3.1 Concept of Q-Class Separation

We proposed a new APUF-based authentication method inspired by [11]. In the following discussion, two types of responses are used. We call the one-bit PUF response as raw response, whereas the data sent from the prover to the verifier during authentication is called prover response. Note that the prover response includes class number as well as prover's ID. In the proposed system, the verifier prepares several classes that are determined by the unreliability of an underlying APUF. Here, the unreliability indicates how much unstable raw responses can be for a fixed challenge. The PUF outputs can be more than ternary depending on the number of classes, i.e., thresholds on the unreliability.

Suppose that we observe 63 PUF outputs for the same challenge. Here, we assume that the output of a strong PUF is 1 bit. Firstly, summation of the 63 outputs is obtained. It is referred to as  $r_s$  hereafter. Then, the CRP is classified by  $r_s$ . For instance, if the verifier prepares 4 classes:

1. Class 1:  $r_s = 0$ ,
2. Class 2:  $1 \leq r_s \leq 32$ ,
3. Class 3:  $33 \leq r_s \leq 62$ ,
4. Class 4:  $r_s = 63$ .

Figure 4 shows an example histogram of  $r_s$  colored by corresponding classes. It is worth noting that the prover is needless to know the classification, i.e., the number of thresholds. Therefore, this authentication system can conceal the accurate unreliability from the attacker.

The proposed method is a certain generalization of the latch-based PUF [11] in the sense that the number of classes could be increased from three to any. The increased number of classes can attribute to improve the security against the ML attacks. The proposed method is advantageous to RG-DTM on the point that it does not need special circuit for measuring delay precisely. Therefore, the proposed method can be combined with the conventional PUFs. Moreover, the number of classes as well as thresholds can be determined after manufacturing.

### 3.2 Q-Class Authentication System

The details of the proposed challenge-response authentication flow are described as follows. Figure 5 illustrates the flow.

1. Several parameters are determined for the verifier: the number of classes  $Q$ , the ranges of each class depending on  $r_s$ , the iteration count of executing PUF  $m_R$  in making CRPs, the number of different challenges required for authentication  $k$ .
2. In the registration phase, the verifier accesses a legitimate PUF to make a table of CRPs and keeps CRPs and ID of the PUF.
3. In the verification phase, the server chooses a master challenge from the stored CRPs, and sends it to the prover (PUF).
4. The prover receives the master challenge, and generates  $k$  challenges from it.

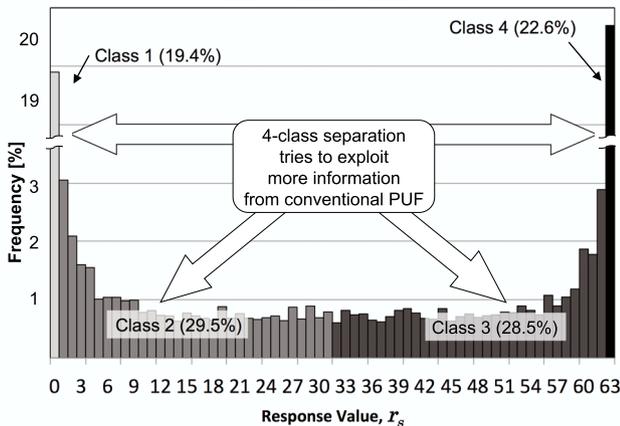


Fig. 4 An example of classification.

5. For each challenge, the prover executes PUF  $m_V$  times to derive the corresponding class by counting the number of 1s in raw responses.
6. The prover sends the prover responses data that contain  $k$  class numbers and ID to the verifier.
7. The verifier compares the prover responses with the stored data in CRPs, and scores the success rate. Then, pass/fail is determined by comparing the success ratio with a threshold.

The verification flow is similar to the conventional PUF-based challenge-response authentication systems except the 5th process in which PUF raw responses are transferred. Accordingly, in the 1st process, the verifier determines the number of classes  $Q$ , the ranges of each class depending on  $r_s$ , and the iteration count of executing PUF  $m_R$  in making CRPs, in addition to the value of  $k$ . In the 5th process, the verifier executes PUF  $m_V$  times to derive the corresponding class by counting the number of 1s in raw response. This process is essential to extend one-bit raw response to the number of  $Q$ .

The ID in the 2nd and 6th steps is an artificial identification number and is not generated by PUF. Verifier uses the ID to select associated CRPs stored in its database and thus makes 1 : 1 (cf. 1 :  $n$ ) authentication.

### 4. Security Evaluations

The security of the proposed authentication scheme is evaluated here. We focus on the authentication operations as described in the previous subsections.

This section is organized as follows. Section 4.1 explains the experimental setup used for all experiments in this paper. Section 4.2 describes preliminary experiments, and Sect. 4.3 explains experiments of  $Q$ -class authentication. Section 4.4 examines the effect of temperature on the accuracy of authentication.

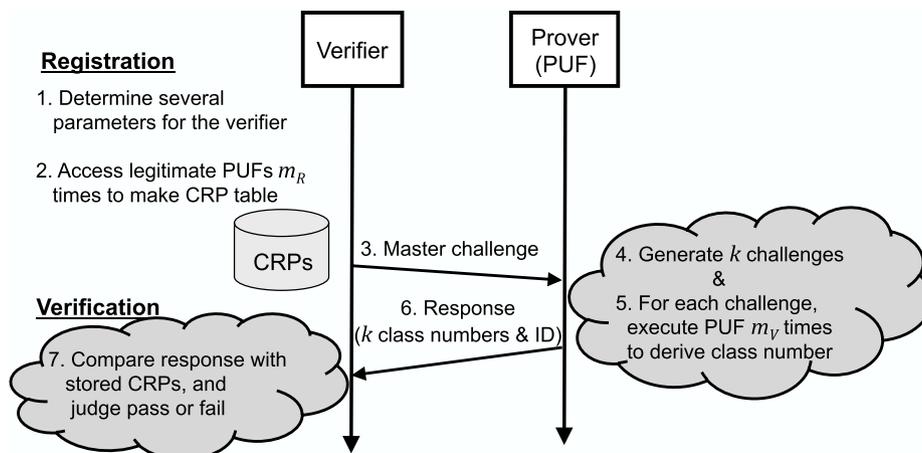


Fig. 5 Overview of new authentication system in the verification phase.

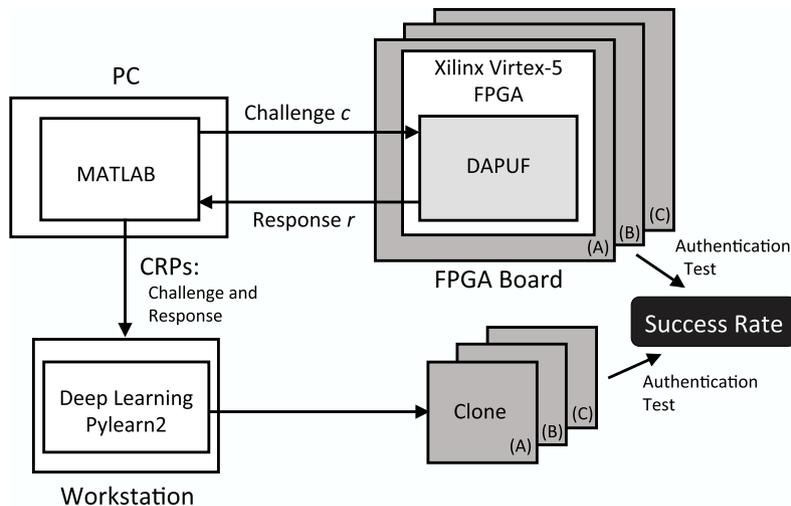


Fig. 6 Environment for calculating success rate.

Table 2 Measurement result for DAPUF ( $kU = 256$ ).

Dataset	2-1 DAPUF			3-1 DAPUF			4-1 DAPUF		
	A	B	C	A	B	C	A	B	C
Correctness [%]	94.7	87.6	88.3	77.0	62.0	73.3	70.0	64.7	62.2
Prediction rate [%]	96.4	84.2	90.0	60.1	60.7	70.2	65.0	58.6	62.9
Uniqueness [%]	73.4			67.3			67.1		
Randomness [bit]	0.09	0.64	0.64	0.89	0.90	0.84	0.80	0.94	0.76
SOM [bit]	-4	8	-5	43	3	7	12	15	-1

#### 4.1 Experimental Setup

Figure 6 shows an experimental environment. Instances of an open-source DAPUF implementation [10] are implemented on a Xilinx Virtex-5 Field-Programmable Gate Array (FPGA) XC5VLX30. CRPs are collected from the three FPGAs, and they are referred to as datasets A, B, and C. The datasets are fed to DL to make clones. Finally, authentication test and measurement of success rate are performed.

The attacker collects CRPs from a legitimate PUF, and uses 50,000 pairs as training data for DL. Pylearn2 is employed for DL [16]. Parameters necessary for Pylearn2 are determined by following the previous work [10].

We also use the idea of challenge vector introduced in [5]. Let  $c_l$  be the  $l$ -th bit of challenge of  $k$  bit length. If  $c_l = 0$ , and  $c_l = 1$  the delay time is represented as  $\delta_l^0$  and  $\delta_l^1$ , respectively.  $\vec{\Phi}$  is so-called challenge vector, and it is expressed as  $\vec{\Phi}(\vec{c}) = (\Phi^1(\vec{c}), \dots, \Phi^d(\vec{c}), 1)$ , where  $\Phi^l(\vec{c}) = \prod_{i=1}^d (1 - 2c_i)$  for  $l = 1, \dots, d$ .

In the training phase, 50,000 CRPs are used to make a clone of a legitimate PUF. Challenges are applied to the above expression to make input data for DL. Depending on the learning data, the clone outputs either raw response or class number.

#### 4.2 Preliminary Experiment

Table 2 shows quantitative metrics, i.e., correctness, predic-

tion rate, uniqueness, randomness, and SOM for the datasets A, B, and C. The prediction rate in Table 2 indicates the ratio of successful predictions by clones over the total predictions. Here, the parameters for measuring the metrics are  $T = 64$ ,  $L = 10,000$ , and  $N = 3$ .

The randomness of 2-1 DAPUF are lower than those of 3-1 and 4-1 DAPUFs. In particular, the randomness of 0.089 bits in the dataset A means the most of the raw responses are either 0 or 1. Randomness of the datasets B and C are better but still low. The results can be explained by unbalanced signal delays caused as a result of layout constraints in FPGA.

The value of  $U$  relates to the total number of raw response needed for authenticating devices appropriately. For instance, when  $U$  is close to 0, the difference of CRPs is small between PUFs. Therefore, if uniqueness is low, more CRPs are needed to distinguish PUFs. In addition, if uniqueness is low, the clone of a legitimate PUF might be able to attack other PUFs.

In the experiment,  $kU = 256$  is used i.e., the number of raw responses  $k = 256/U$ . The parameter is determined in order to set effective response length to be 256 bits considering bits lost by  $U < 1$ . SOM is calculated based on the above parameter. As shown in the SOM row in the table, cloning is successful when 2-1 DAPUF or 4-1 DAPUF are used. In contrast, the clones can be distinguished from the legitimate prover when 3-1 is used.

#### 4.3 Experiments on Q-Class Authentication

The following experiments to evaluate the security of the

**Table 3** Parameter settings used for experiments on 2-, 3-, 4-, and 5-class authentication ( $m_R = 63$ ).

$m_V$	Range of raw response value for each class number									
	$Q = 2$				$Q = 3$					
	Class 1		Class 2		Class 1		Class 2		Class 3	
1	$r_s = 0$		$r_s = 1$		-		-		-	
3	$0 \leq r_s \leq 1$		$2 \leq r_s \leq 3$		$r_s = 0$		$1 \leq r_s \leq 2$		$r_s = 3$	
7	$0 \leq r_s \leq 3$		$4 \leq r_s \leq 7$		$r_s = 0$		$1 \leq r_s \leq 6$		$r_s = 7$	
15	$0 \leq r_s \leq 7$		$8 \leq r_s \leq 15$		$r_s = 0$		$1 \leq r_s \leq 14$		$r_s = 15$	
31	$0 \leq r_s \leq 15$		$16 \leq r_s \leq 31$		$r_s = 0$		$1 \leq r_s \leq 30$		$r_s = 31$	
63	$0 \leq r_s \leq 31$		$32 \leq r_s \leq 63$		$r_s = 0$		$1 \leq r_s \leq 62$		$r_s = 63$	

$m_V$	Range of raw response value for each class number										
	$Q = 4$				$Q = 5$						
	Class 1	Class 2	Class 3		Class 4	Class 1	Class 2	Class 3		Class 4	Class 5
1	-										
3	$r_s = 0$	$r_s = 1$	$r_s = 2$		$r_s = 3$	-	-	-		-	-
7	$r_s = 0$	$1 \leq r_s \leq 3$	$4 \leq r_s \leq 6$		$r_s = 7$	$r_s = 0$	$1 \leq r_s \leq 2$	$3 \leq r_s \leq 4$		$5 \leq r_s \leq 6$	$r_s = 7$
15	$r_s = 0$	$1 \leq r_s \leq 7$	$8 \leq r_s \leq 14$		$r_s = 15$	$r_s = 0$	$1 \leq r_s \leq 5$	$6 \leq r_s \leq 9$		$10 \leq r_s \leq 14$	$r_s = 15$
31	$r_s = 0$	$1 \leq r_s \leq 15$	$16 \leq r_s \leq 30$		$r_s = 31$	$r_s = 0$	$1 \leq r_s \leq 10$	$11 \leq r_s \leq 20$		$21 \leq r_s \leq 30$	$r_s = 31$
63	$r_s = 0$	$1 \leq r_s \leq 31$	$32 \leq r_s \leq 62$		$r_s = 63$	$r_s = 0$	$1 \leq r_s \leq 21$	$22 \leq r_s \leq 41$		$42 \leq r_s \leq 62$	$r_s = 63$

authentication system are conducted using the following parameters. In the registration phase, 10,000 different challenges are used to make CRPs consisting of class numbers. The thresholds of classification are set as summarized in Table 3. In the verification phase, 10,000 challenges are sent from the verifier to the prover. For each challenge, the prover iterates PUF operations  $m_V$  times, and answers the corresponding class based on the number of 1s in the raw response value  $r_s$ .

An attacker creates a clone as described in Sect. 4.1. The quality of clones is evaluated with success rate that is derived by counting the number of class numbers, which is the same as the registered class. The success rate is obtained by dividing the number of prover's correct answers by 10,000 CRPs.

Cloning attack is considered and evaluated because that should be the strongest attack for the proposed system. Another potential threat is a spoofing attack in which an attacker replaces a public ID for 1 : 1 authentication thereby converting a legitimate PUF into another one. The attack succeeds when a PUF is coincidentally accepted as another PUF. The resistance against such attack can be evaluated by  $1 - U$  when  $Q = 2$ . As shown in Table 2, the predictability by clones are 84%, 60%, 58% for 2-, 3-, and 4-1 DAPUFs, respectively. They are much higher than corresponding  $1 - U$  that are 27%, 33%, 33%, respectively. The result indicates that the spoofing attack is less effective compared to cloning attack.

Therefore, the experimental results are evaluated using success rate. Let  $A^P$  and  $A^C$  be the success rate of a prover, and success rate from an attacker. The difference of success rate is defined as

$$\mathcal{D} = A_{Q, m_V}^P - \max_{1 \leq m_V \leq 63} A_{Q, m_V}^C. \quad (8)$$

The maximum  $A_{Q, m_V}^C$  is used considering the most suc-

cessful clones for strict security evaluation. As  $\mathcal{D}$  becomes higher, authentication errors i.e., false acceptance or false rejection become smaller. Because success rate indicates the average rate of correct responses, the number of correct responses goes up and down. In general, when the threshold is set high, false rejection increases. In contrast, false acceptance increases when the threshold is set low. If  $\mathcal{D}$  is large, it is possible to set a threshold that might not occur authentication error.

Figure 7 shows the success rates, and Table 4 summarizes the success rate difference  $\mathcal{D}$  between the legitimate prover (PUF) and the attacker (clone) for the datasets A, B, and C from DAPUF.

When  $Q = 2$ ,  $Q$ -class authentication system is the same as that using conventional DAPUF. Therefore, success rate is almost flat even when  $m_V$  is increased as shown in Fig. 7.

In 3-class authentication, there is a bias in the total number of raw responses in each class when using 3-1 DAPUF and 4-1 DAPUF. That is because 3-1 DAPUF and 4-1 DAPUF have numerous unstable raw responses. DL predicts efficiently when the ratio between the classes is biased. The success rates increase by  $m_V$  with a few exceptions. When using 2-1 DAPUF,  $\mathcal{D}$  is larger compared to that of  $Q = 2$ . However, in the case of 3-1 DAPUF and 4-1 DAPUF,  $\mathcal{D}$  does not become larger as expected, since the clone's success rate is higher.

For  $Q = 4$ ,  $\mathcal{D}$  is improved in comparison to that with  $Q = 2$ . When using 2-1 DAPUF and 3-1 DAPUF,  $\mathcal{D}$  increases slightly as shown in Table 4. Furthermore, the differences of the success rates  $\mathcal{D}$  are 41%, 44%, and 36% for datasets A, B, and C, respectively, when using 4-1 DAPUF. In addition,  $\mathcal{D}$  increases by 10% compared to that with  $Q = 2$ , because the number of raw responses are uniform in each class. In other words, the chance to get correct response becomes low if the number of classes increased.

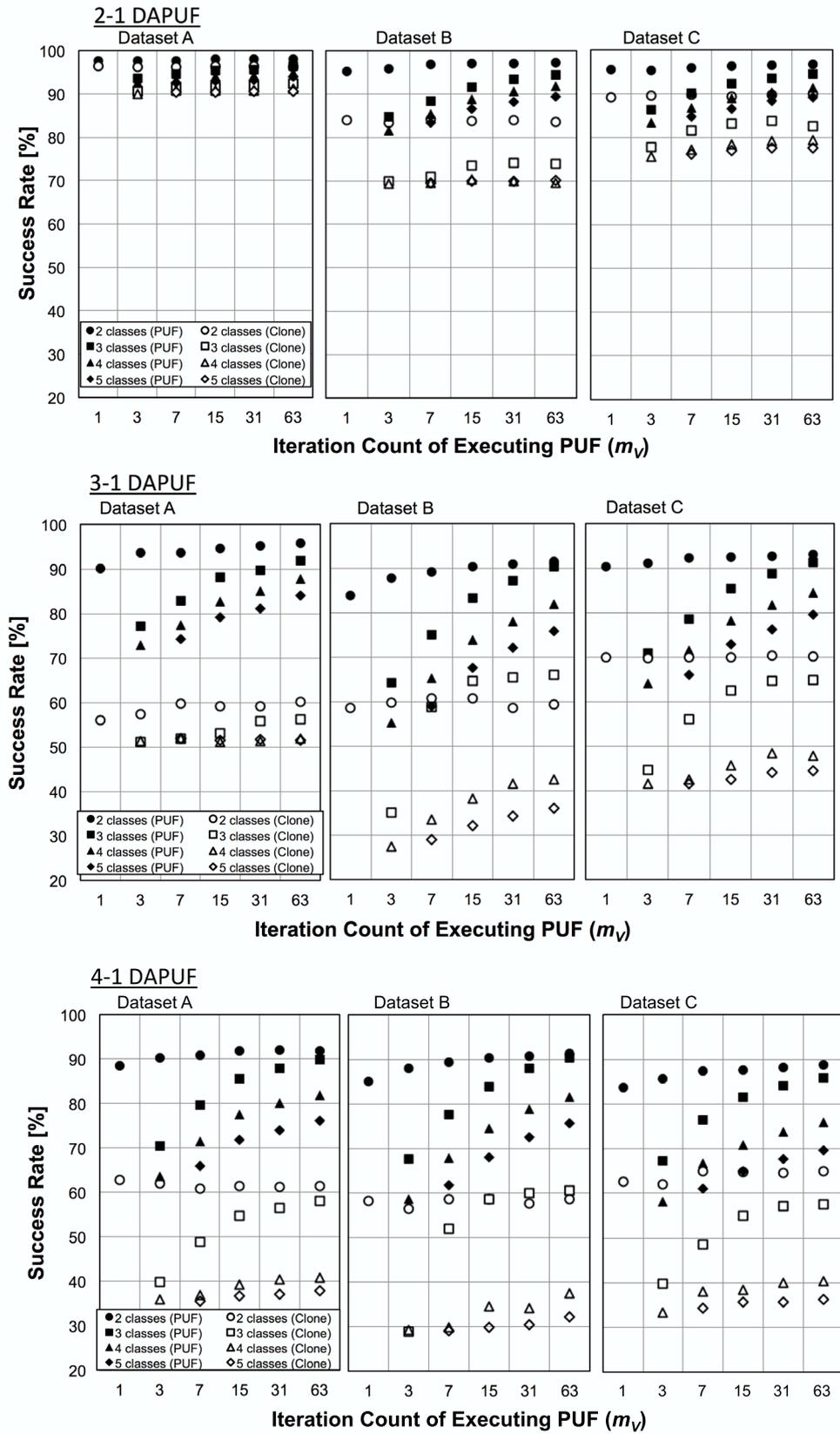


Fig. 7 Success rate of Q-class authentication using 2-1, 3-1, and 4-1 DAPUF for different iteration counts,  $m_V$  of executing PUF in verification ( $m_R = 63$ ).

**Table 4** The difference of success rate  $\mathcal{D}$  ( $m_V=63$ ).

Dataset ( $Q$ )	2-1 DAPUF	3-1 DAPUF	4-1 DAPUF
A (2)	1.6	35.5	29.3
A (3)	4.0	35.5	31.7
A (4)	3.6	35.8	40.9
A (5)	3.3	32.2	38.3
B (2)	13.1	30.8	32.6
B (3)	20.3	24.3	29.8
B (4)	21.4	39.4	44.1
B (5)	19.4	39.8	43.5
C (2)	6.8	22.8	23.7
C (3)	10.6	26.4	28.4
C (4)	12.1	36.1	35.5
C (5)	11.8	35.2	33.3

For  $Q = 5$ , both the first and second terms of  $\mathcal{D}$ , i.e.,  $A_{Q,m_V}^P$  and  $\max_{1 \leq m_V \leq 63} A_{Q,m_V}^C$  become smaller because more misclassifications are observed. The first term decreases more rapidly, and thus  $\mathcal{D}$  is smaller compared to that with  $Q = 4$ . In summary, improvement in resistance against ML attack comes at the cost of misclassification. In some cases, the degradation of  $A_{Q,m_V}^P$  is larger than the degradation of the  $\max_{1 \leq m_V \leq 63} A_{Q,m_V}^C$  as with  $Q = 5$ . Therefore, it is important to choose an appropriate  $Q$  in order to maximize the performance of the authentication system.

When increasing the value of  $Q$ , it becomes more difficult to break the authentication system as shown in Fig. 7. However, when the value of  $Q$  is increased too much, the unreliability of the raw responses would be known to the attacker. Such unreliability can be used for an attack in [7]. Altogether, there is a strong relationship between the number of  $Q$  and the clone's ability.

Figure 7 shows the clone accuracy does not increase significantly by the value of  $m_V$ .

As Fig. 7 shows, the resistance against DL-attack does not improve even if  $Q$  is increased when using 2-1 DAPUF. For this reason, the prover's advantage difference is not higher. It is reported that 3-1 DAPUF is more suitable for challenge-response authentication than 2- or 4-1 DAPUFs in [10]. In case of 3-1 DAPUF, there is a case in which  $Q$  does not make significant change (see dataset A in Table 4). In contrast,  $\mathcal{D}$  of 4-1 DAPUF are improved by 10% in all the datasets when  $Q$  is changed from 2 to 4. This is presumably because the correctness of 4-1 DAPUF is lower than 2- or 3-1 DAPUF. To sum up,  $Q$ -class authentication is tended to suitable for unstable PUF.

#### 4.4 Change of Temperature

DAPUF outputs a lot of unstable raw responses. Accordingly, it is necessary to evaluate performance under various environments. One of the factors of authentication accuracy is authentication environment, i.e., temperature of FPGA boards. This experiment uses 4-1 DAPUF that becomes the most unstable in the case of  $m_V = 63$ . The success rates of provers are measured under different temperatures.

**Table 5** The success rate of environmental change ( $m_V = 63$ ).

class ( $Q$ )	temp. [ $^{\circ}\text{C}$ ]	prover			clone		
		A	B	C	A	B	C
2	5	86.22	83.79	83.29	62.13	58.62	64.72
	room	91.92	91.2	88.72	60.19	59.34	70.06
	50	86.35	88.31	87.07	62.20	58.32	65.16
	70	61.78	72.89	66.31	42.81	48.2	62.45
3	5	81.71	82.59	79.88	56.89	60.31	56.81
	room	89.82	90.32	85.90	56.13	66.04	64.88
	50	84.44	87.50	82.77	56.96	60.76	57.43
	70	56.20	72.64	62.27	18.87	54.43	54.60
4	5	69.07	67.20	65.29	40.77	36.83	39.27
	room	81.75	81.54	75.83	51.84	42.46	47.87
	50	71.50	75.98	71.28	40.26	37.32	39.11
	70	36.62	54.07	42.46	18.76	28.28	32.19
5	5	62.51	59.66	58.52	37.39	32.05	34.75
	room	76.23	75.58	69.62	37.90	32.08	36.30
	50	64.75	69.05	64.3	37.55	32.32	36.28
	70	30.22	46.85	34.41	19.24	22.84	25.29

Table 5 shows the experimental result when changing temperature of FPGA. The room temperature is around from  $24^{\circ}\text{C}$  to  $28^{\circ}\text{C}$ . Then, the boards are exposed to environments with  $5^{\circ}\text{C}$ ,  $50^{\circ}\text{C}$ , and  $70^{\circ}\text{C}$ . The environment of less than  $5^{\circ}\text{C}$  and over  $50^{\circ}\text{C}$  is realized by using a cold storage chamber or a heat gun, respectively. Meanwhile the temperature is monitored using a thermocouple. Note that the temperatures measured on the surface of the FPGA chip are  $0$ – $6^{\circ}\text{C}$ ,  $49$ – $58^{\circ}\text{C}$ , and  $69$ – $77^{\circ}\text{C}$ , respectively in reality. When the temperature is changed, the success rate decreases, because the number of unstable raw responses increases.

The FPGA boards do not work due to communication failure from  $-17^{\circ}\text{C}$  to  $-10^{\circ}\text{C}$ . In other words, the raw responses could not be obtained since the communication of FPGA boards was not performed. When the temperature is in the range between  $5^{\circ}\text{C}$  and  $50^{\circ}\text{C}$ , the success rates of provers are higher than those of clones, thus the verifier can set a threshold to distinguish the prover from the clone. When the temperature is higher than  $70^{\circ}\text{C}$ , PUFs sometimes did not operate normally. If the temperature is over  $80^{\circ}\text{C}$ , PUFs did not work, and the raw responses become all the same even if the given challenges are random.

Among the success rates of clone A for 2- and 3-class at  $5^{\circ}\text{C}$  or  $50^{\circ}\text{C}$  and clone B for 5-class at  $50^{\circ}\text{C}$ , no significant difference is observed. The clone's success rate does not increase by changing temperature, thus  $\mathcal{D}$  at room temperature does not decrease any further.

Furthermore, the success rates of the prover decreases when the temperature of FPGA boards are changed. Although, the prover's success rates are higher than those of clones. Therefore, it would not become a serious threat to authentication.

## 5. Conclusions and Future Works

This study proposed a new authentication system that uses the  $Q$ -class separation of raw response values. The experimental evaluation of  $Q$ -class authentication using DAPUF (for  $Q = 3, 4, 5$ ) revealed that the proposed authentication

system enhanced ML tolerance compared to the conventional scheme. In particular, when using 4-1 DAPUF, the difference between the prover and the clone increases more effectively. The authentication system accuracy might be unaffected by temperature changes from 5 to 50°C.

In future work, we will evaluate various PUF primitives that are considered difficult to use in a conventional authentication system because of the unreliable raw responses.

## Acknowledgements

We are grateful to the associate editor and the anonymous reviewers for various constructive comments that helped significantly improve the presentation of the paper. This paper is based on results obtained from a project commissioned by the New Energy and Industrial Technology Development Organization (NEDO).

## References

- [1] R. Pappu, "Physical one-way functions," PhD Thesis, Massachusetts Institute of Technology, 2001
- [2] R. Pappu, B. Recht, J. Taylor, and N. Gershenfeld, "Physical one-way functions," *Science*, vol.297, no.5589, pp.2026–2030, 2002.
- [3] B. Gassend, D. Clarke, M. Van Dijk, and S. Devadas, "Silicon physical random functions," *Proc. 9th ACM Conference on Computer and Communications Security*, pp.148–160, ACM, 2002.
- [4] D. Lim, "Extracting secret keys from integrated circuits," Master's thesis, Massachusetts Institute of Technology, 2004.
- [5] U. Rührmair, F. Sehnke, J. Sölter, G. Dror, S. Devadas, and J. Schmidhuber, "Modeling attacks on physical unclonable functions," *Proc. 17th ACM Conference on Computer and Communications Security*, pp.237–249, ACM, 2010.
- [6] G.E. Suh and S. Devadas, "Physical unclonable functions for device authentication and secret key generation," *Proc. 44th Annual Design Automation Conference*, pp.9–14, ACM, 2007.
- [7] G.T. Becker, "The gap between promise and reality: On the insecurity of XOR arbiter PUFs," *International Workshop on Cryptographic Hardware and Embedded Systems (CHES)*, pp.535–555, Springer, 2015.
- [8] J. Delvaux and I. Verbauwhede, "Side channel modeling attacks on 65 nm arbiter PUFs exploiting CMOS device noise," *Hardware-Oriented Security and Trust (HOST), 2013 IEEE International Symposium on*, pp.137–142, IEEE, 2013.
- [9] T. Machida, D. Yamamoto, M. Iwamoto, and K. Sakiyama, "A new mode of operation for arbiter PUF to improve uniqueness on FPGA," *Computer Science and Information Systems (FedCSIS), 2014 Federated Conference on*, pp.871–878, 2014.
- [10] R. Yashiro, T. Machida, M. Iwamoto, and K. Sakiyama, "Deep-learning-based security evaluation on authentication systems using arbiter PUF and its variants," *International Workshop on Security*, pp.267–285, Springer, 2016.
- [11] D. Yamamoto, K. Sakiyama, M. Iwamoto, K. Ohta, M. Takenaka, and K. Itoh, "variety enhancement of PUF responses using the locations of random outputting RS latches," *J. Cryptographic Engineering*, vol.3, no.4, pp.197–211, 2013.
- [12] Y. Hori, T. Yoshida, T. Katashita, and A. Satoh, "Quantitative and statistical performance evaluation of arbiter physical unclonable functions on FPGAs," *Proc. 2010 International Conference on Reconfigurable Computing and FPGAs, ReConFig 2010*, pp.298–303, 2010.
- [13] A. Maiti, G. Vikash, and S. Patrick, "A systematic method to evaluate and compare the performance of physical unclonable functions," *Embedded Systems Design with FPGAs*, pp.245–267, 2013.
- [14] K. Furuhashi, M. Shiozaki, A. Fukushima, T. Murayama, and T. Fujino, "The arbiter-PUF with high uniqueness utilizing novel arbiter circuit with delay-time measurement," *2011 IEEE International Symposium of Circuits and Systems (ISCAS)*, pp.2325–2328, IEEE, 2011.
- [15] T. Machida, D. Yamamoto, M. Iwamoto, and K. Sakiyama, "A new arbiter PUF for enhancing unpredictability on FPGA," *The Scientific World Journal*, vol.2015, pp.1–13, 2015.
- [16] I.J. Goodfellow, D. Warde-Farley, P. Lamblin, V. Dumoulin, M. Mirza, R. Pascanu, J. Bergstra, F. Bastien, and Y. Bengio, "Pylearn2: A Machine Learning Research Library," *arXiv preprint arXiv:1308.4214*, 2013.



**Risa Yashiro** received the B.E. degree from Tokai University, Japan, in 2015. She is a master's student in Graduate School of The University of Electro-Communications, Japan since 2015. She is a student member of IEICE.



**Takeshi Sugawara** received B.E., M.Sc., and Ph.D. degrees from Tohoku University, Japan, in 2006, 2008, and 2011, respectively. In 2011, he joined Mitsubishi Electric Corporation. He is currently an Associate Professor at The University of Electro-Communications, Japan since 2017. His research interests include cryptography, anti-tamper design, and embedded systems security.



**Mitsugu Iwamoto** received the B.E., M.E., and Ph.D. degrees from the University of Tokyo, Tokyo, Japan, in 1999, 2001, and 2004, respectively. In 2004, he joined the University of Electro-Communications, where he is currently an Associate Professor of Department of Informatics. His research interests include information theory, information security, and cryptography. He is a member of IEICE, IEEE, and IACR.



**Kazuo Sakiyama** received the B.E. and M.E. degrees from Osaka University in 1994 and 1996, respectively, the M.S. degree from UCLA in 2003, and the Ph.D. degree in electrical engineering from KU Leuven in 2007. From 1996 to 2004, he was with the Semiconductor and IC Division, Hitachi, Ltd. He has been Professor at The University of Electro-Communications, Tokyo since 2013. He is a member of IACR, IPSJ and IEEE.