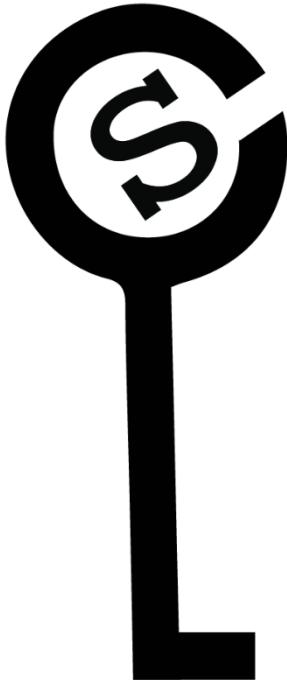


The Interpose PUF: Secure PUF Design against State-of-the-art Machine Learning Attacks

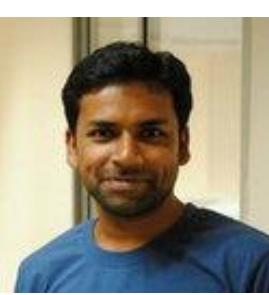


Phuong Ha Nguyen,

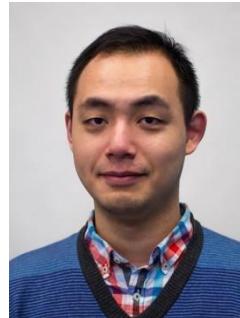
**Durga P. Sahoo, Kaleel Mahmood, Chenglu Jin,
Ulrich Rührmair and Marten van Dijk**



Ha



Durga



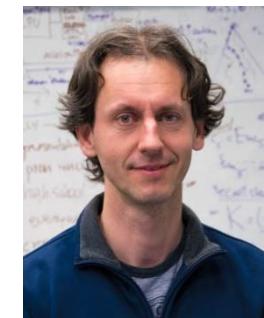
Chenglu



Kaleel



Uli



Marten

CHES 2019

UCONN

Content

1. Concept - Overview - Motivation
2. APUF- XOR APUF -iPUF
3. Machine Learning based Modeling Attacks on PUFs
4. Classical Machine Learning Attacks on APUFs
5. Classical Machine Learning Attacks on XOR APUFs
6. Interpose PUF (iPUF) – Resistance against Classical Machine Learning Attacks
7. Short-term Reliability
8. Reliability based Modeling Attacks on APUFs, XOR APUFs and iPUF's resistance against reliability based modeling attacks.



1. Concept - Overview - Motivation



Concept - Overview – Motivation



Concept - Overview – Motivation

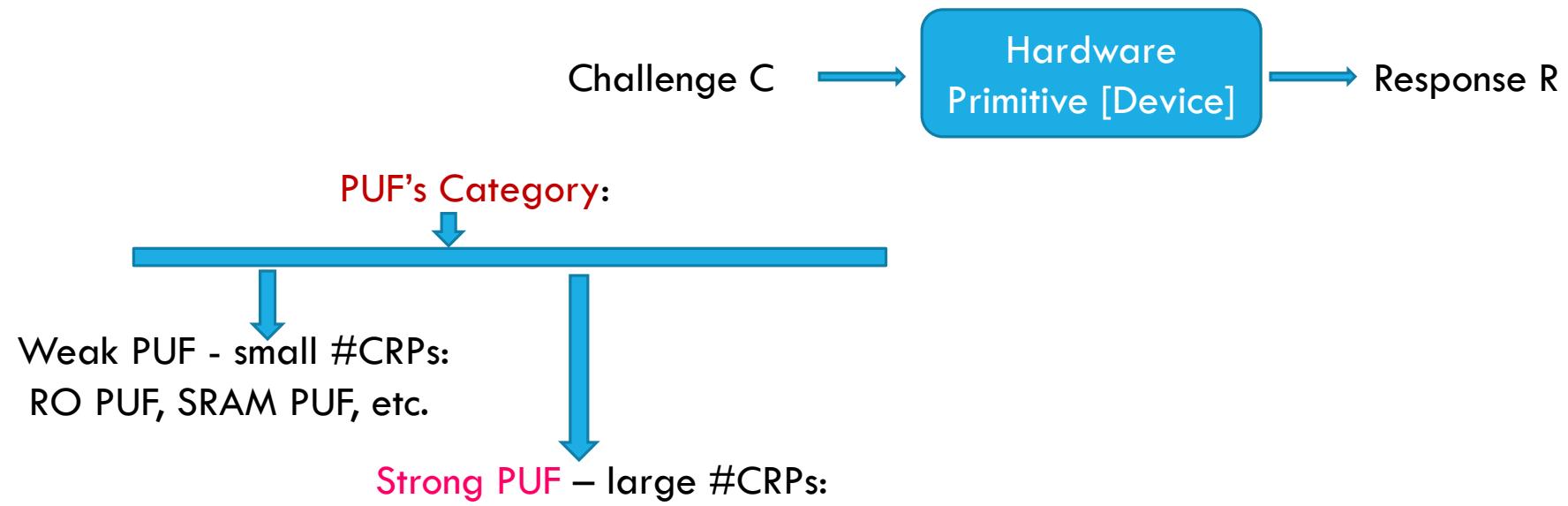


Nature: process variation – physically unclonability - unique

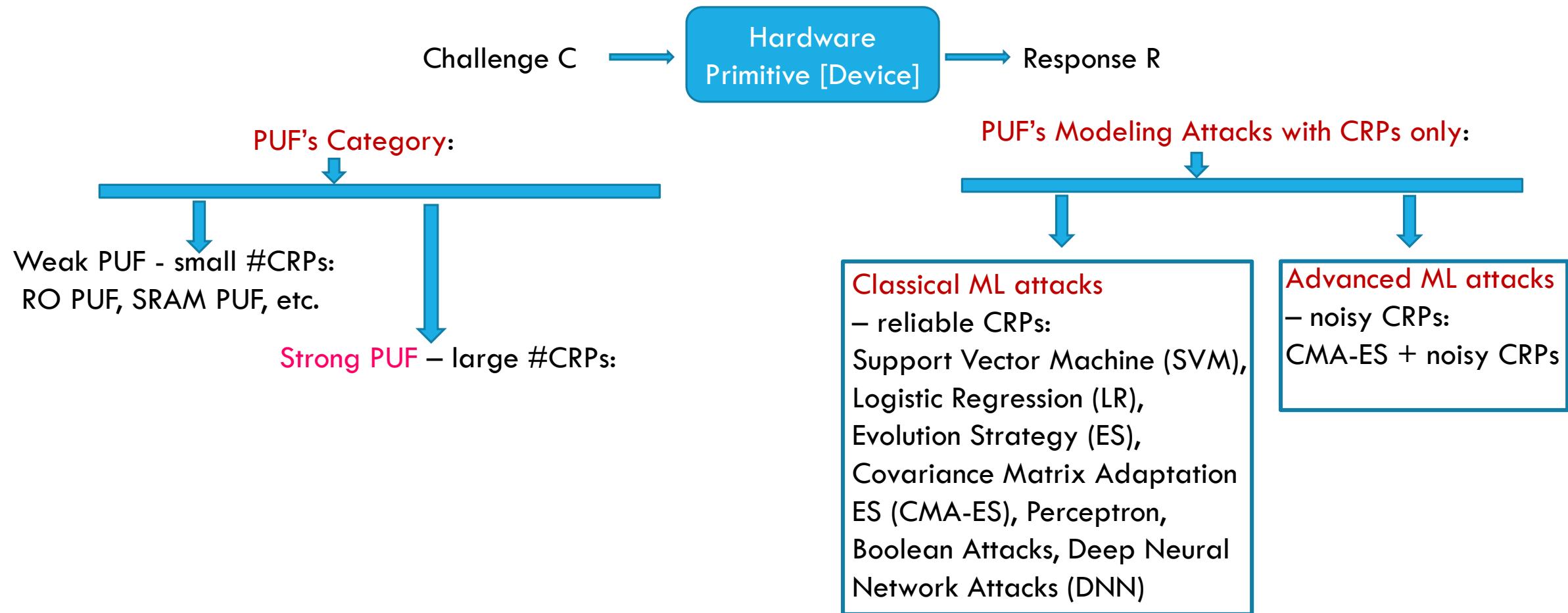
Application: device identification, authentication
and crypto key generation



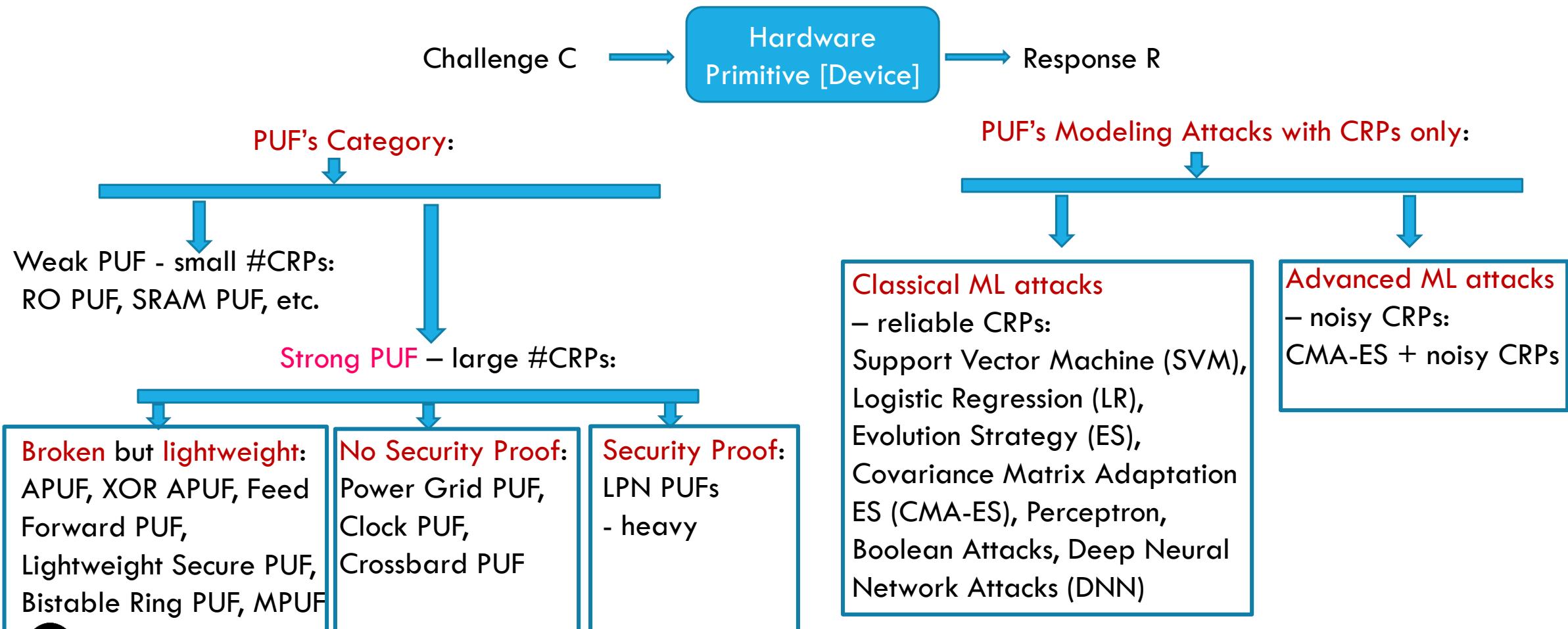
Concept - Overview – Motivation



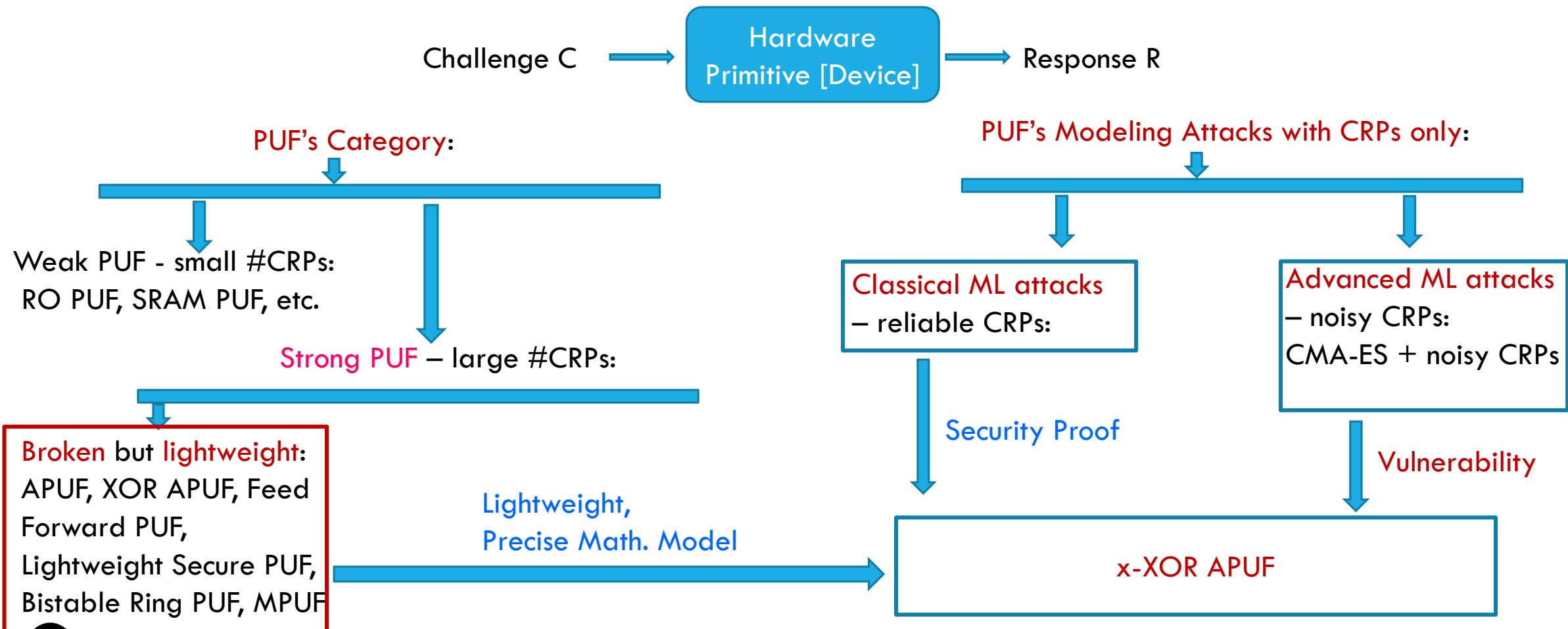
Concept - Overview – Motivation



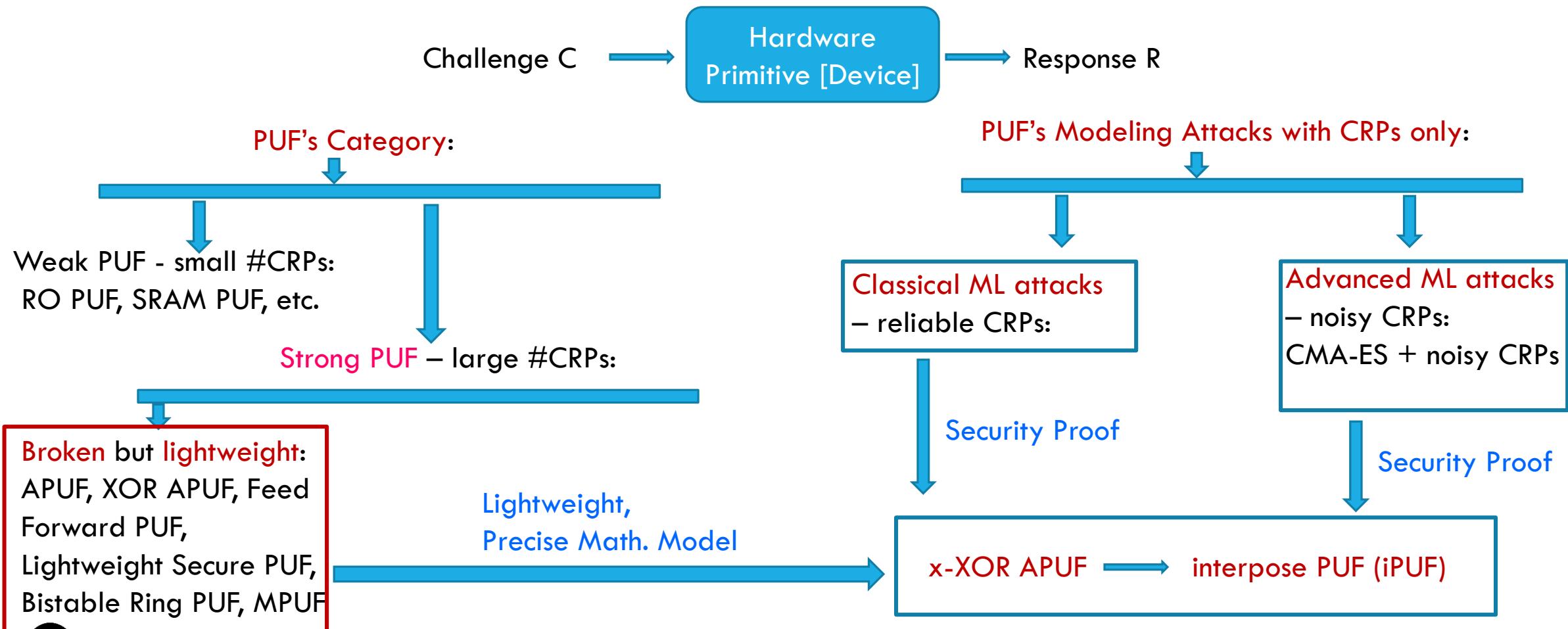
Concept - Overview – Motivation



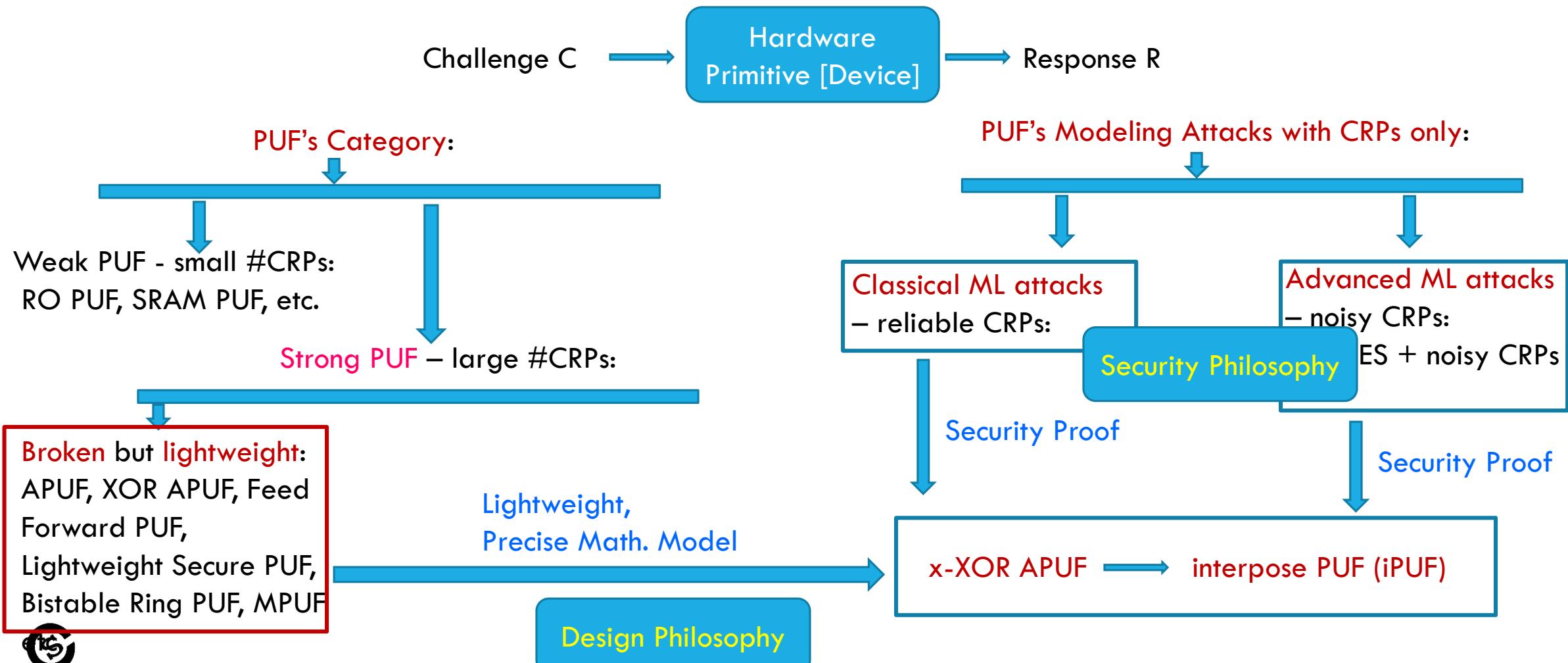
Concept - Overview – Motivation



Concept - Overview – Motivation



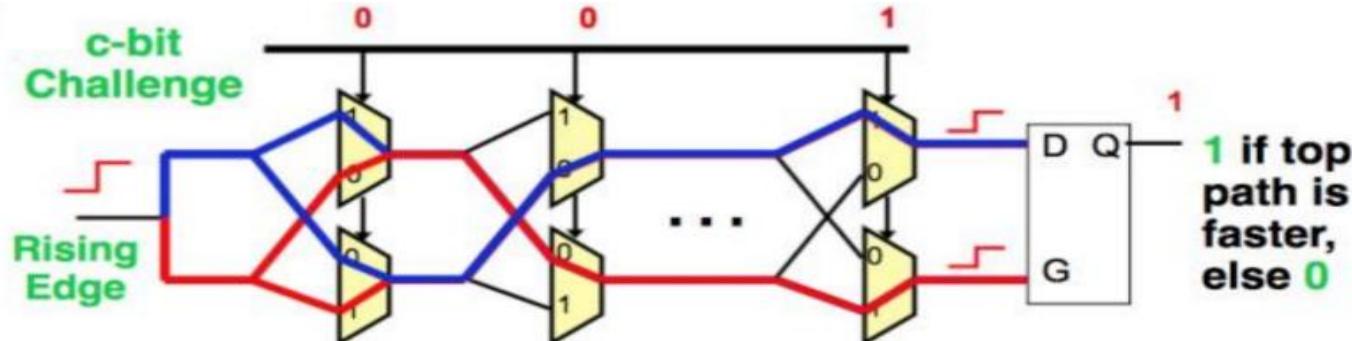
Concept - Overview – Motivation



2. APUF- XOR APUF -iPUF



APUF, XOR APUF and iPUF

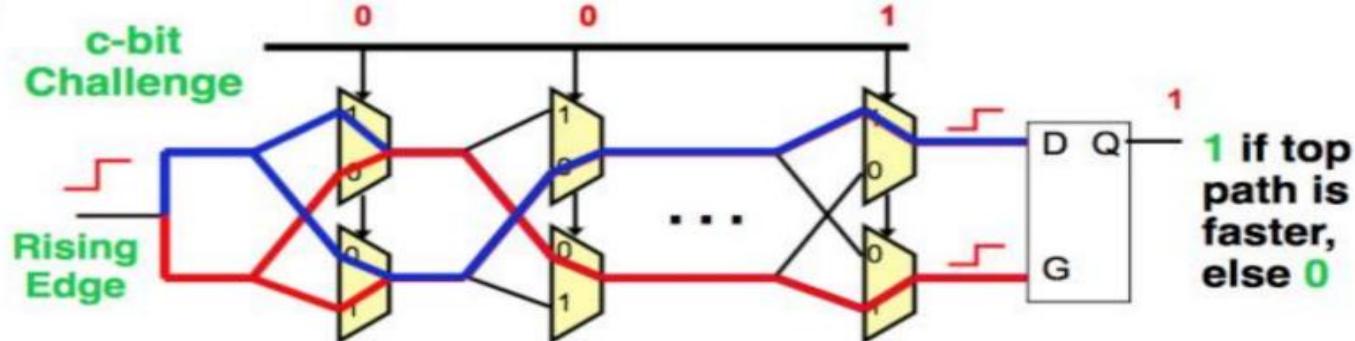


Arbiter PUF (APUF) [1]

- Extremely lightweight and large number of CRPs i.e, 2^n CRPs
- Process variation makes PUF's output not always reliable
- Not secure against modeling attacks

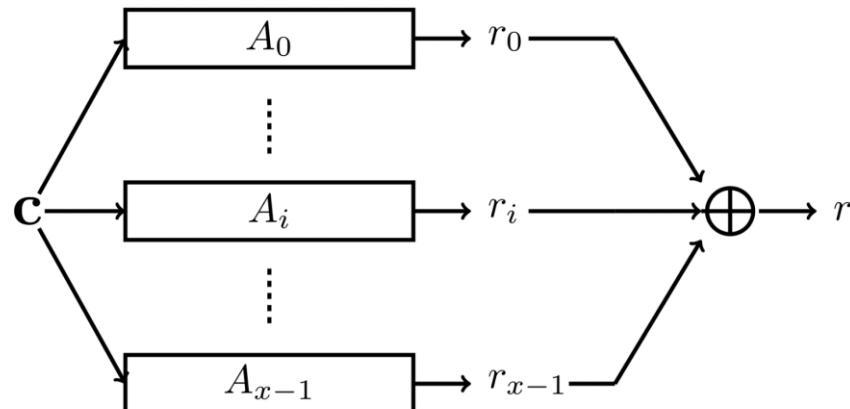


APUF, XOR APUF and iPUF

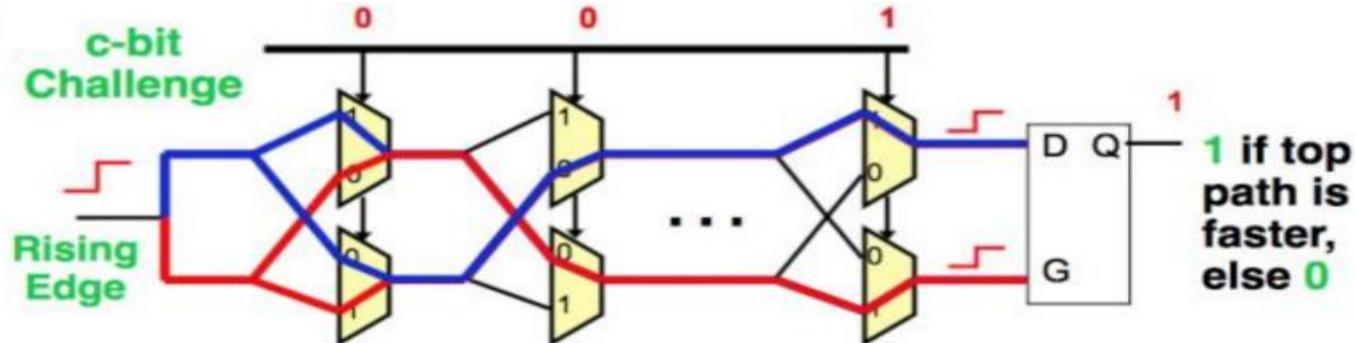


Arbiter PUF (APUF)

x-XOR APUF

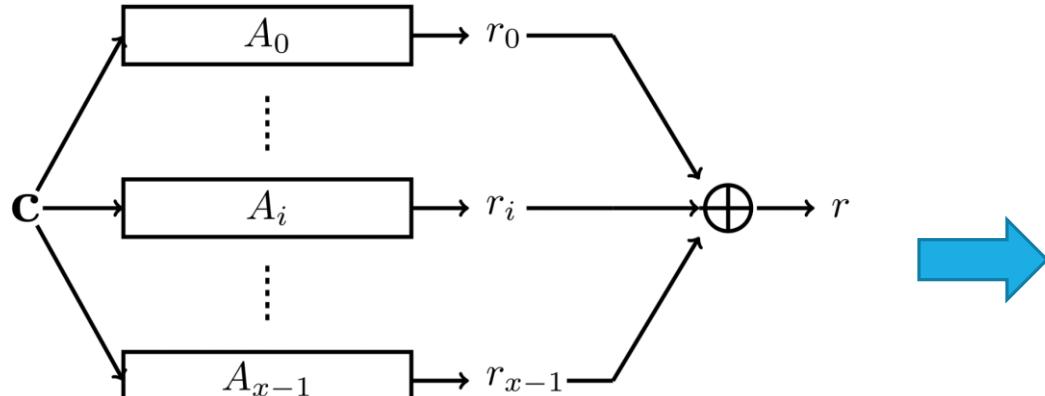


APUF, XOR APUF and iPUF

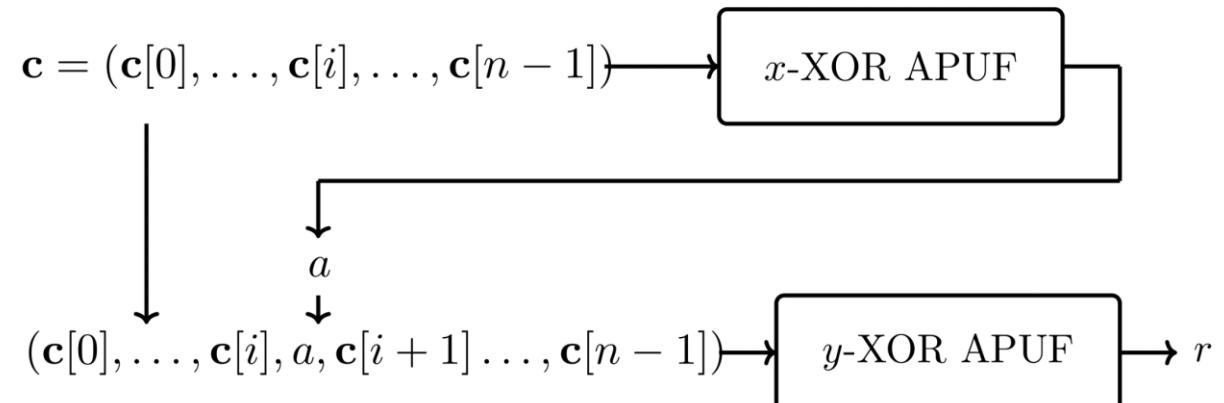


Arbiter PUF (APUF)

x-XOR APUF

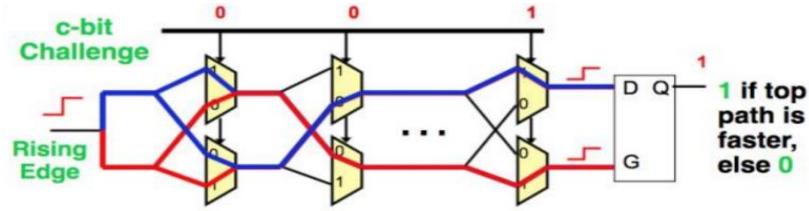


Interpose PUF (iPUF)

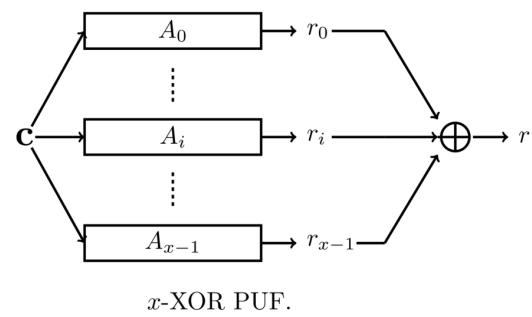


APUF, XOR APUF and iPUF

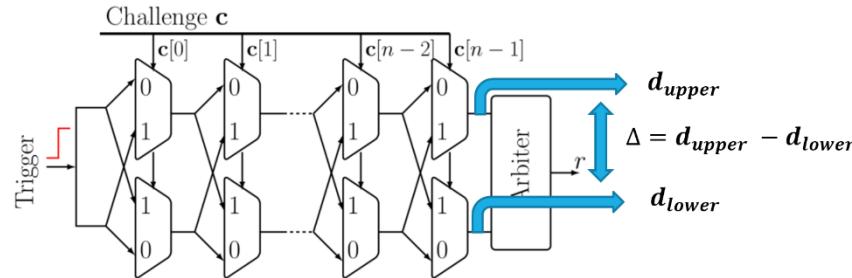
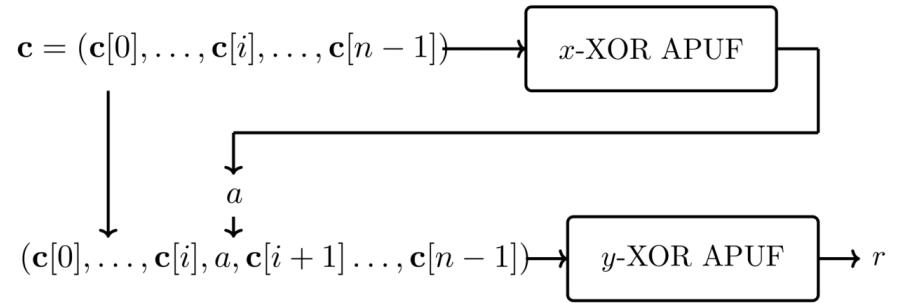
Arbiter PUF (APUF)



x-XOR Arbiter PUF



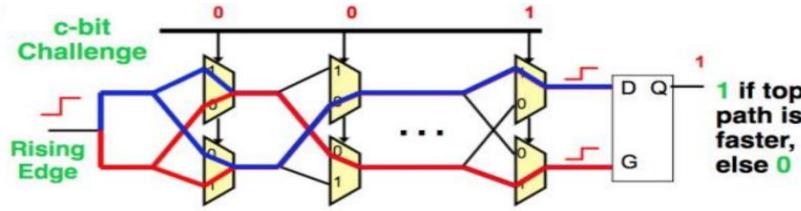
Interpose PUF (iPUF)



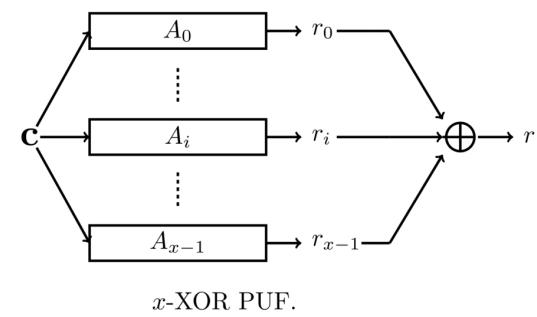
A n -stage classic Arbiter PUF with challenge $c \in \{0, 1\}^n$.

APUF, XOR APUF and iPUF

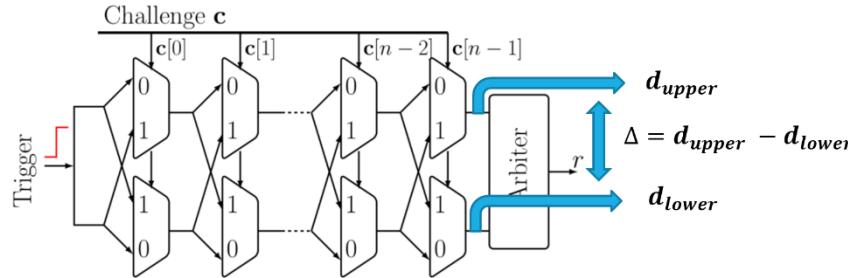
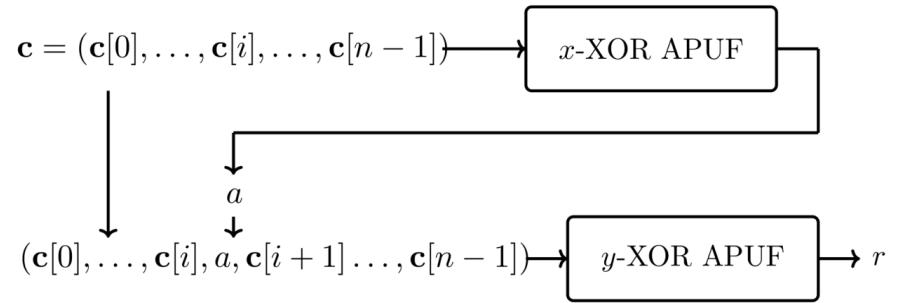
Arbiter PUF (APUF)



x-XOR Arbiter PUF



Interpose PUF (iPUF)



A n -stage classic Arbiter PUF with challenge $c \in \{0, 1\}^n$.

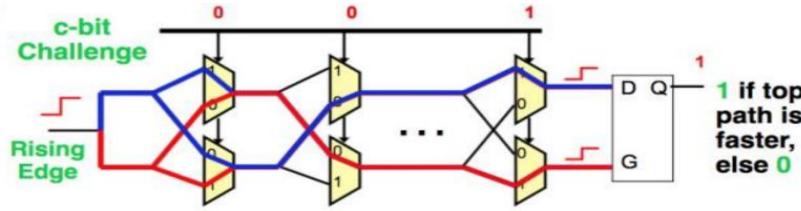
- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = w \cdot \Phi$
- w : unique for any APUF instance
- Φ is the parity vector



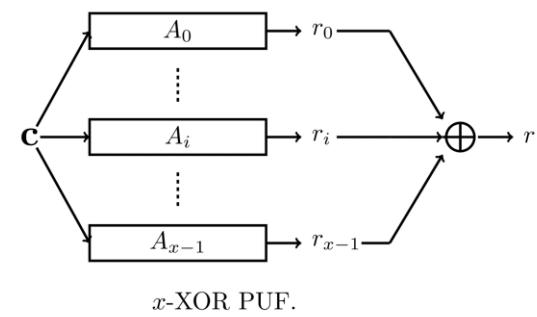
Precise linear model + CRPs + ML
= practically and softwarely clonable

APUF, XOR APUF and iPUF

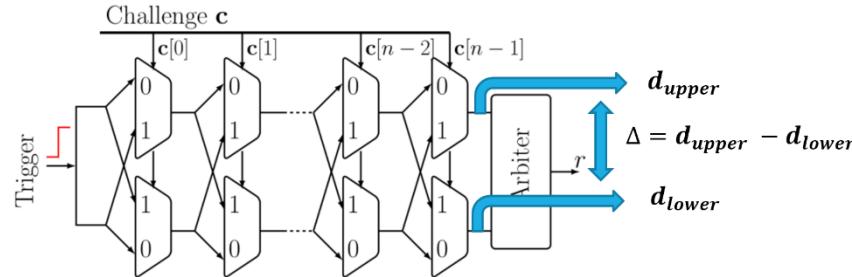
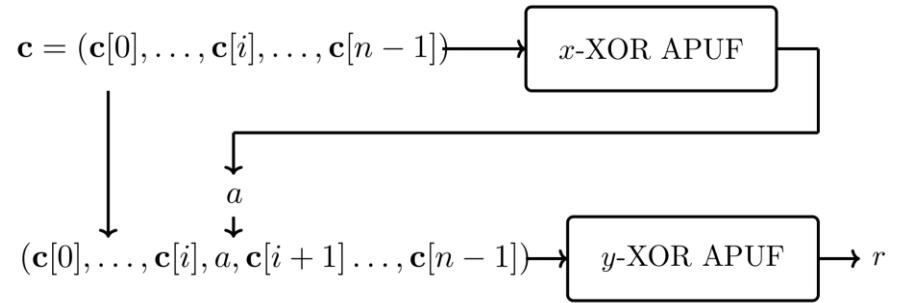
Arbiter PUF (APUF)



x-XOR Arbiter PUF



Interpose PUF (iPUF)



A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.

- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- \mathbf{w} : unique for any APUF instance
- Φ is the parity vector



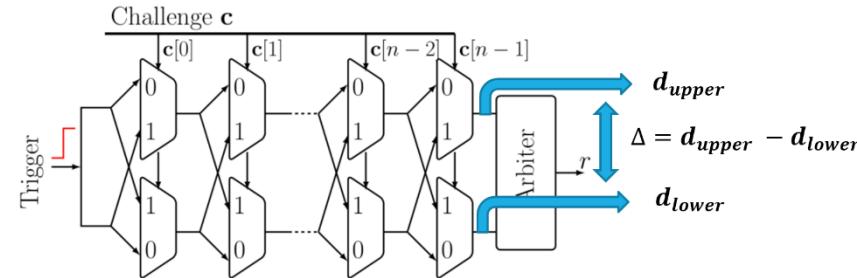
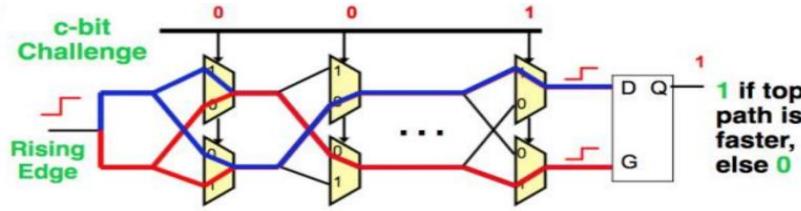
Precise linear model + CRPs + ML
= practically and softwarely clonable

Precise non-linear model + CRPs + classical ML
= impractically softwarely clonable

$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x w_i^T \Phi\right) [2]$$

APUF, XOR APUF and iPUF

Arbiter PUF (APUF)

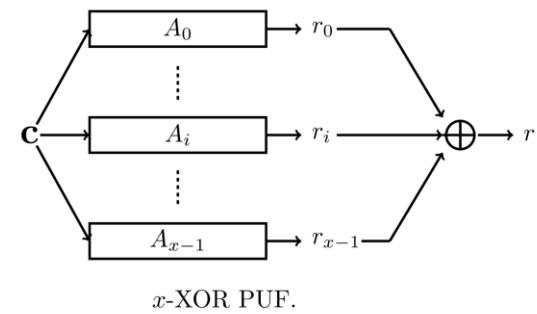


- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- \mathbf{w} : unique for any APUF instance
- Φ is the parity vector

$$\Phi[i] = \prod_{j=i, \dots, n-1} (1 - \mathbf{c}[j]), i = 0, \dots, n-1, \Phi[n] = 1$$

Precise linear model + CRPs + ML
= practically and softwarely clonable

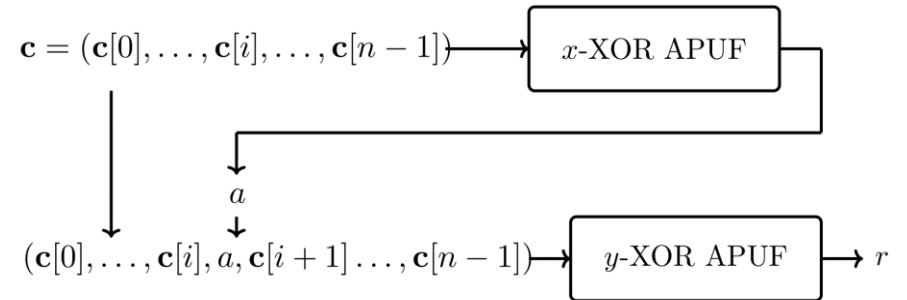
x-XOR Arbiter PUF



$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x \mathbf{w}_i^T \Phi\right)$$

Precise non-linear model + CRPs + classical ML
= impractically softwarely clonable

Interpose PUF (iPUF)

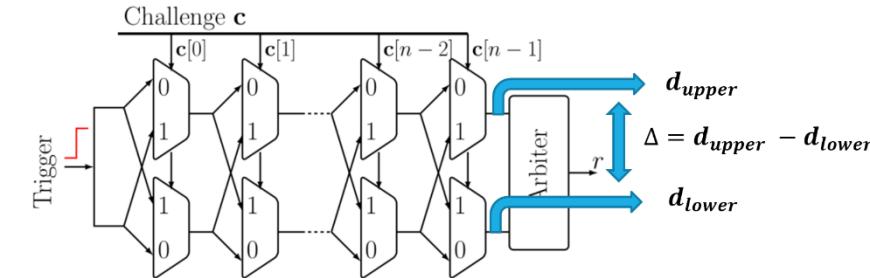
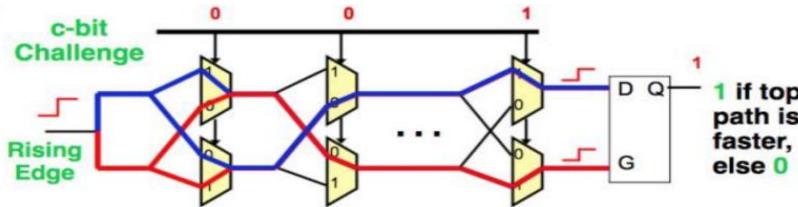


$$(x, y) - IPUF \approx \left(y + \frac{x}{2}\right) - XOR PUF$$

if a is inserted at the middle

APUF, XOR APUF and iPUF

Arbiter PUF (APUF)

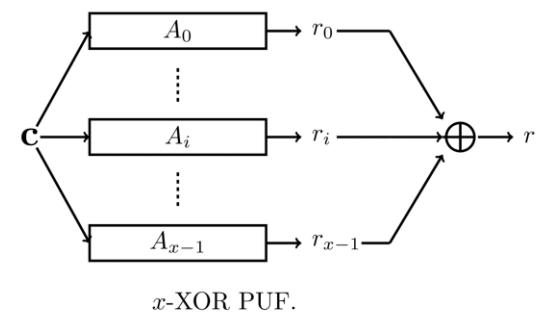


- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- \mathbf{w} : unique for any APUF instance
- Φ is the parity vector



Precise linear model + CRPs + ML
= practically and softwarely clonable

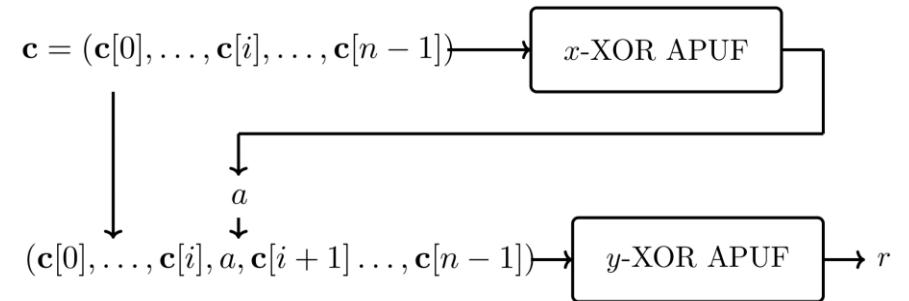
x-XOR Arbiter PUF



$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x \mathbf{w}_i^T \Phi\right)$$

Precise non-linear model + CRPs + classical ML
= impractically softwarely clonable

Interpose PUF (iPUF)



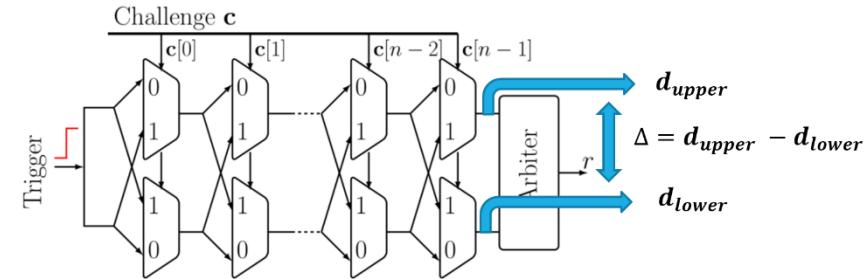
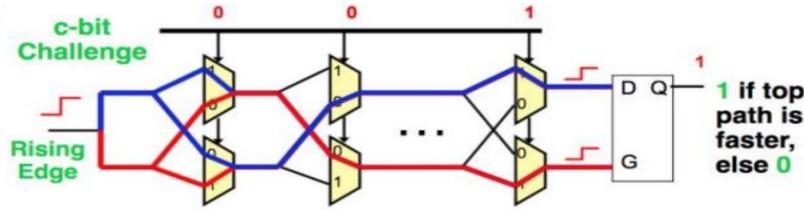
$$(x, y) - IPUF \approx \left(y + \frac{x}{2}\right) - XOR PUF$$

if a is inserted at the middle

Precise non-linear model + CRPs + classical ML
= impractically softwarely clonable

APUF, XOR APUF and iPUF

Arbiter PUF (APUF)



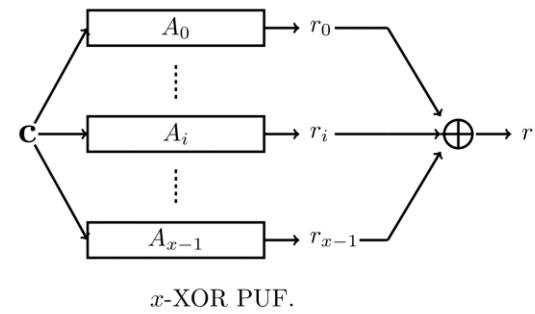
A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0,1\}^n$.

- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- \mathbf{w} : unique for any APUF instance
- Φ is the parity vector

$$\Phi[i] = \prod_{j=i, \dots, n-1} (1 - \mathbf{c}[j]), i = 0, \dots, n-1, \Phi[n] = 1$$

Precise linear model + CRPs + ML
= practically and softwarely clonable

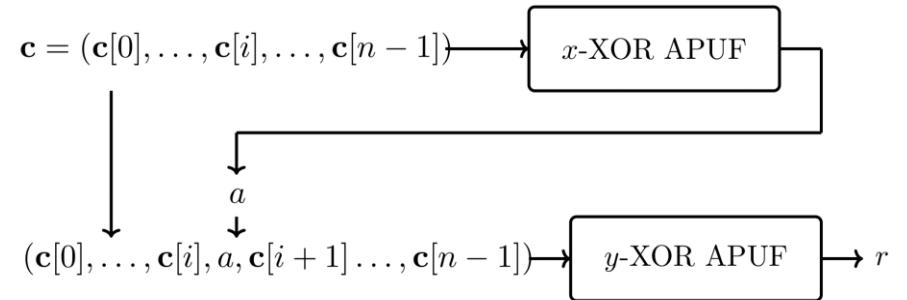
x-XOR Arbiter PUF



$$r_{XOR\ APUF} = sign\left(\prod_{i=1}^x w_i^T \Phi\right)$$

Precise non-linear model + CRPs + classical ML
= impractically softwarely clonable

Interpose PUF (iPUF)



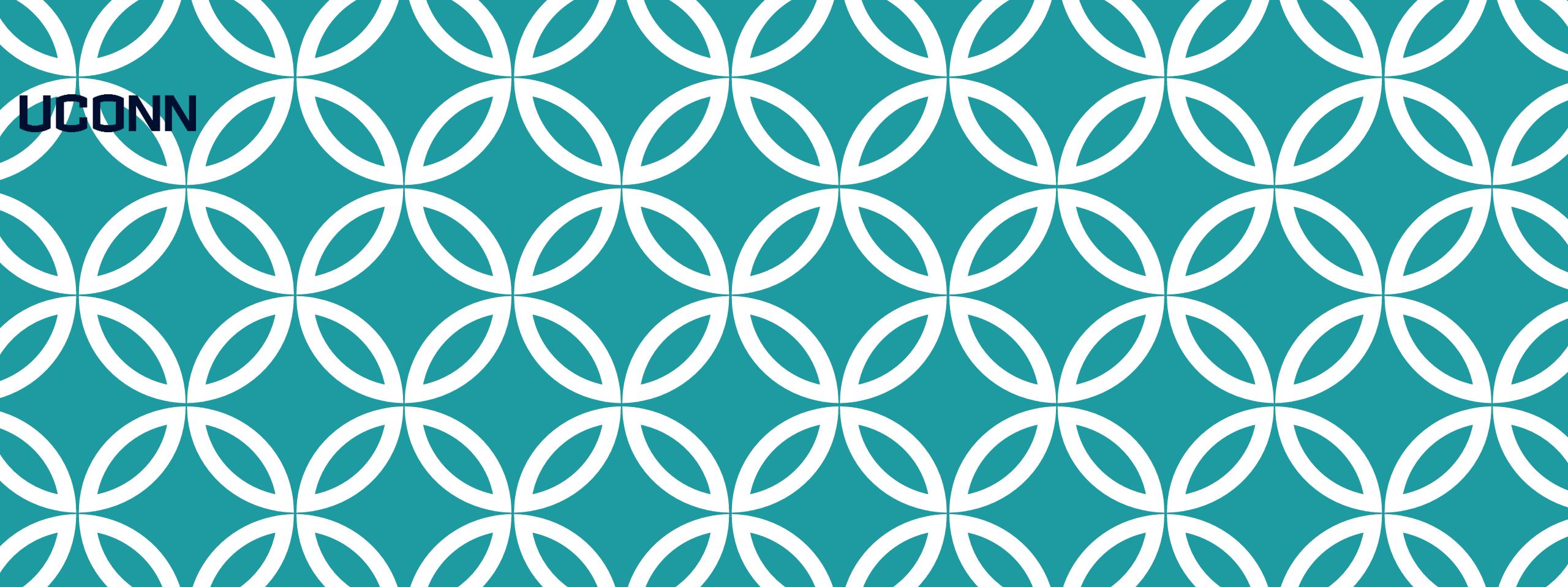
$$(x, y) - IPUF \approx \left(y + \frac{x}{2}\right) - XOR\ PUF$$

if a is inserted at the middle

Precise non-linear model + CRPs + classical ML
= impractically softwarely clonable

XOR APUF is not Secure against noisy CRPs + CMA-ES [Advanced ML]! (CHES2015) why?

Why not for iPUF?



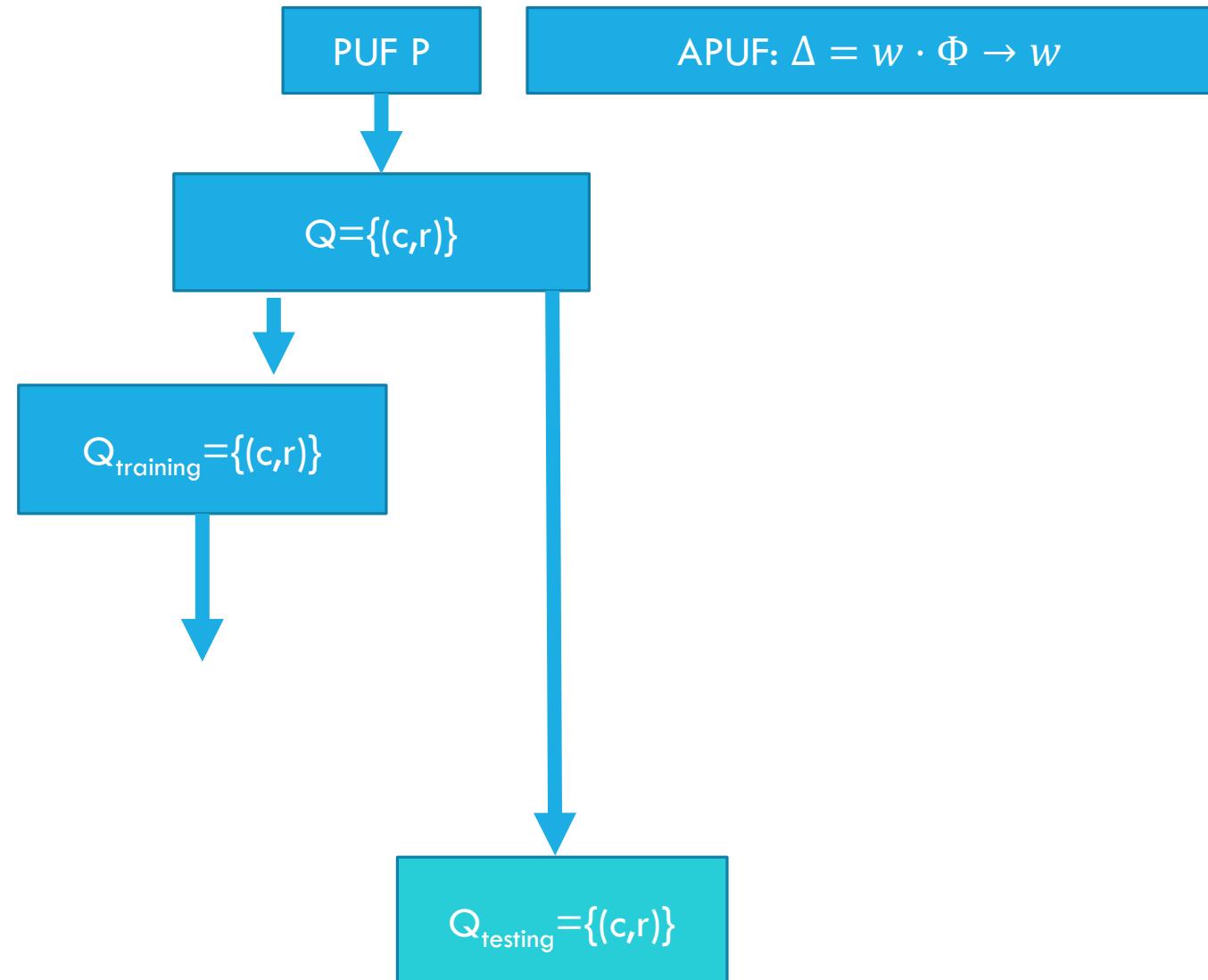
3. Machine Learning based Modeling Attacks on PUFs



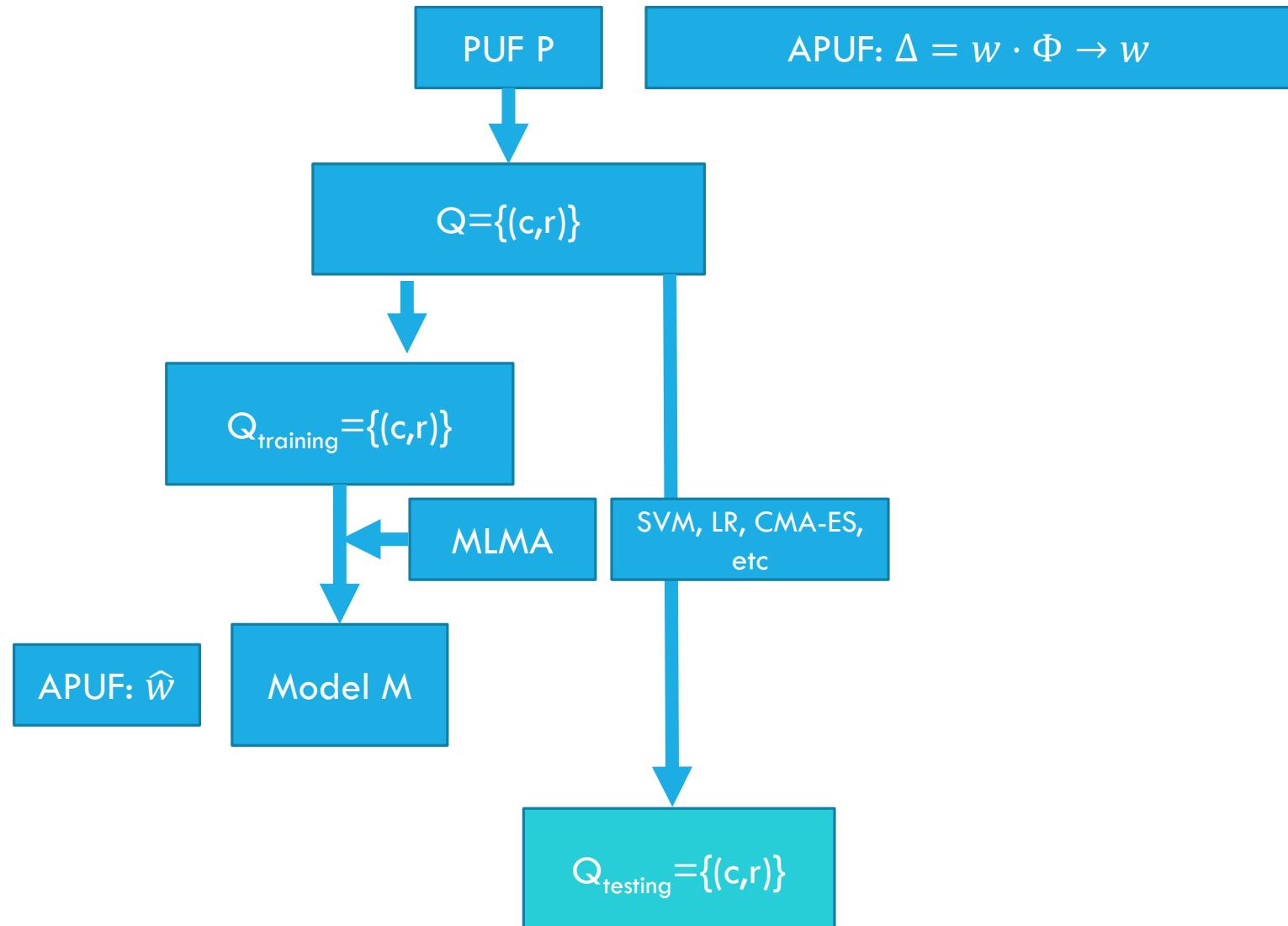
Basic Steps in MLMA



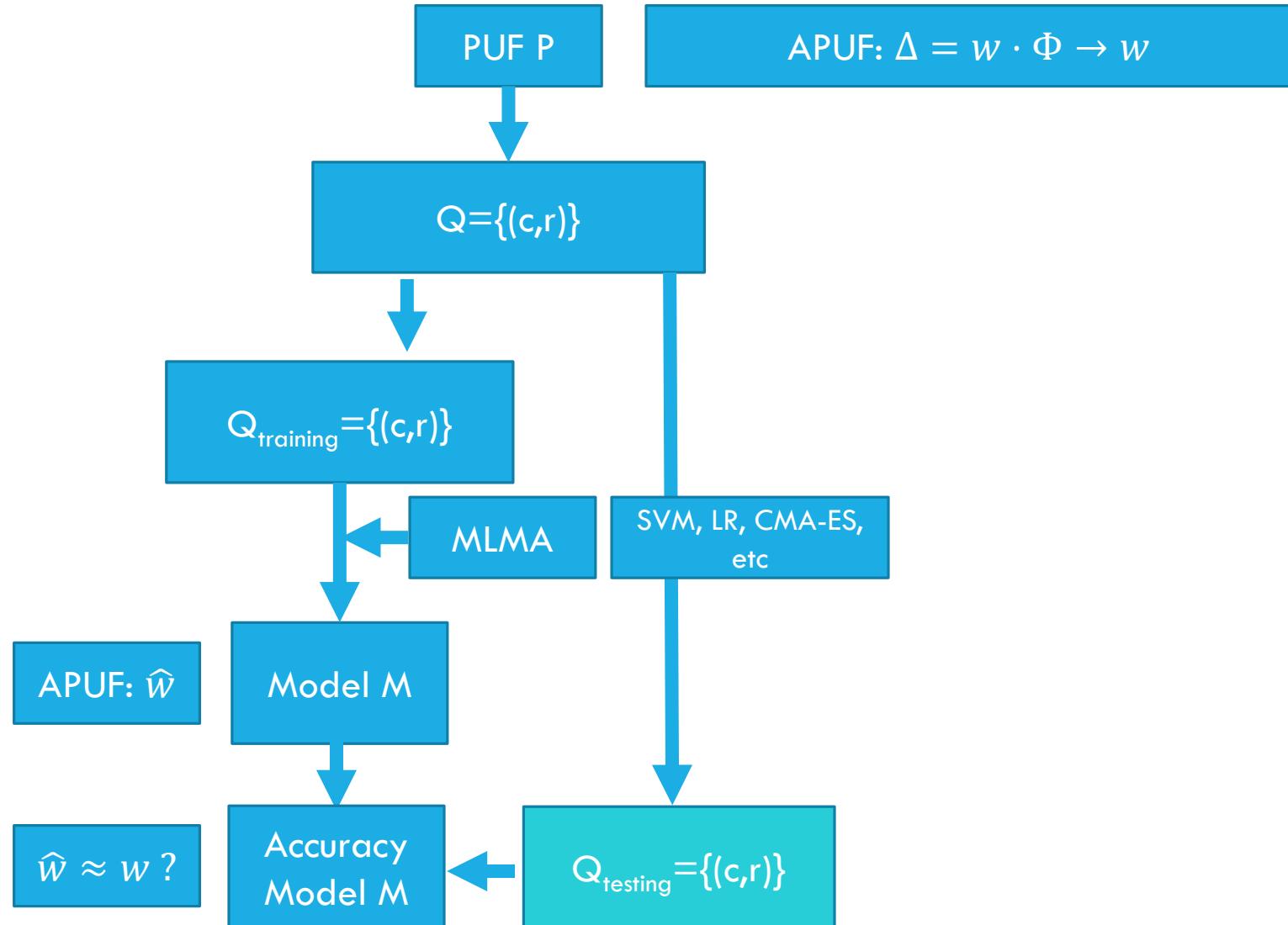
Basic Steps in MLMA



Basic Steps in MLMA



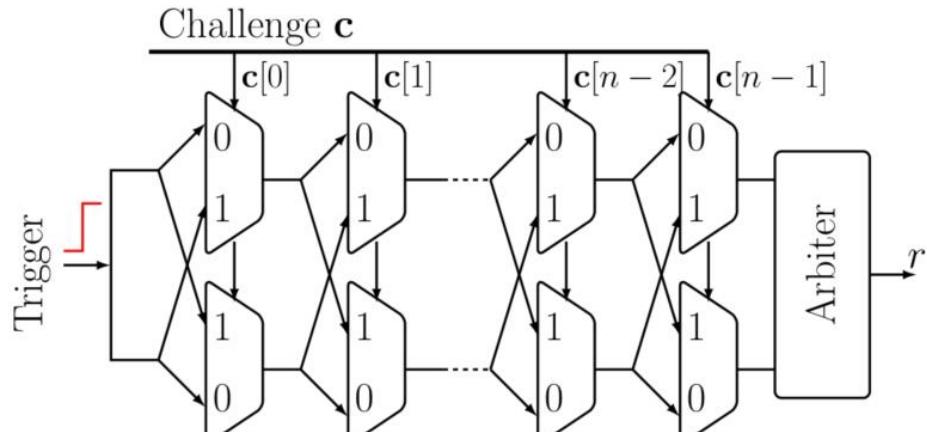
Basic Steps in MLMA



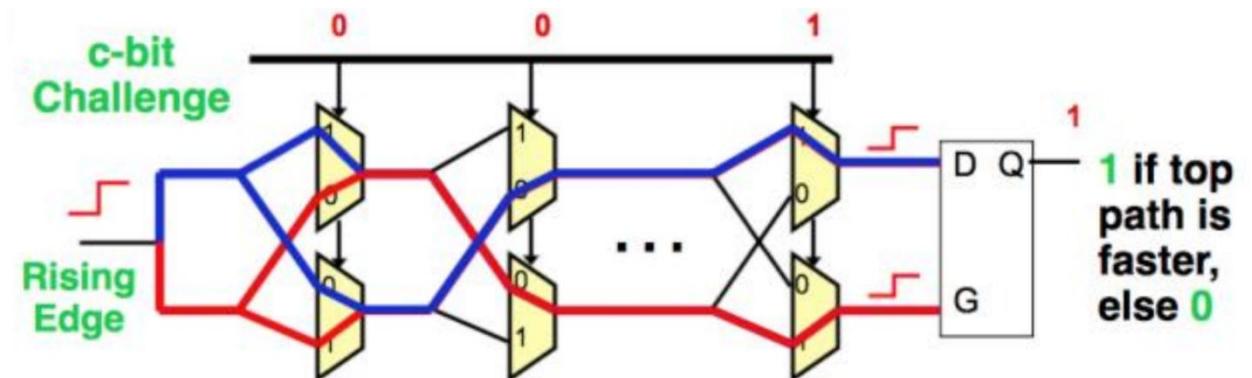
4. Classical Machine Learning Attacks on APUFs



Arbiter PUF



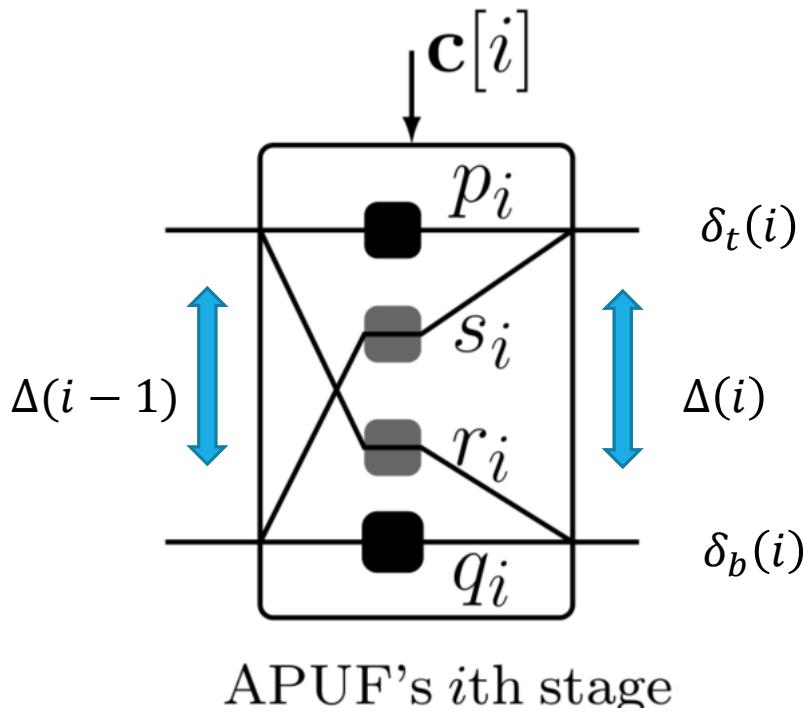
A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.



[10]

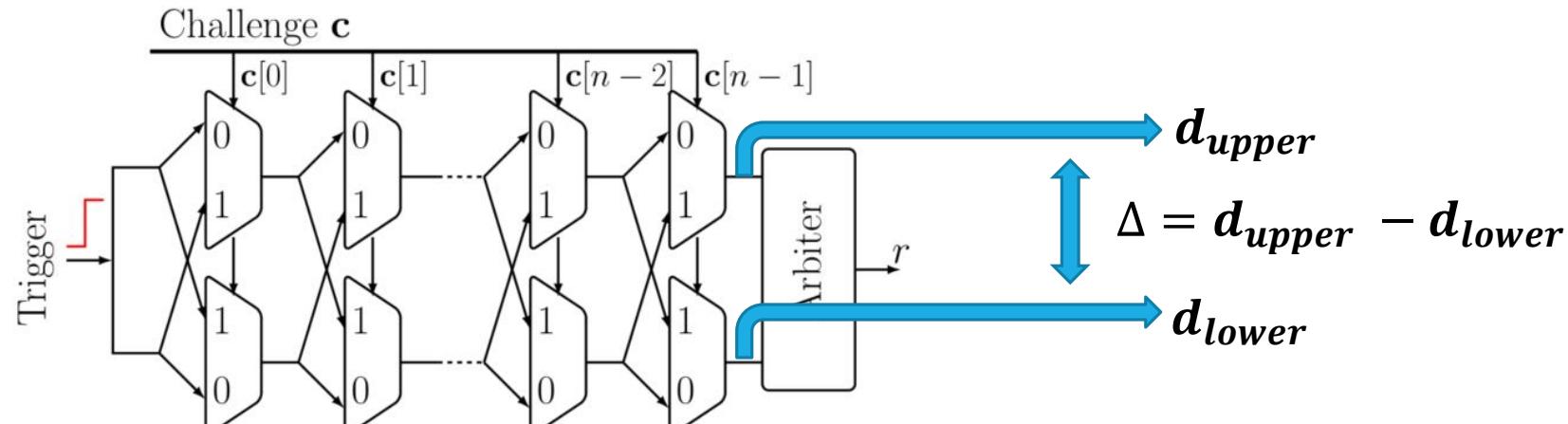
- Blaise Gassend, Dwaine E. Clarke, Marten van Dijk, Srinivas Devadas:
Silicon physical random functions. ACM Conference on Computer and Communications Security 2002: 148-160

Delays in APUFs



$$\begin{aligned}\Delta(i) &= \delta_t(i) - \delta_b(i) \\&= \Delta(i-1)\mathbf{c}[i] + \alpha_i\mathbf{c}[i] + \beta_i \\ \alpha_i &= \frac{p_i - q_i + r_i - s_i}{2} \\ \beta_i &= \frac{p_i - q_i - r_i + s_i}{2} \\ \mathbf{c}[i] &\in \{0, 1\}\end{aligned}$$

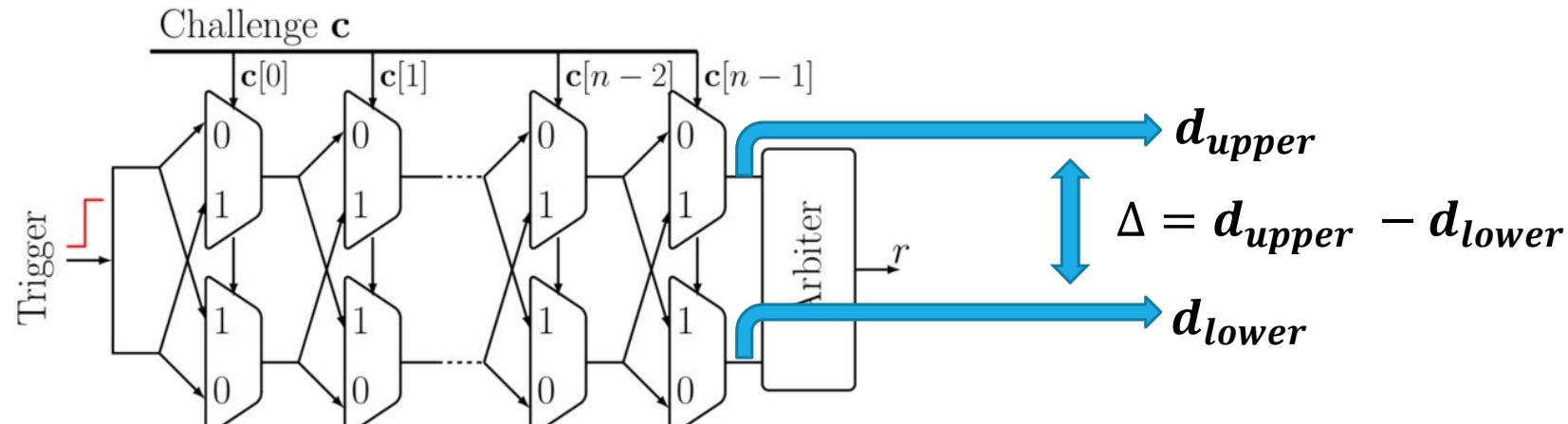
APUF linear delay model



A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.



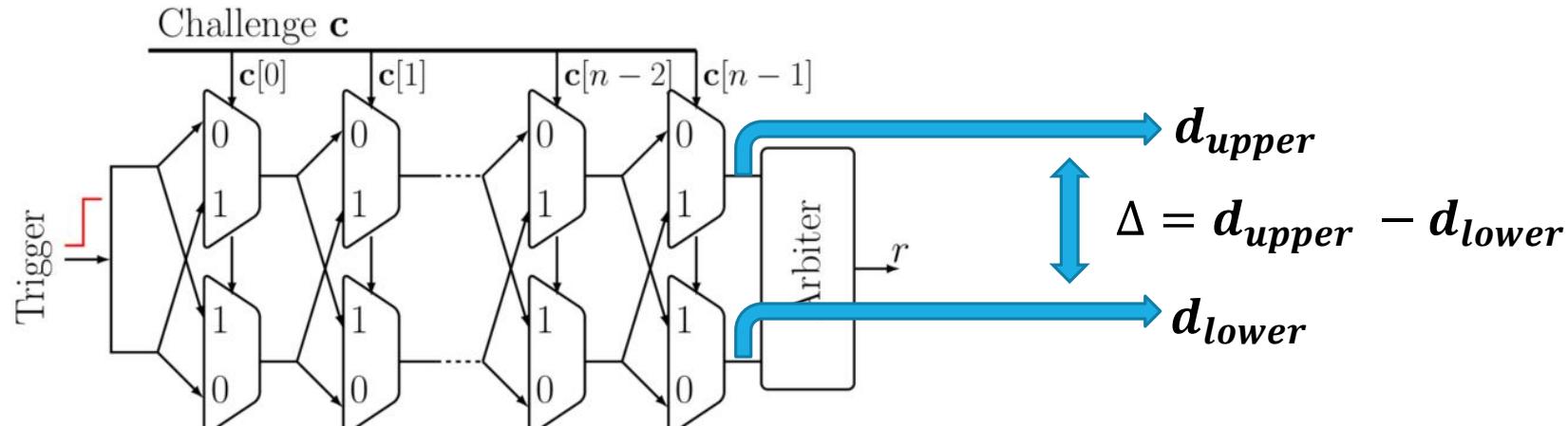
APUF linear delay model



A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.

- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$

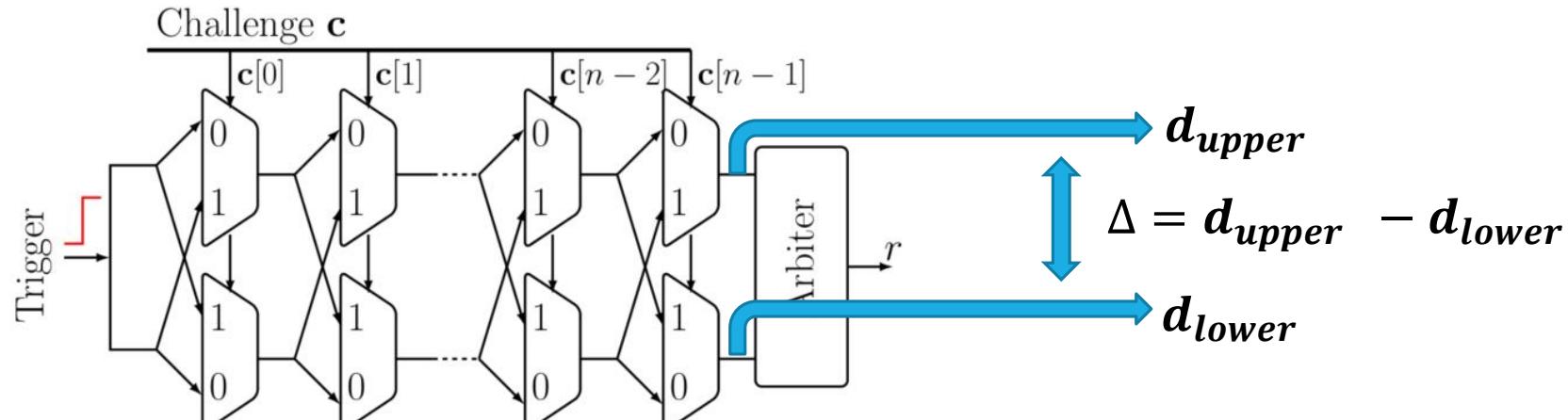
APUF linear delay model



A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.

- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- \mathbf{w} is the weight vector, unique for each APUF instance and represents the delay of a given APUF instance

APUF linear delay model



A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.

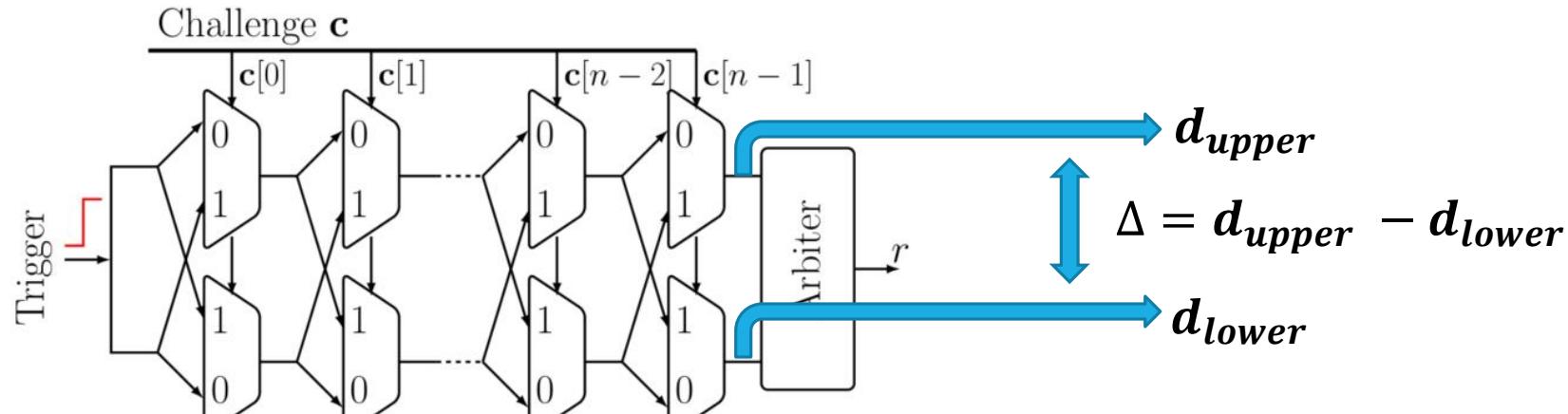
- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- $\mathbf{w}, \mathbf{w}[i] = \alpha_{i+1} + \beta_i, i = 1, \dots, n-1, \mathbf{w}[0] = \alpha_0, \mathbf{w}[n] = \beta_n$
- Φ is the parity vector

$$\Phi[i] = \prod_{j=i, \dots, n-1} (1 - c[j]), i = 0, \dots, n-1$$

$$\Phi[n] = 1$$



APUF linear delay model



A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0,1\}^n$.

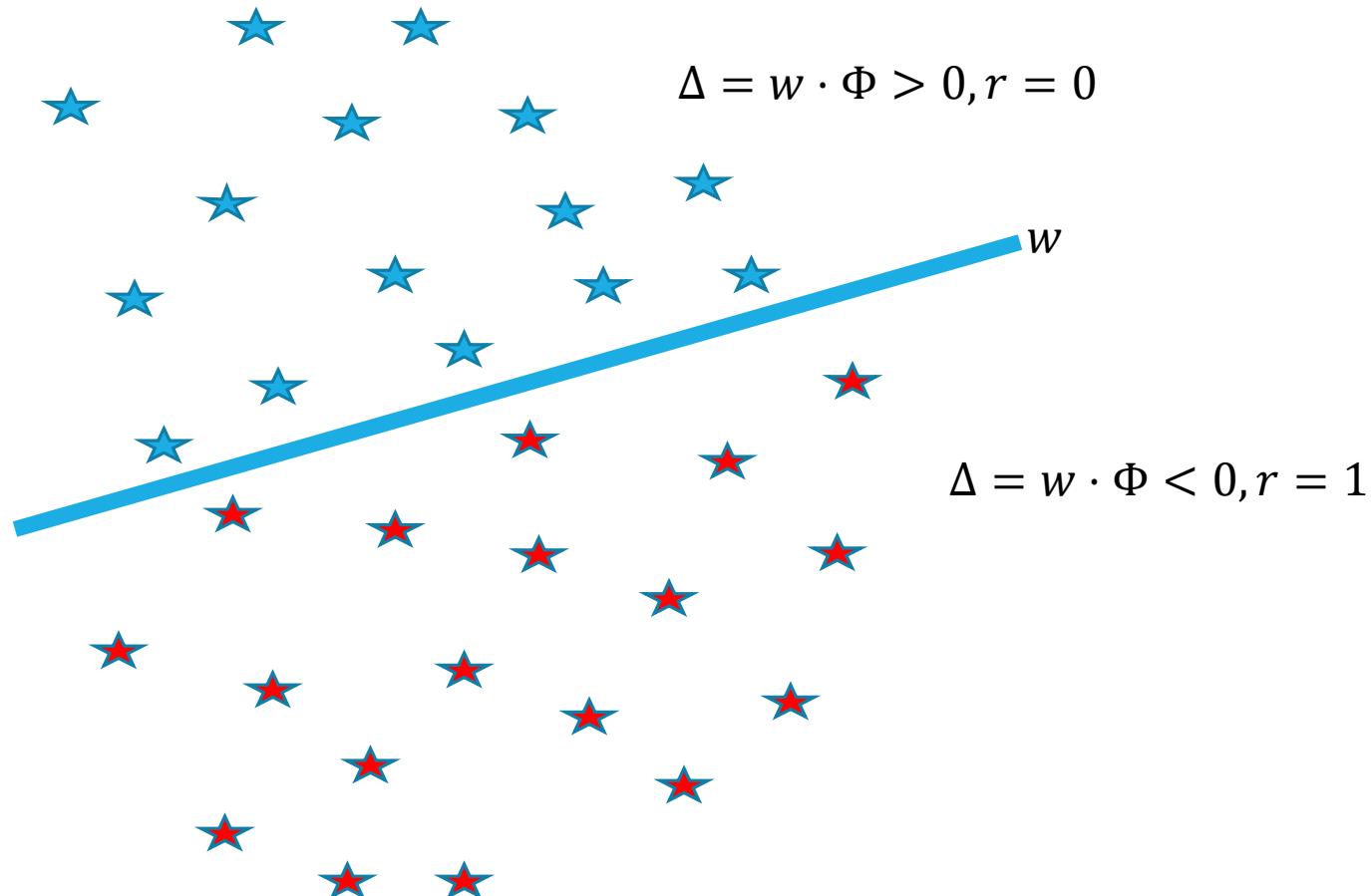
- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$
- \mathbf{w} is the weight vector, unique for each APUF instance and represents the delay of a given APUF instance
- Φ is the parity vector

$$\Phi[i] = \prod_{j=i, \dots, n-1} (1 - c[j]), i = 0, \dots, n-1$$

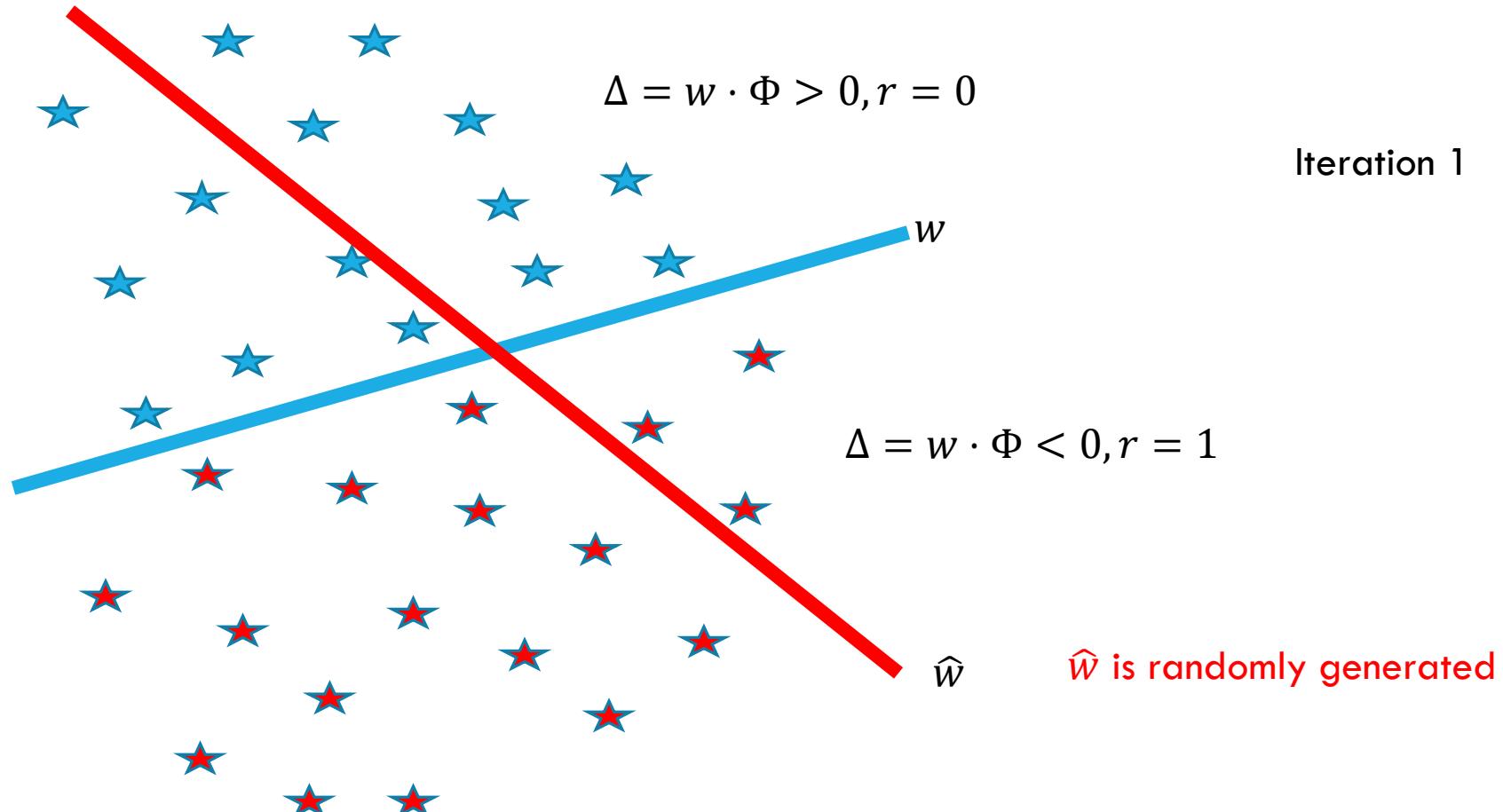
$$\Phi[n] = 1$$

Security: it can be mathematically modeled if the adversary can access a set of CRPs of its by using machine learning techniques

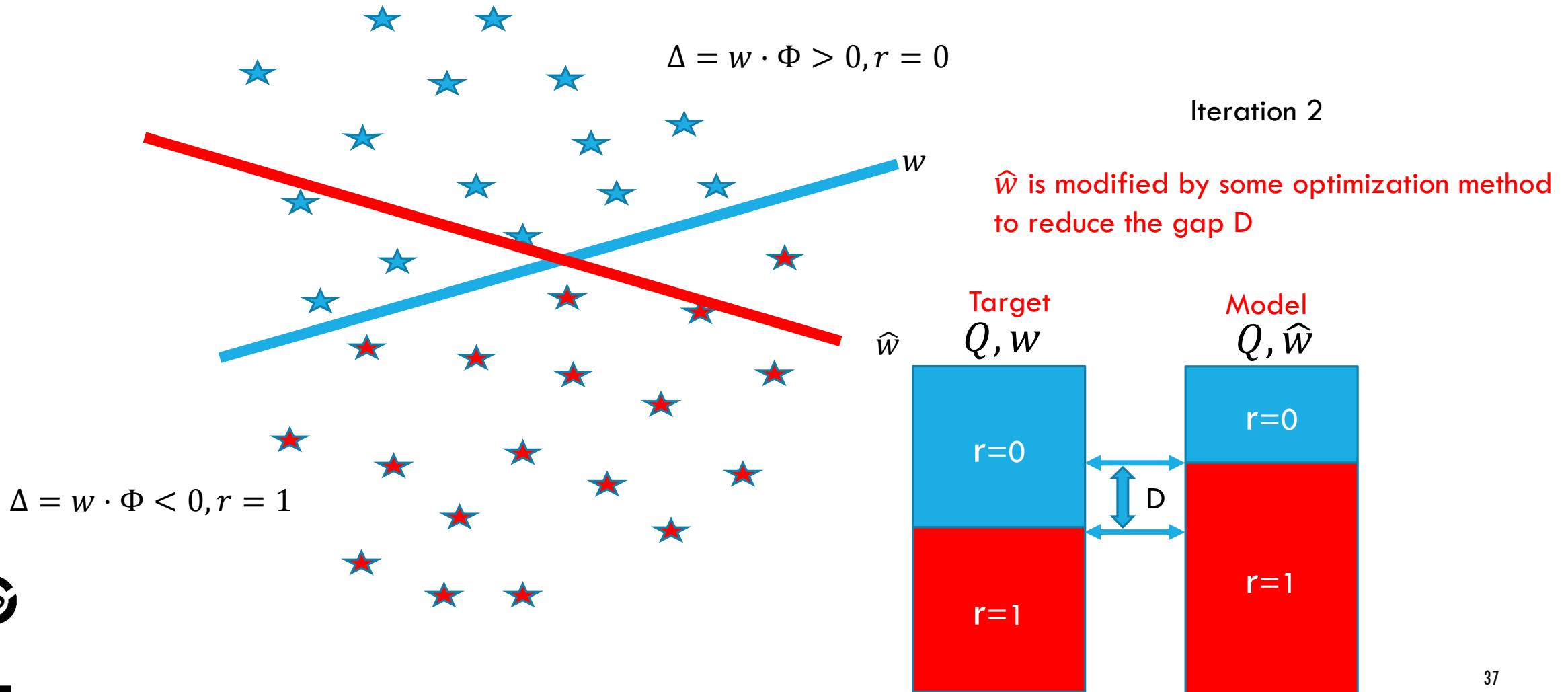
Machine Learning based modeling attack on APUF



Machine Learning based modeling attack on APUF



Machine Learning based modeling attack on APUF



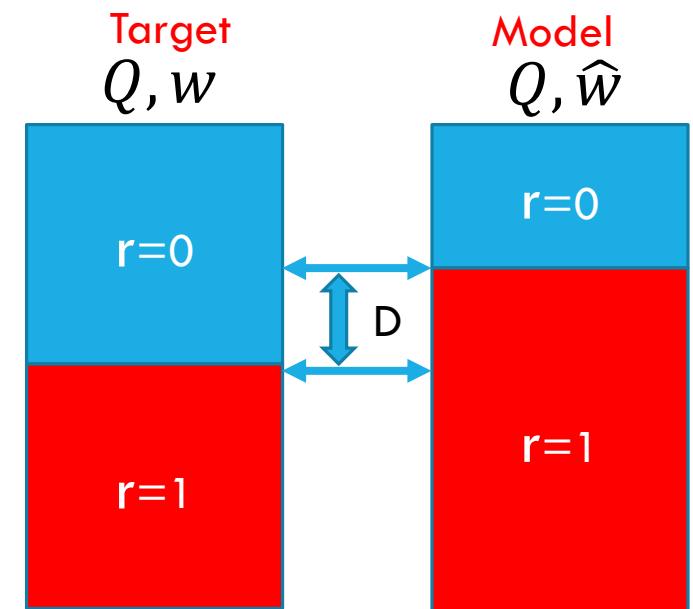
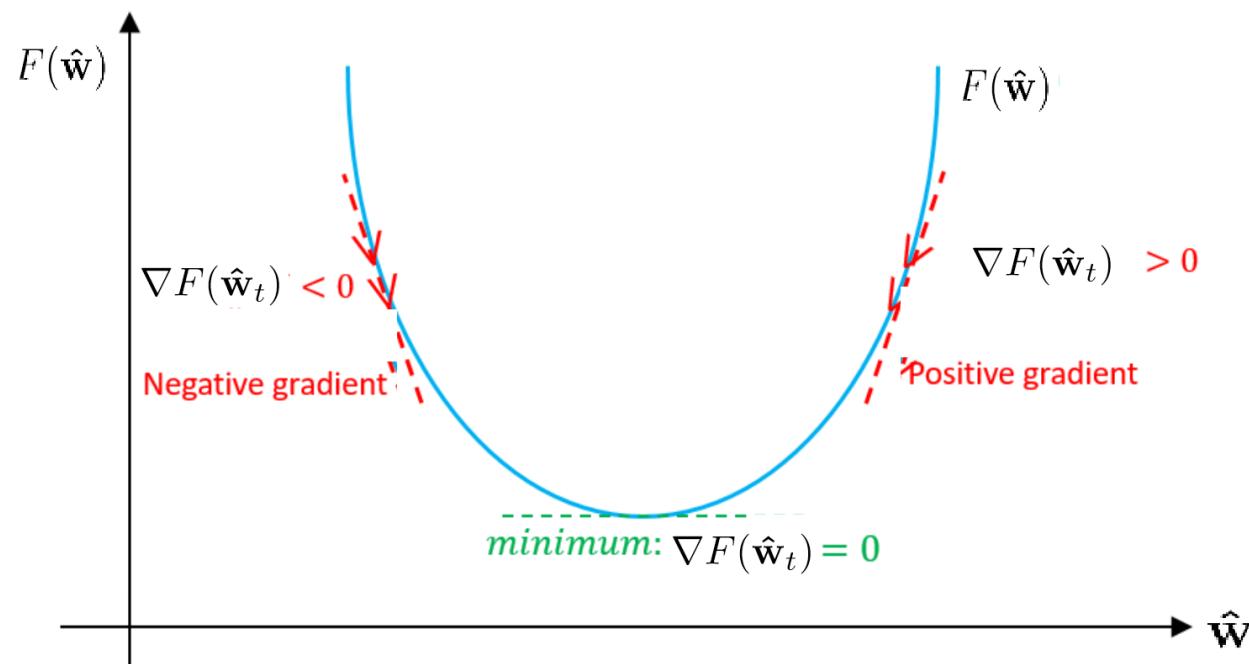
Machine Learning based modeling attack on APUF

$$\min_{\hat{\mathbf{w}}} F(\hat{\mathbf{w}}) = \frac{1}{N} \sum_{i=1}^N L(\hat{\mathbf{w}}, (\mathbf{c}_i, r_i)),$$

where L is the loss function and $\{(\mathbf{c}_i, r_i)_{i=1,\dots,N}\}$ is the set of N CRPs.

$\hat{\mathbf{w}}$ is modified by some optimization method to reduce the gap D

To find $\hat{\mathbf{w}}$, gradient based optimization methodology can be applied, e.g.,



Machine Learning based modeling attack on APUF

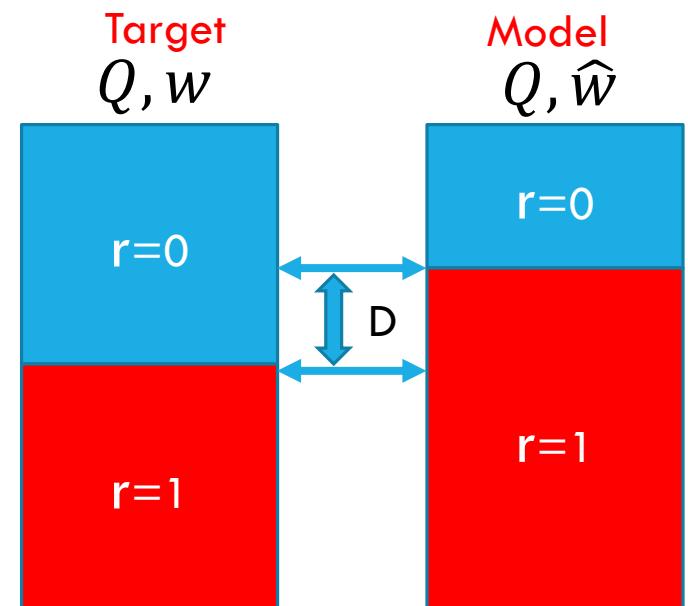
$$\min_{\hat{\mathbf{w}}} F(\hat{\mathbf{w}}) = \frac{1}{N} \sum_{i=1}^N L(\hat{\mathbf{w}}, (\mathbf{c}_i, r_i)),$$

$\hat{\mathbf{w}}$ is modified by some optimization method to reduce the gap D

where L is the loss function and $\{(\mathbf{c}_i, r_i)_{i=1,\dots,N}\}$ is the set of N CRPs.

To find $\hat{\mathbf{w}}$, gradient based optimization methodology can be applied, e.g.,

$$\begin{aligned}\hat{\mathbf{w}}_{t+1} &= \hat{\mathbf{w}}_t - \eta \nabla F(\hat{\mathbf{w}}_t) \\ \nabla F(\hat{\mathbf{w}}_t) &= \left(\frac{\partial F}{\partial \hat{\mathbf{w}}_t[0]}, \frac{\partial F}{\partial \hat{\mathbf{w}}_t[1]}, \dots, \frac{\partial F}{\partial \hat{\mathbf{w}}_t[n-1]} \right)\end{aligned}$$



$$\min_{\hat{\mathbf{w}}} F(\hat{\mathbf{w}}) = \frac{1}{N} \sum_{i=1}^N L(\hat{\mathbf{w}}, (\mathbf{c}_i, r_i)),$$

$\hat{\mathbf{w}}$ is modified by some optimization method to reduce the gap D

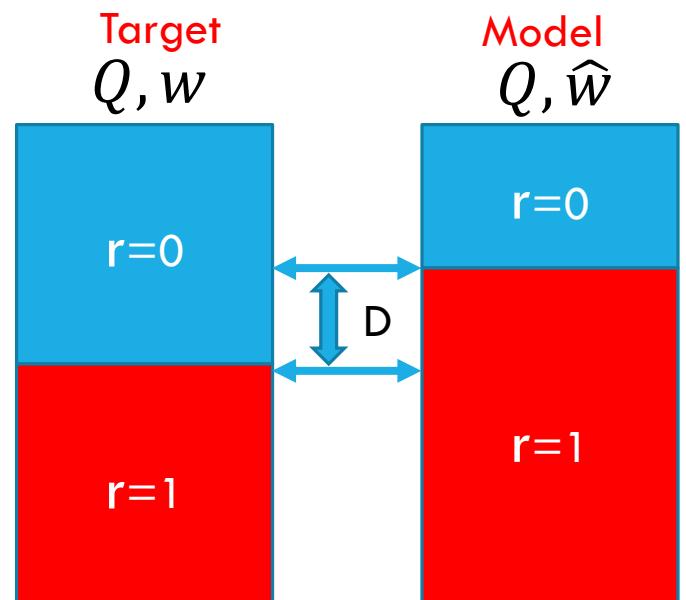
where L is the loss function and $\{(\mathbf{c}_i, r_i)_{i=1,\dots,N}\}$ is the set of N CRPs.

To find $\hat{\mathbf{w}}$, gradient based optimization methodology can be applied, e.g.,

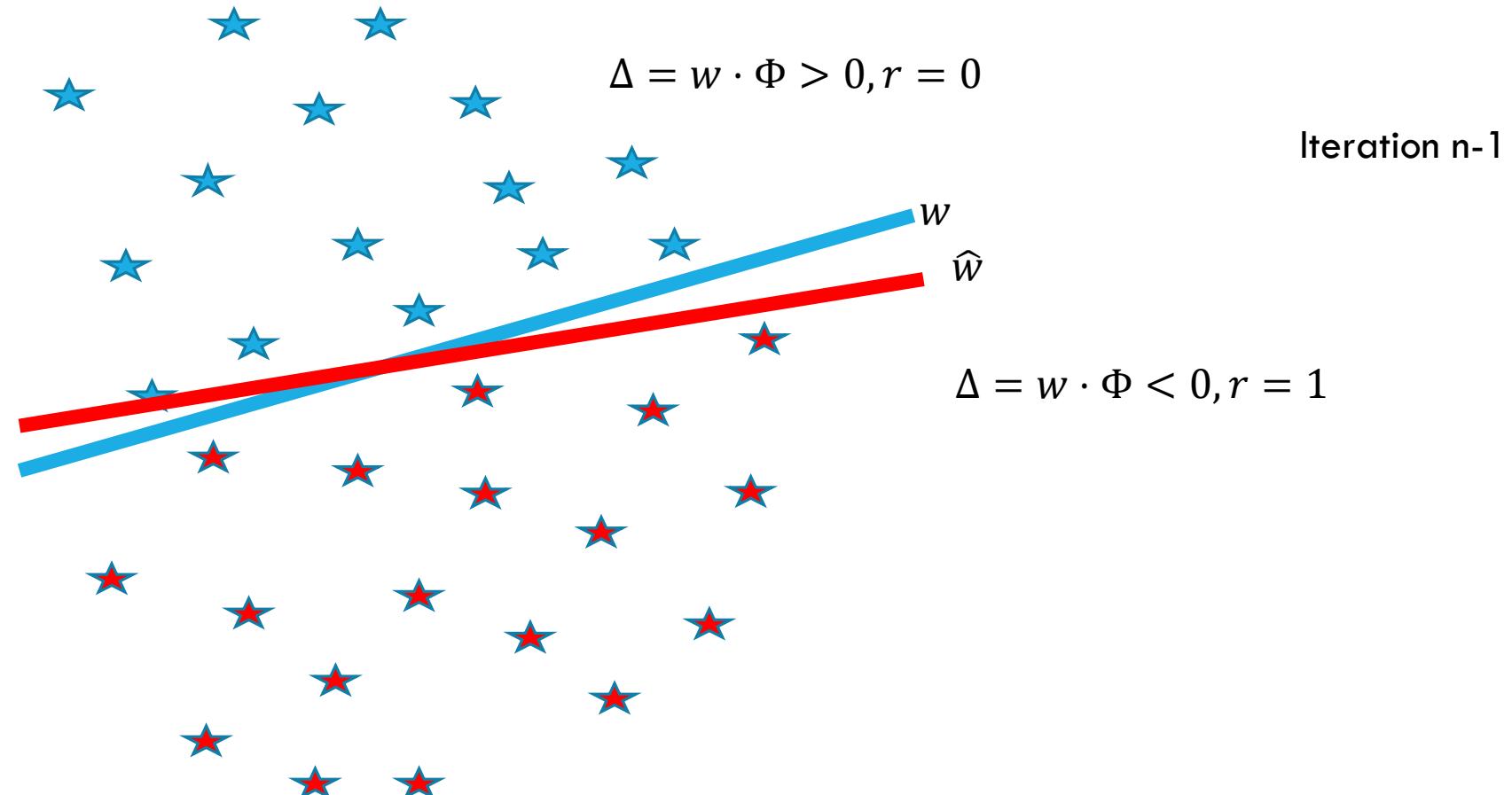
$$\begin{aligned}\hat{\mathbf{w}}_{t+1} &= \hat{\mathbf{w}}_t - \eta \nabla F(\hat{\mathbf{w}}_t) \\ \nabla F(\hat{\mathbf{w}}_t) &= \left(\frac{\partial F}{\partial \hat{\mathbf{w}}_t[0]}, \frac{\partial F}{\partial \hat{\mathbf{w}}_t[1]}, \dots, \frac{\partial F}{\partial \hat{\mathbf{w}}_t[n-1]} \right)\end{aligned}$$

For example, $L(\hat{\mathbf{w}}, (\mathbf{c}_i, r_i)) = \frac{1}{2} \|\hat{\mathbf{w}} \Phi(\mathbf{c}_i) - r_i\|^2$ then

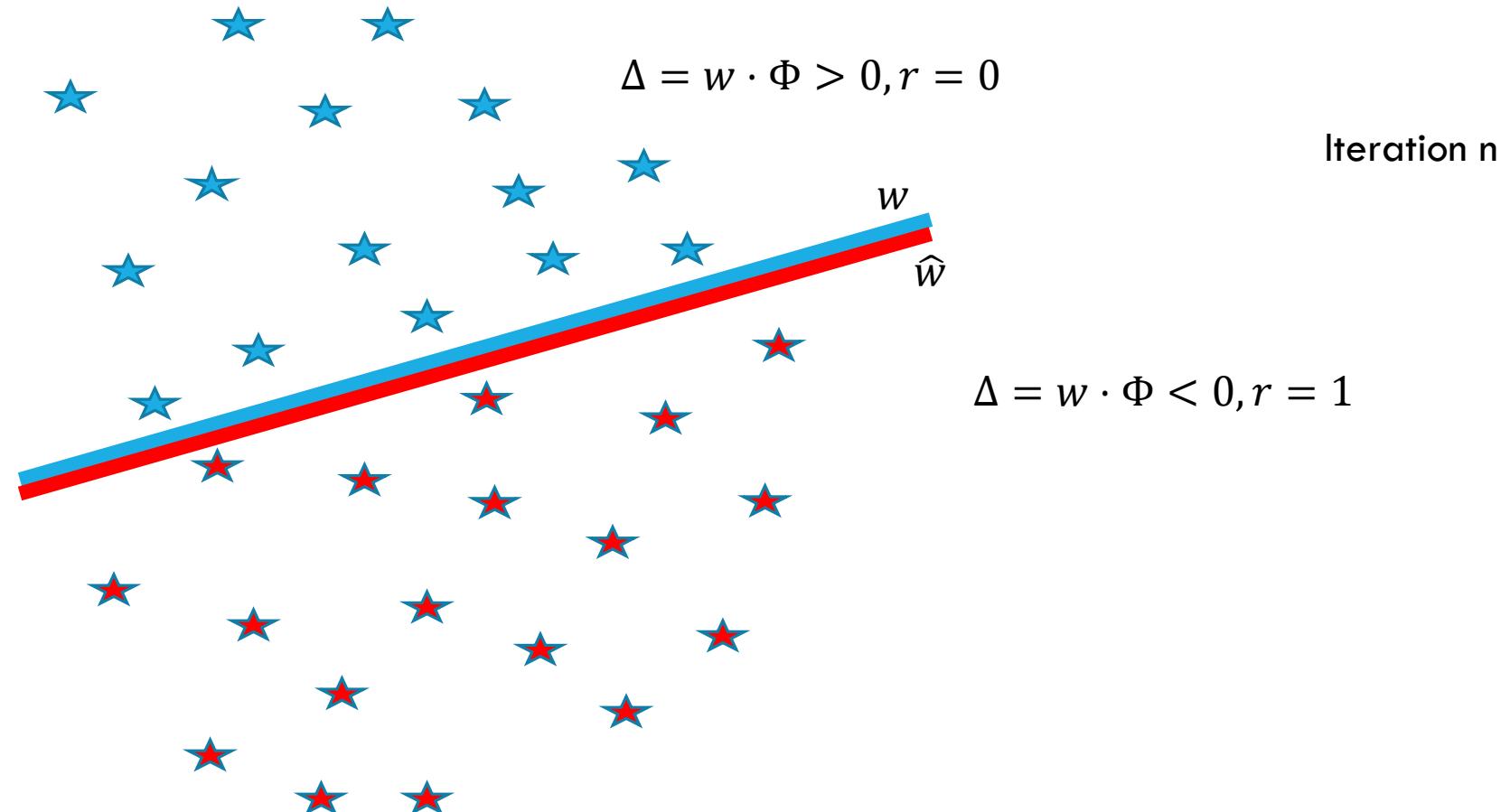
$$\frac{\partial F}{\partial \hat{\mathbf{w}}_t[k]} = \sum_{i=1}^N (\hat{\mathbf{w}}[k] \Phi(\mathbf{c}_i))[k] - r_i) \times \Phi(\mathbf{c}_i)[k], \quad k = 0, \dots, n-1.$$



Machine Learning based modeling attack on APUF



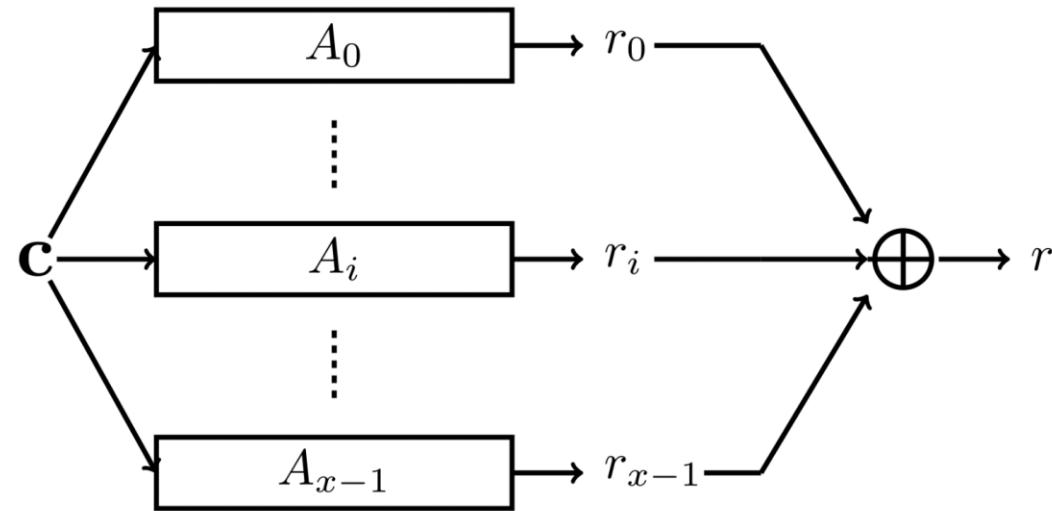
Machine Learning based modeling attack on APUF



5. Classical Machine Learning Attacks on XOR APUFs



x-XOR APUF

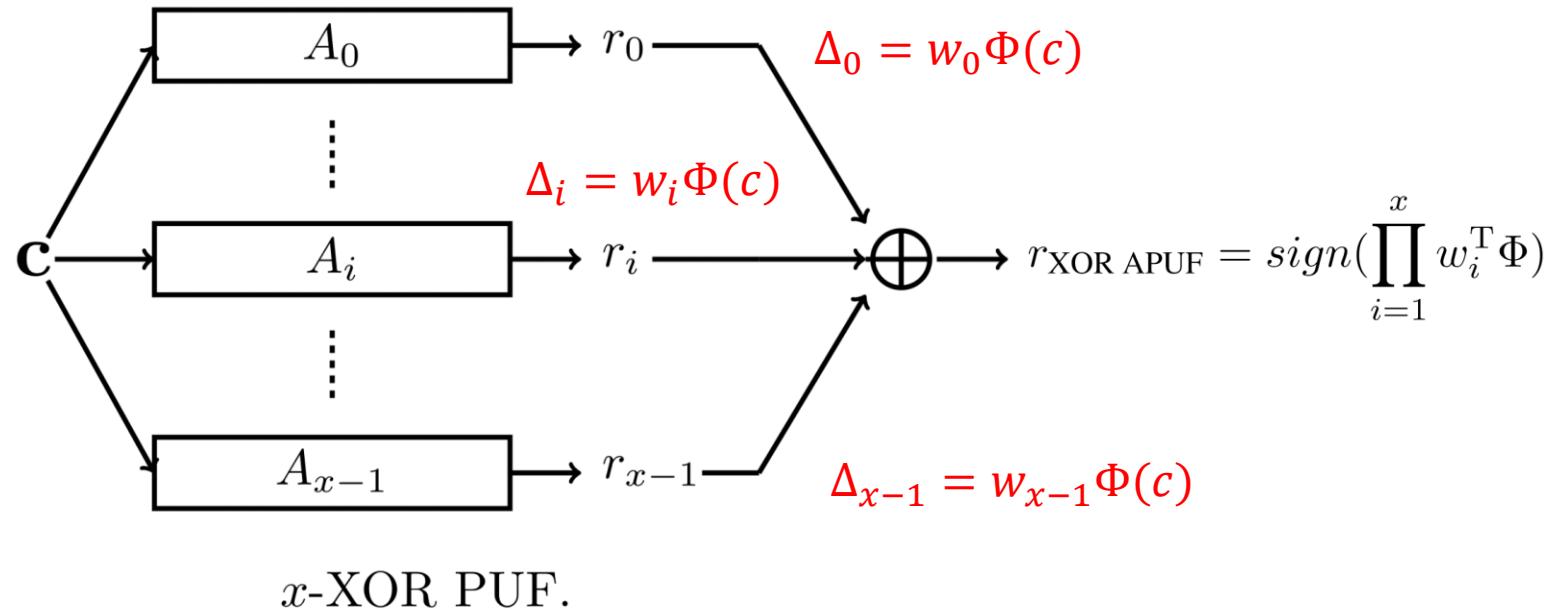


x-XOR PUF.

Physical unclonable functions for device authentication and secret key generation. DAC2007, G. E. Suh and S. Devadas



x-XOR APUF

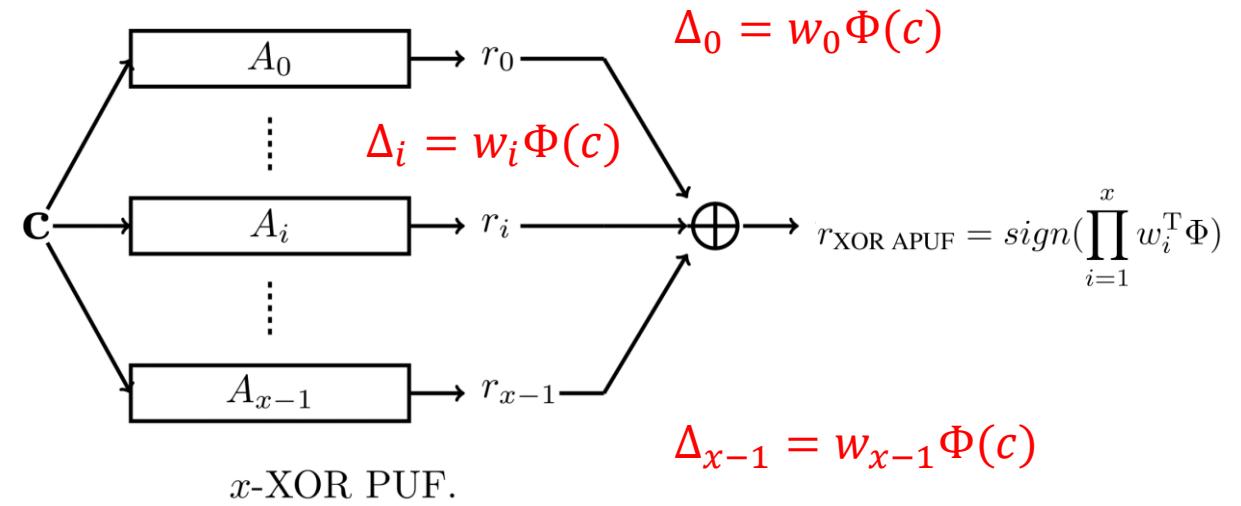


x -XOR PUF.

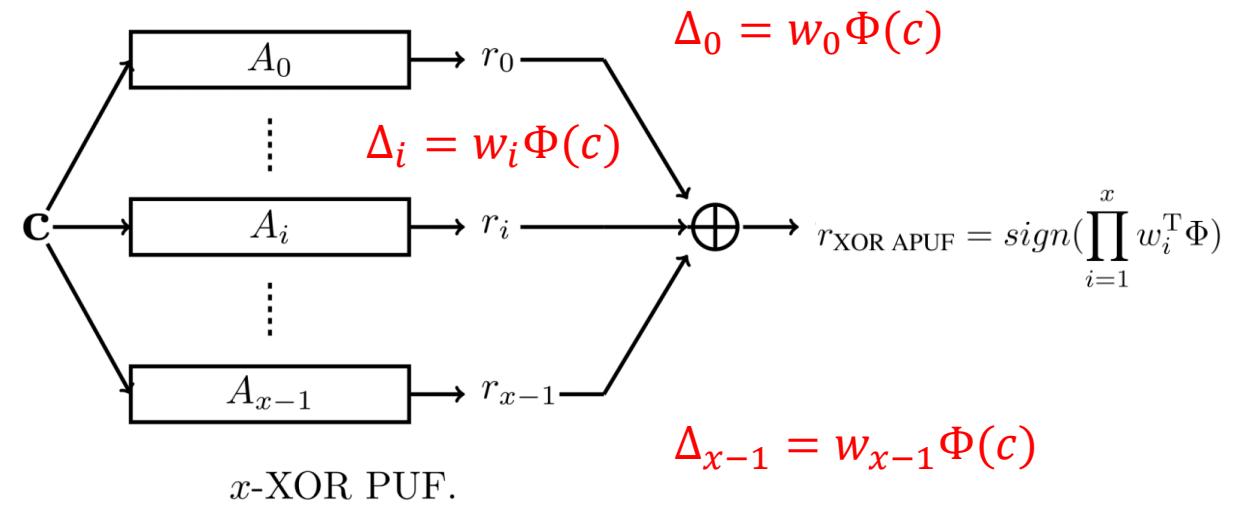
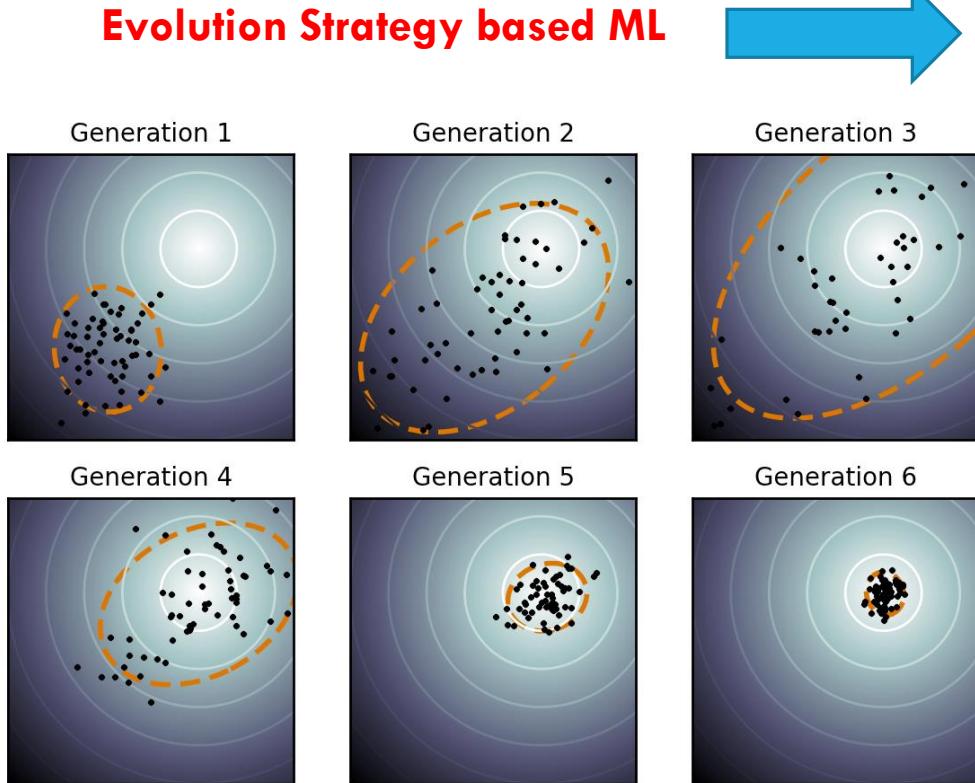
Physical unclonable functions for device authentication and secret key generation. DAC2007, G. E. Suh and S. Devadas



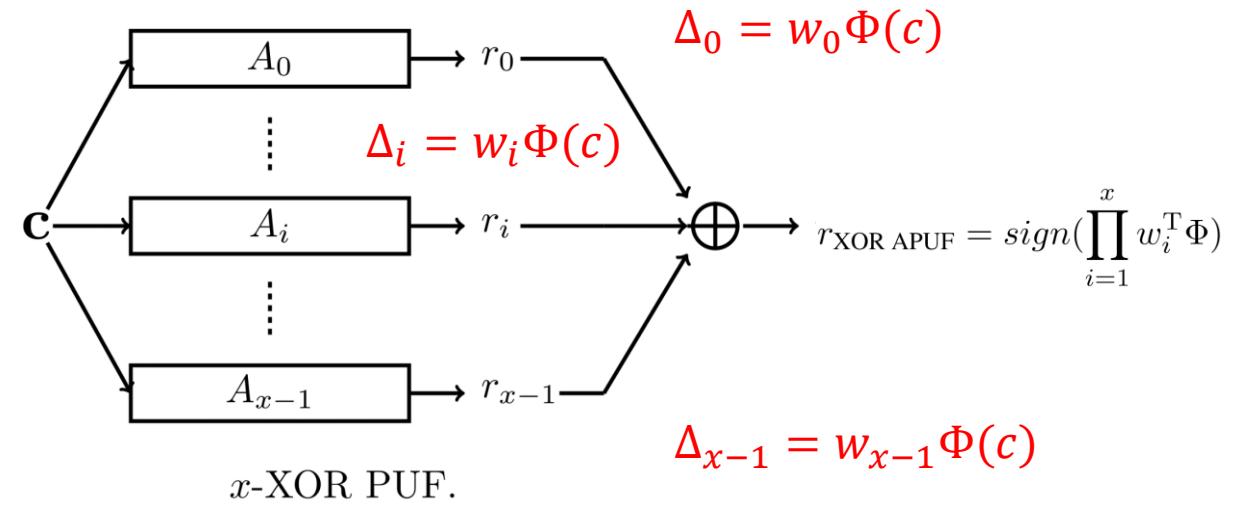
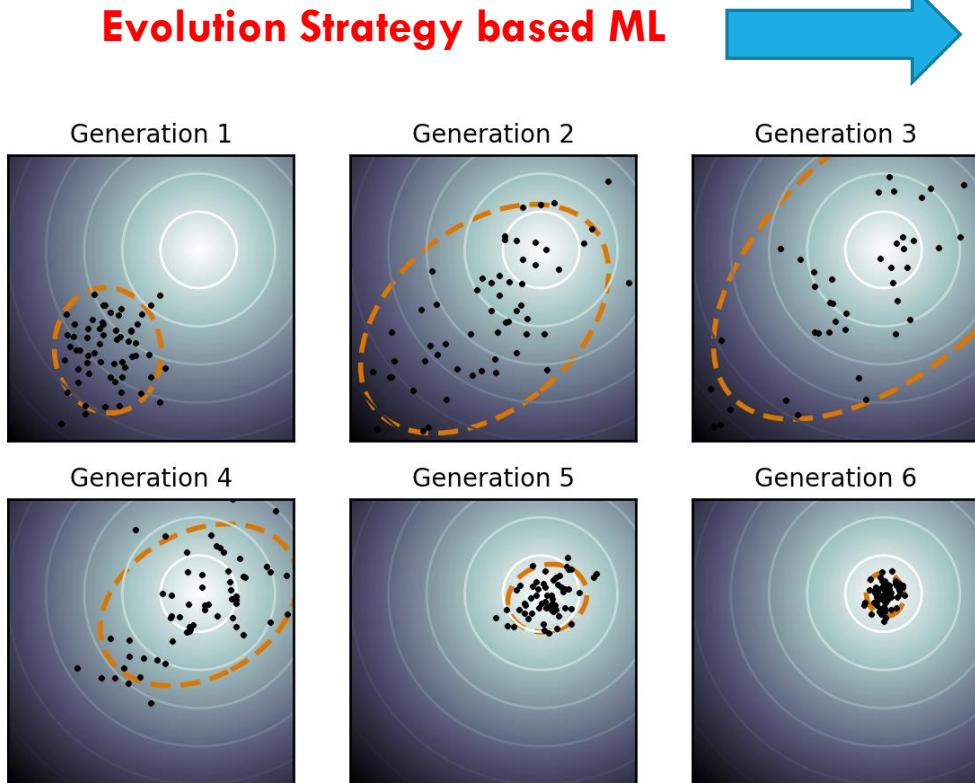
Evolution Strategy based ML



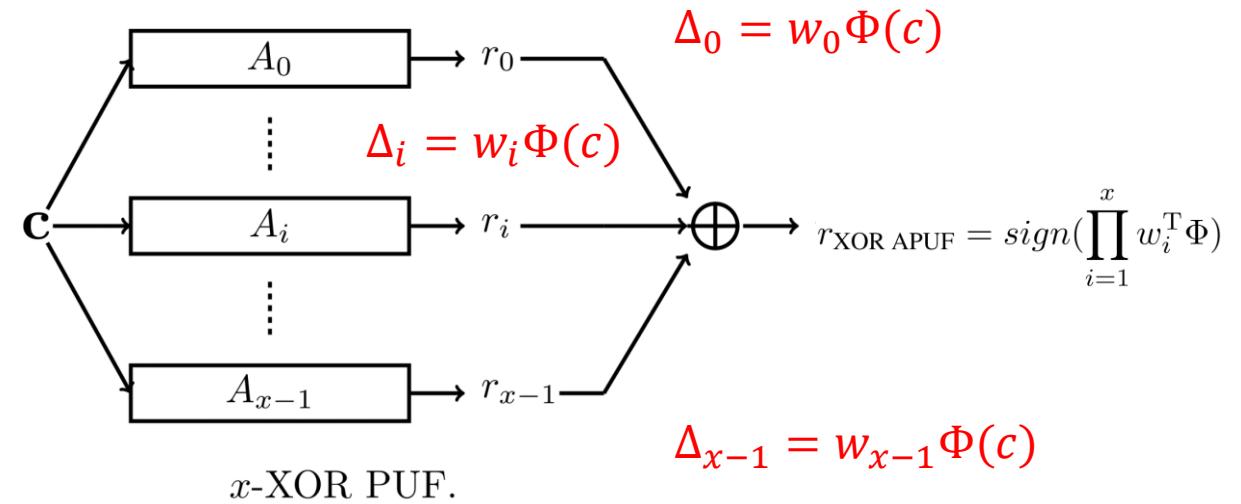
Attacks on x-XOR APUF: Evolution Strategy



Attacks on x-XOR APUF: Evolution Strategy



1. Use the **precise model of XOR APUF**
 2. **Not use gradient for optimization process.**
- Hence, it is **not an optimal method for finding \hat{w}**

1. Evolution Strategy based ML**Logistic Regression**

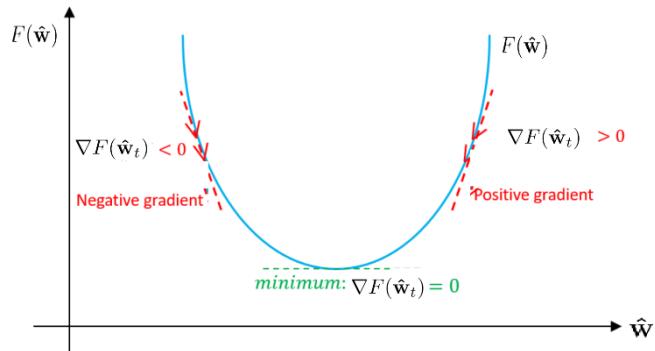
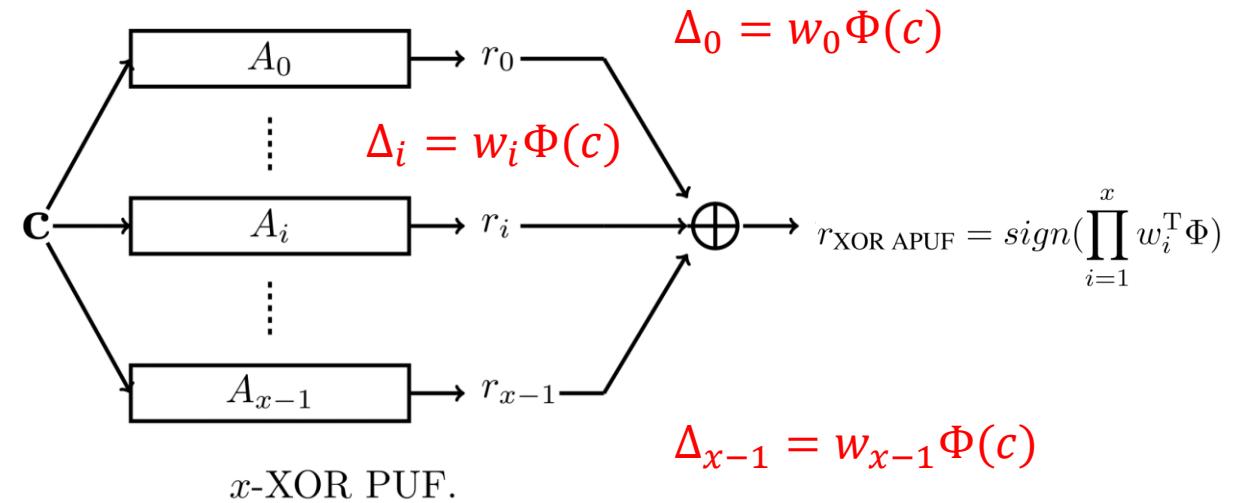
1. Evolution Strategy based ML

Logistic Regression

$$r_{\text{XOR APUF}} = \text{sign}(\prod_{i=1}^x w_i \Phi(\mathbf{c}))$$

$$\approx \sigma(\prod_{i=1}^x w_i \Phi(\mathbf{c}))$$

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



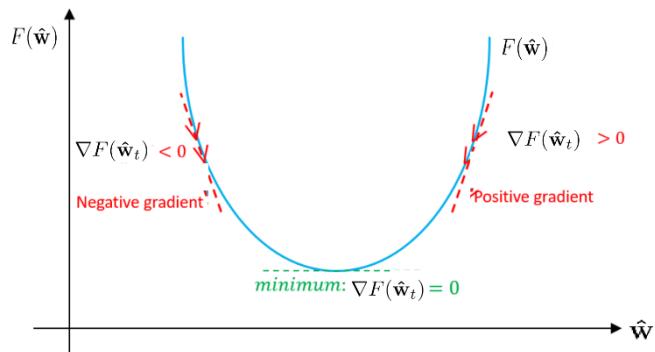
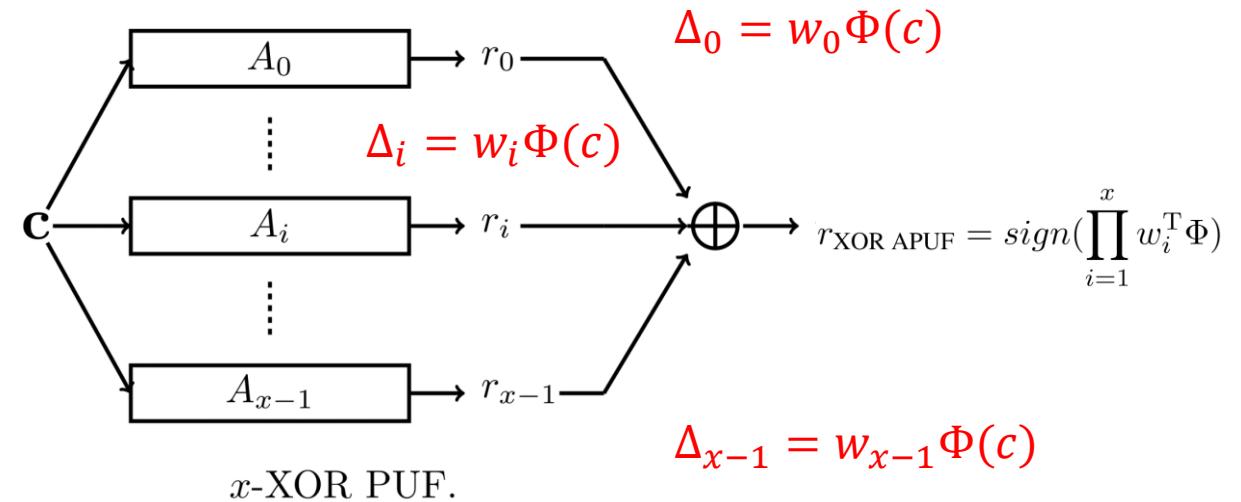
1. Evolution Strategy based ML

Logistic Regression

$$r_{\text{XOR APUF}} = \text{sign}(\prod_{i=1}^x w_i \Phi(c))$$

$$\approx \sigma(\prod_{i=1}^x w_i \Phi(c))$$

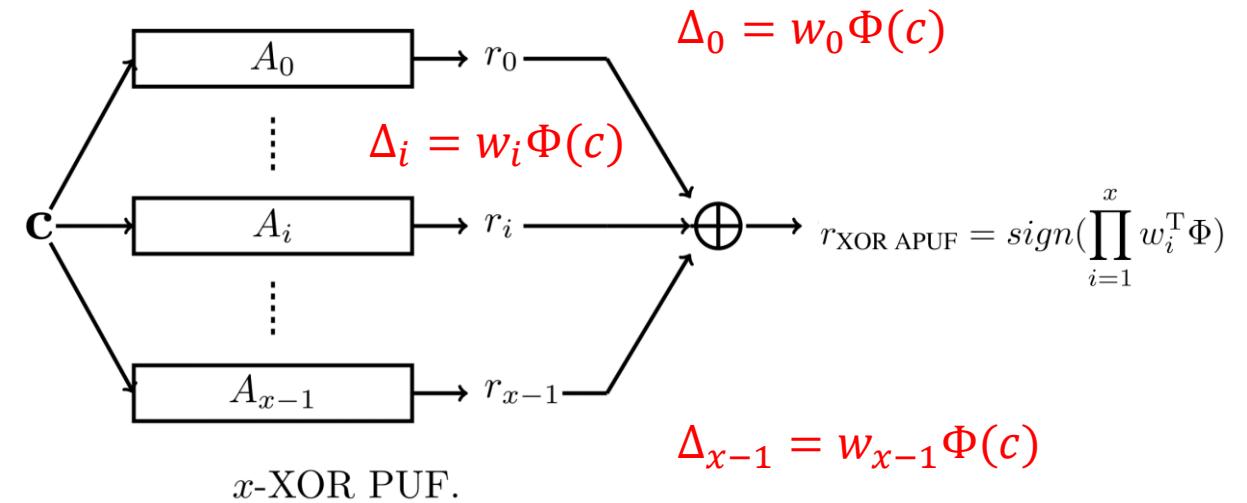
$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$



1. Use the **precise** model of XOR APUF
 2. Use **gradient** for optimization process.
- Hence, it is **an optimal** method for finding \hat{w}

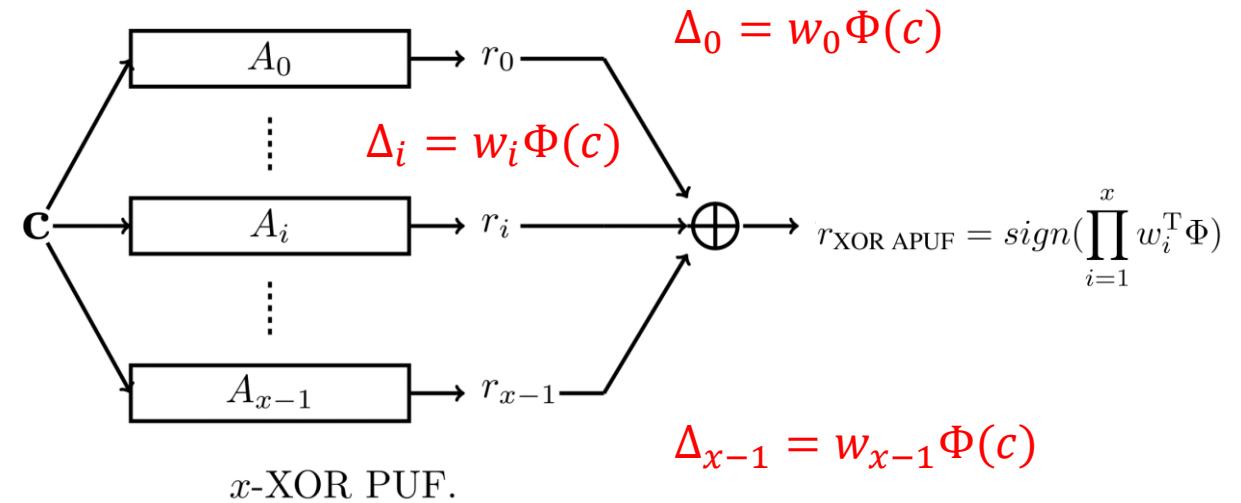
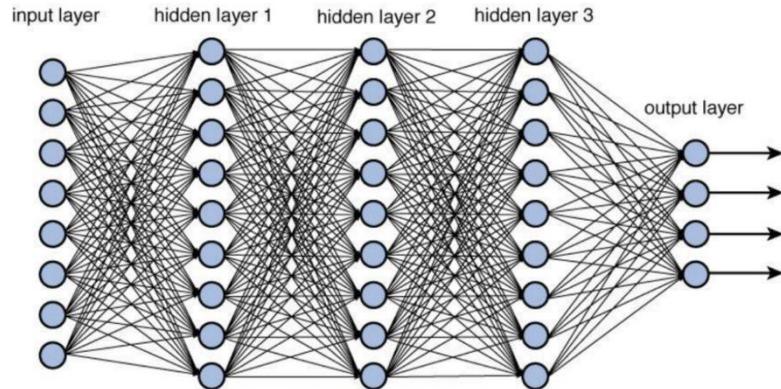
1. Evolution Strategy based ML
2. Logistic Regression based ML

Deep Neural Network



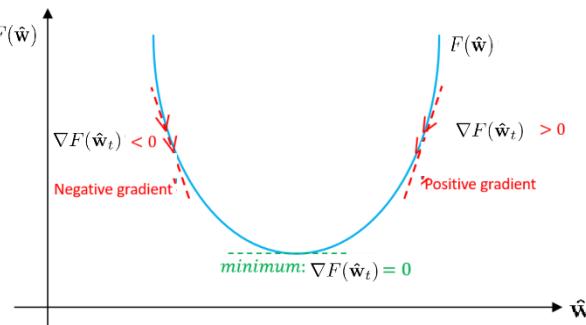
1. Evolution Strategy based ML
2. Logistic Regression based ML

Deep Neural Network



$$r_{\text{XOR APUF}} = \text{sign}(\prod_{i=1}^x w_i \Phi(\mathbf{c}))$$

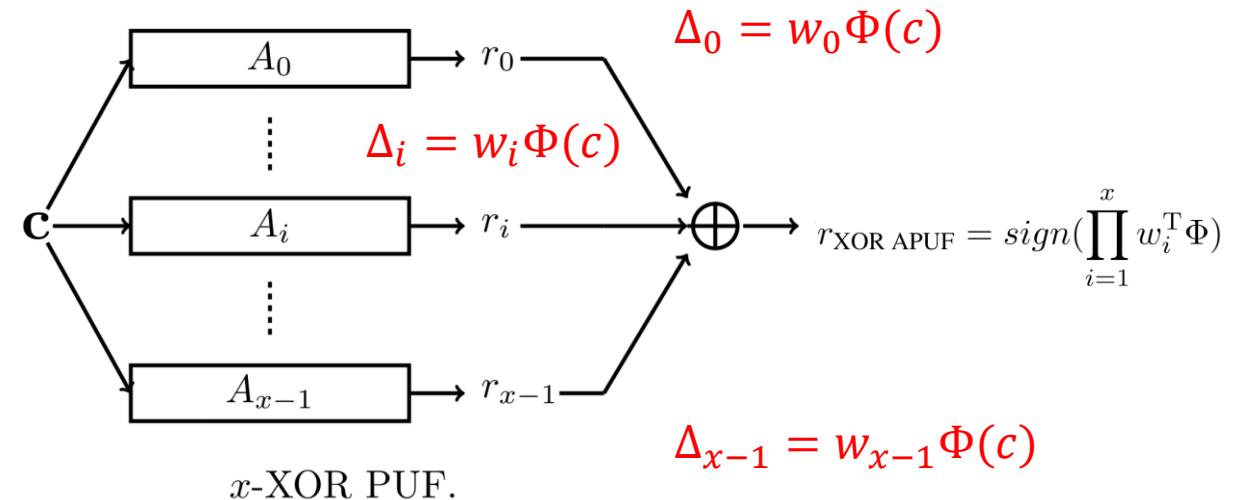
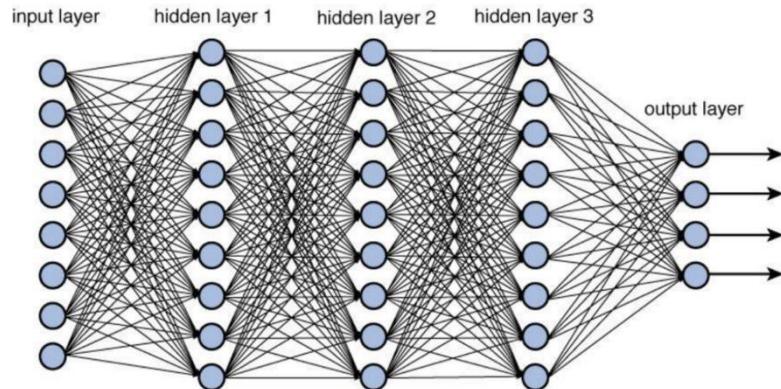
$$\approx \text{DNN}_{\mathbf{w}}(\Phi(\mathbf{c}))$$



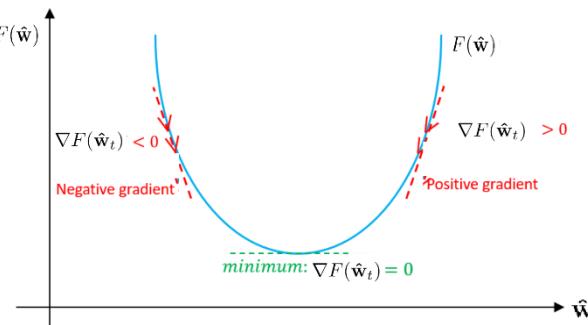
Attacks on x-XOR APUF: Deep Neural Network

1. Evolution Strategy based ML
2. Logistic Regression based ML

Deep Neural Network



$$r_{\text{XOR APUF}} = \text{sign}(\prod_{i=1}^x w_i \Phi(c)) \approx \text{DNN}_w(\Phi(c))$$

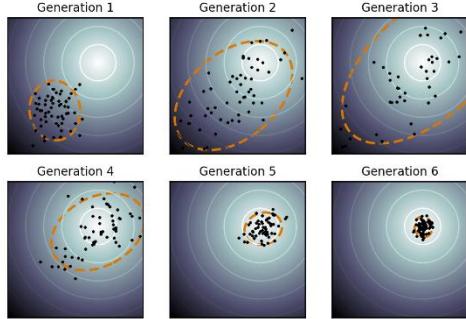


1. Use non-precise model of XOR APUF
 2. Use gradient for optimization process.
- Hence, it is **not optimal** method for finding \hat{w}



Attacks on x-XOR APUF: Comparison

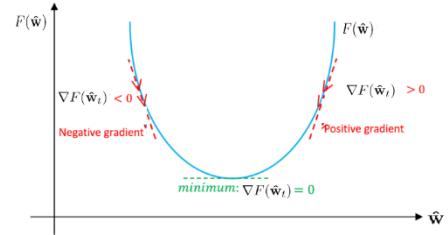
Evolution Strategy based ML



$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x w_i^T \Phi\right)$$

1. Use **precise model of XOR APUF**
2. **Not use gradient** for optimization process.
Hence, it is **not optimal** method for finding \hat{w}

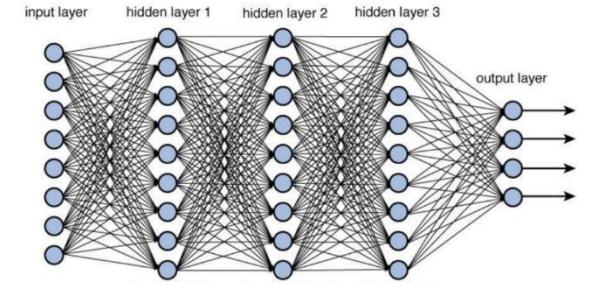
Logistic Regression based ML



$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x w_i^T \Phi\right)$$

1. Use **precise model of XOR APUF**
2. **Use gradient** for optimization process.
Hence, it is **an optimal** method for finding \hat{w}

Deep Neural Network based ML

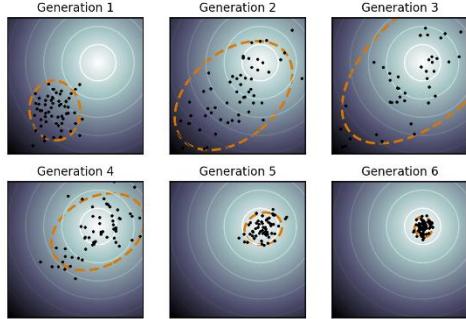


$$\begin{aligned} r_{\text{XOR APUF}} &= \text{sign}\left(\prod_{i=1}^x w_i \Phi(\mathbf{c})\right) \\ &\approx \text{DNN}_{\mathbf{w}}(\Phi(\mathbf{c})) \end{aligned}$$

1. Use **non-precise model of XOR APUF**
2. **Use gradient** for optimization process.
Hence, it is **not optimal** method for finding \hat{w}

Attacks on x-XOR APUF: Comparison

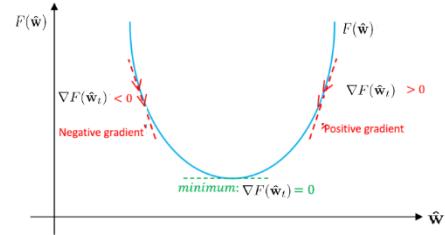
Evolution Strategy based ML



$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x w_i^T \Phi\right)$$

1. Use **precise model of XOR APUF**
2. **Not use gradient** for optimization process.
Hence, it is **not optimal** method for finding \hat{w}

Logistic Regression based ML



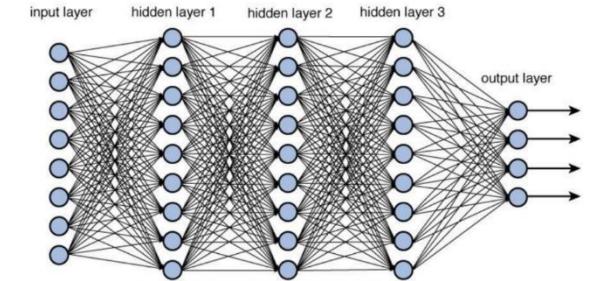
$$r_{\text{XOR APUF}} = \text{sign}\left(\prod_{i=1}^x w_i^T \Phi\right)$$

1. Use **precise model of XOR APUF**
2. **Use gradient** for optimization process.
Hence, it is **an optimal** method for finding \hat{w}



Best Attack

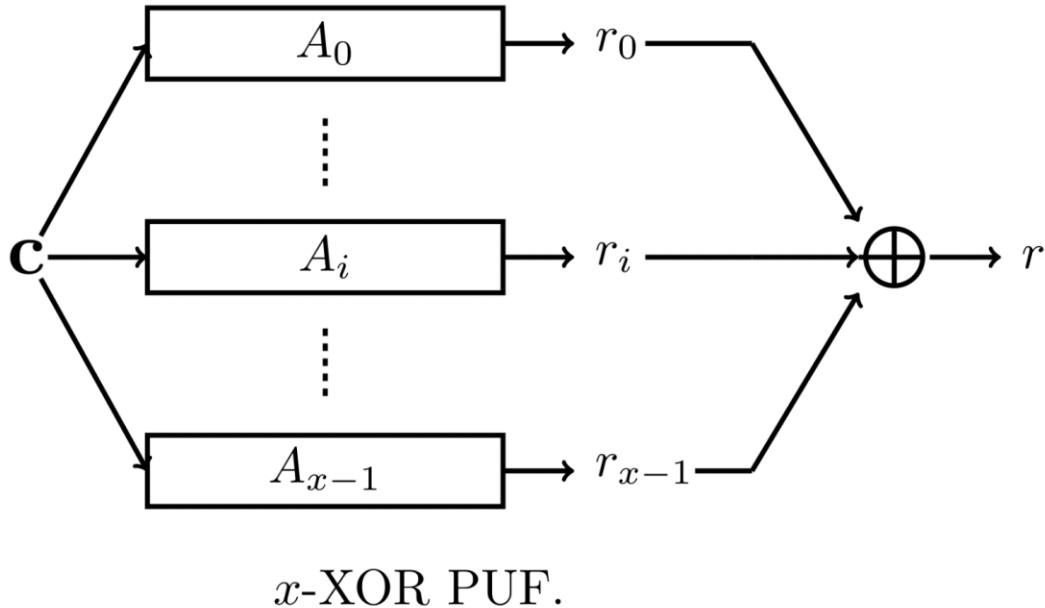
Deep Neural Network based ML



$$\begin{aligned} r_{\text{XOR APUF}} &= \text{sign}\left(\prod_{i=1}^x w_i \Phi(\mathbf{c})\right) \\ &\approx \text{DNN}_{\mathbf{w}}(\Phi(\mathbf{c})) \end{aligned}$$

1. Use **non-precise model of XOR APUF**
2. **Use gradient** for optimization process.
Hence, it is **not optimal** method for finding \hat{w}

x-XOR APUF



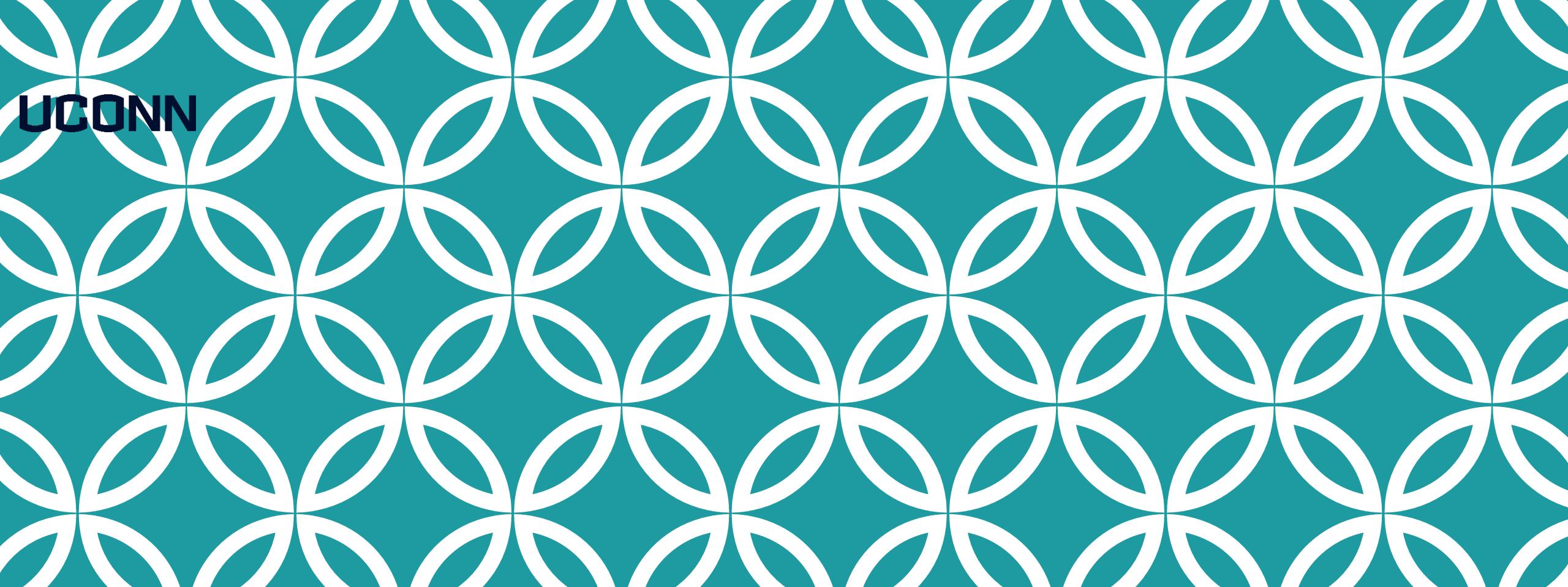
Security: when the number of APUFs is large enough (64-bit 10-XOR APUF), then the CRP-based modeling attacks are not practical anymore.

As shown in [6], in case of Logistic Regression based modelling attacks on n -bit x -XOR PUF has , the number of iterations for the training phase is: $N_{iteration} = O\left(\frac{(n+1)^x}{N_{data} \cdot x!}\right)$ and the total time complexity is:

$$T = N_{iteration} \cdot (n + 1) \cdot x \cdot N_{data}$$

Physical unclonable functions for device authentication and secret key generation. DAC2007, G. E. Suh and S. Devadas

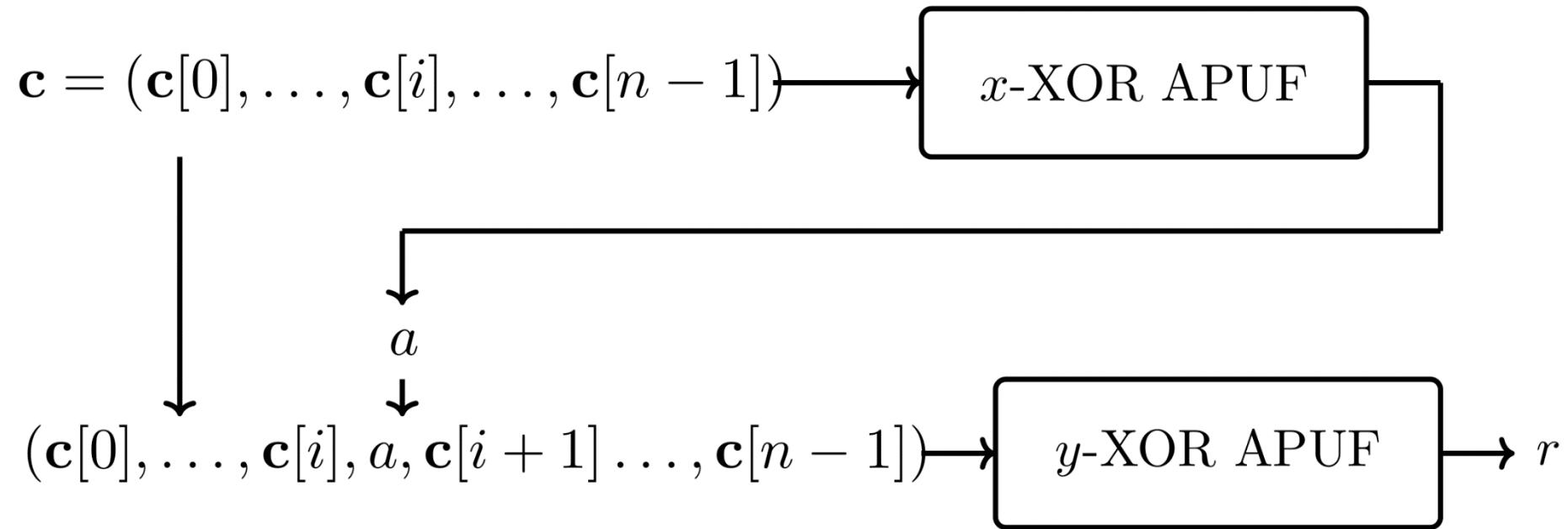


A decorative background pattern consisting of a grid of white, stylized flower or leaf motifs on a teal-colored hexagonal grid.

6. Interpose PUF (iPUF) and Classical Modeling Attacks



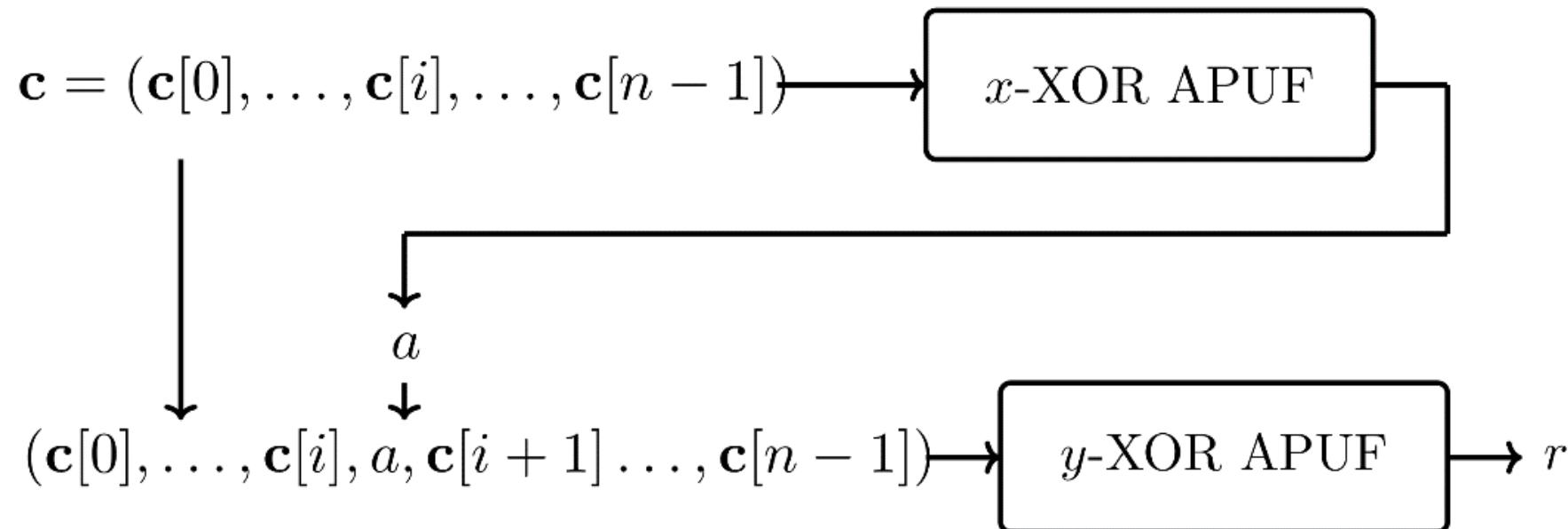
5. Interposed PUF (iPUF)



Traditional
Modeling Attacks

$$(x, y) - IPUF \approx \left(y + \frac{x}{2} \right) - XOR PUF$$

if a is inserted at the middle



Road Map

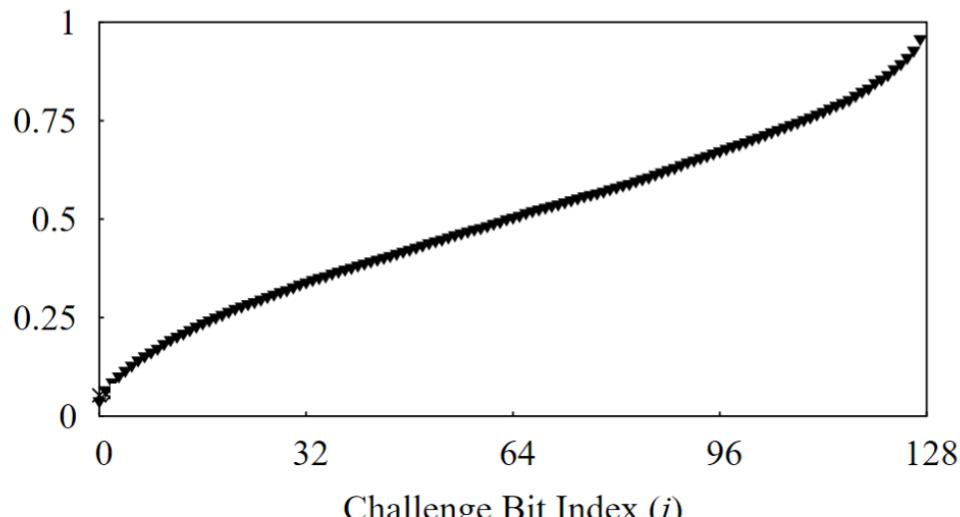
- Strict Avalanche Criterion (SAC) Property of APUF
- Influence of each APUF in XOR APUF



Strict Avalanche Criterion (SAC) Property



Strict Avalanche Criterion (SAC) Property

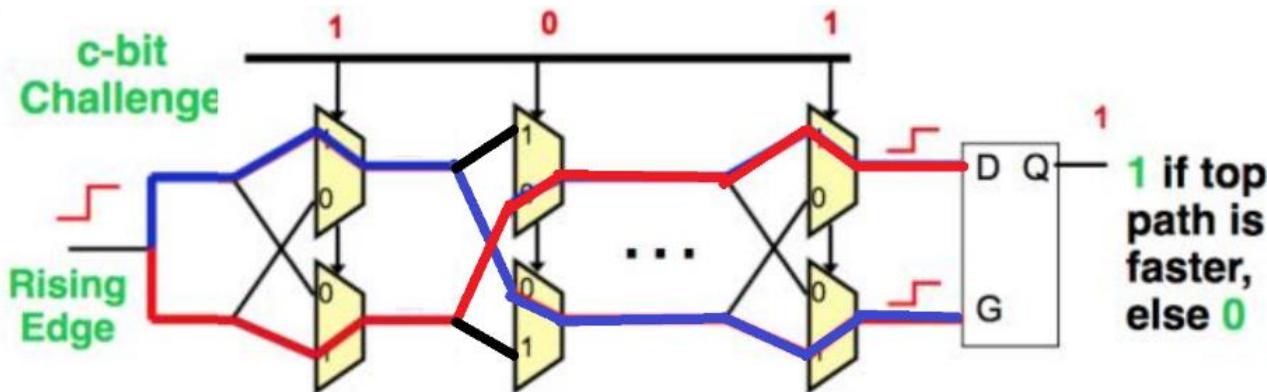
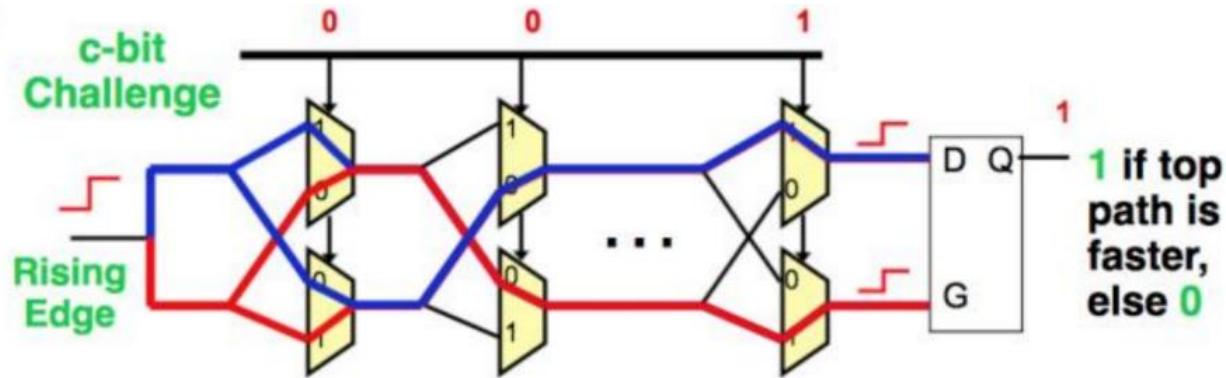


(b) 128-bit APUF

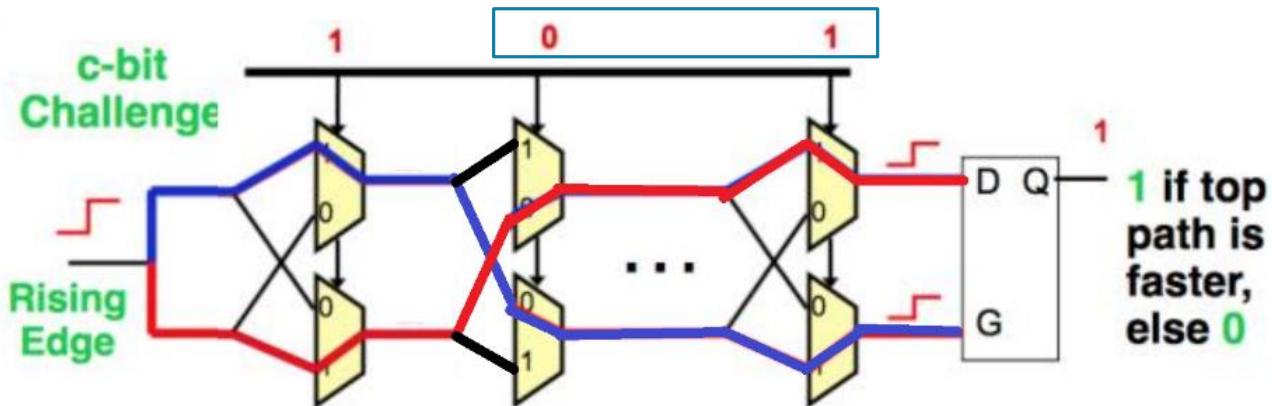
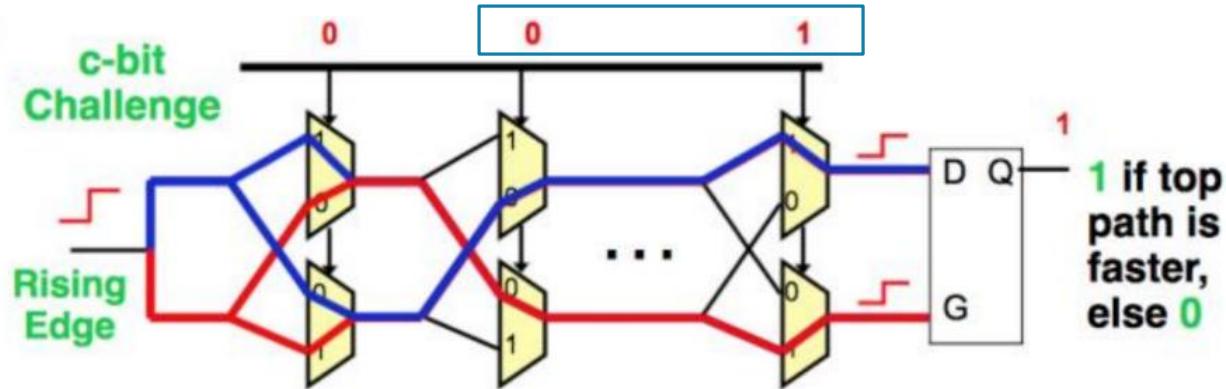
$$\Pr_{\mathbf{c}}(r_0 \neq r_1) \approx \frac{n-i}{n}$$

$$\Pr_{\mathbf{c}}(r_0 \neq r_1) = \frac{1}{2}, i = \frac{n}{2}$$

Position of Challenge bit of APUF

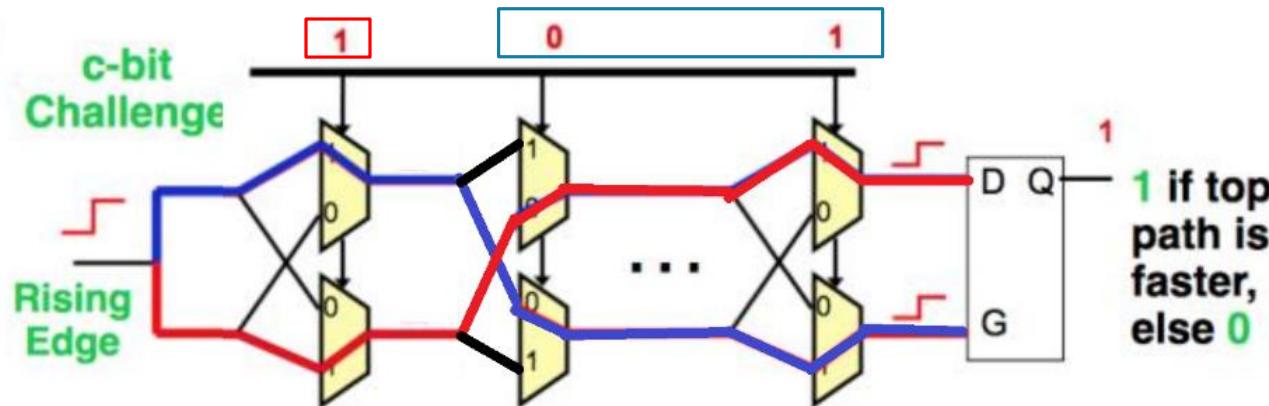
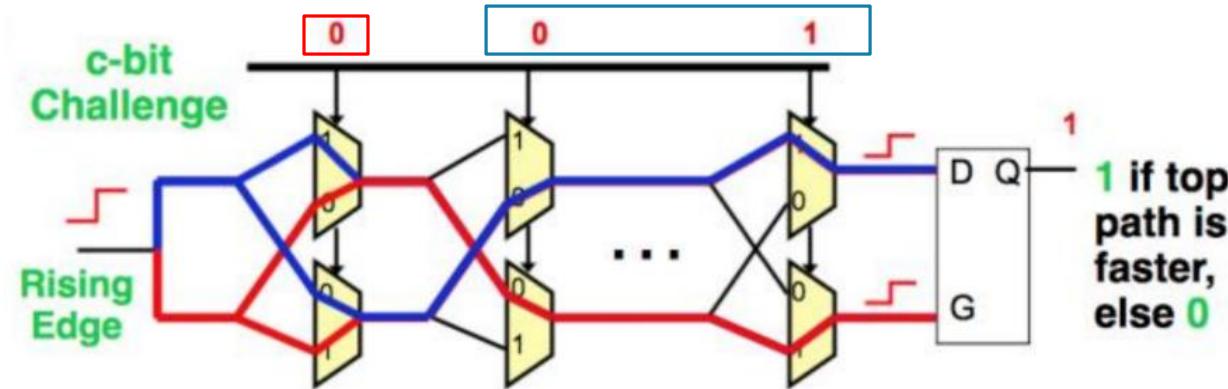


Position of Challenge bit of APUF



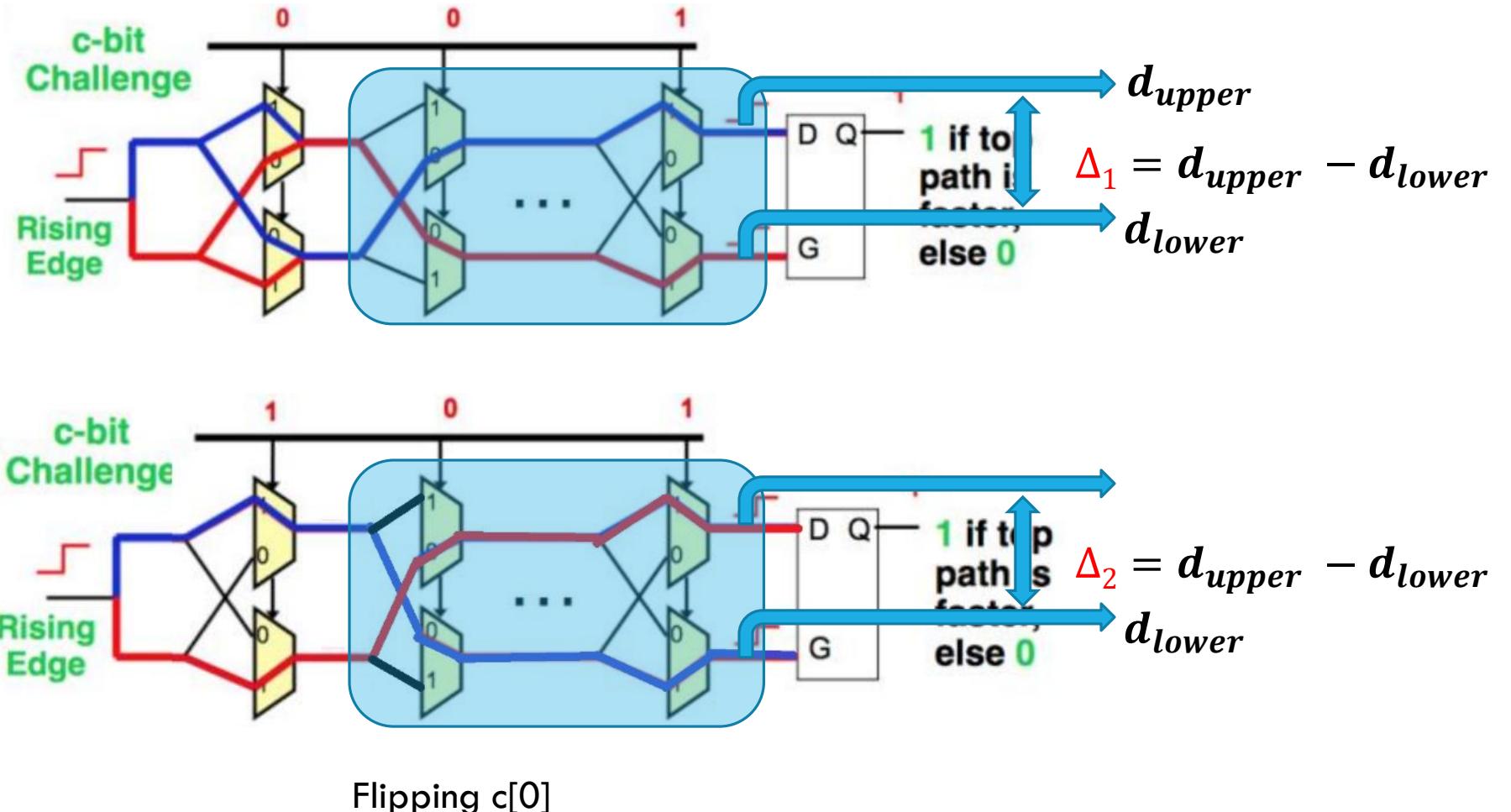
Flipping $c[0]$

Position of Challenge bit of APUF

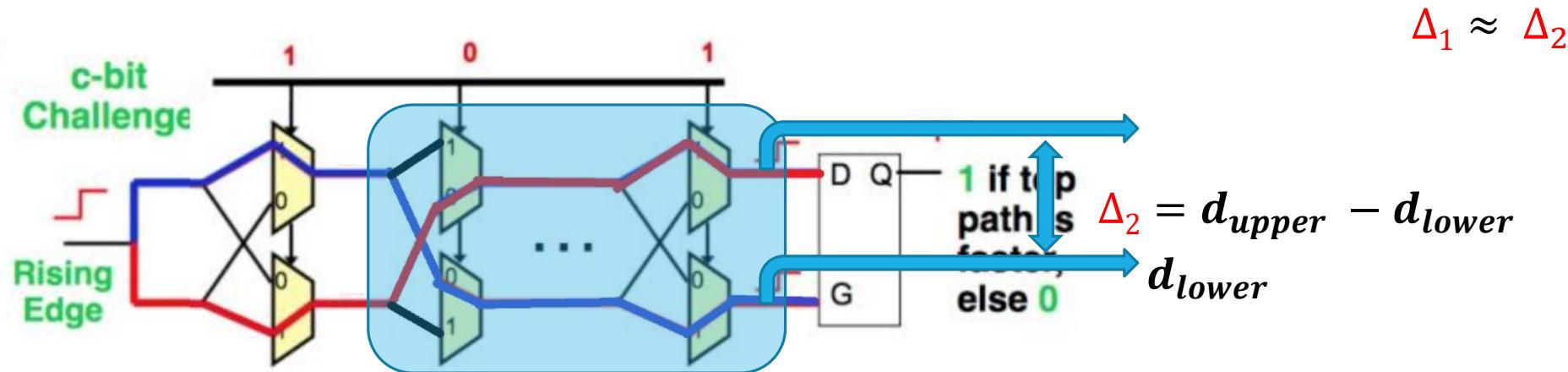
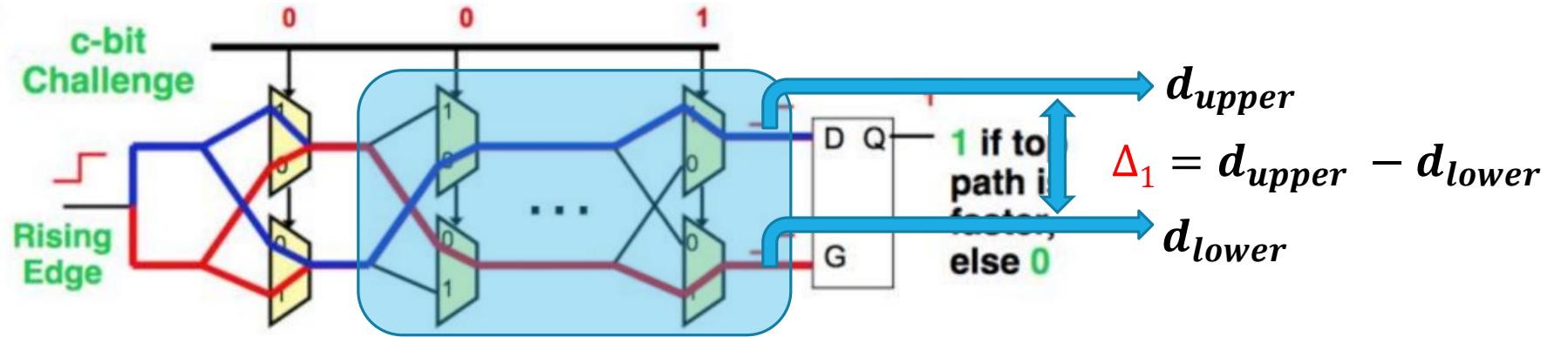


Flipping $c[0]$

Position of Challenge bit of APUF

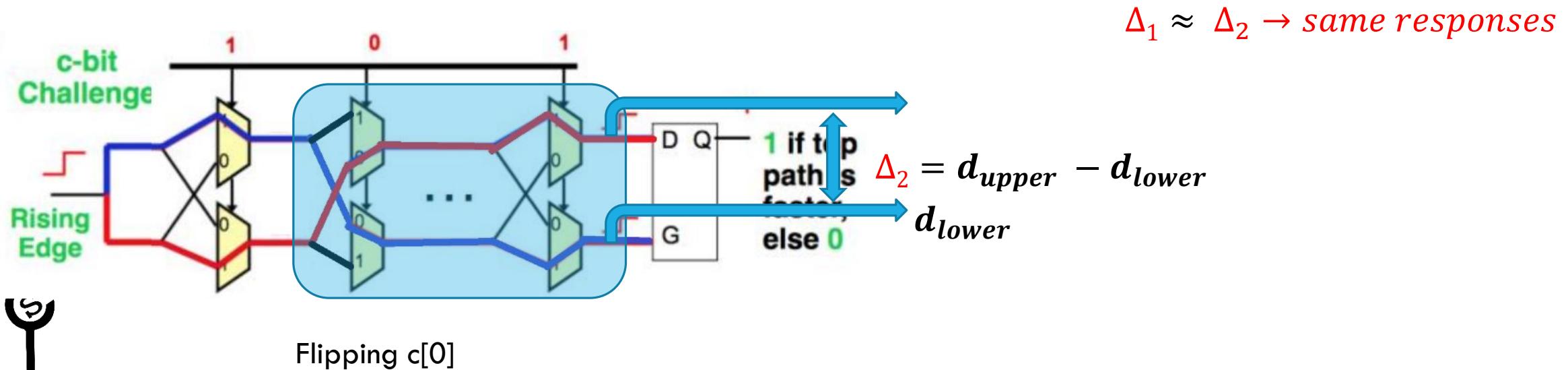
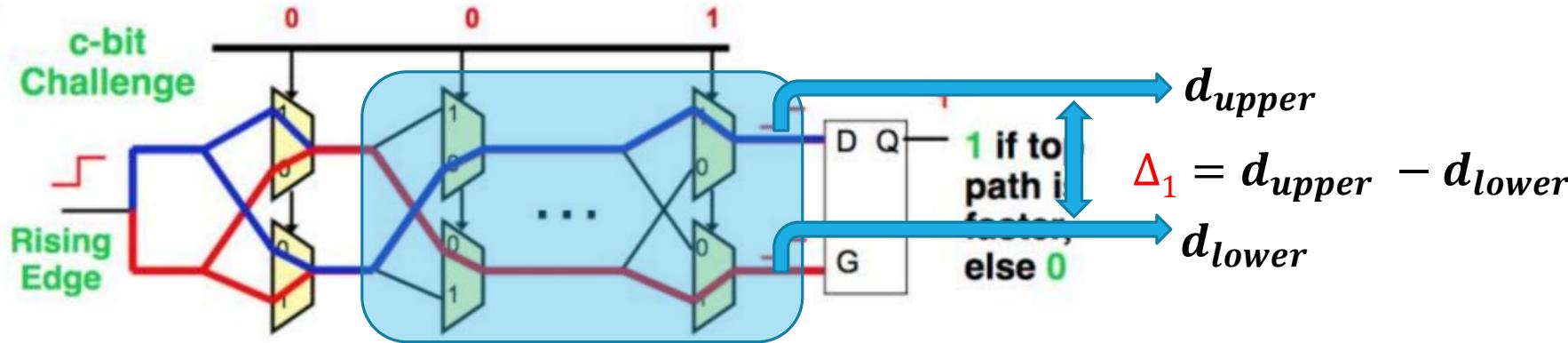


Position of Challenge bit of APUF

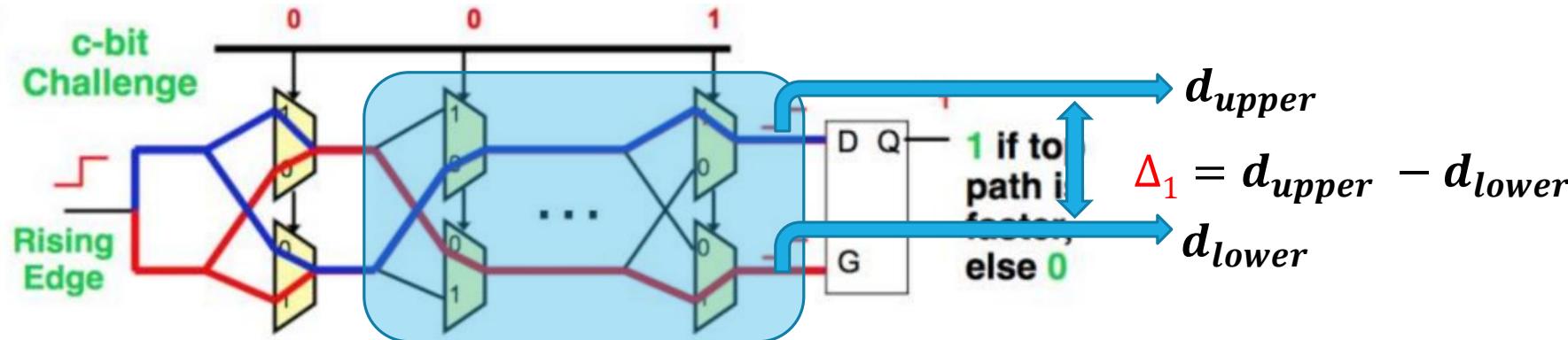


Flipping $c[0]$

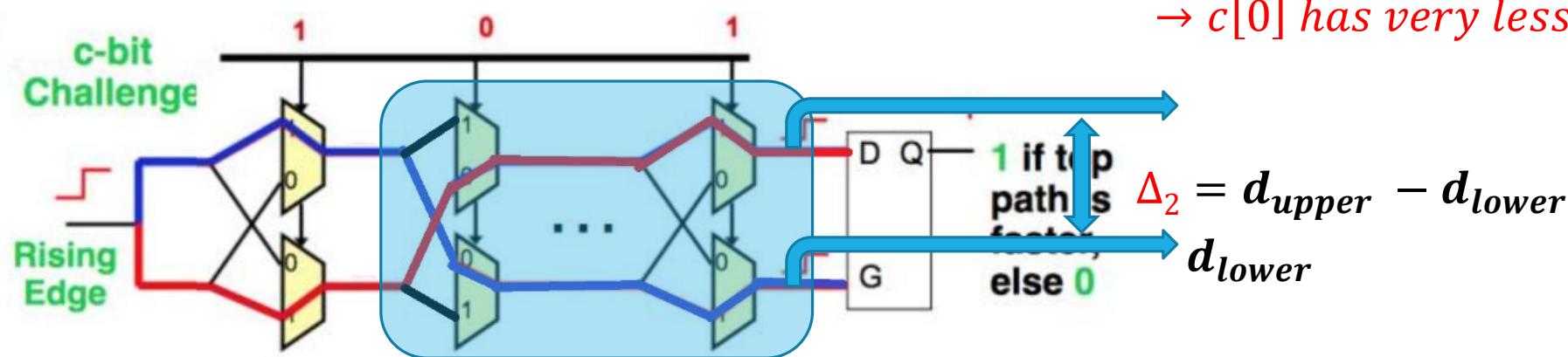
Position of Challenge bit of APUF



Position of Challenge bit of APUF

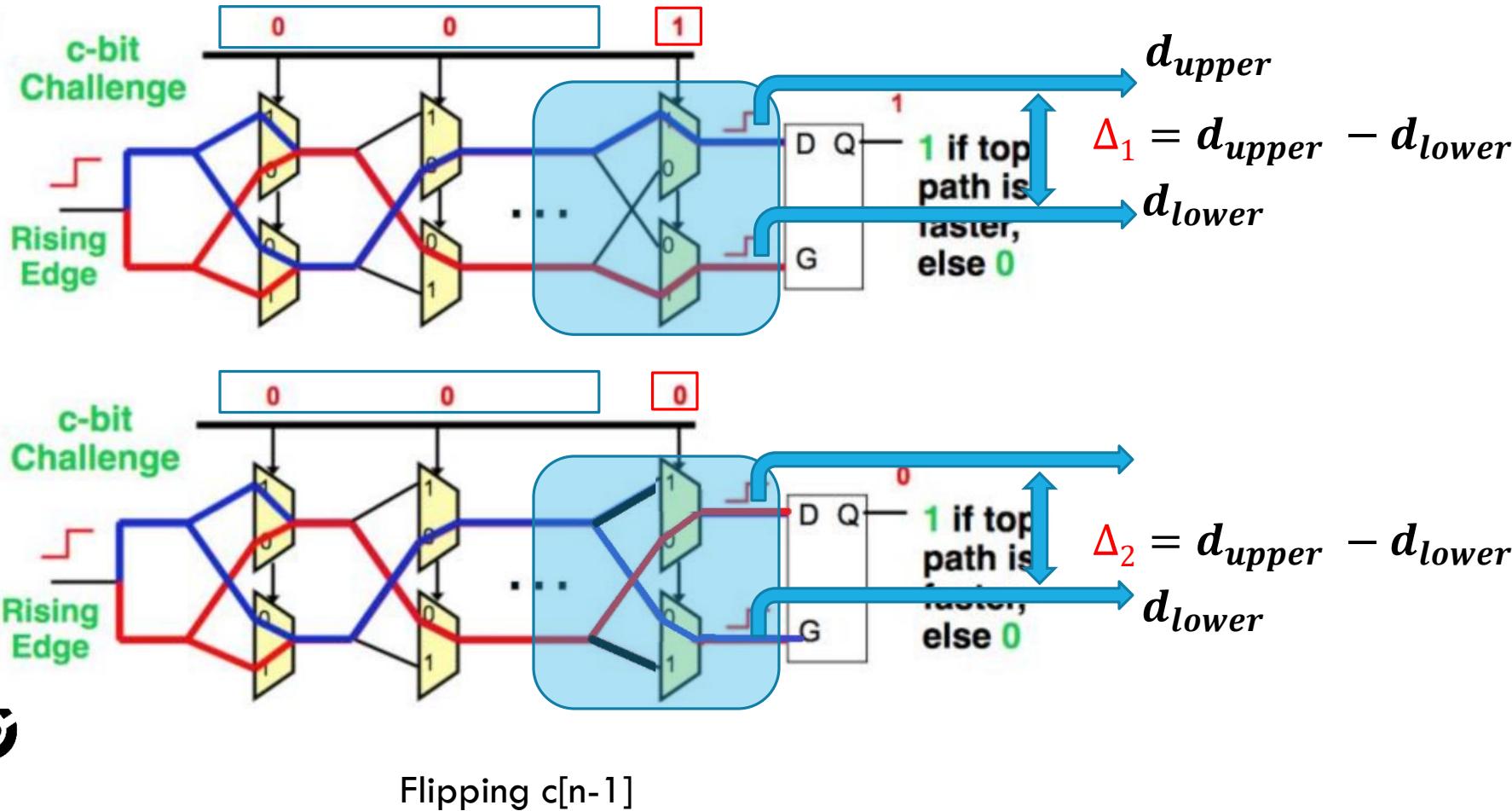


$\Delta_1 \approx \Delta_2 \rightarrow \text{same responses}$
 $\rightarrow c[0] \text{ has very less influence on APUF's output}$

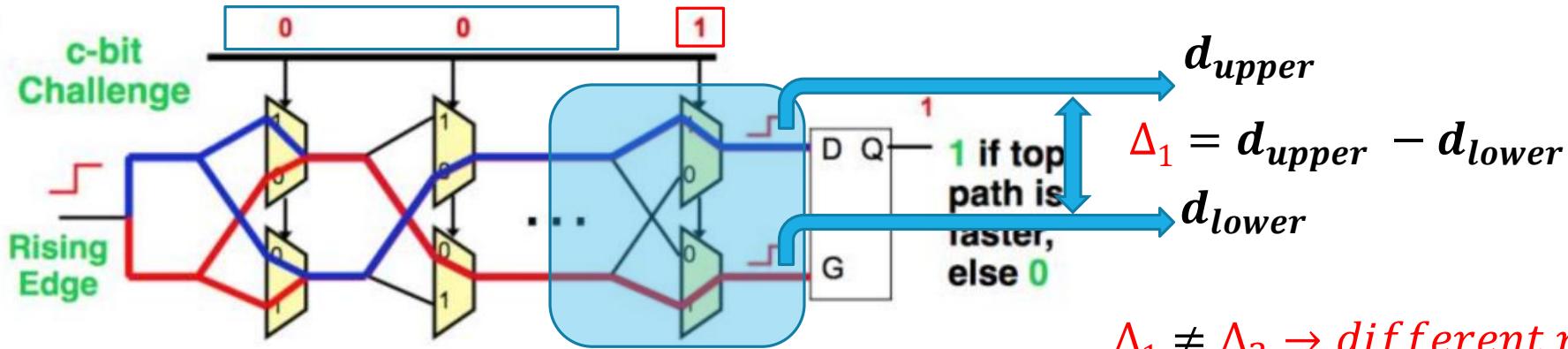


Flipping $c[0]$

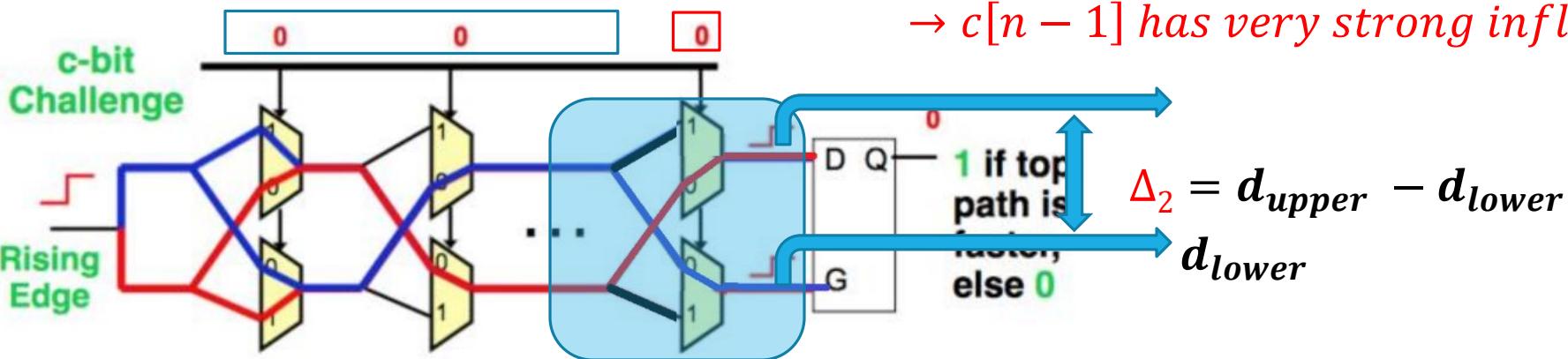
Position of Challenge bit of APUF



Position of Challenge bit of APUF

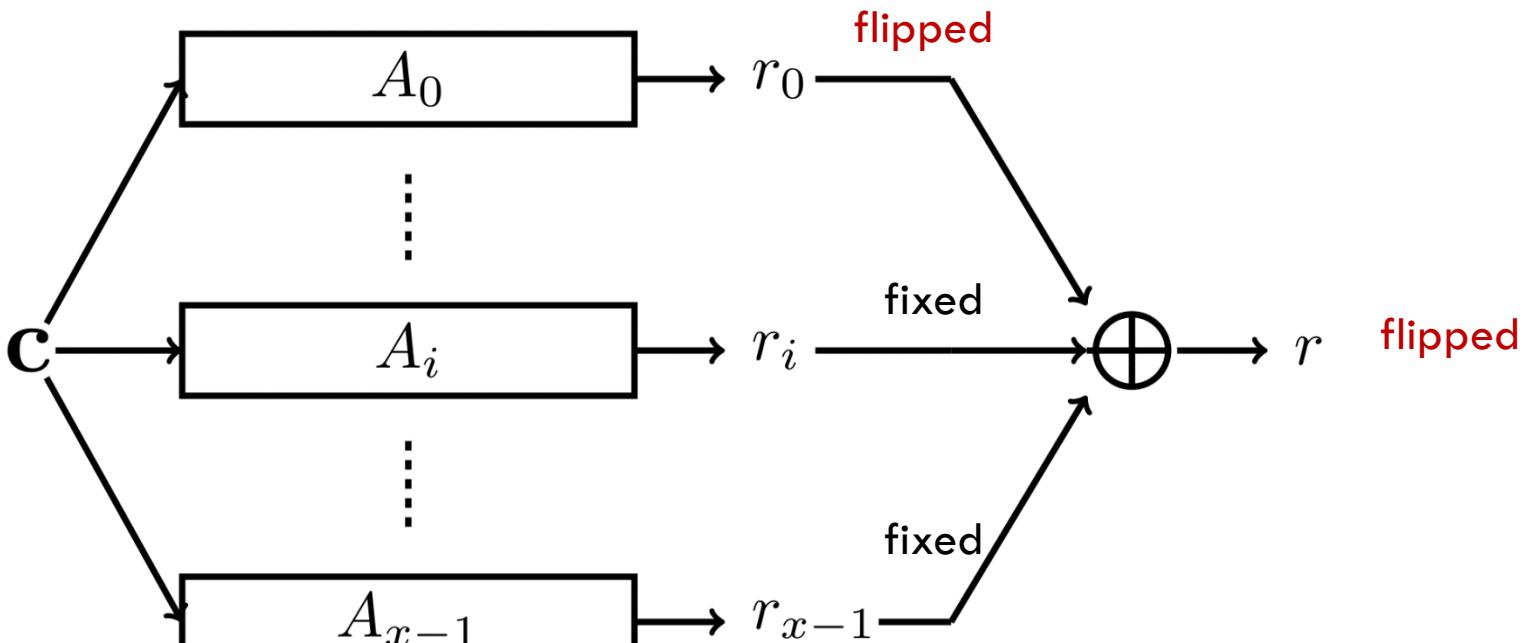


$\Delta_1 \neq \Delta_2 \rightarrow$ different responses
 $\rightarrow c[n-1]$ has very strong influence on APUF's output



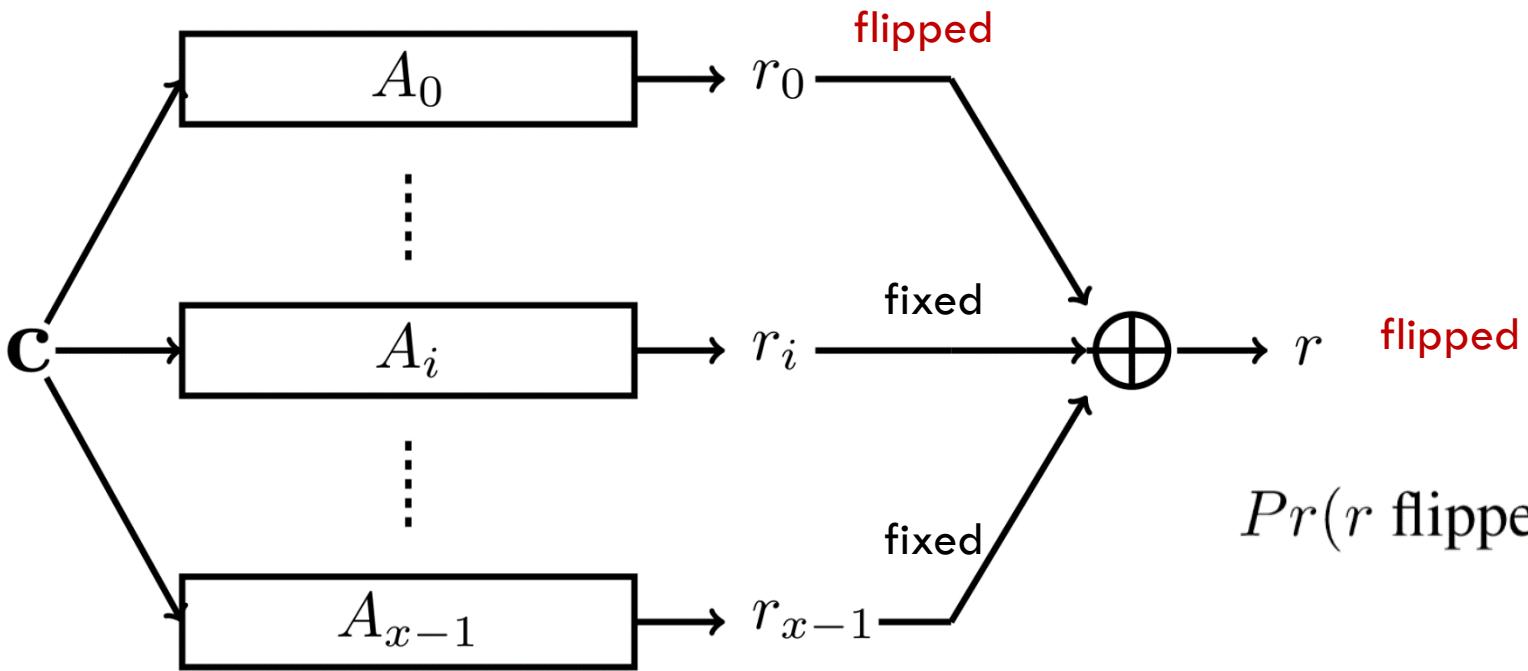
Flipping $c[n-1]$





x -XOR PUF.

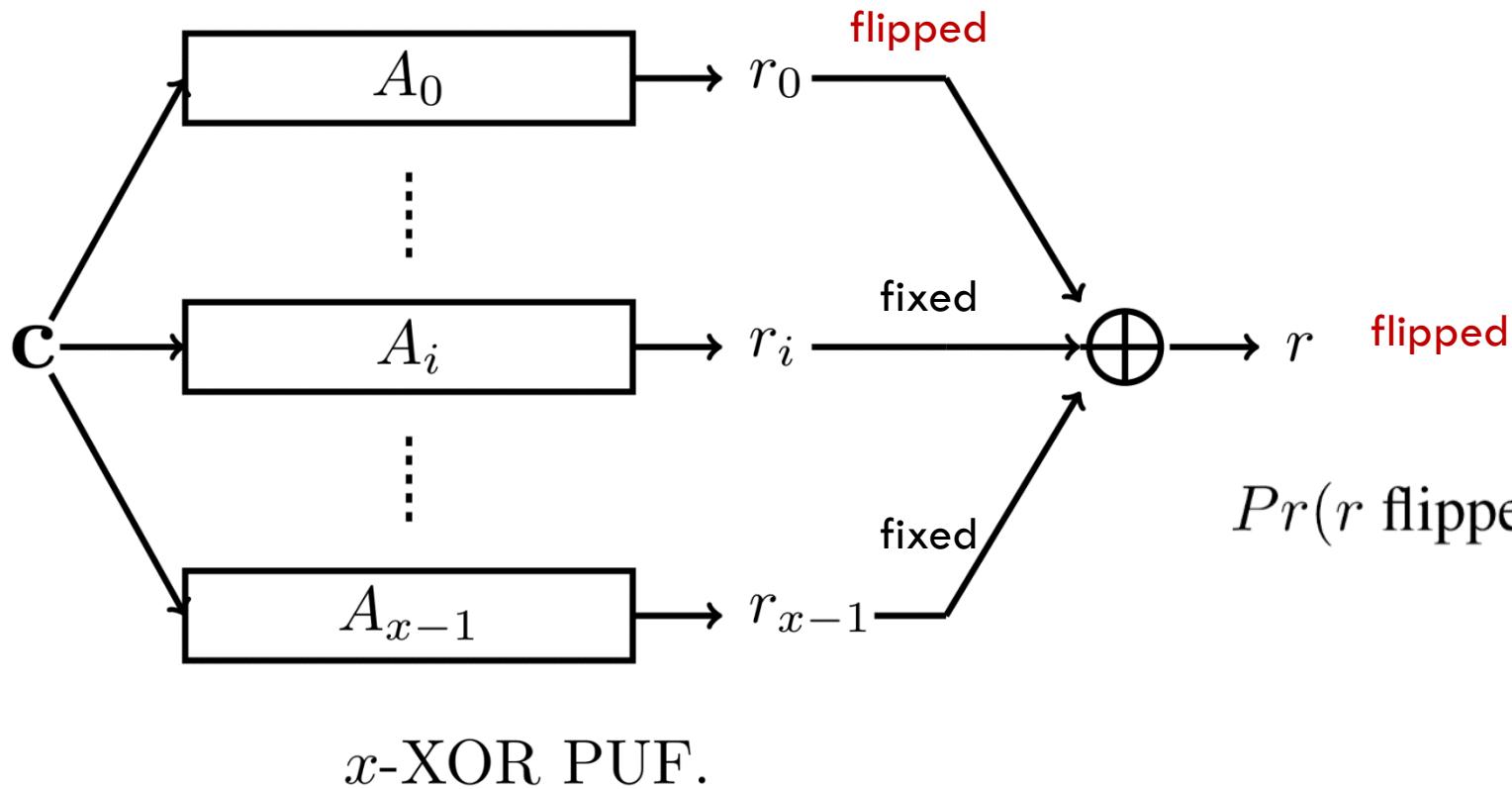




$$\Pr(r \text{ flipped} | r_i \text{ flipped}, r_{j \neq i} \text{ fixed}) = 1$$

x -XOR PUF.





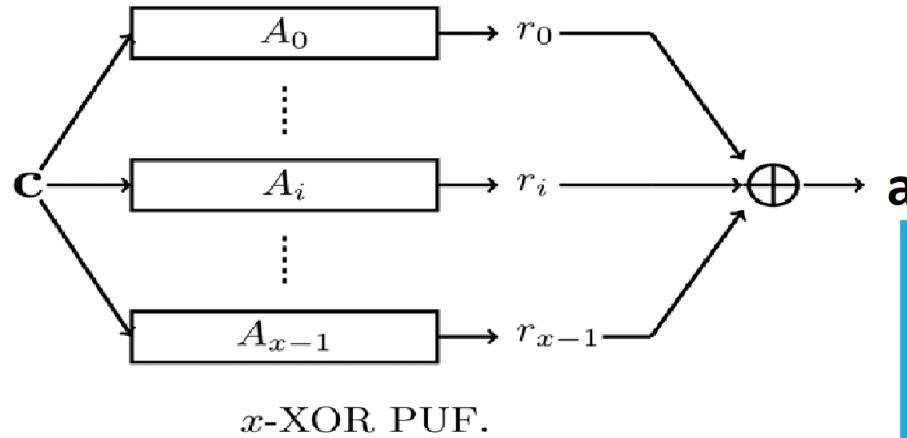
$$\Pr(r \text{ flipped} | r_i \text{ flipped}, r_{j \neq i} \text{ fixed}) = 1$$

All APUFs' outputs contribute to the XOR APUF's output



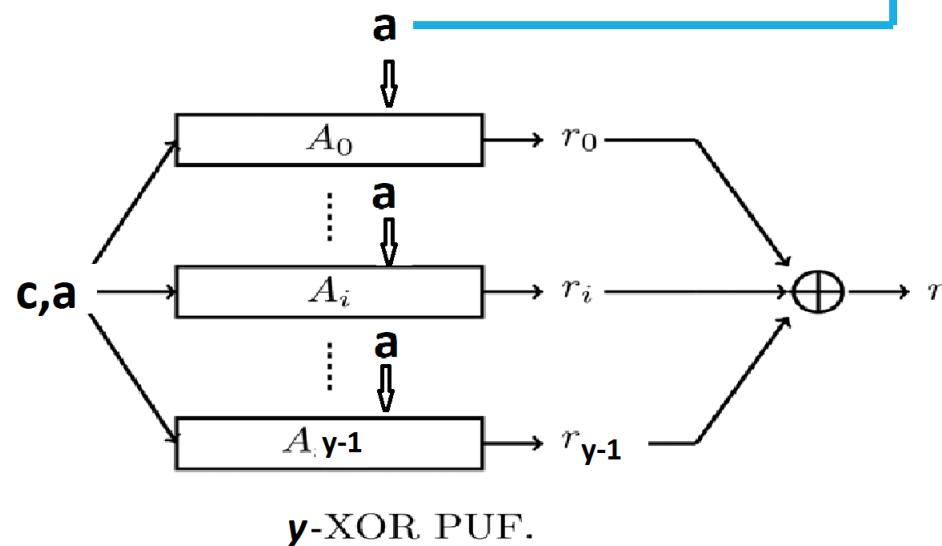
XOR APUF and iPUF

(x,y)-iPUF



- All APUFs' outputs in x -XOR APUF contribute to a

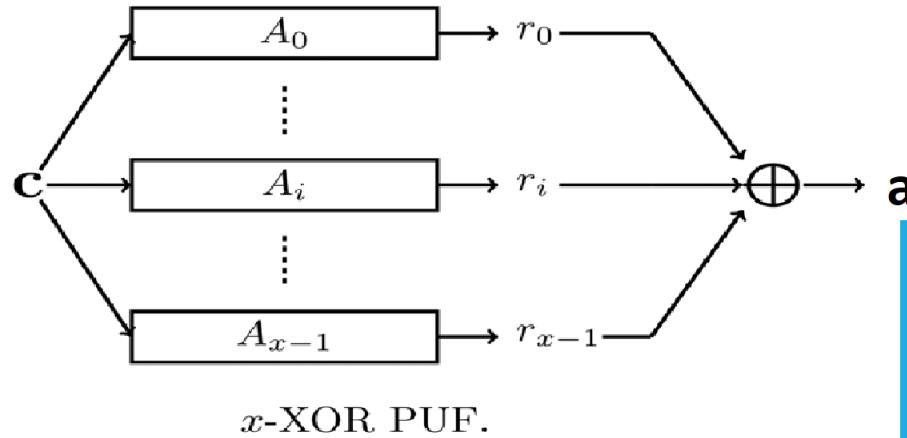
x -XOR PUF.



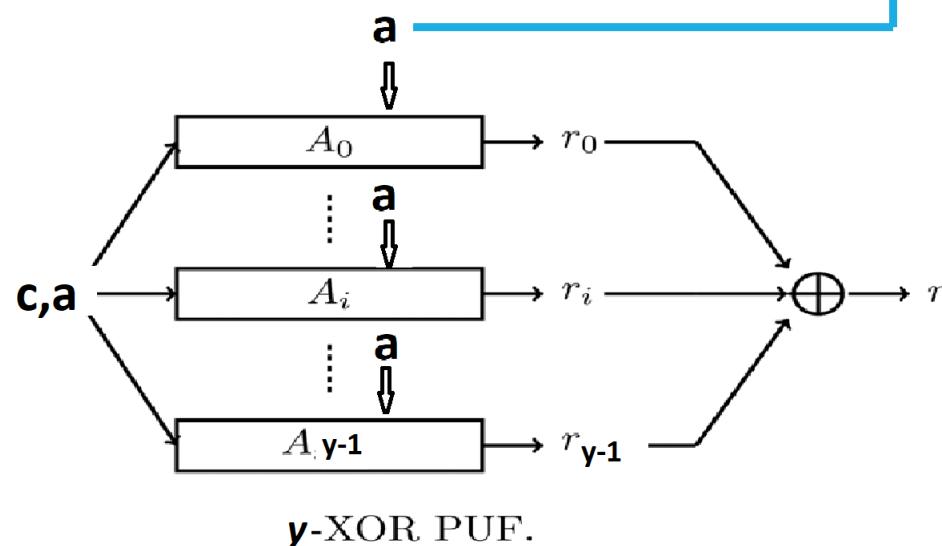
y -XOR PUF.

XOR APUF and iPUF

(x,y)-iPUF

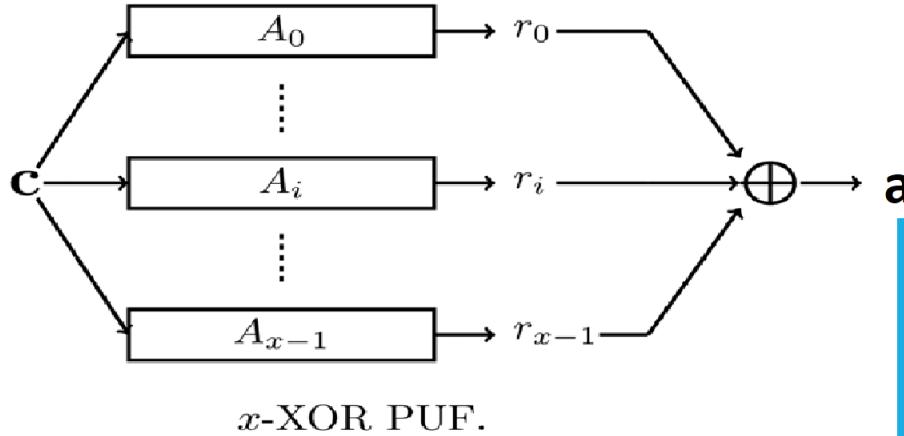


- All APUFs' outputs in x -XOR APUF contribute to a
- All APUF's outputs in y -XOR APUF contribute to r

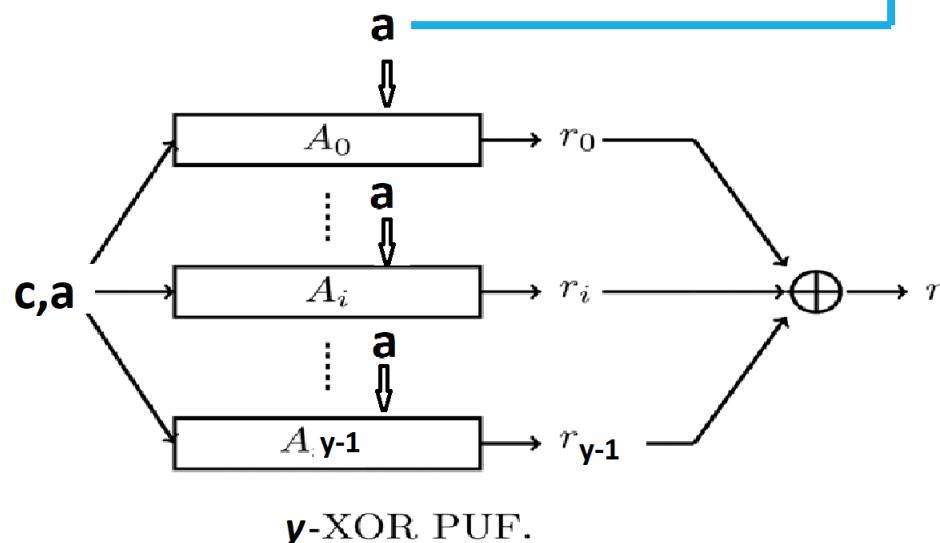


XOR APUF and iPUF

(x,y)-iPUF

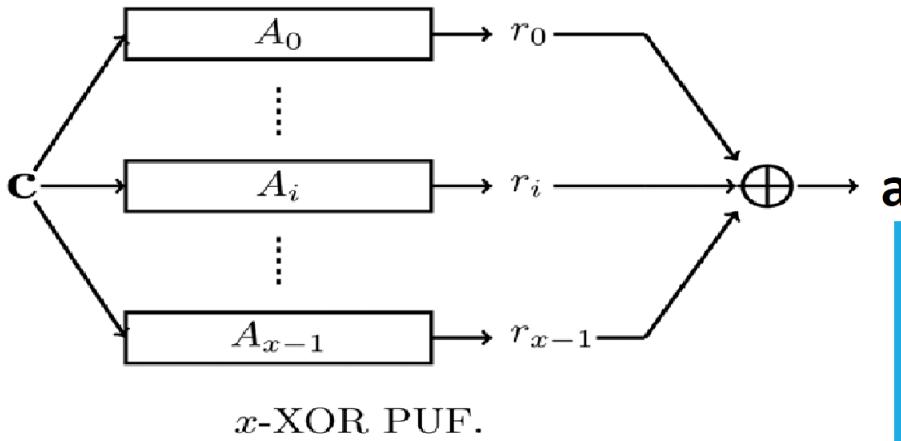


- All APUFs' outputs in x -XOR APUF contribute to a
- All APUF's outputs in y -XOR APUF contribute to r
- For a given challenge c ,
 - If flipping a makes iPUF's output r not flip due to SAC of APUF, then (x,y) -iPUF = y -XOR APUF
 - because only y APUFs at y -XOR APUF contribute to r .

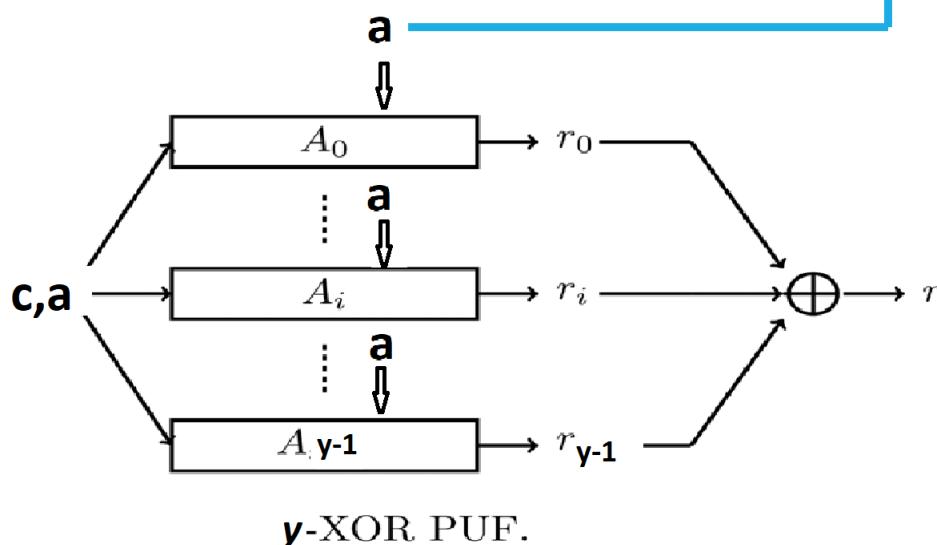


XOR APUF and iPUF

(x,y) -iPUF

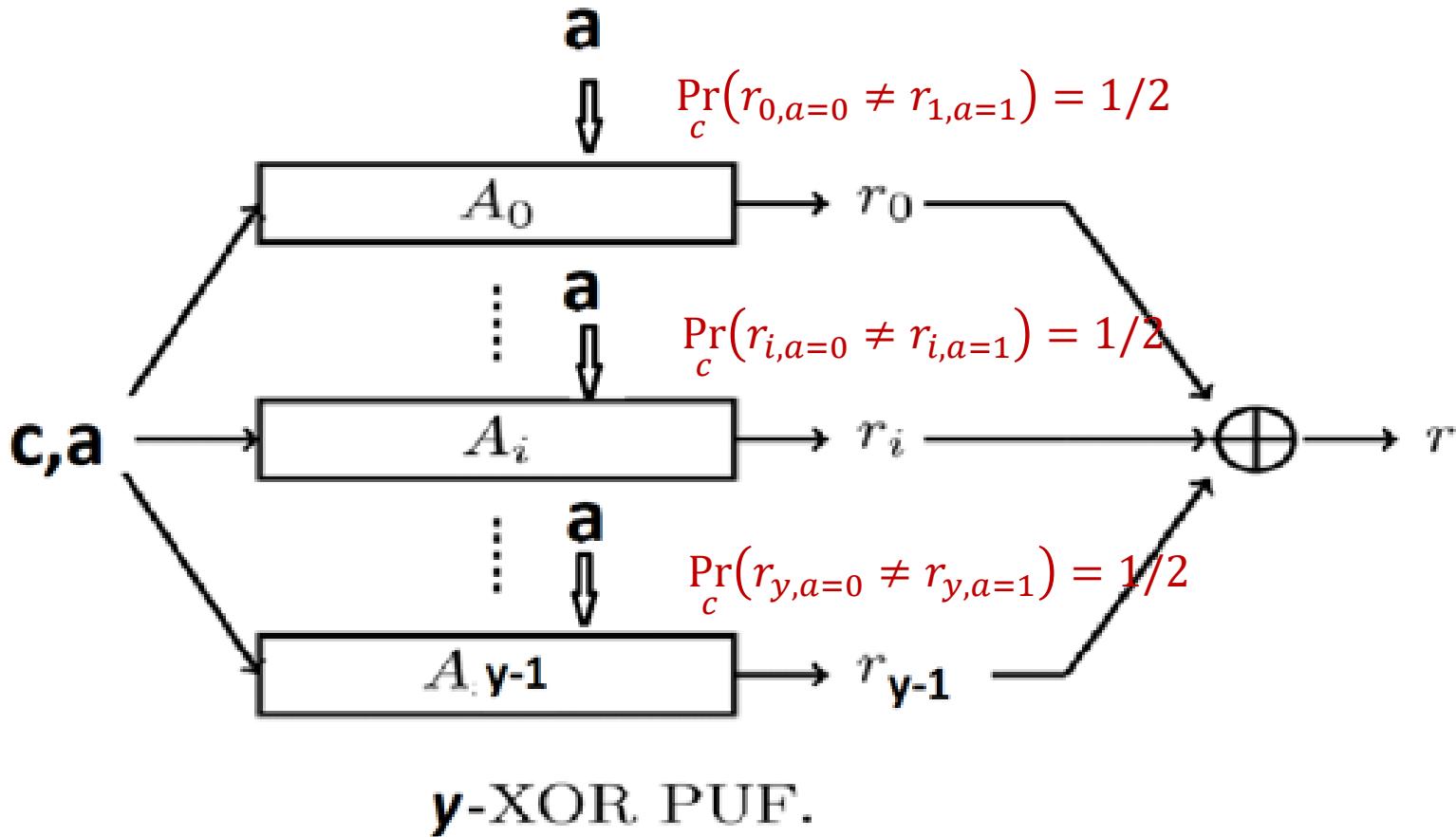


- All APUFs' outputs in x -XOR APUF contribute to a
- All APUF's outputs in y -XOR APUF contribute to r
- For a given challenge c ,
If flipping a makes iPUF's output r not flip due to SAC of APUF,
then (x,y) -iPUF = y -XOR APUF
because only y APUFs at y -XOR APUF contribute to r .



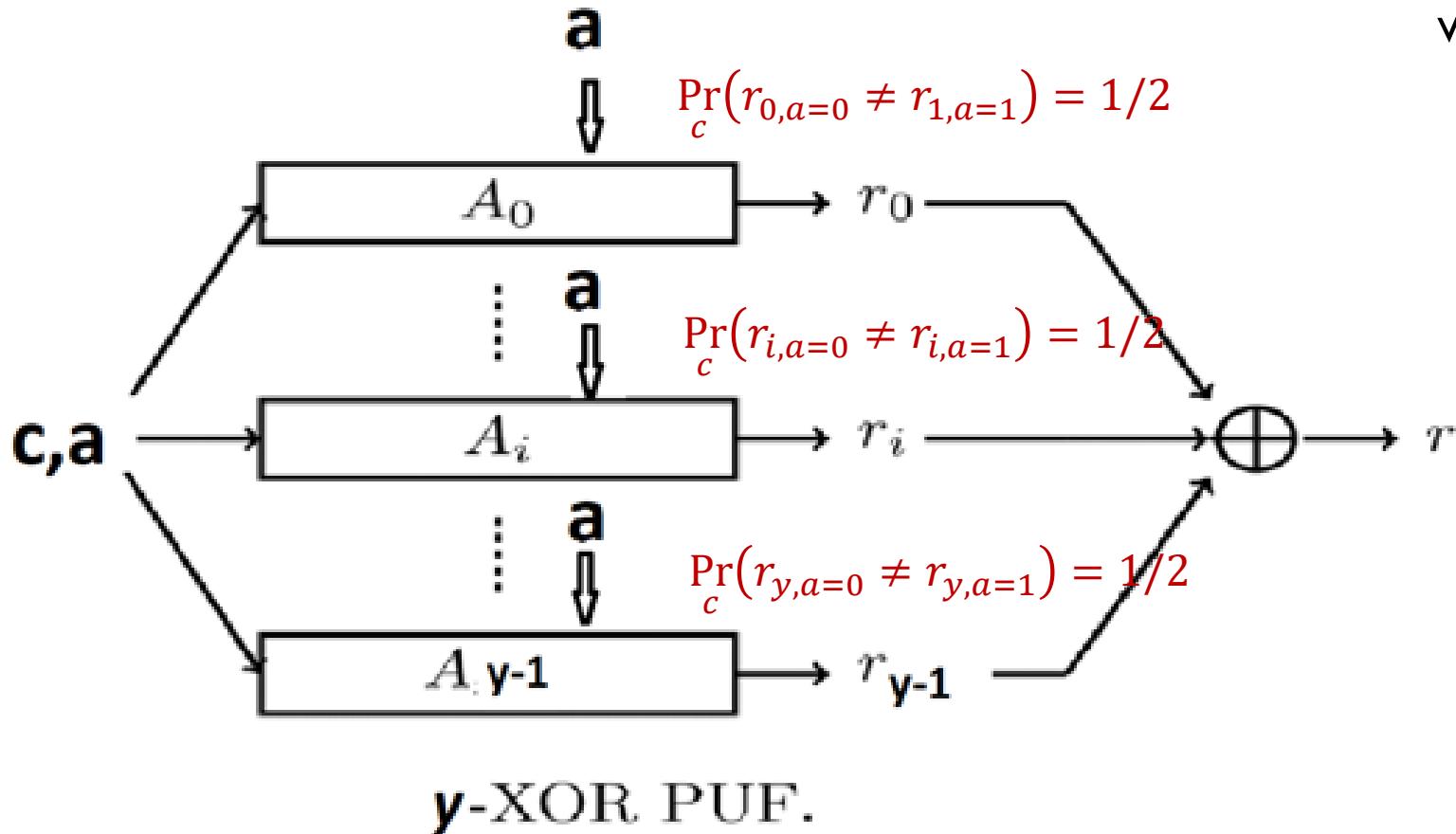
- For a given challenge c ,
If flipping a makes iPUF's output r flip due to SAC of APUF,
then (x,y) -iPUF = $(x+y)$ -XOR APUF
because all $x+y$ APUFs at x and y -XOR APUFs contribute to r

XOR APUF and iPUF

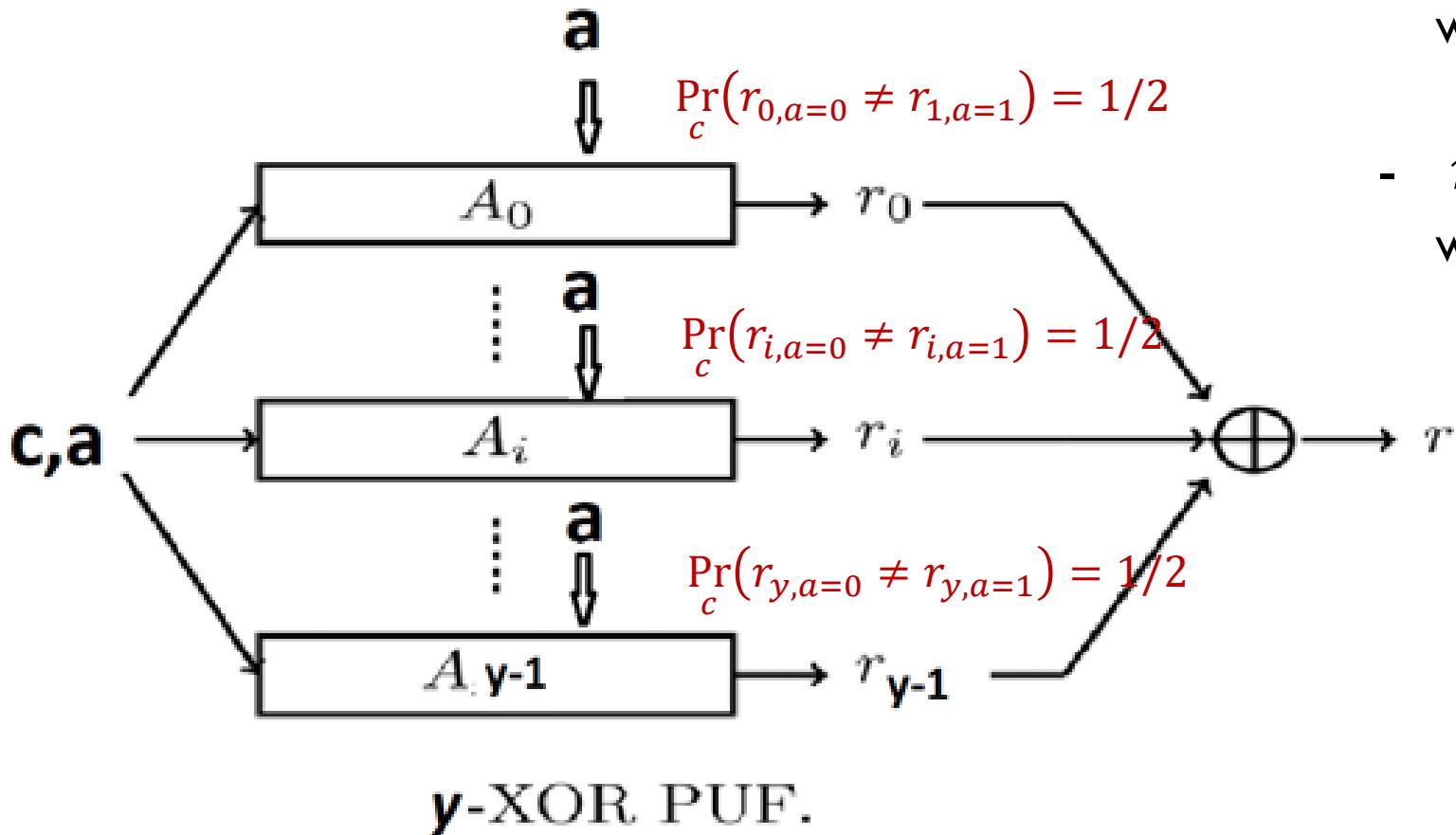


XOR APUF and iPUF

- r flips when odd # flipped r_i when a flips

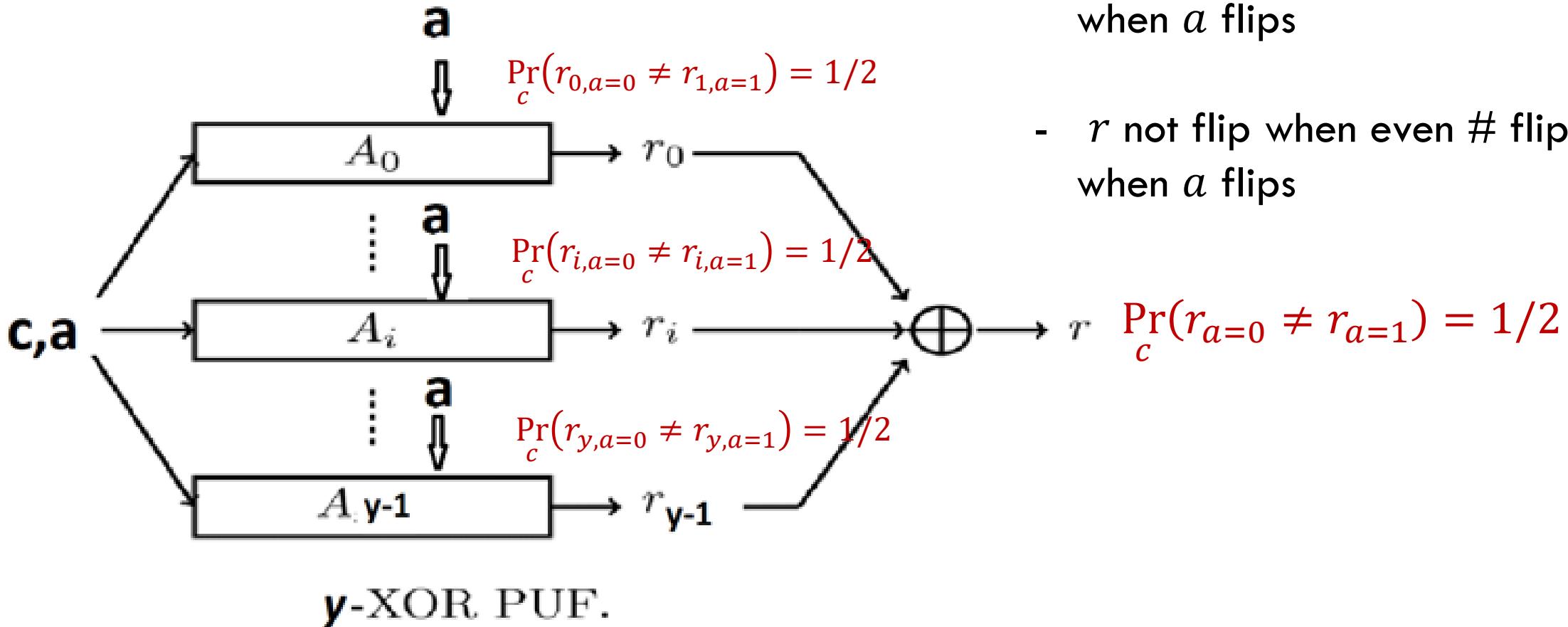


XOR APUF and iPUF



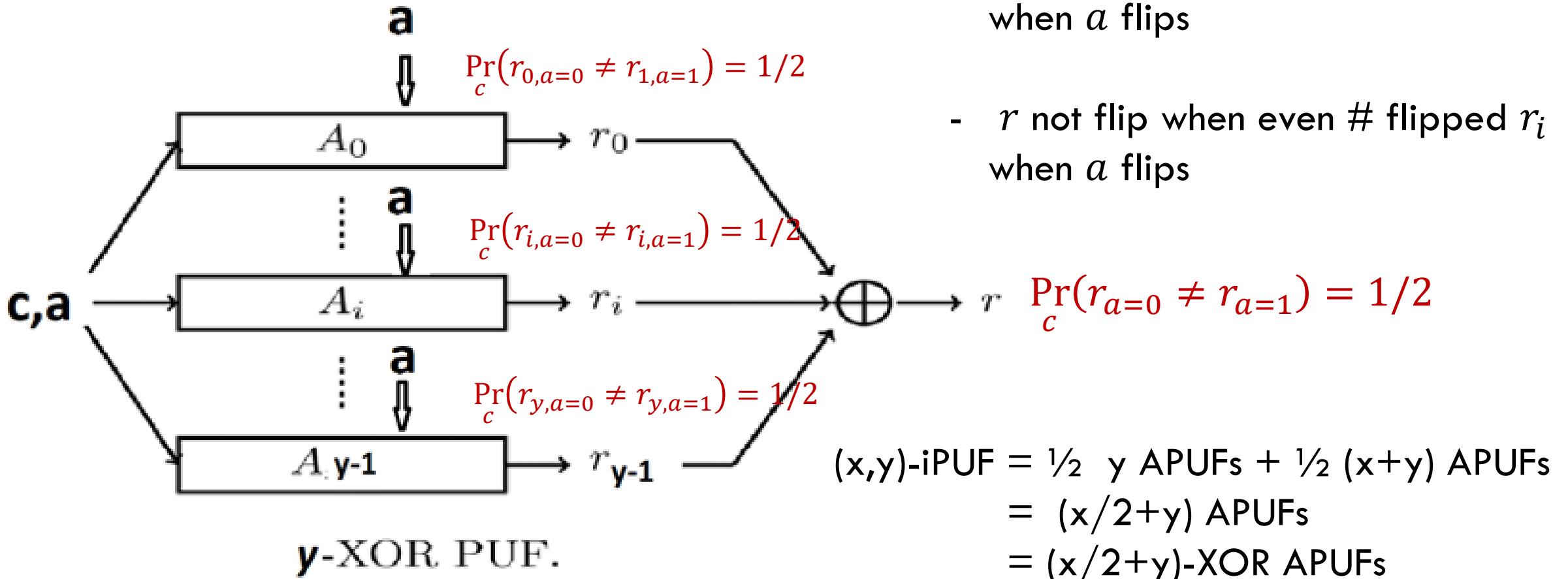
- r flips when odd # flipped r_i when a flips
- r not flip when even # flipped r_i when a flips

XOR APUF and iPUF



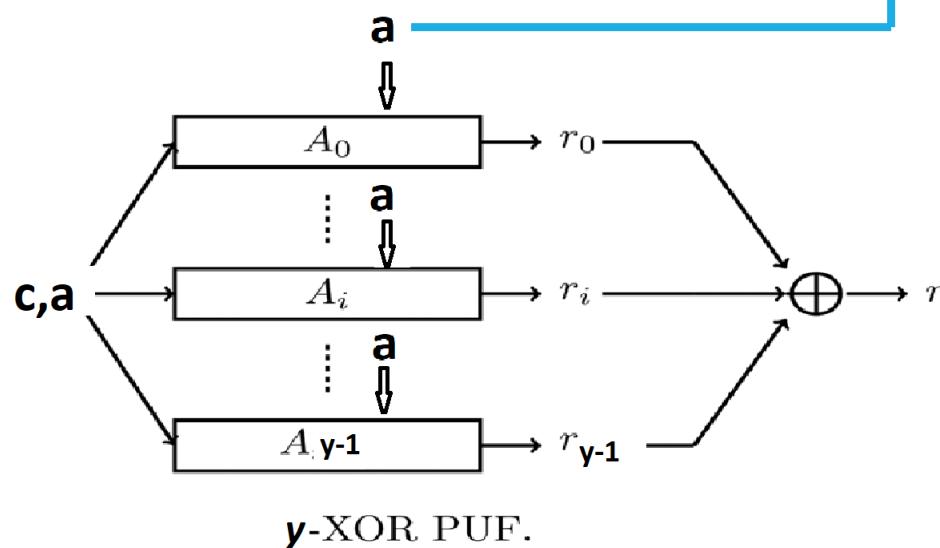
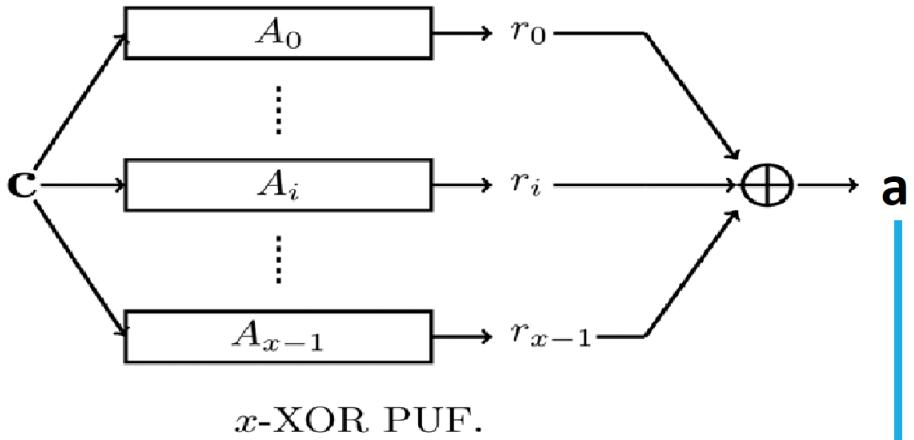
- r flips when odd # flipped r_i when a flips
- r not flip when even # flipped r_i when a flips

XOR APUF and iPUF



XOR APUF and iPUF

(x,y)-iPUF

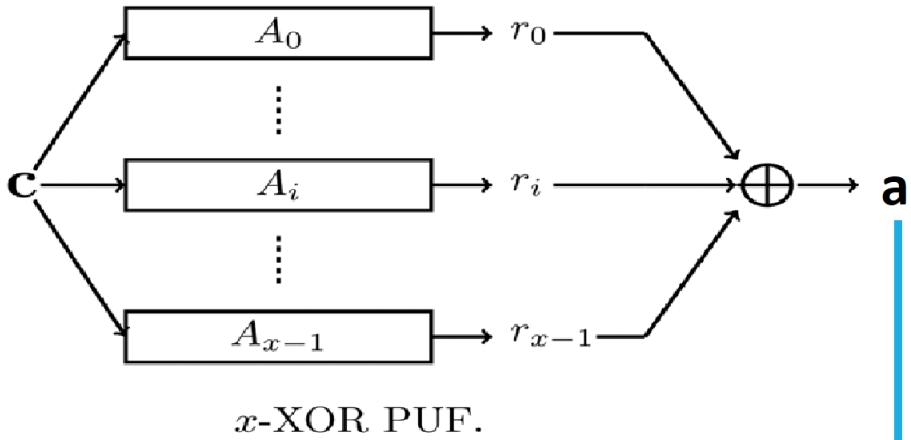


- All APUFs' outputs in x -XOR APUF contribute to a
- All APUF's outputs in y -XOR APUF contribute to r
- For a given challenge c ,
 - If flipping a makes iPUF's output r not flip due to SAC of APUF, then (x,y) -iPUF = y -XOR APUF because only y APUFs at y -XOR APUF contribute to r .
- For a given challenge c ,
 - If flipping a makes iPUF's output r flip due to SAC of APUF, then (x,y) -iPUF = $(x+y)$ -XOR APUF because all $x+y$ APUFs at x and y -XOR APUFs contribute to r
- If a is interposed at the middle ($=n/2$), then (x,y) -iPUF = $(x/2+y)$ -XOR APUF

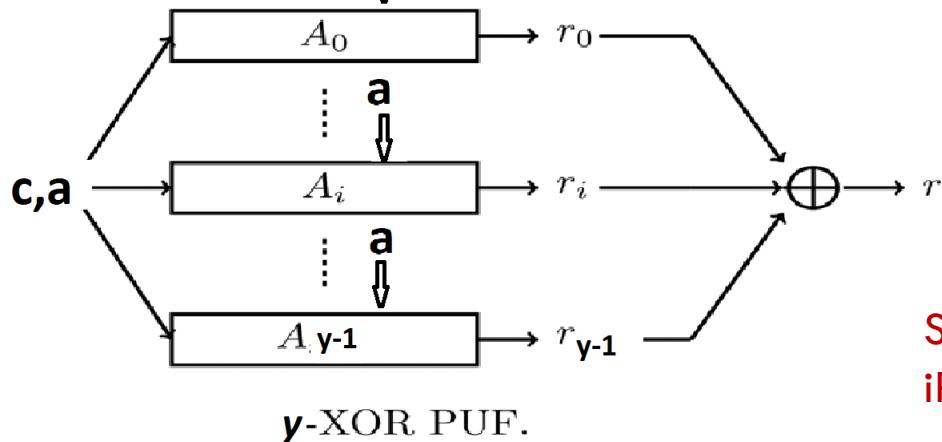


XOR APUF and iPUF

(x,y)-iPUF



- All APUFs' outputs in x -XOR APUF contribute to a
- All APUF's outputs in y -XOR APUF contribute to r
- For a given challenge c ,
 - If flipping a makes iPUF's output r not flip due to SAC of APUF, then (x,y) -iPUF = y -XOR APUF because only y APUFs at y -XOR APUF contribute to r .
- For a given challenge c ,
 - If flipping a makes iPUF's output r flip due to SAC of APUF, then (x,y) -iPUF = $(x+y)$ -XOR APUF because all $x+y$ APUFs at x and y -XOR APUFs contribute to r
- If a is interposed at the middle ($=n/2$), then (x,y) -iPUF = $(x/2+y)$ -XOR APUF

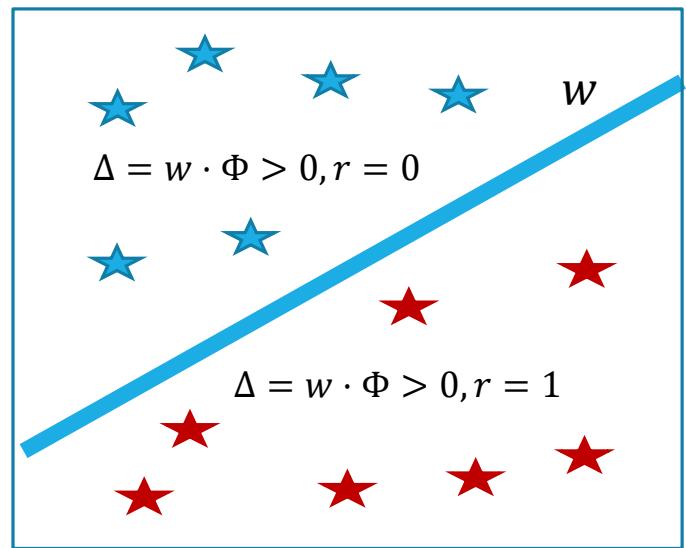
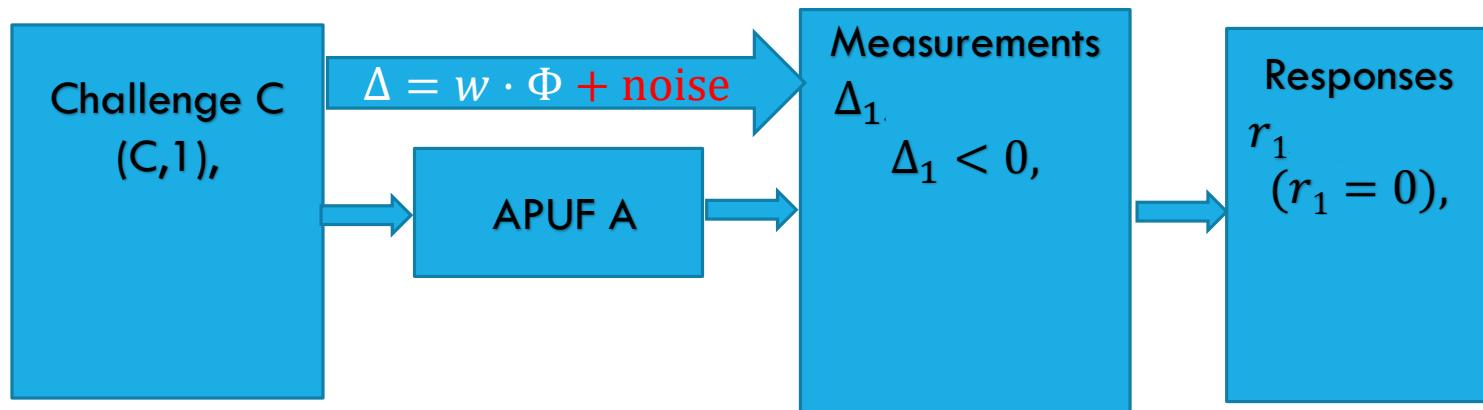


Since XOR APUF is secure against classical ML attack, iPUF is also secure against classical ML attacks

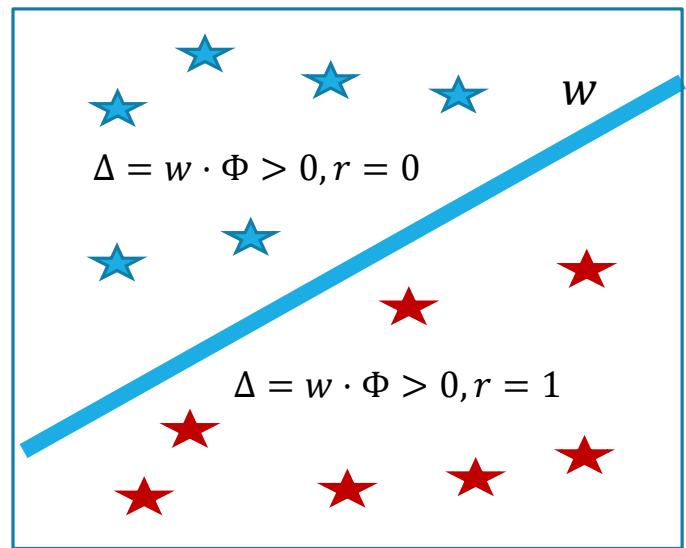
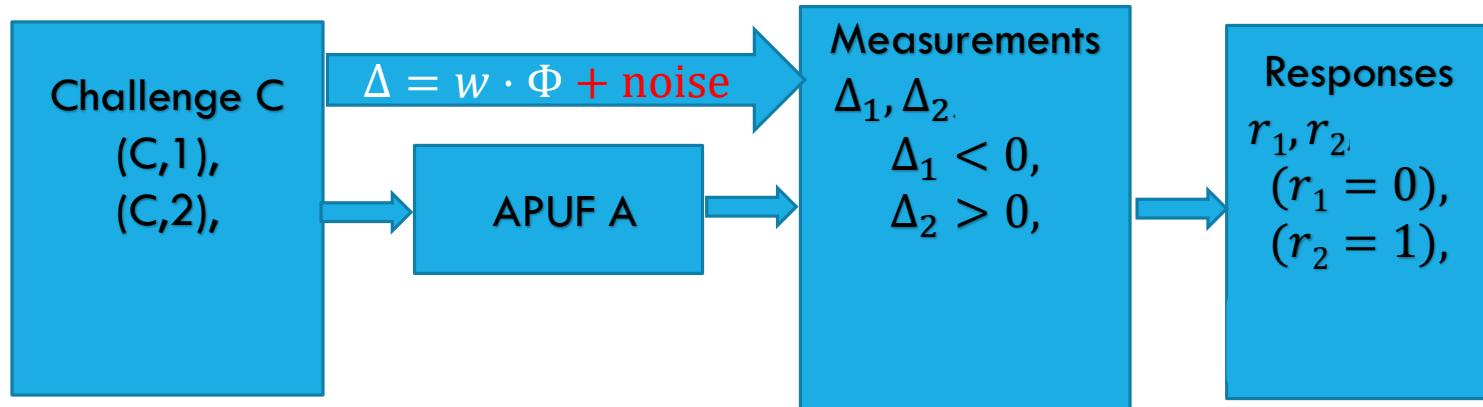
7. Short-term Reliability



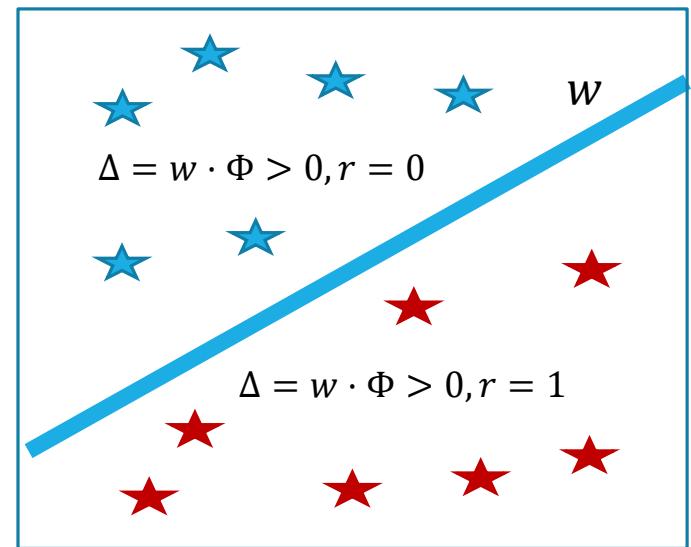
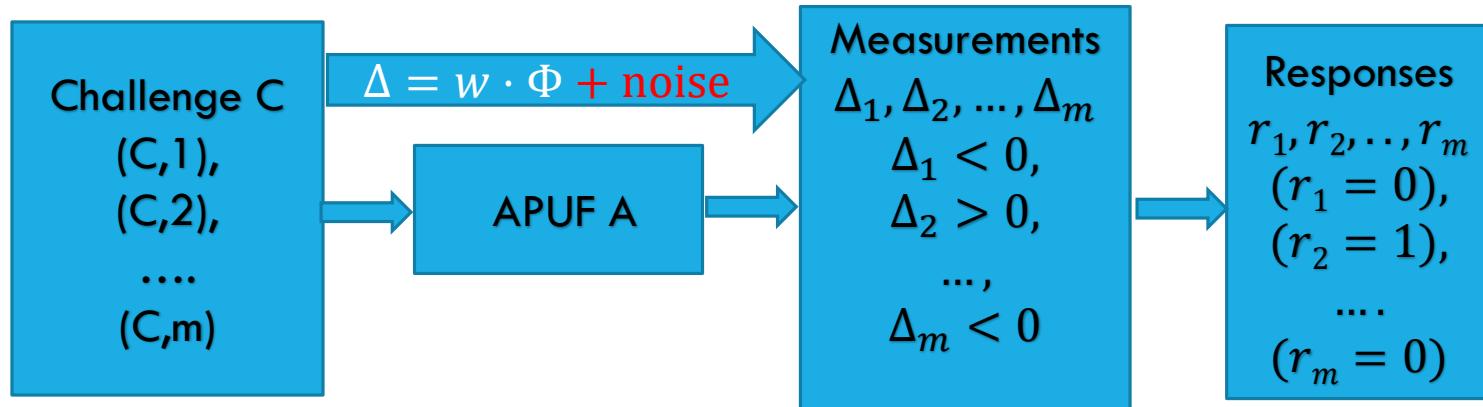
Arbiter: Repeatability – short-term Reliability



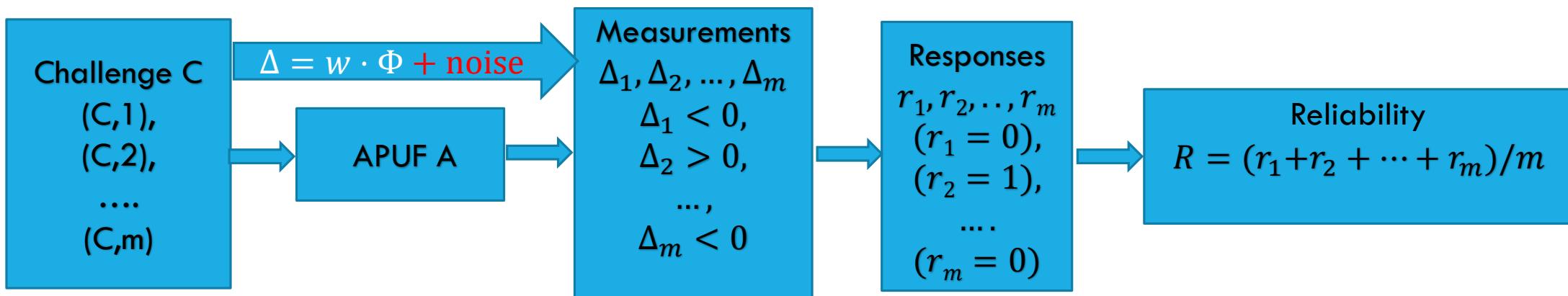
Arbiter: Repeatability – short-term Reliability



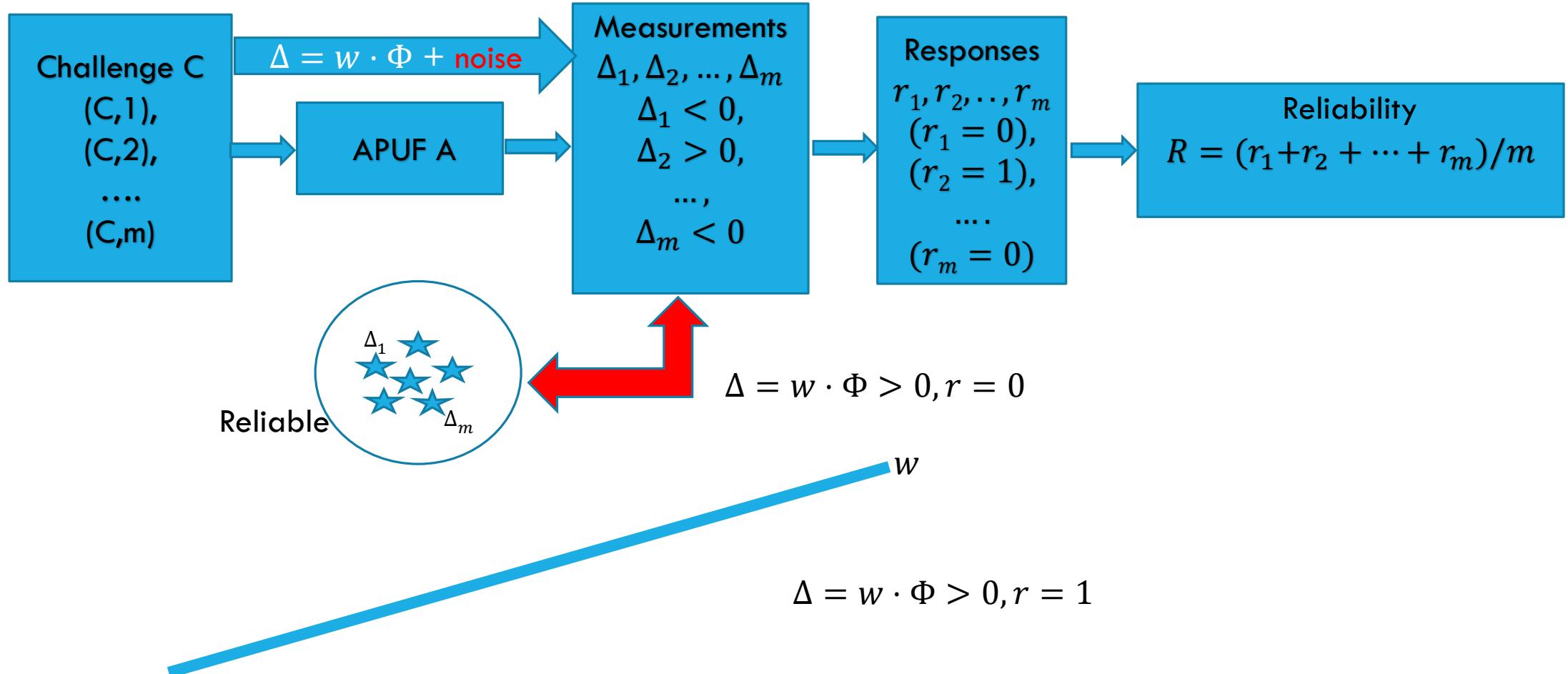
Arbiter: Repeatability – short-term Reliability



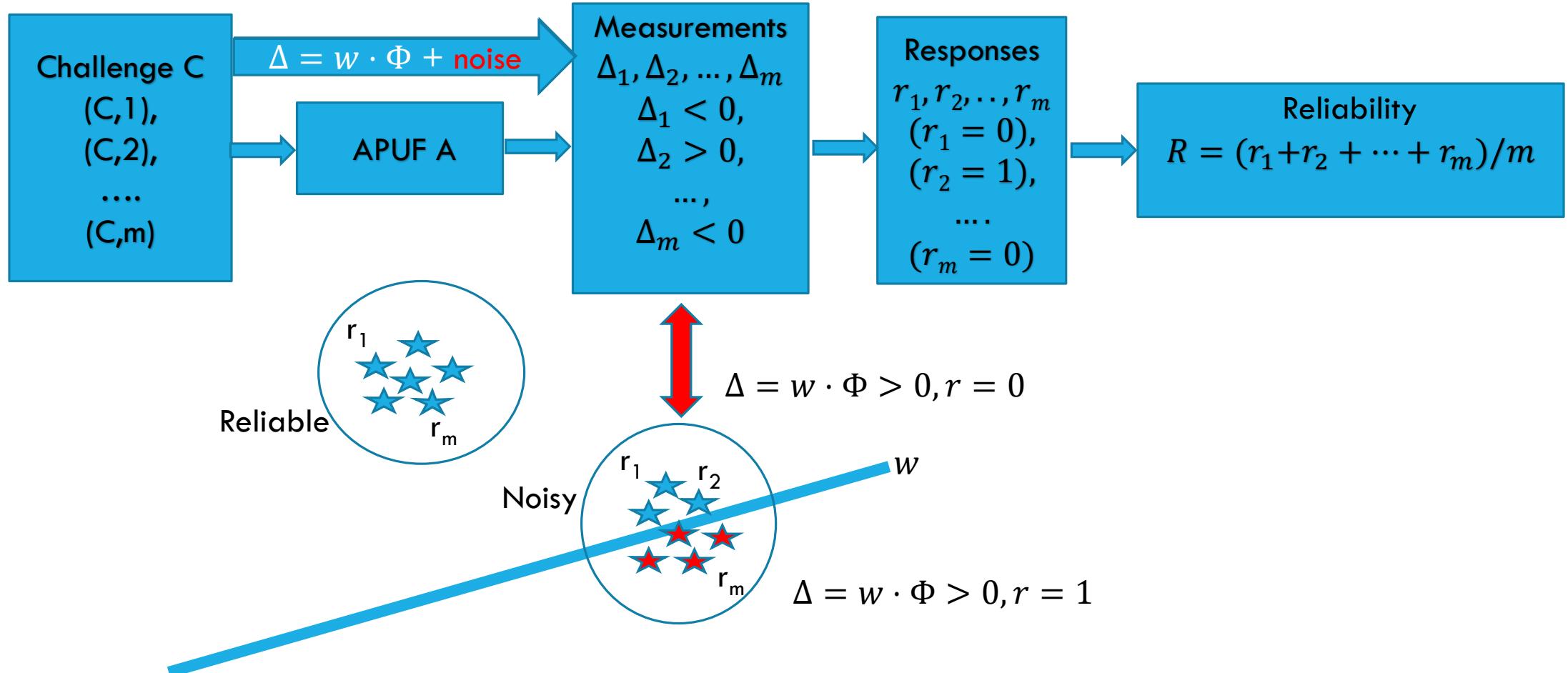
Arbiter: Repeatability – short-term Reliability



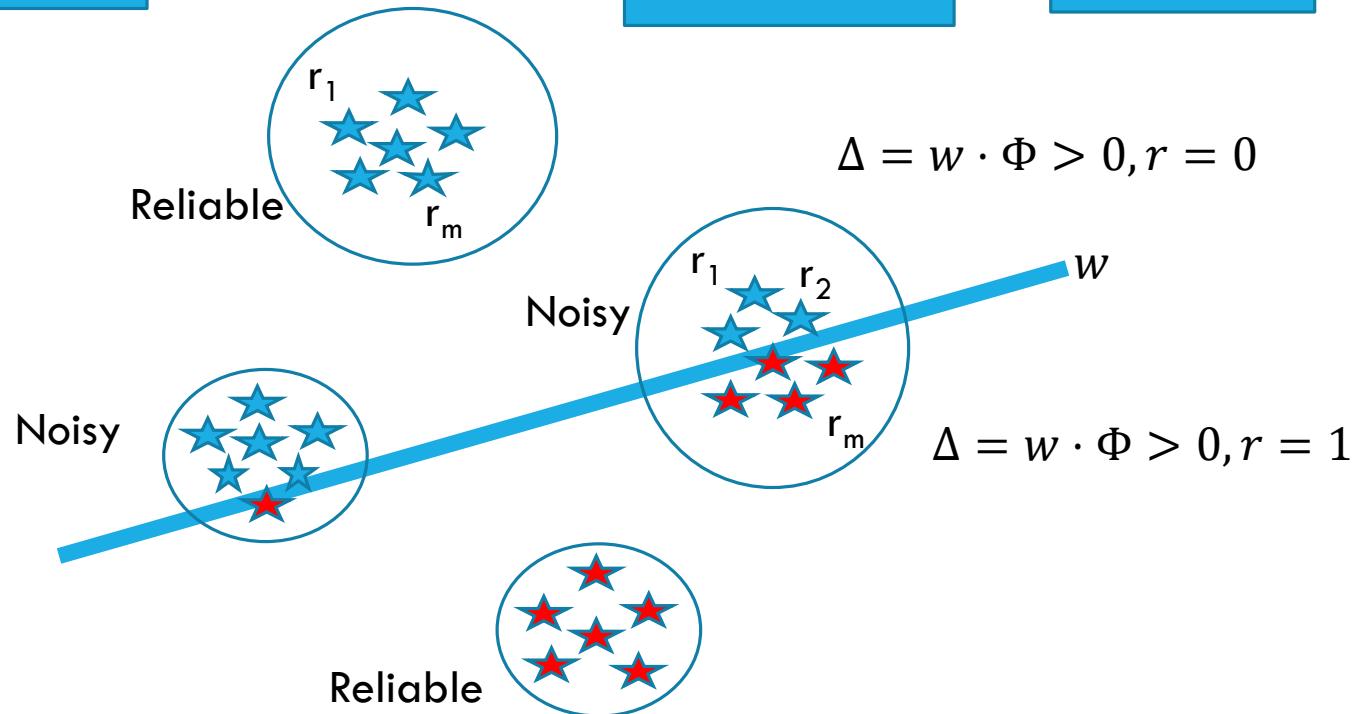
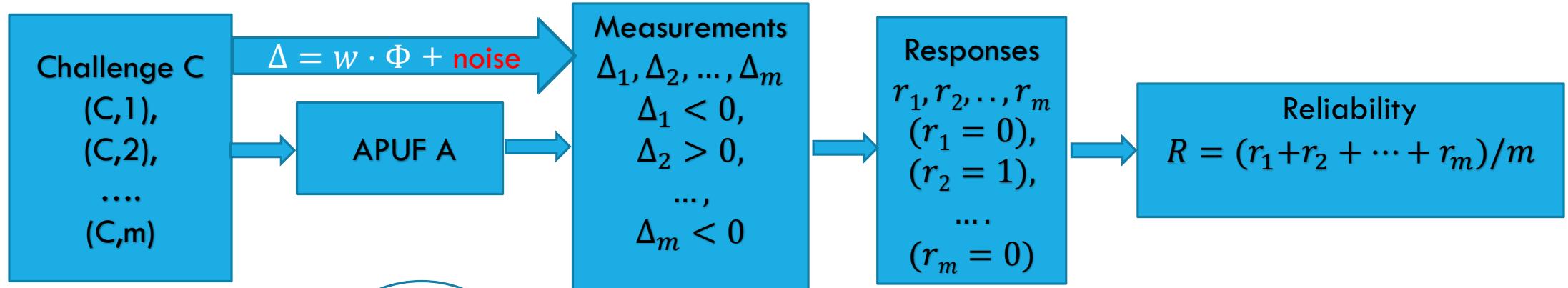
Arbiter: Repeatability – short-term Reliability



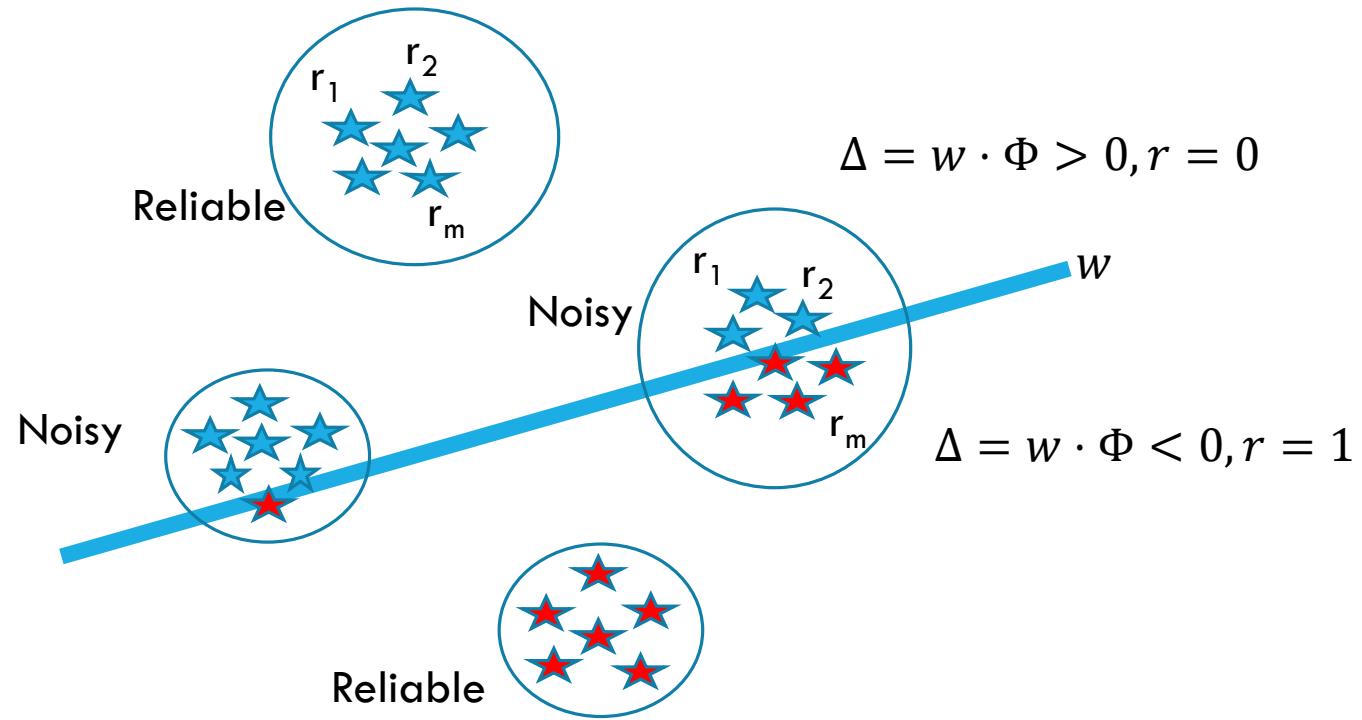
Arbiter: Repeatability – short-term Reliability



Arbiter: Repeatability – short-term Reliability



Arbiter: Repeatability – short-term Reliability



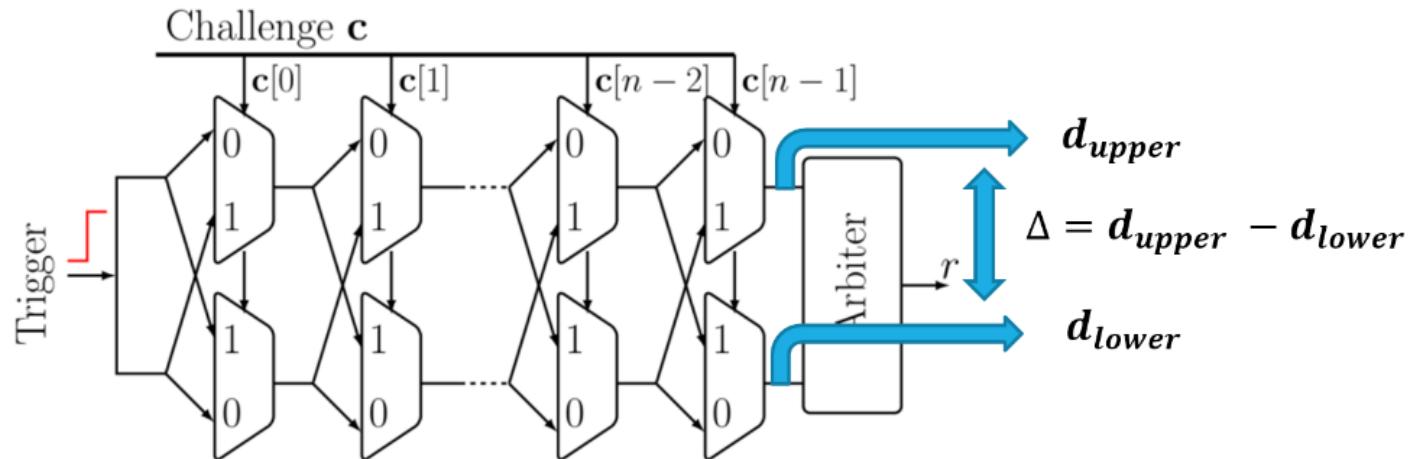
Reliability of C and w are related

CMA-ES + many pairs of (c, R) : Reliability based modeling attack

The Gap Between Promise and Reality: On the Insecurity of XOR Arbiter PUFs CHES, 2015, Georg T. Becker

8. Reliability based Modeling Attacks





A n -stage classic Arbiter PUF with challenge $\mathbf{c} \in \{0, 1\}^n$.

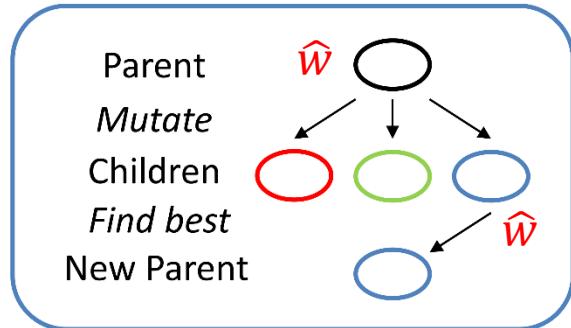
- $\Delta > 0 \rightarrow r = 1$. Otherwise $r = 0$
- $\Delta = d_{upper} - d_{lower} = \mathbf{w} \cdot \Phi$



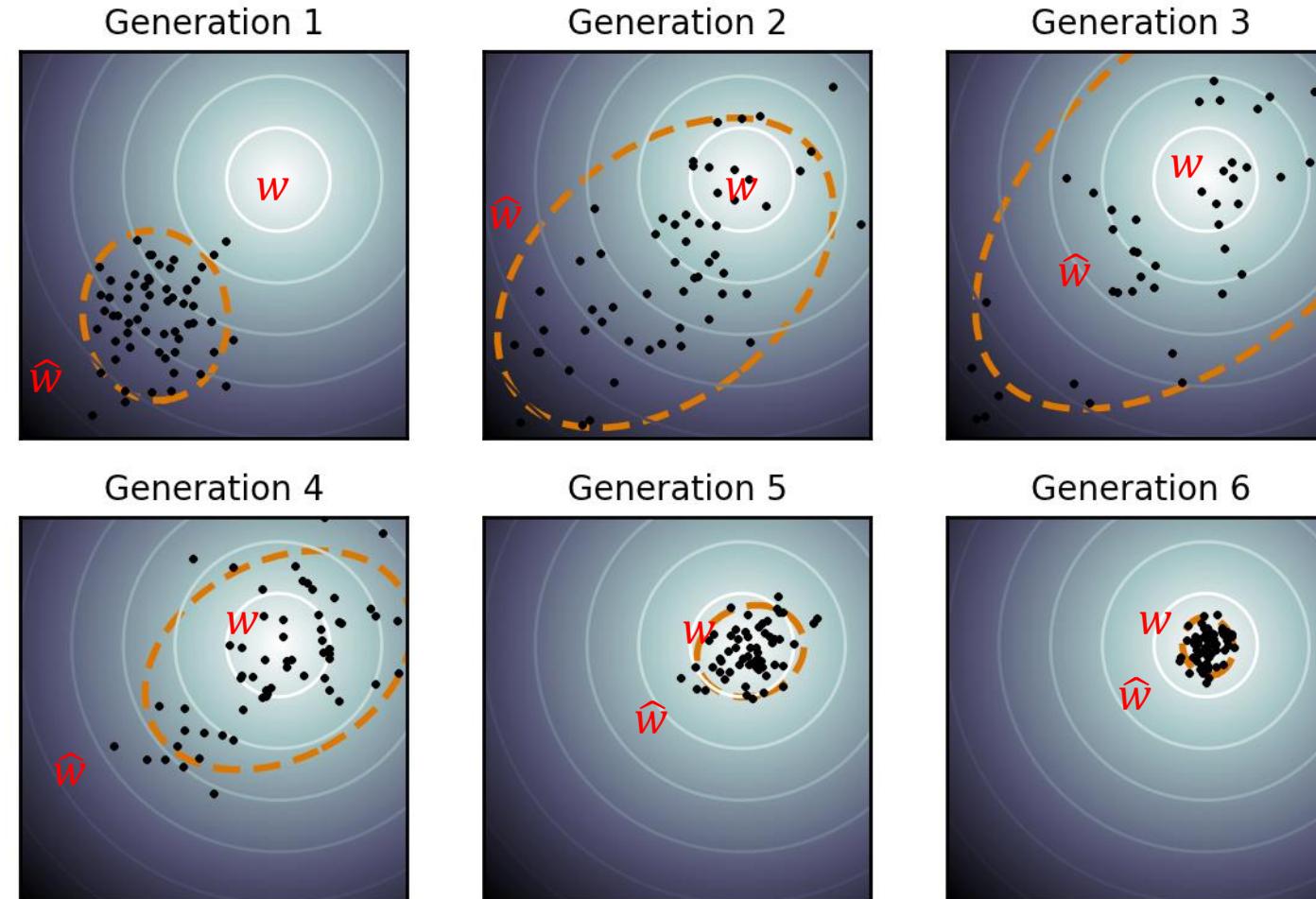
Covariance Matrix Adaptation Evolution Strategy (CMA-ES) Algorithm

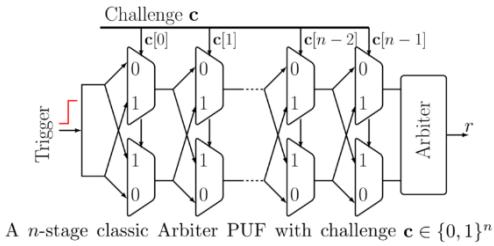
w : target

\hat{w} : estimator or model

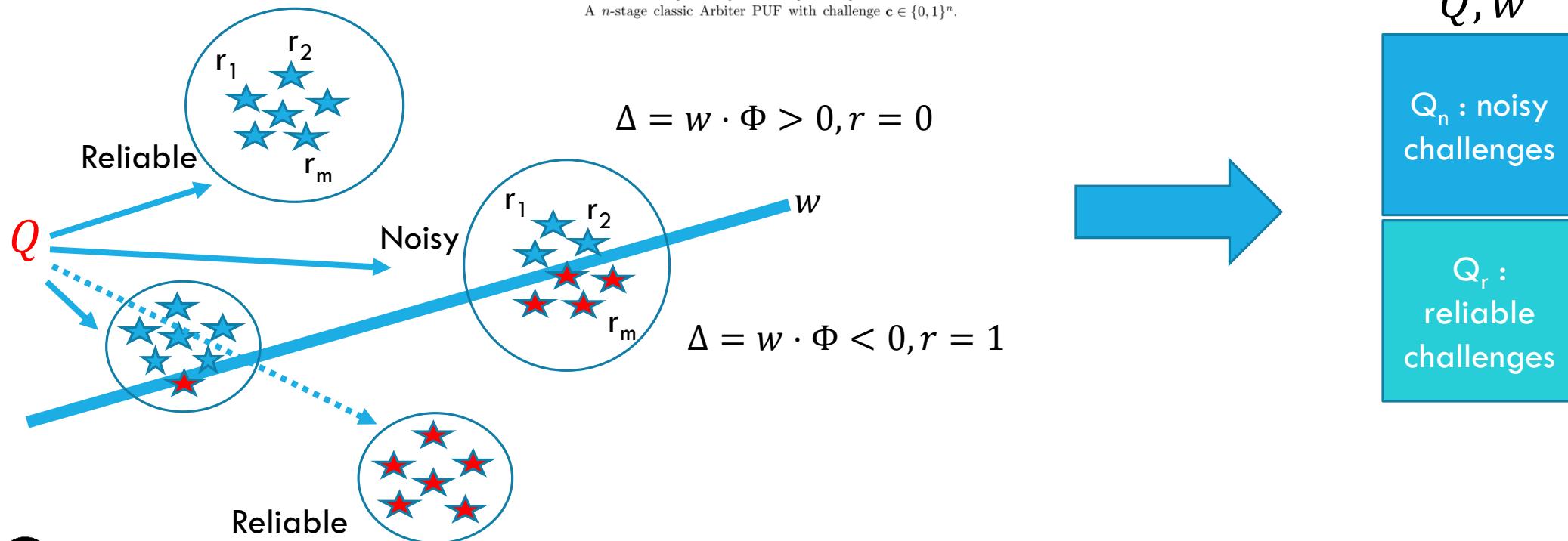


[4]

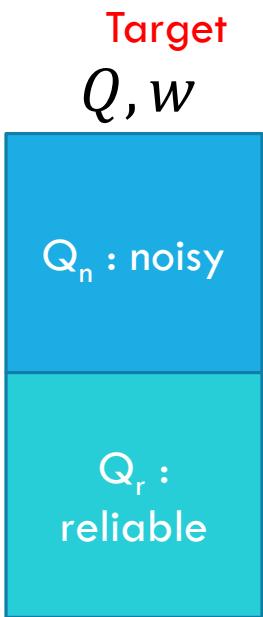




Q : set of CRPs ,
 w : APUF



Iteration 1



Model

Q, ϵ_1, \hat{w}_1

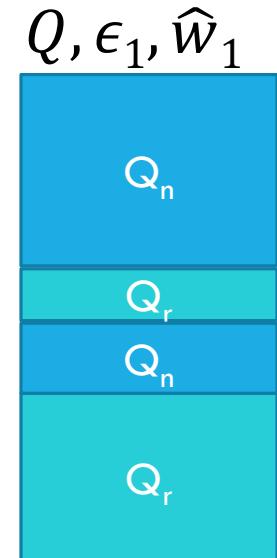
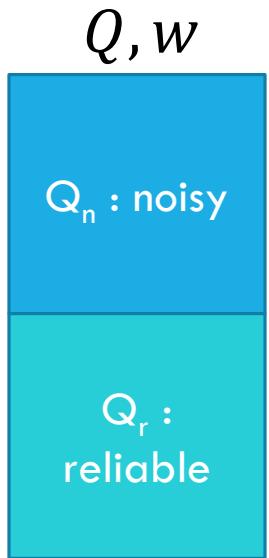
$Q \rightarrow \text{challenge } c \rightarrow \Phi(c) \rightarrow \Delta = \hat{w}_1 \cdot \Phi(c)$

$|\Delta| \leq \epsilon_1 \rightarrow \text{challenge } c \text{ is noisy} \longrightarrow Q_n : \text{noisy}$

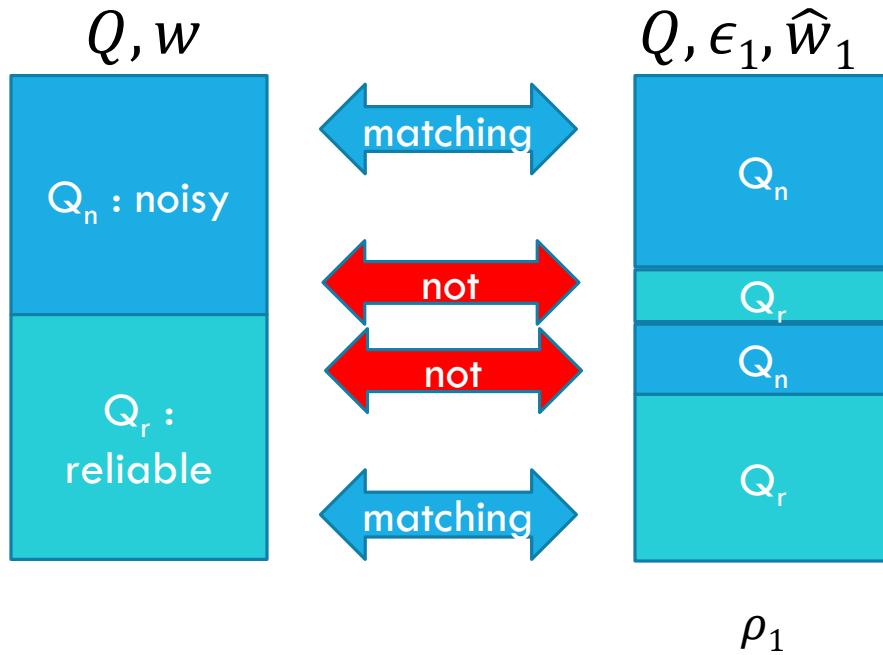
$|\Delta| > \epsilon_1 \rightarrow \text{challenge } c \text{ is reliable} \longrightarrow Q_r : \text{reliable}$



Iteration 1



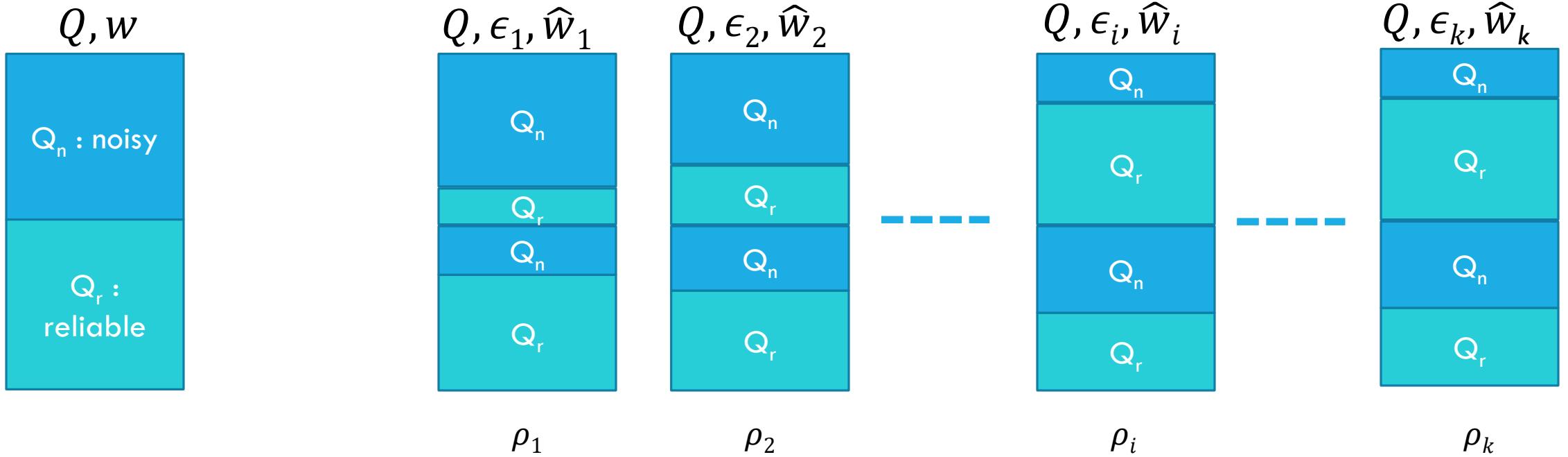
Iteration 1

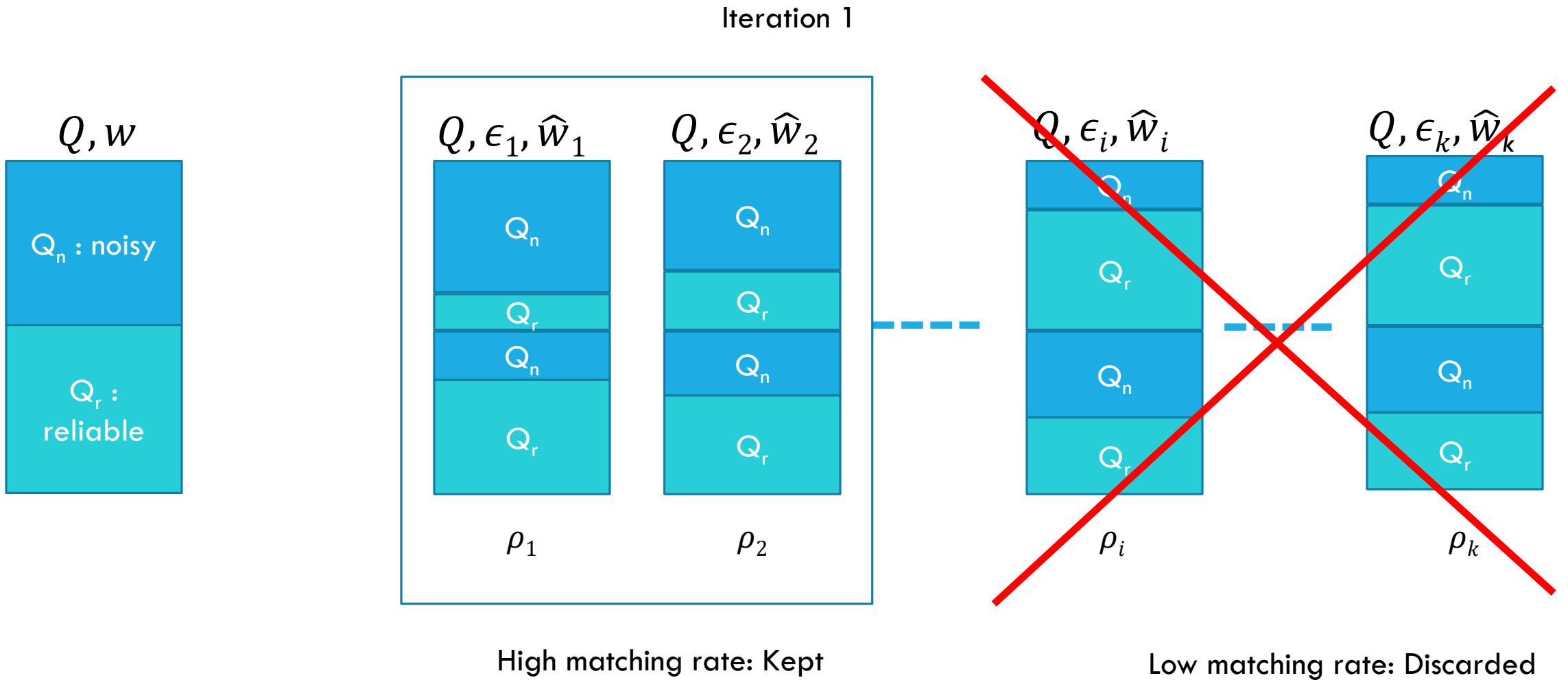


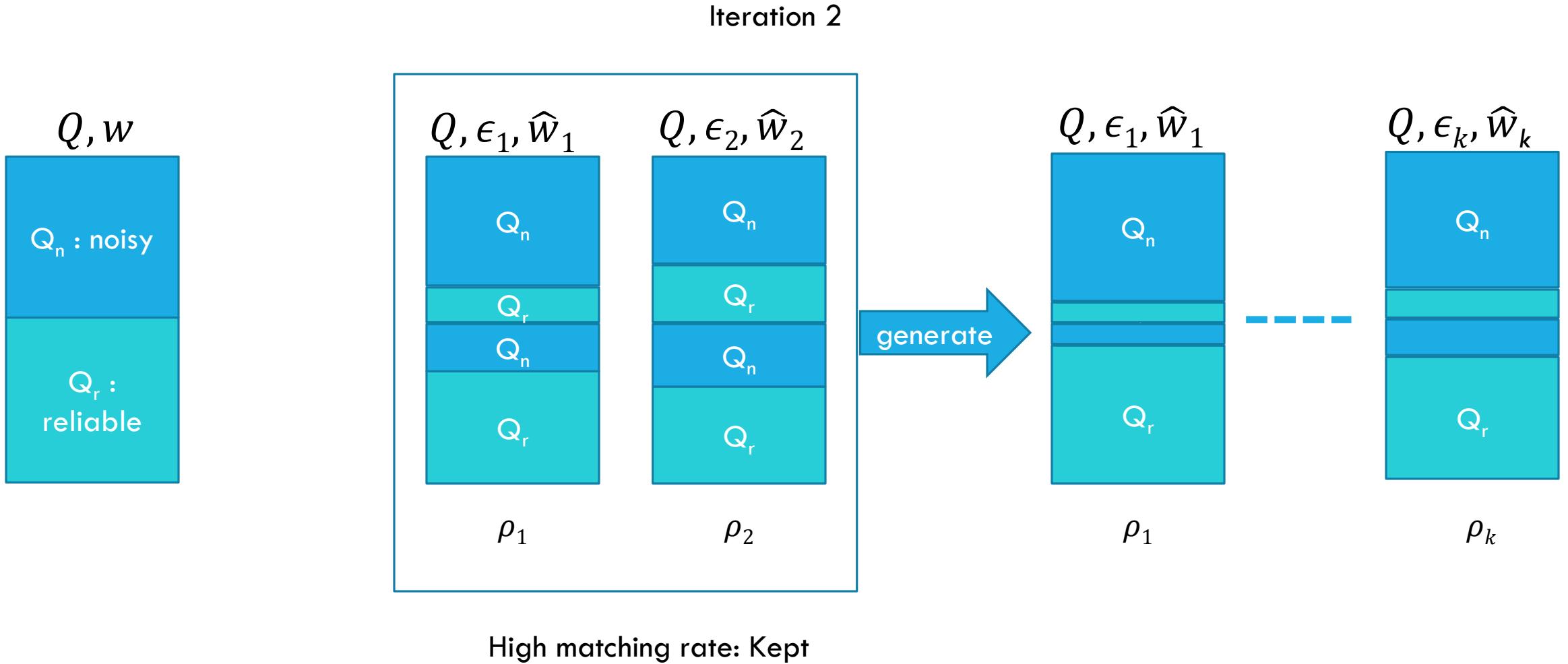
Compute the matching rate ρ_1 between $[Q, w]$ and $[Q, \epsilon_1, \hat{w}_1]$

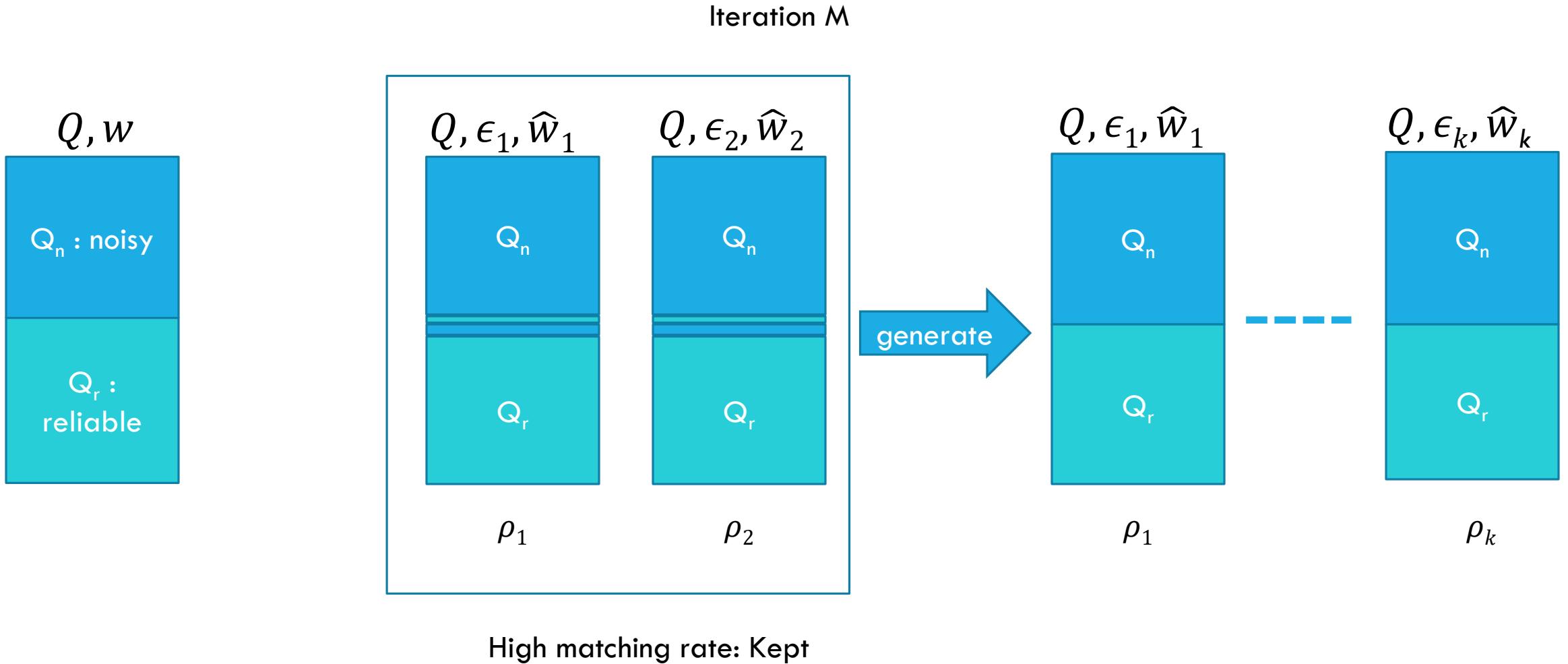


Iteration 1

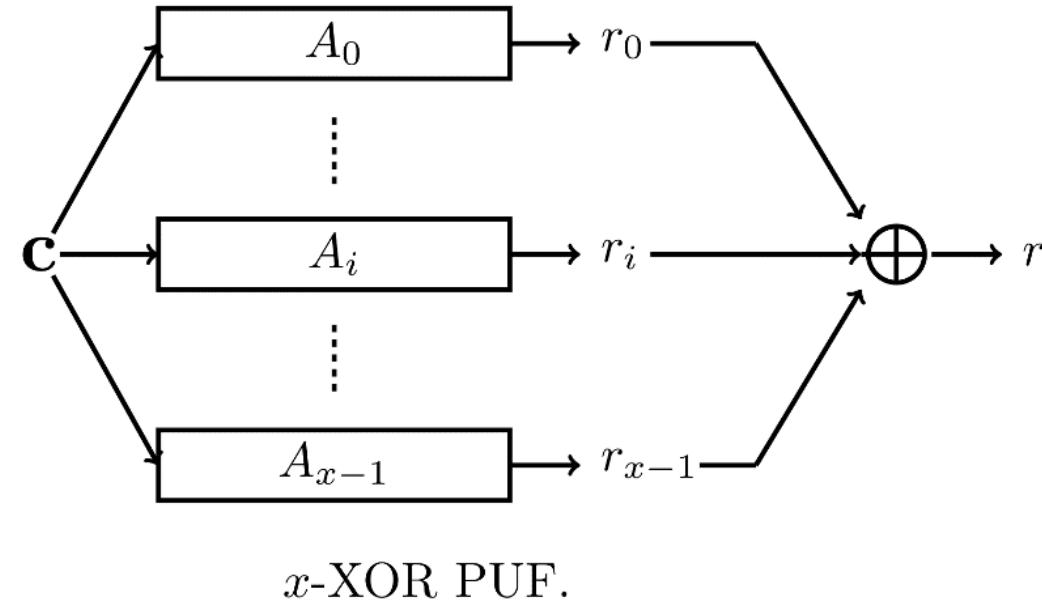




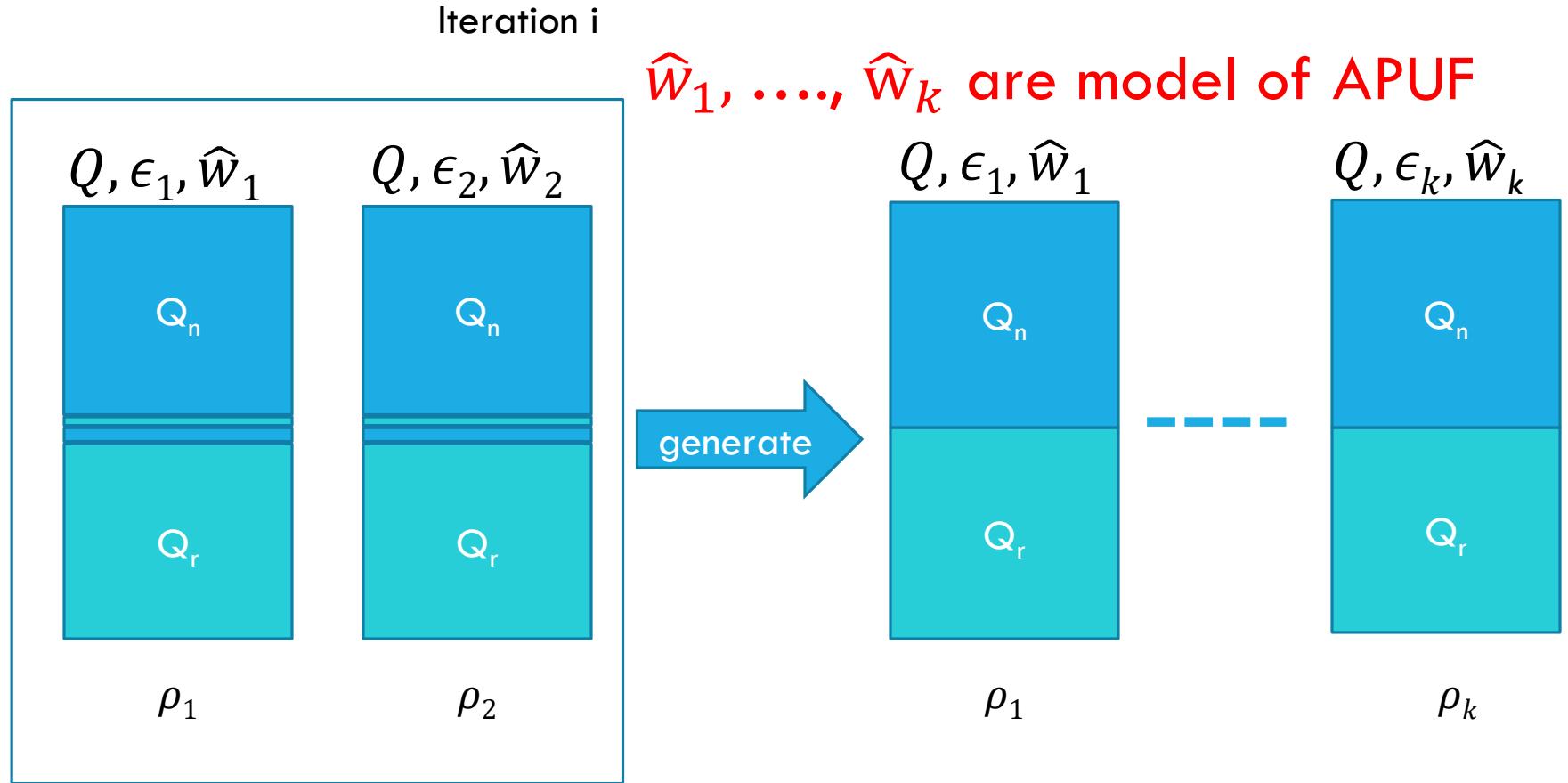
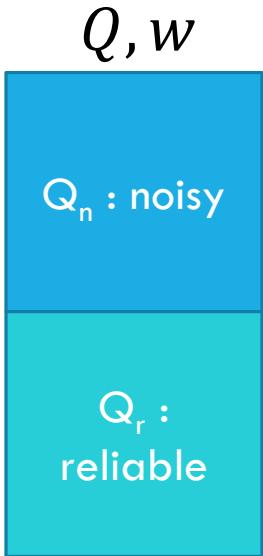




χ -XOR APUF



APUF $\rightarrow Q$

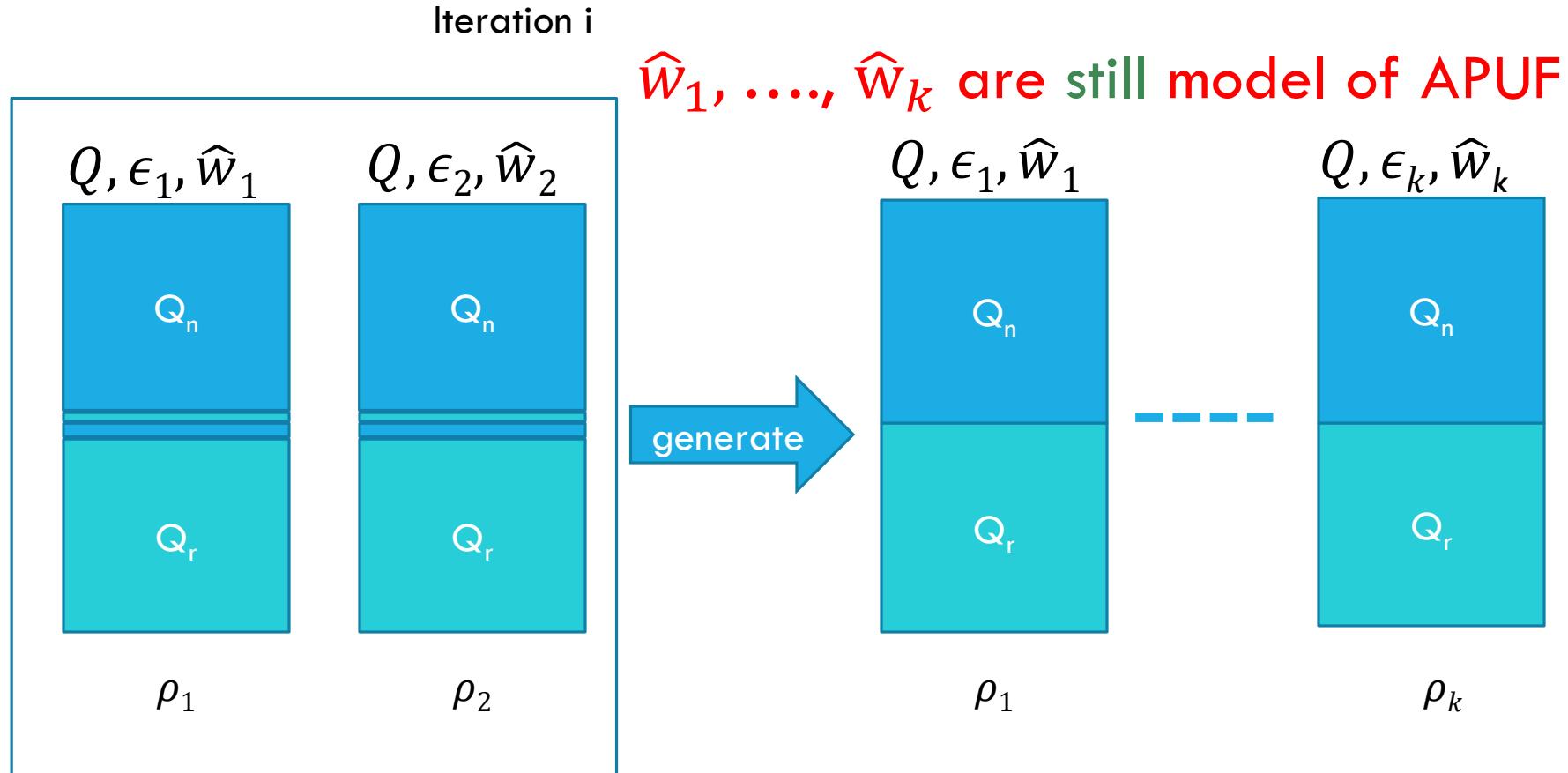


High matching rate: Kept

CMA-ES attack on APUF



XOR APUF $\rightarrow Q$

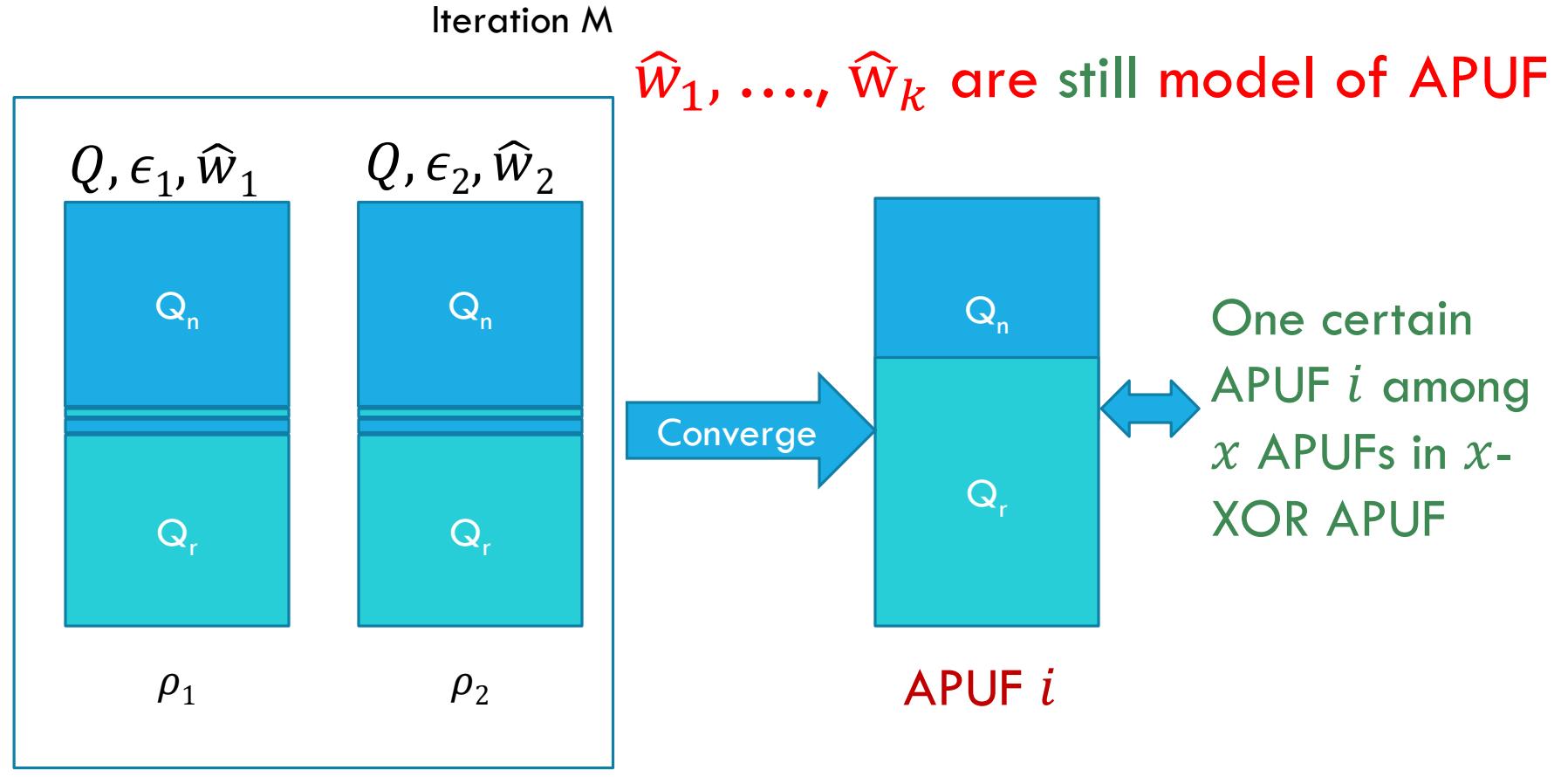
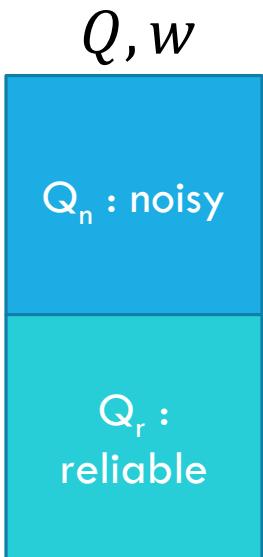


High matching rate: Kept

CMA-ES attack on XOR APUF



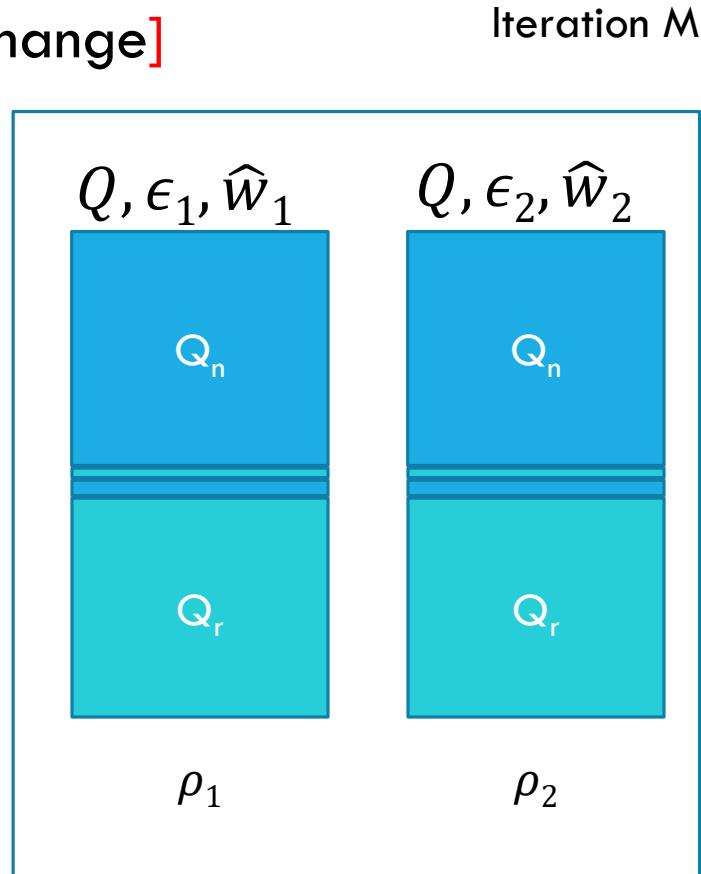
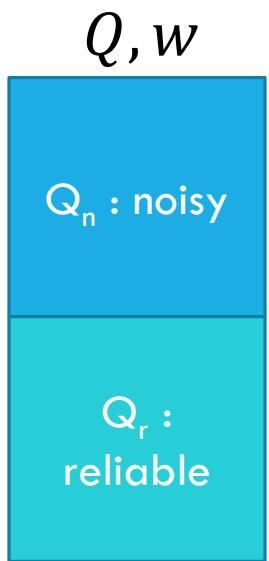
XOR APUF $\rightarrow Q$



CMA-ES attack on XOR APUF

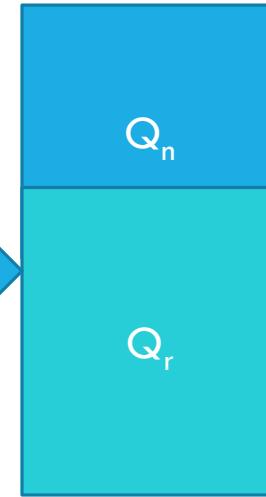


XOR APUF $\rightarrow Q$ [Change]



$\hat{w}_1, \dots, \hat{w}_k$ are still model of APUF

Converge



One certain
APUF j among
 x APUFs in x -
XOR APUF

High matching rate: Kept

CMA-ES attack on XOR APUF



Understanding Reliability based modeling attack on XOR PUF

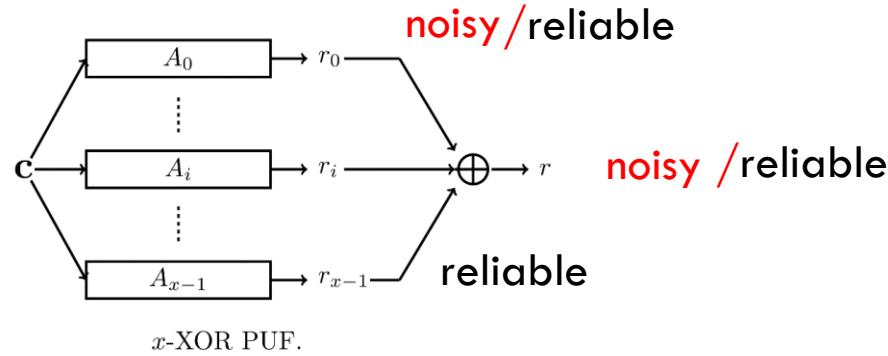
- Question 1: How does the attack on XOR PUF work?
- Question 2: How can we make the attack on XOR PUF fail?



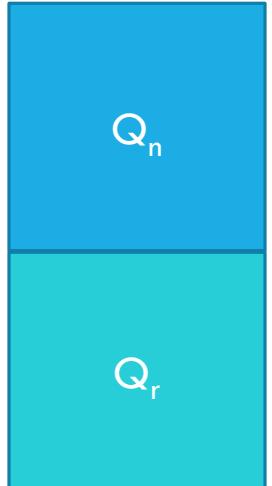
Question 1: How does the attack on XOR PUF work?



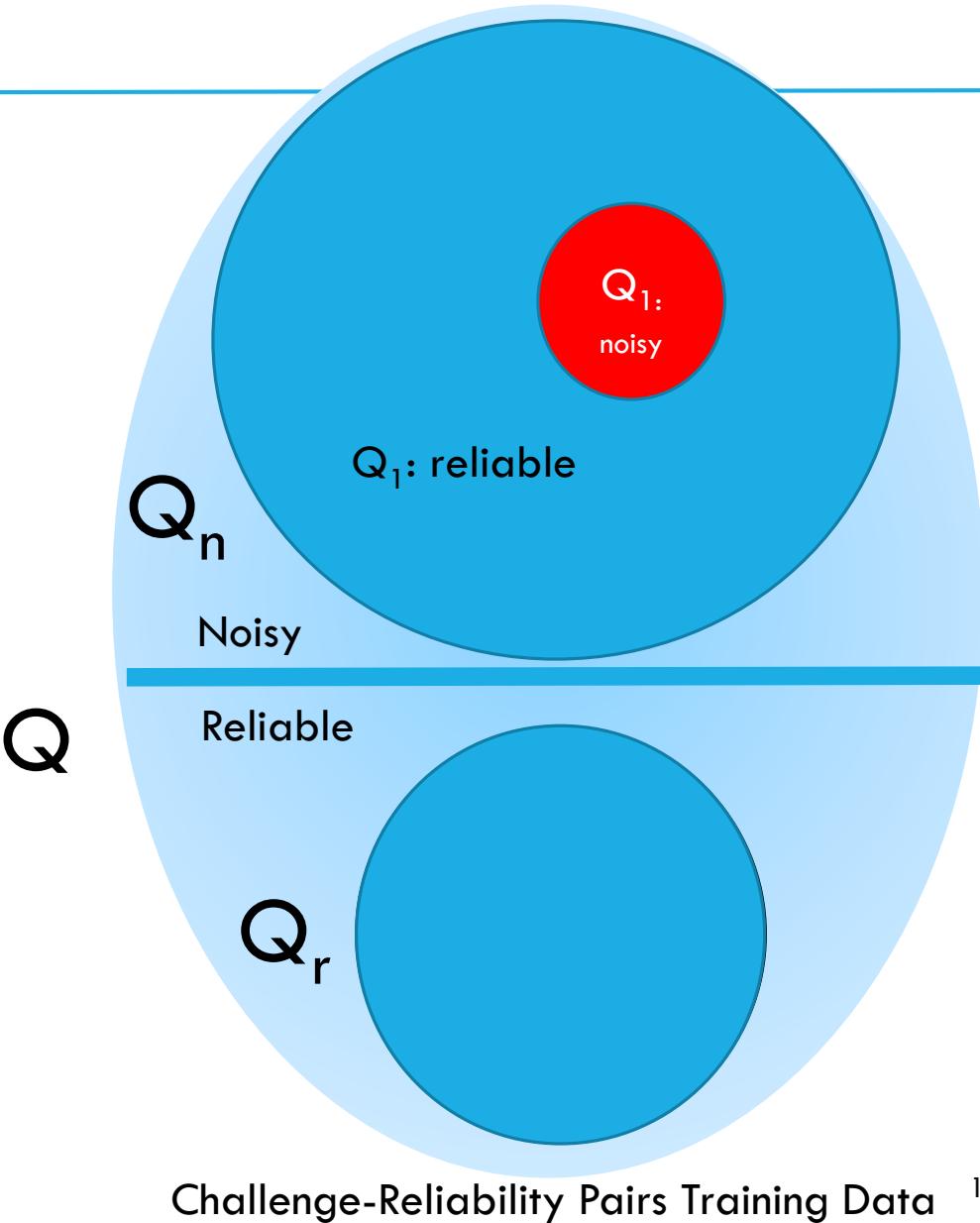
The noisy and reliable challenges in XOR PUF



XOR APUF

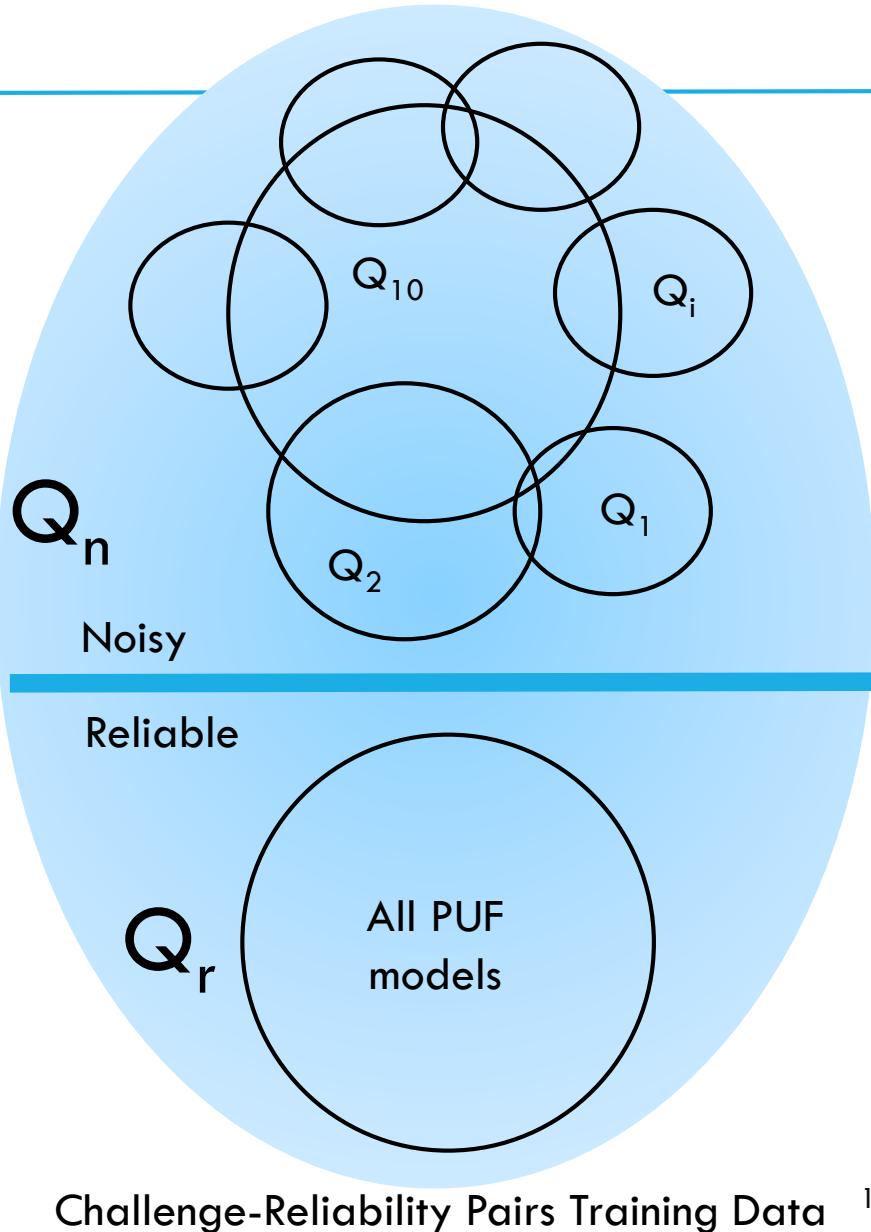
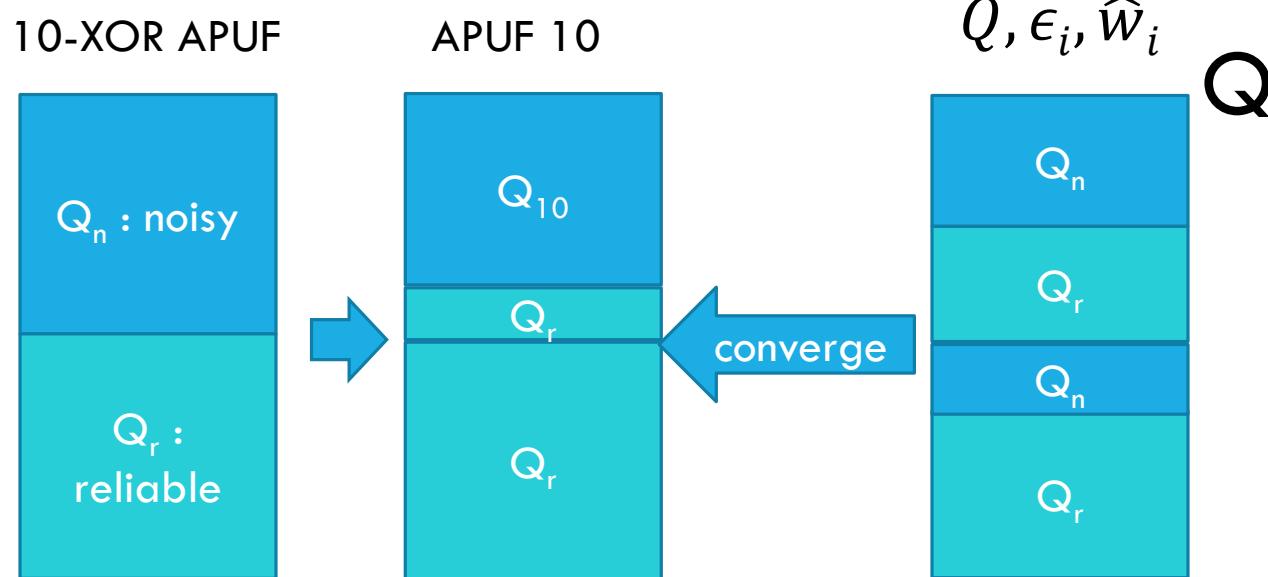


APUF 1



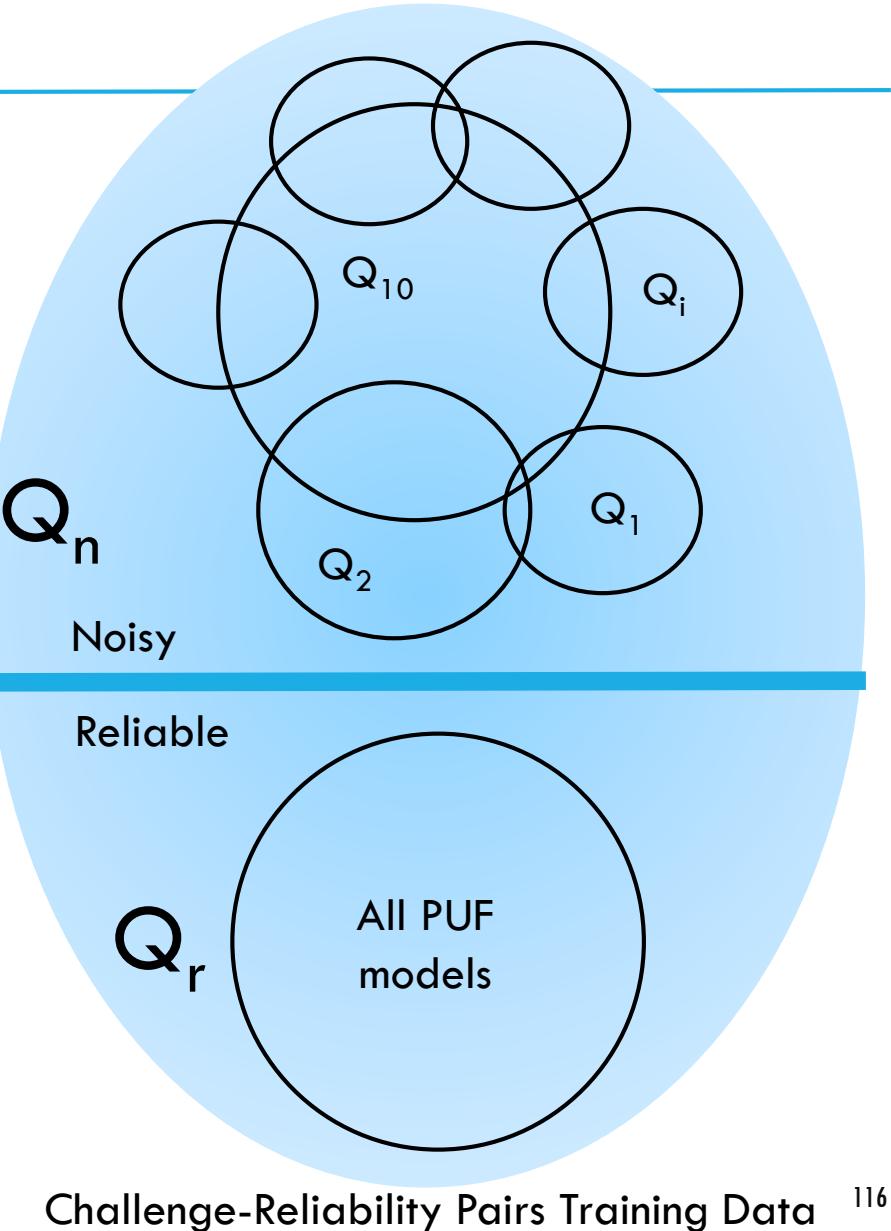
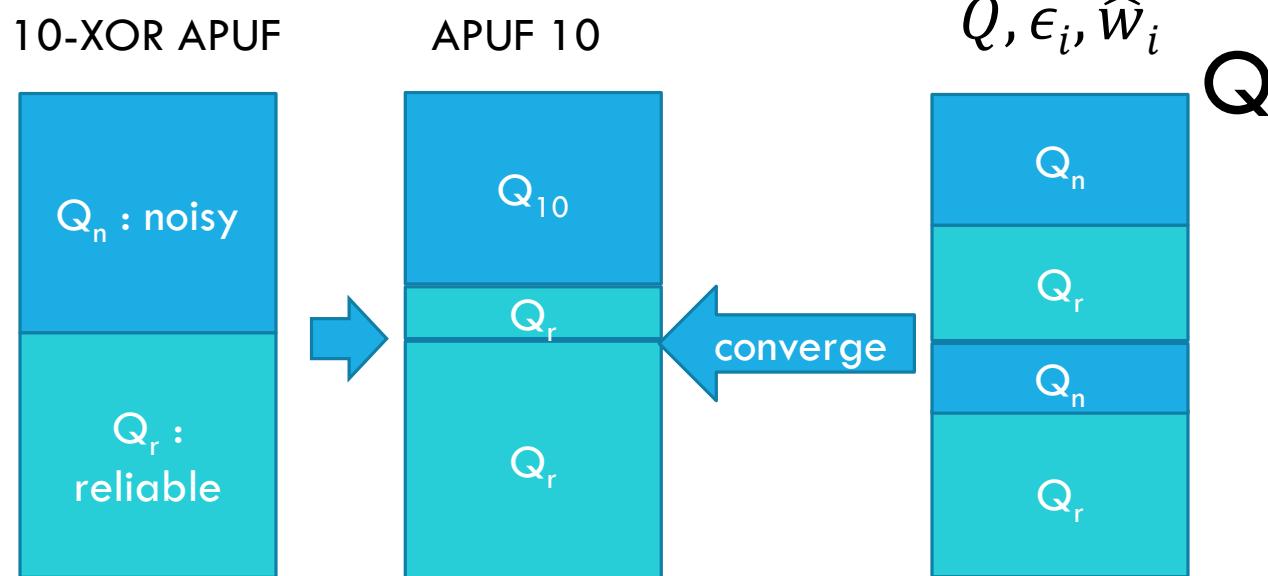
Key idea of the attack on XOR PUF

- (1) All the models \hat{w}_i in CMA-ES are **models of APUF**



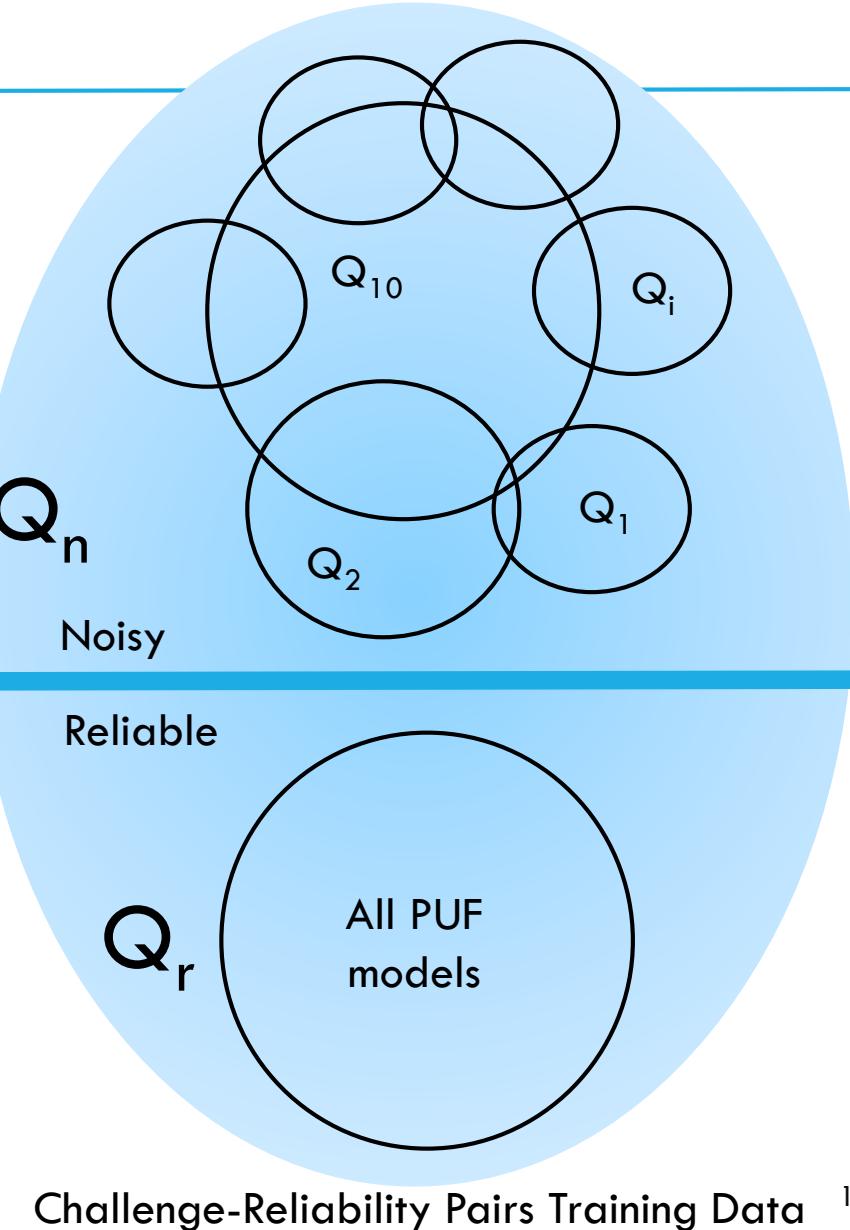
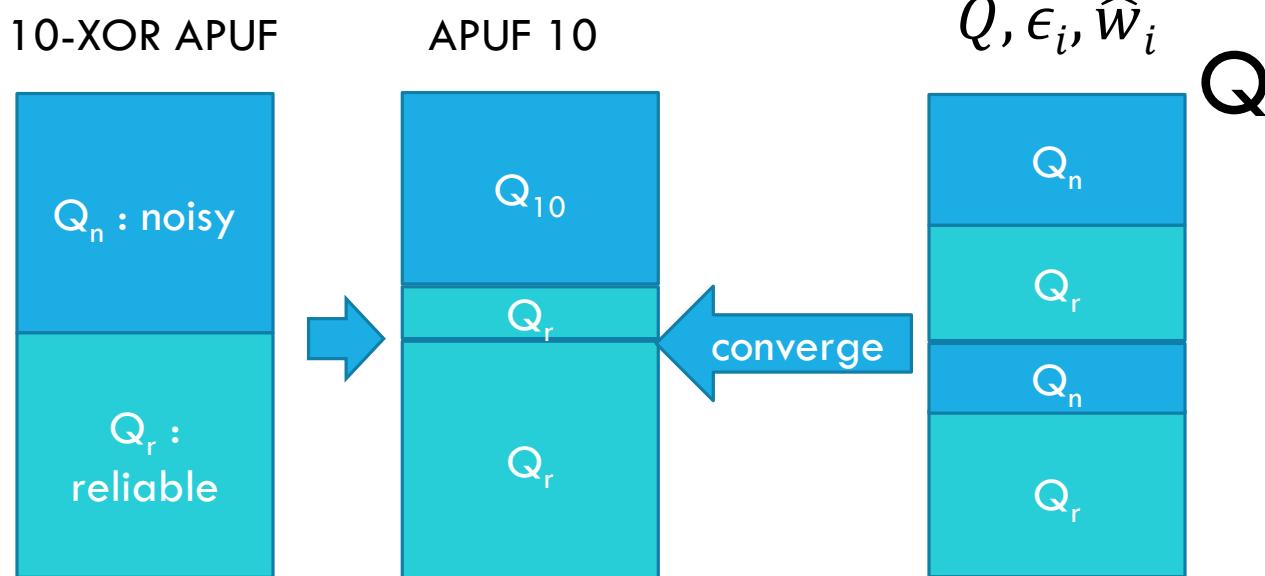
Key idea of the attack on XOR PUF

- (1) All the models \hat{w}_i in CMA-ES are **models of APUF**
- (2) \hat{w}_i can **only** converge to an APUF instance



Key idea of the attack on XOR PUF

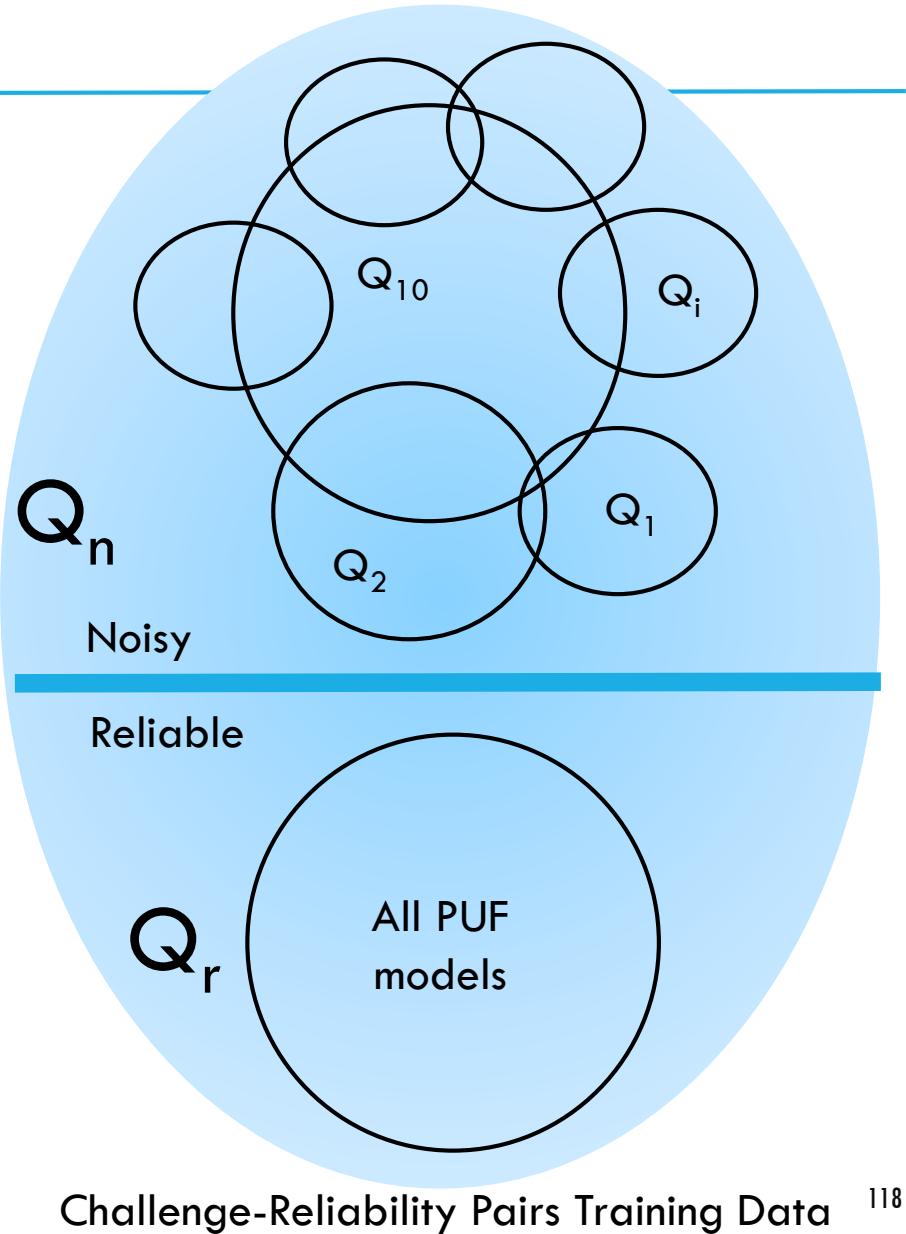
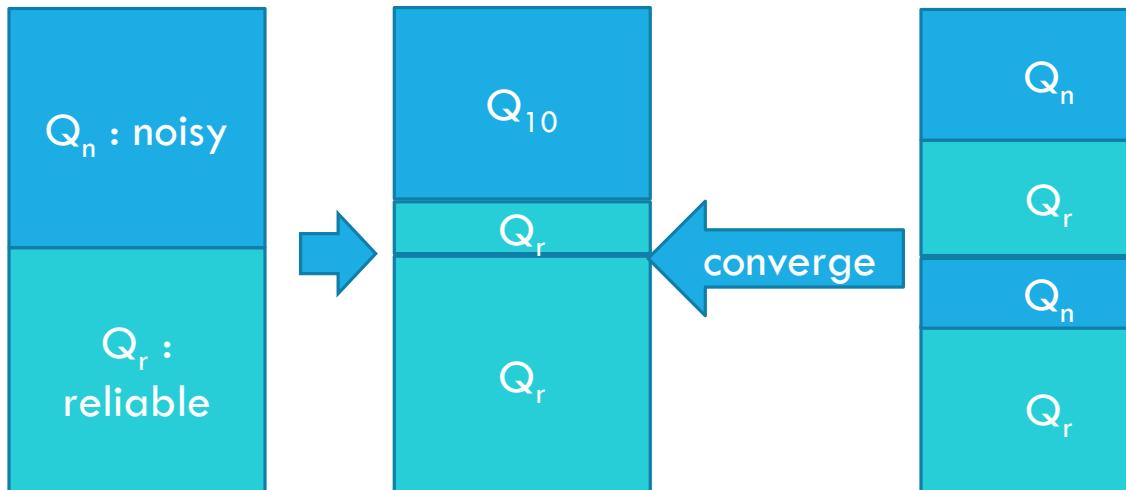
- (1) All the models \hat{w}_i in CMA-ES are **models of APUF**
- (2) \hat{w}_i can **only** converge to an APUF instance
- (3) CMA ES **maximizes the matching Q** of \hat{w}_i and Q of XOR APUF



Key idea of the attack on XOR PUF

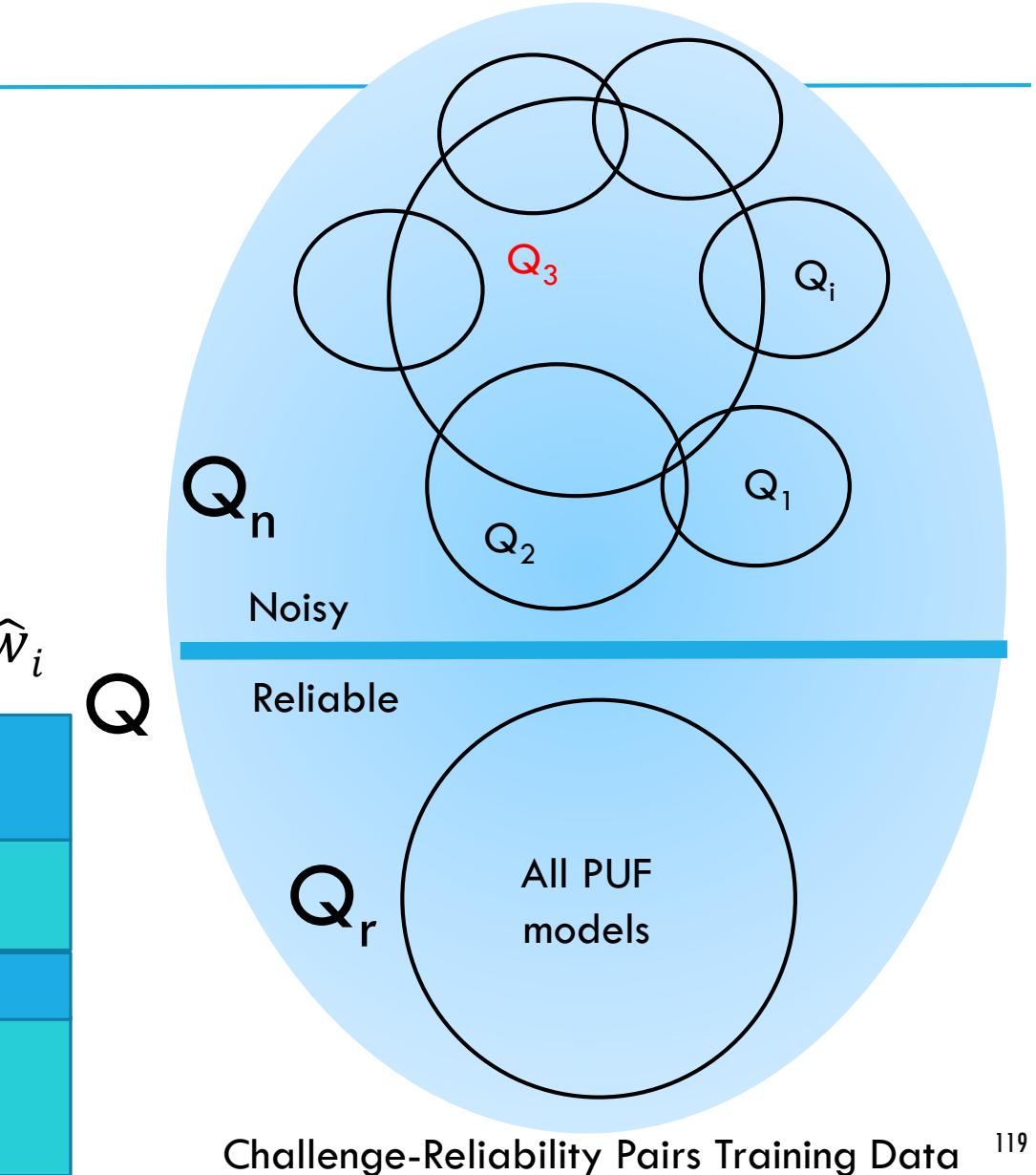
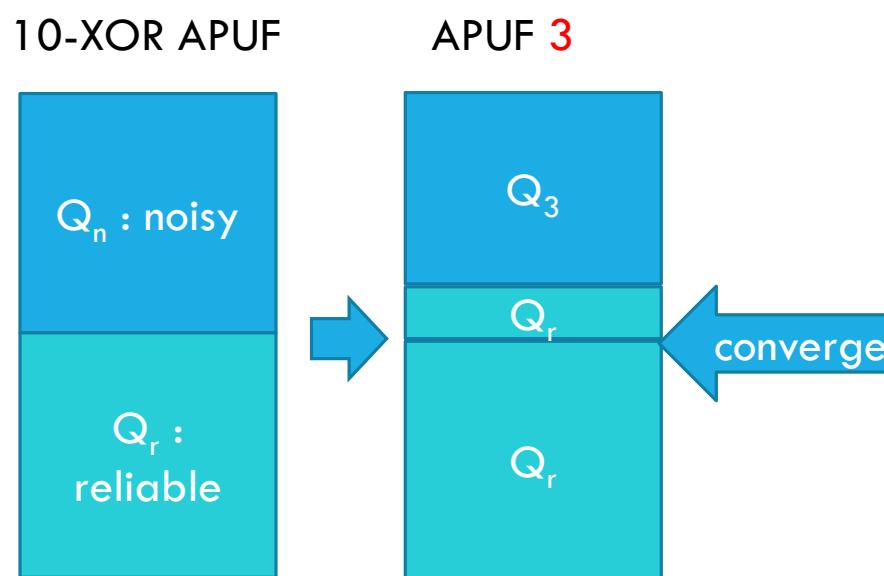
- (1) All the models \hat{w}_i in CMA-ES are **models of APUF**
- (2) \hat{w}_i can **only** converge to an APUF instance
- (3) CMA ES **maximizes the matching Q** of \hat{w}_i and Q of XOR APUF
- (1)+(2)+(3) CMA ES forces \hat{w}_i converges to APUF 10 because **Q of APUF 10 is the representative of Q of XOR APUF.**

10-XOR APUF APUF 10



Key idea of the attack on XOR PUF [5]

- Changing Q makes Q_3 largest

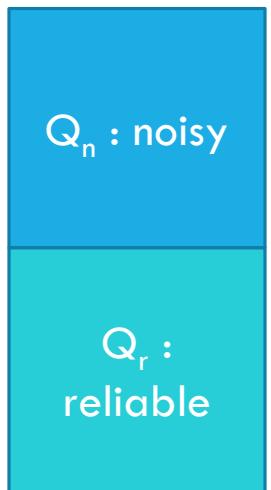


Key idea of the attack on XOR PUF

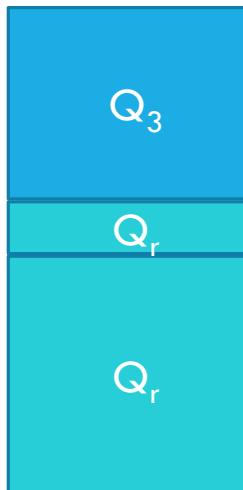
- Changing Q makes Q_3 largest

Keep change Q and apply CMA-ES attack on Q
to get models of all APUF instances

10-XOR APUF

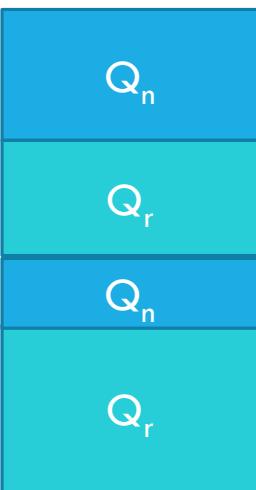


APUF 3

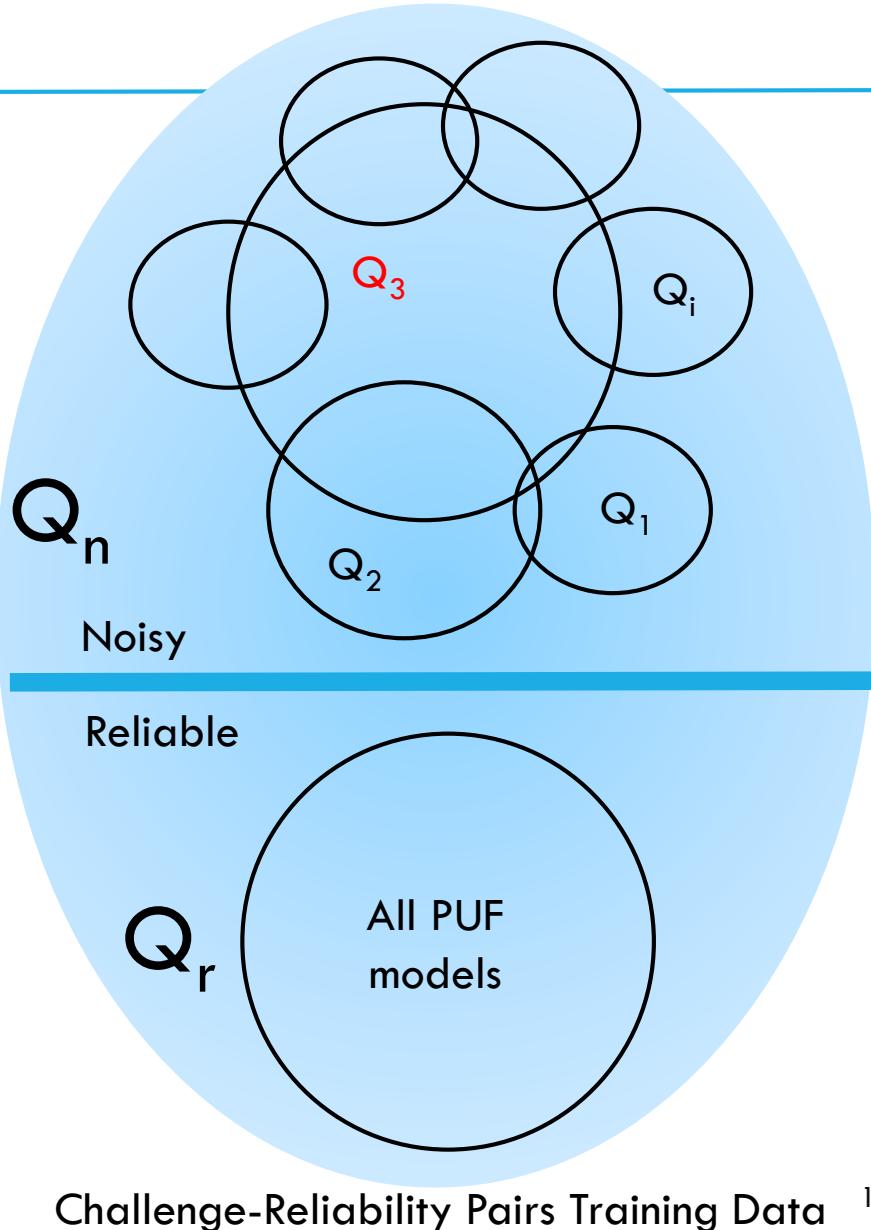


converge

Q, ϵ_i, \hat{W}_i



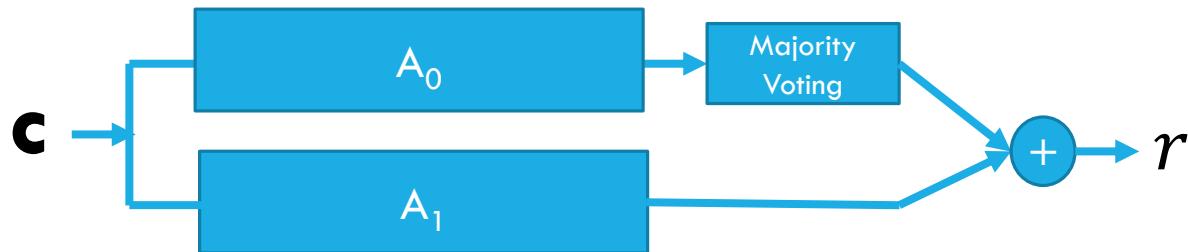
Q



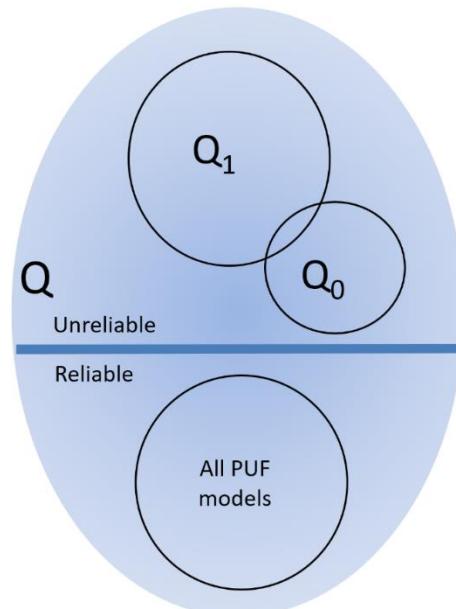
Question 2: How to make the attack on XOR PUF fail?



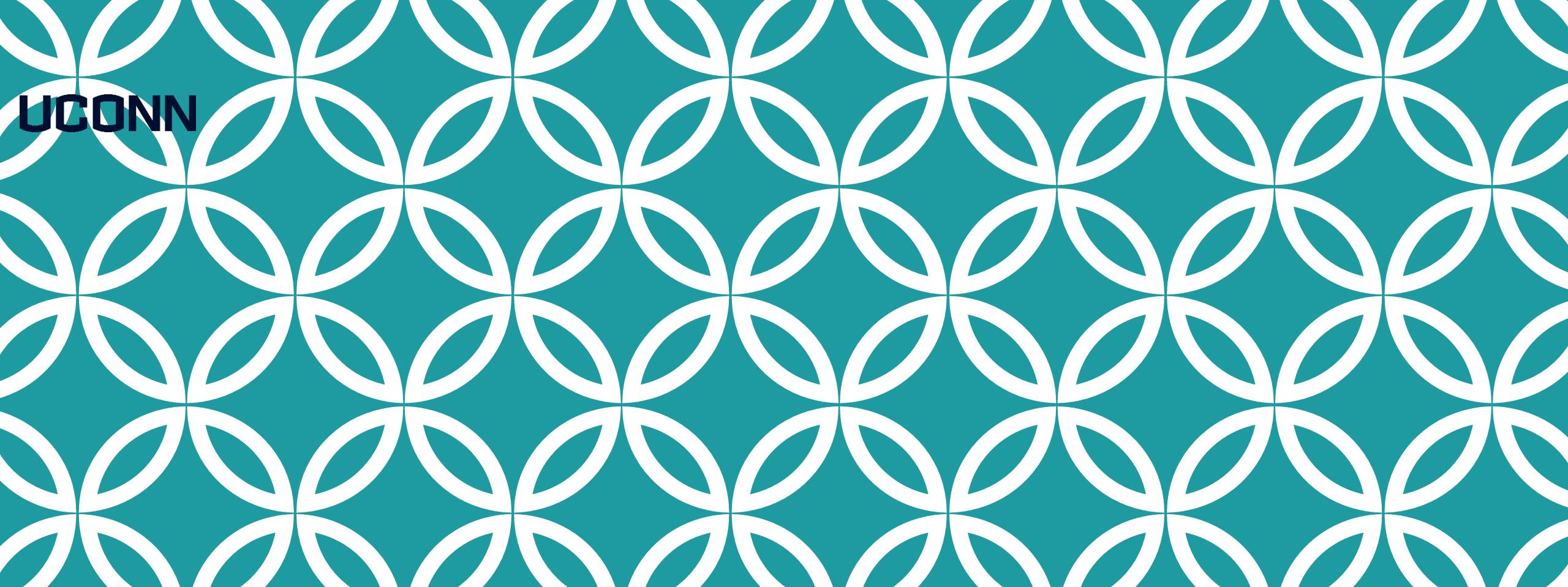
Attack fails



2-XOR APUF with majority voting circuit at A_0



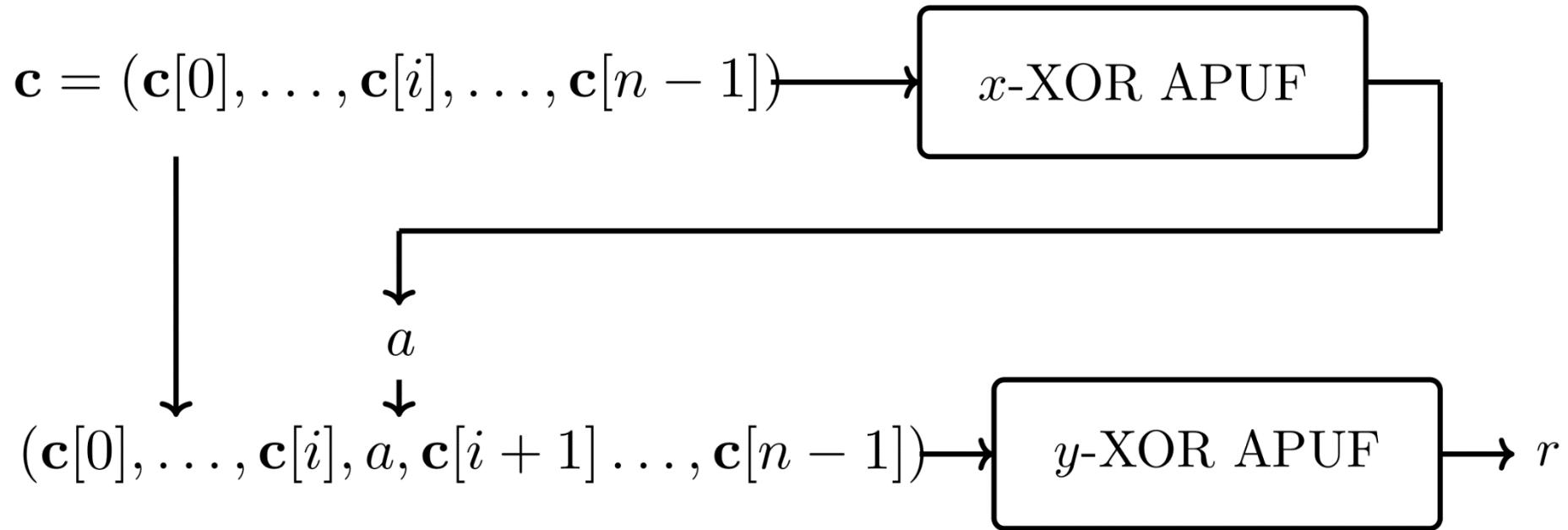
CMA ES never converges to APUF A_0 and always converges to APUF A_1 when majority voting mechanism in use.



9. Becker Attack Resistance of iPUF

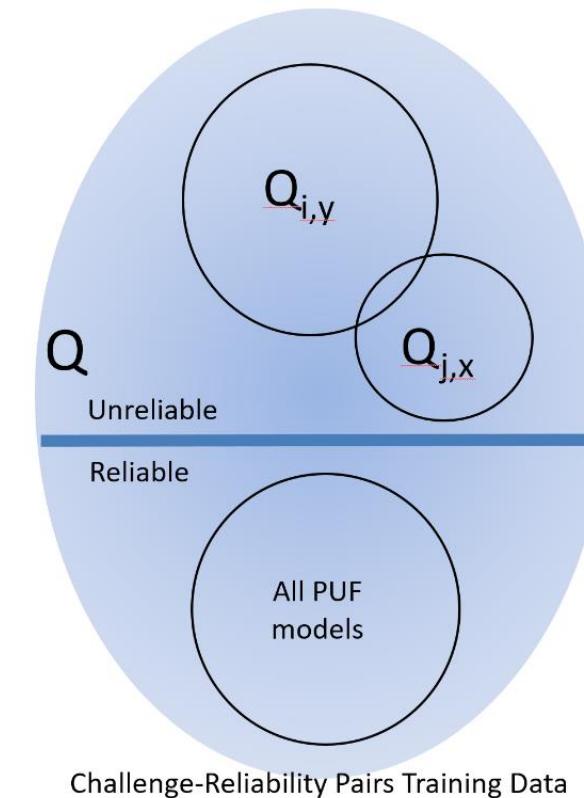
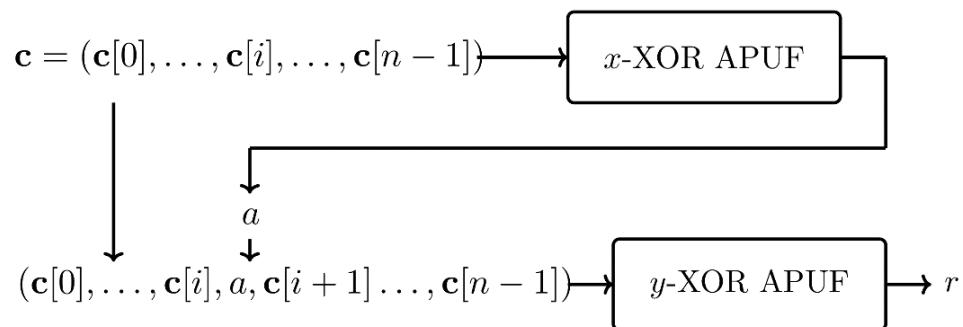


Interposed PUF (iPUF)

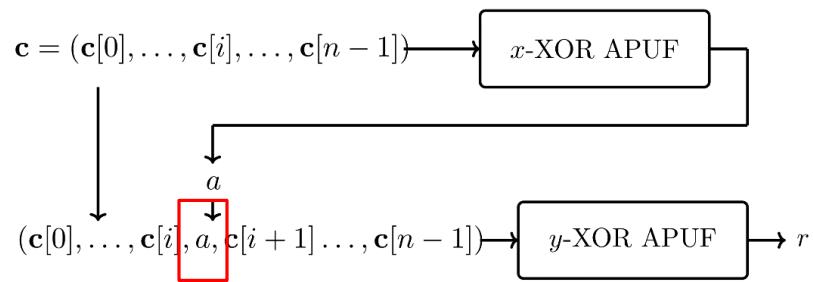
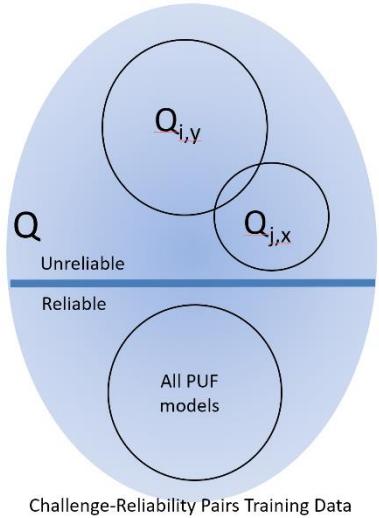


Security of iPUF wrt Reliability-based modeling attack

- Reason 1: the information of APUF instances in x-XOR PUF presented at the iPUF output is less compared to APUF instances in y-XOR PUF. Thus, the reliability based modeling attack never converges to any APUF instance in x-XOR PUF



- **Reason 2:** to attack APUFs at x-XOR APUF, the adversary needs to compute Δ . But compute Δ is infeasible because the output of x-XOR PUF (a) is not known.



Q, ϵ_1, \hat{w}_1

Cannot compute $\Phi(c)$ or Δ

$Q \rightarrow challenge\ c \rightarrow \Phi(c) \rightarrow \Delta = \hat{w}_1 \cdot \Phi(c)$

$|\Delta| \leq \epsilon_1 \rightarrow challenge\ c\ is\ noisy \longrightarrow Q_n : noisy$

$|\Delta| > \epsilon_1 \rightarrow challenge\ c\ is\ reliable \longrightarrow Q_r : reliable$

$$\Delta = \Delta(n-1) = \mathbf{w} \cdot \boldsymbol{\Phi}^T$$

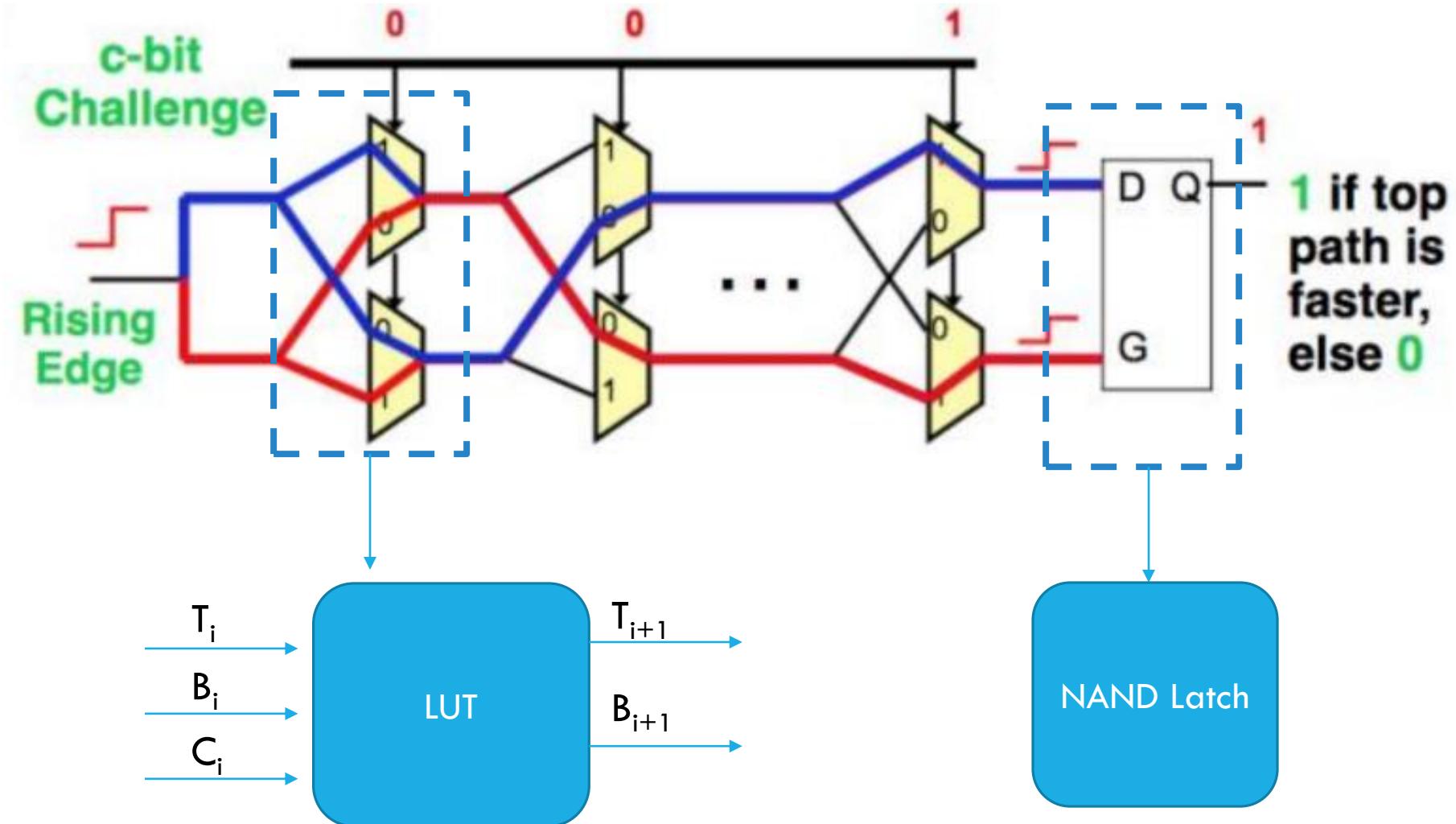
$$\boldsymbol{\Phi}[i] = \prod_{j=i, \dots, n-1} (1 - c[j]), i = 0, \dots, n-1$$

$$\boldsymbol{\Phi}[n] = 1$$

10. Implementation

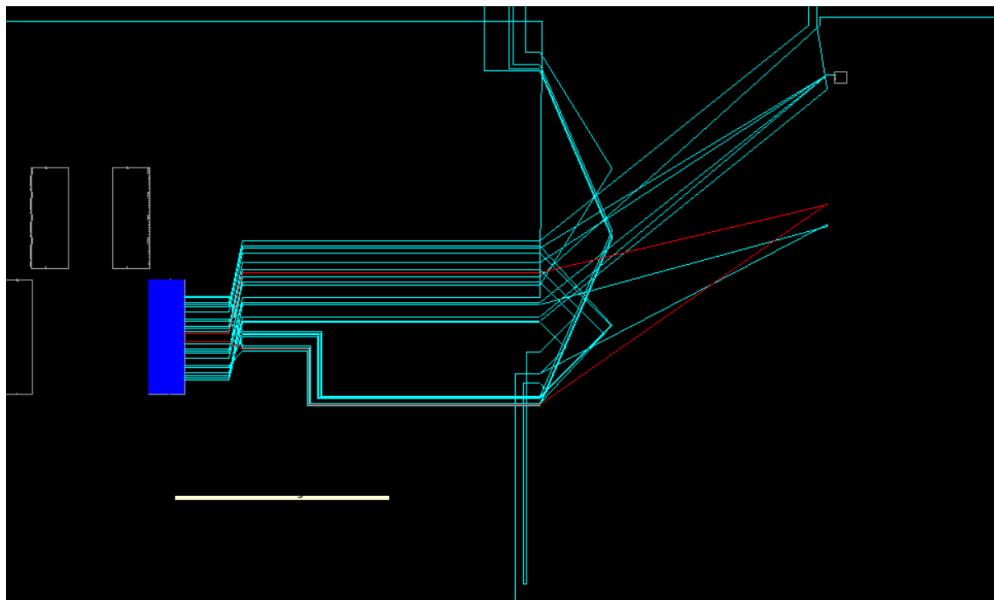


APUF Design

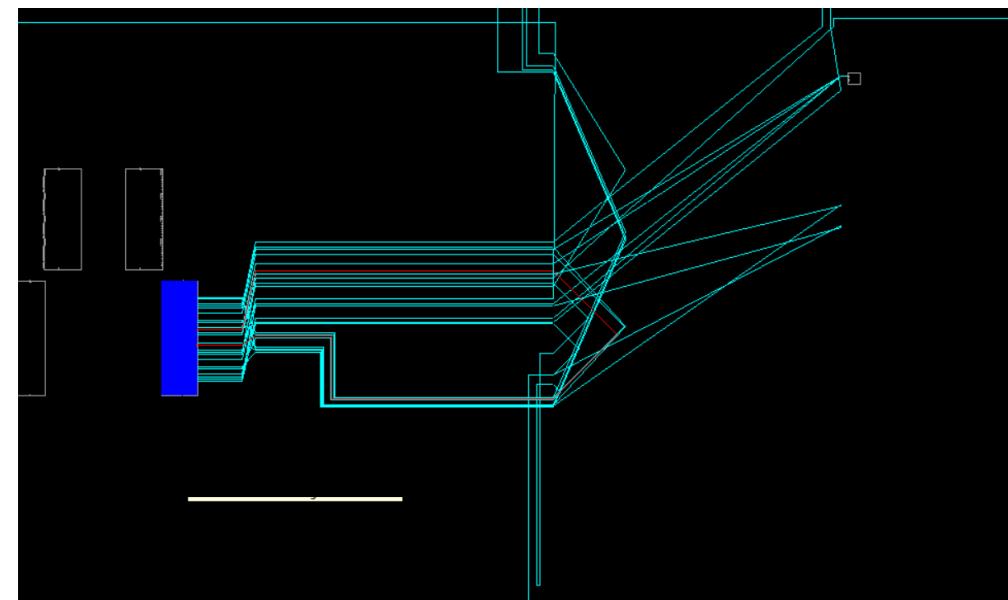


Implementing APUFs on FPGAs

- The APUFs on FPGAs have very poor uniqueness (same APUF designs on an FPGA will create almost the same APUF instances)
 - Process variation in FPGAs does not play a big role in creating an APUF instance.
 - Routing delay variation in FPGAs, which is out of the control of an FPGA user, largely influences the created APUF instance.



Top Path



Bottom Path

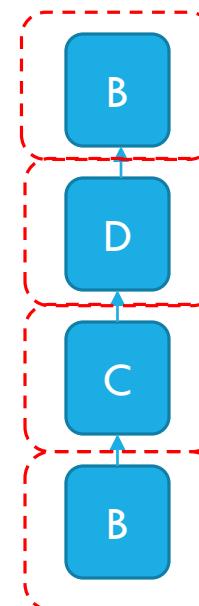
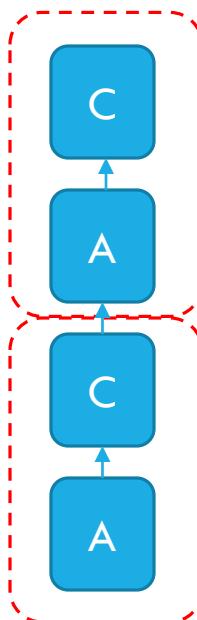
Solutions

- Vary the APUF design slightly in every instance
 - 1. Programmable Delay Line (PDL) – based APUF
 - Recent study shows a bias in the responses of PDL-based APUFs.
 - 2. Vary the placement of switches in different APUF instances
 - Our choice
 - Place the switches in one column of an FPGA and vary the selection of the LUTs in the column



Placement of Switches

- Pattern Placement
- Place the switches based on a pre-defined pattern (e.g., AC)
- Random Placement
- Place the switches in one randomly-selected LUT in each slice

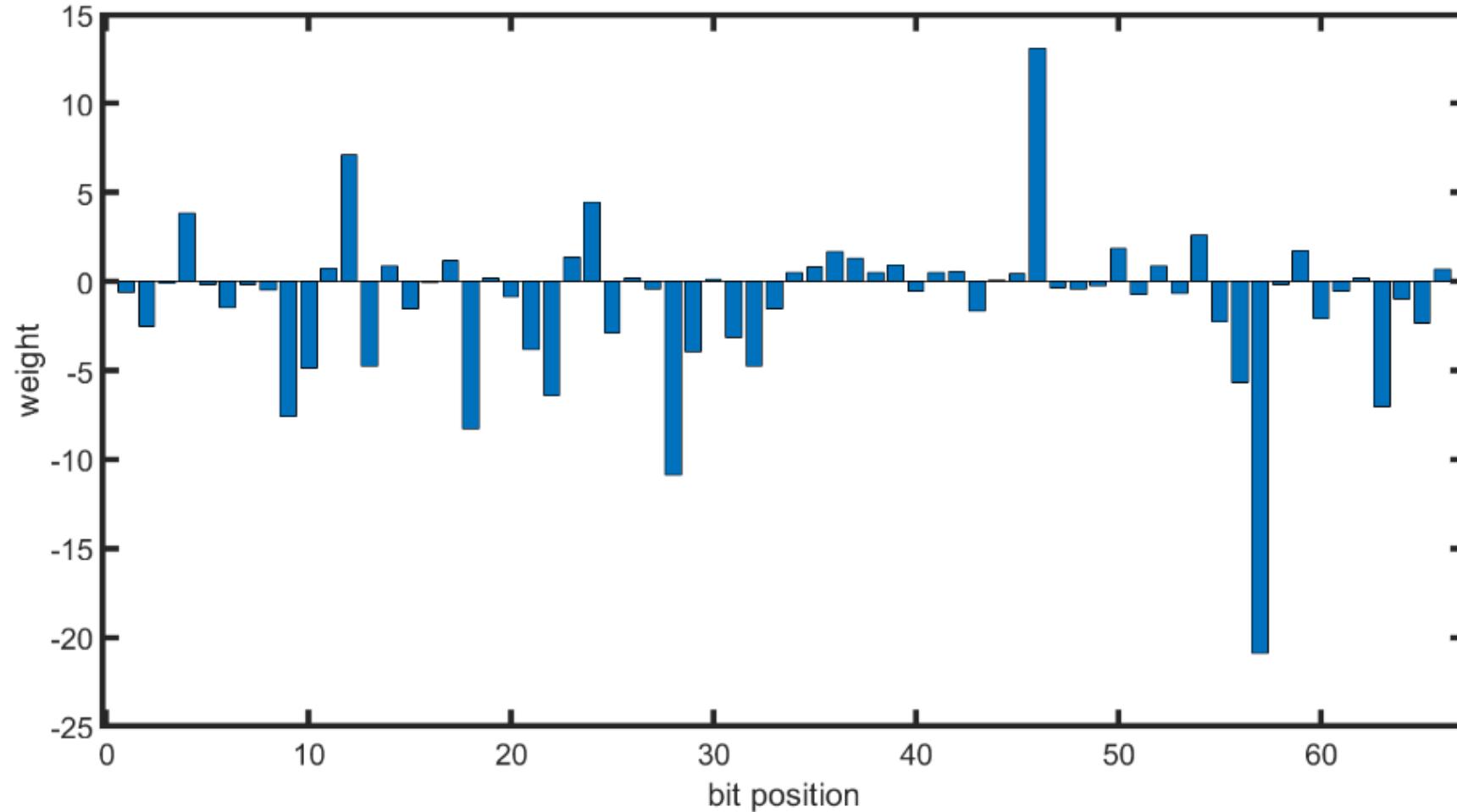


Placement of Switches

- Pattern Placement **Good Implementation**
 - Place the switches based on a pre-defined pattern (e.g., AC)
 - Characterize the uniformity, uniqueness, and reliability before using the PUF instance.
 - Pros:
 - The delay between each switch is more equally distributed.
 - The standard deviation of weights ranges from 1.11 to 3.77
 - Cons:
 - Limited choices of patterns
 - One may not be able to build a sufficiently large XOR APUF.
- Random Placement **Bad Implementation**
 - Place the switches in one randomly-selected LUT in each slice
 - Characterize the uniformity, uniqueness, and reliability before using the PUF instance.
 - Pros:
 - Lots of choices of placement
 - A lot more unique PUF instances
 - Cons:
 - Some delays may get significantly larger or smaller than the others.
 - It can sometimes lead to a successful attack on the PUF.
 - The standard deviation of weights ranges from 4.23 to 7.48



Distribution of the weights of a badly implemented APUF

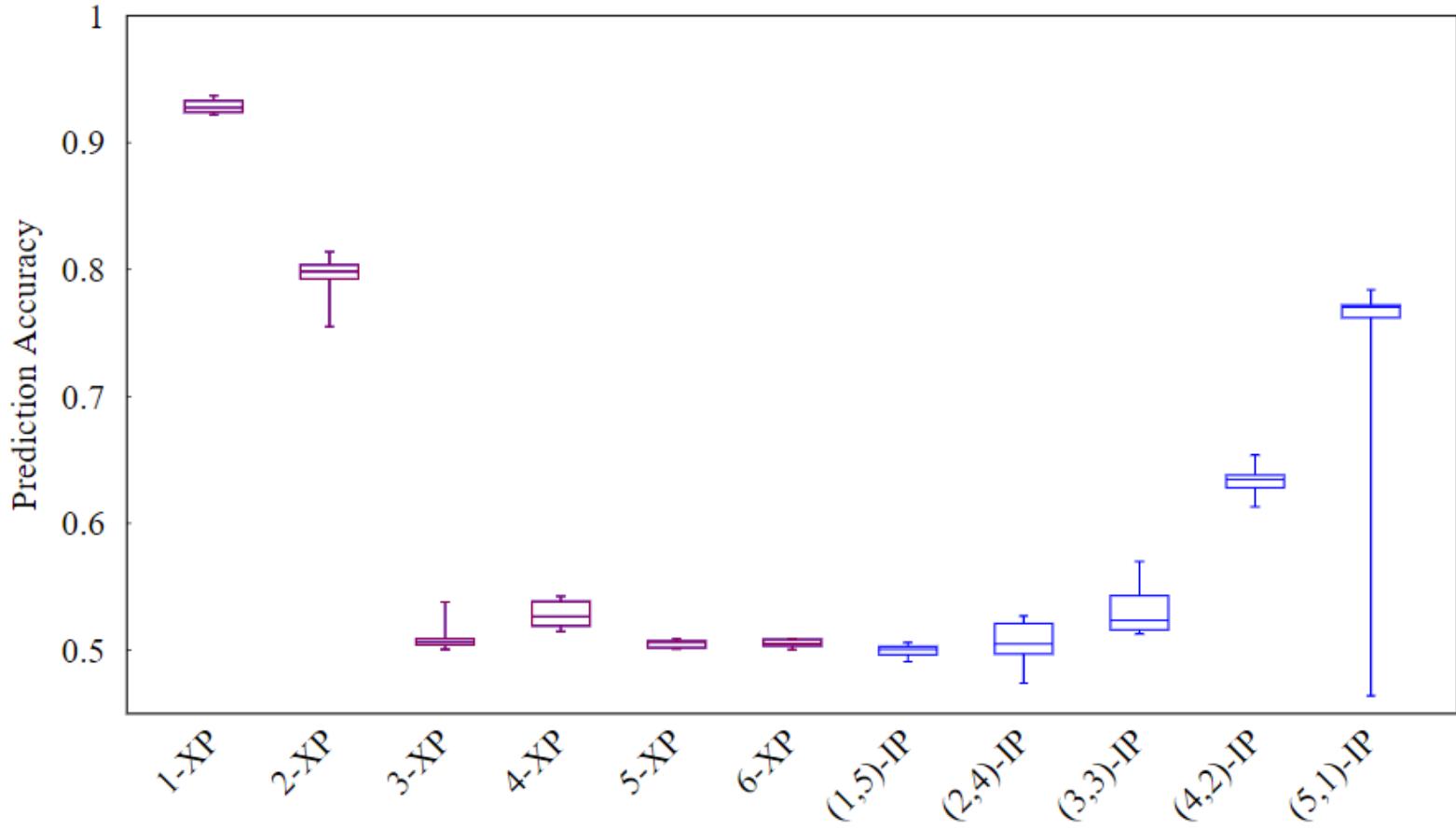


Reliability-based attack on XOR APUFs on FPGA

	#CRPs used in one attack	Overall Noise Rate	Average Prediction Accuracy
2-XOR	50,000	1.44%	98.22%
3-XOR	90,000	2.38%	96.38%
4-XOR	140,000	2.92%	96.15%
5-XOR	200,000	3.80%	96.62%
6-XOR	260,000	4.47%	91.58%*



Model Complexity Comparison between iPUF and XOR APUF



$$(x, y) - IPUF \approx \left(y + \frac{x}{2} \right) - XOR PUF$$



Reliability-based Attack on iPUF

- iPUF is secure against reliability-based attack by its design, so we used an (1,1)-iPUF in the experiments
- An (1,1)-iPUF with random placement (bad implementation) **can be** attacked by the reliability-based attack
 - The model of the upper layer APUF was discovered, then an attacker can carry on to attack the lower layer.
- An (1,1)-iPUF with pattern placement (good implementation) **cannot be** attacked by the reliability-based attack
 - No correct model was built



Open Source

- All the PUF simulation, attacks, and FPGA implementation can be found at

https://github.com/scluconn/DA_PUF_Library



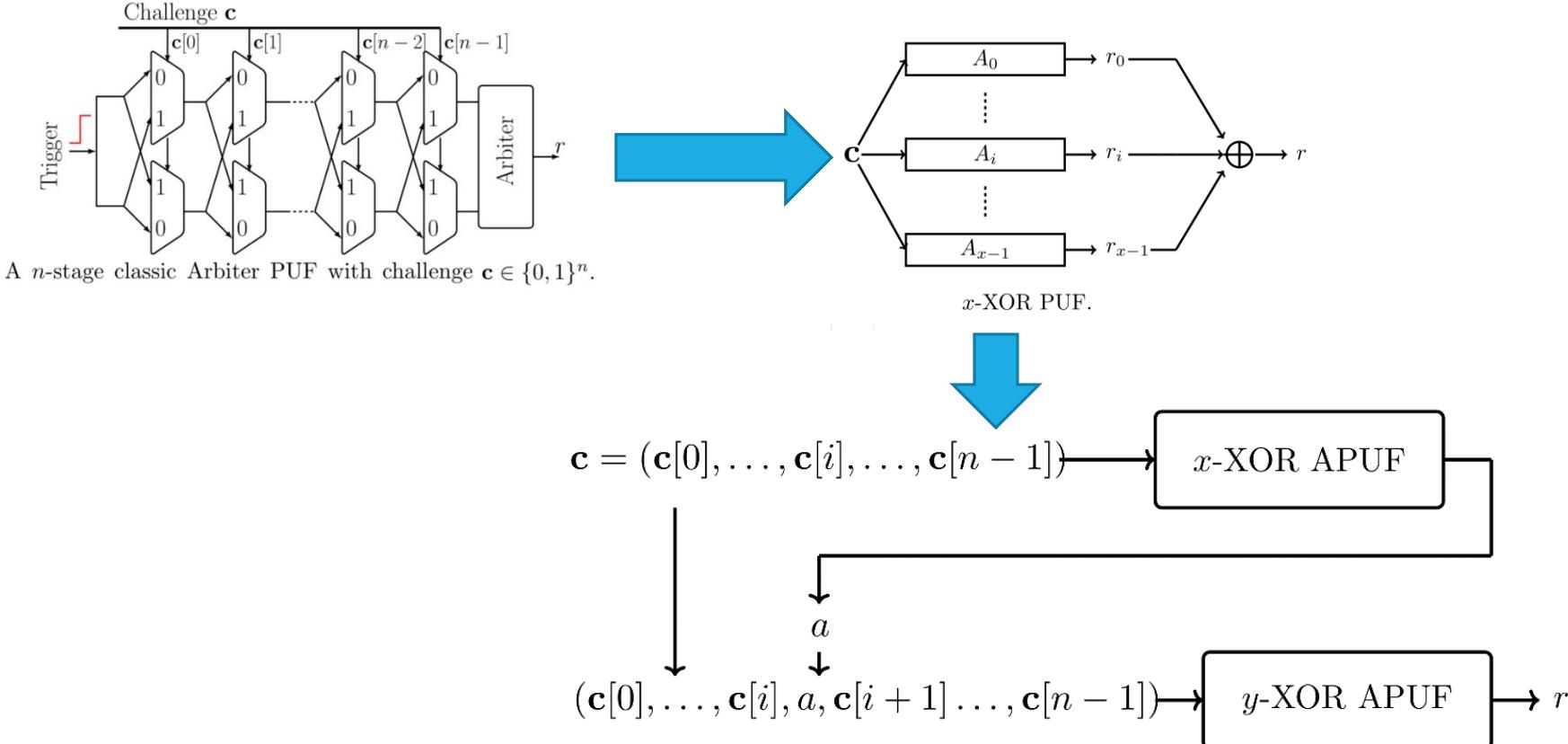
Other Contributions

- Enhanced Reliability based Modeling Attacks on APUF and XOR APUFs
- Prove why Logistic Regression on XOR APUF is the best attack and why Logistic Regression on iPUF is not applicable
- Implementation of iPUF and of modeling attacks on APUFs, XOR APUF, and iPUF
https://github.com/scluconn/DA_PUF_Library/



Conclusion

- We explain how the reliability-based modeling attack on XOR PUF works
- We propose a new lightweight PUF design (iPUF) which is secure against the state-of-the art of modelling attacks.



Literature

1. <https://slideplayer.com/slide/3927633/>
2. Cryptanalysis of electrical PUFs via machine learning algorithms – Master Thesis of Jan Solter
3. The Gap Between Promise and Reality: On the Insecurity of XOR Arbiter PUFs
CHES, September 16 th , 2015, Georg T. Becker
4. <https://en.wikipedia.org/wiki/CMA-ES>

Thank you for your attention!
and any questions?

