


# Watermarking PRFs from Lattices via Extractable PRFs

Sam Kim and David J. Wu

# Watermarking Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMWW19]

```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```



Embed a “mark” within a program



```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```

If mark is removed, then program is destroyed


Two main algorithms (simplified):

- $\text{Mark}(C) \rightarrow C'$ : Takes a circuit  $C$  and outputs a marked circuit  $C'$
- $\text{Verify}(C') \rightarrow \{0,1\}$ : Tests whether a circuit  $C'$  is marked or not

# Watermarking Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```



```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```

Embed

Notion extend to setting where watermark can be any string

If mark is removed, then program is destroyed

Two main algorithms (simplified).

- $\text{Mark}(C) \rightarrow C'$ : Takes a circuit  $C$  and outputs a marked circuit  $C'$
- $\text{Verify}(C') \rightarrow \{0,1\}$ : Tests whether a circuit  $C'$  is marked or not

# Watermarking Programs


[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```

Mark



```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```



**Functionality-preserving:** On input a circuit  $C$ , the Mark algorithm outputs a circuit  $C'$  where


$$C(x) = C'(x)$$

on almost all inputs  $x$

# Watermarking Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```



```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk = _mm_aesenc_si128(*blk, sched[j]);
}
```

**Unremovability:** Given a marked program  $C'$ , no efficient adversary can construct a circuit  $C^*$  where

- $C^*(x) = C'(x)$  on almost all inputs  $x$
- The circuit  $C^*$  is unmarked:  $\text{Verify}(C^*) = 0$

# Watermarking Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk =
}
```



```
static void AES_enc_blk(block *blk, const AES_KEY *key) {
    unsigned j, rnds = ROUNDS(key);
    const __m128i *sched = ((__m128i *) (key->rd_key));
    *blk = _mm_xor_si128(*blk, sched[0]);
    for (j = 1; j < rnds; ++j) {
        *blk = _mm_aesenc_si128(*blk, sched[j]);
    }
    *blk =
}
```

Adversary is very powerful: sees the code of the marked program  $C'$  and has complete flexibility in crafting  $C^*$


**Unremovability:** Given a marked program  $C'$ , no efficient adversary can construct a circuit  $C^*$  where

- $C^*(x) = C'(x)$  on almost all inputs  $x$
- The circuit  $C^*$  is unmarked:  $\text{Verify}(C^*) = 0$

# Watermarking Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

```
static void AES_enc_blk(block *blk, const AES_KEY *key) {  
    unsigned j, rnds = ROUNDS(key);  
    const __m128i *sched = ((__m128i *) (key->rd_key));  
    *blk = _mm_xor_si128(*blk, sched[0]);  
    for (j = 1; j < rnds; ++j) {  
        *blk = _mm_aesenc_si128(*blk, sched[j]);  
    }  
    *blk = _mm_aesenc_si128(*blk, sched[j]);  
}
```



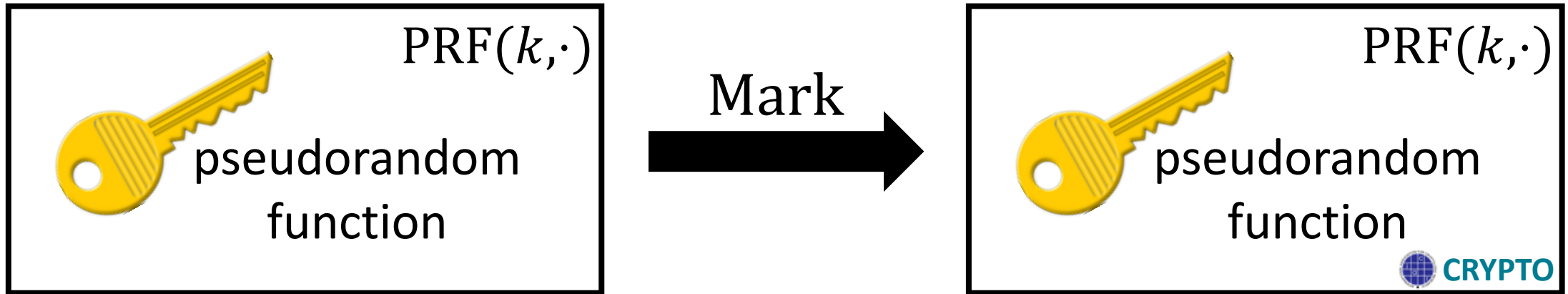
```
static void AES_enc_blk(block *blk, const AES_KEY *key) {  
    unsigned j, rnds = ROUNDS(key);  
    const __m128i *sched = ((__m128i *) (key->rd_key));  
    *blk = _mm_xor_si128(*blk, sched[0]);  
    for (j = 1; j < rnds; ++j) {  
        *blk = _mm_aesenc_si128(*blk, sched[j]);  
    }  
    *blk = _mm_aesenc_si128(*blk, sched[j]);  
}
```

Learning the original  
(unmarked) function gives a  
way to remove the watermark

- Notion only achievable for functions that are not learnable
- Focus has been on cryptographic functions

# Watermarking Cryptographic Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

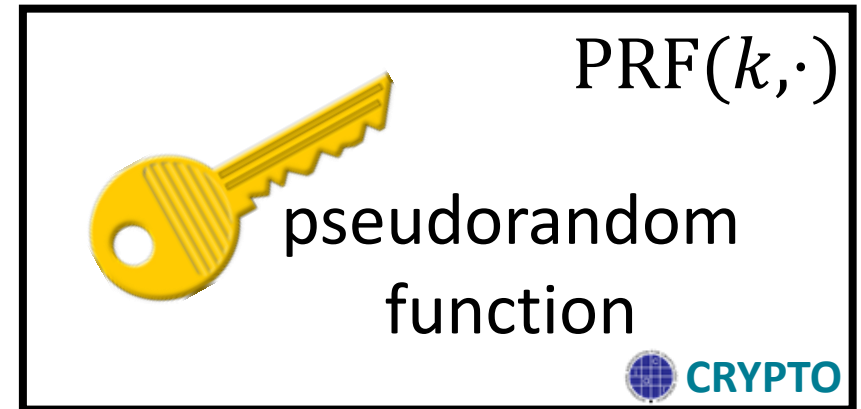
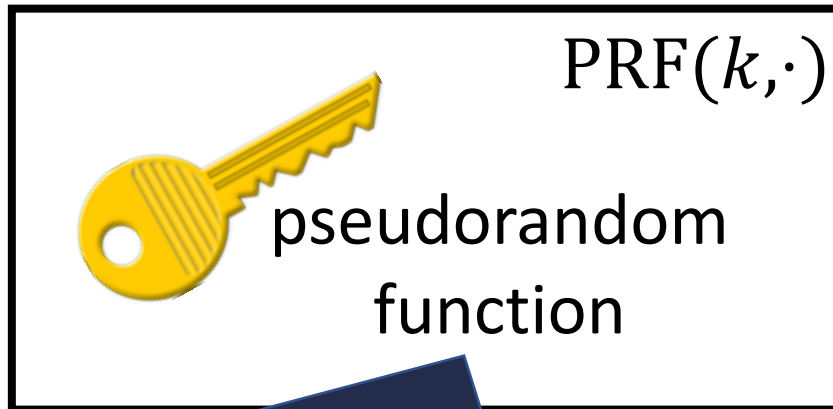


- Focus of this work: watermarking PRFs [CHNVW16, BLW17, KW17, QWZ18]



# Watermarking Cryptographic Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]

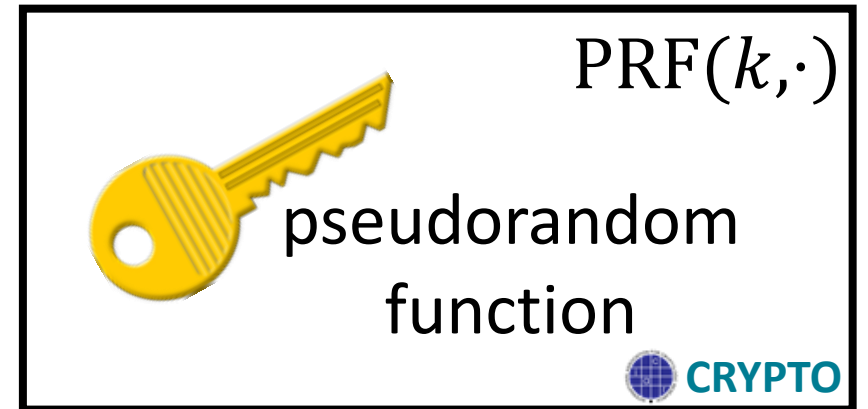
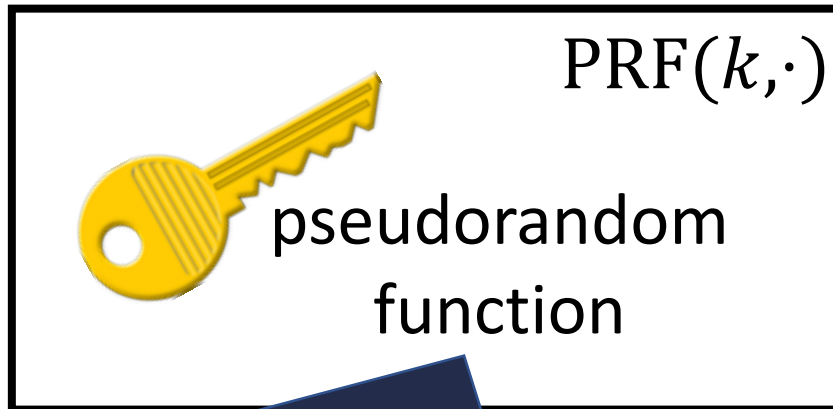


A function whose input-output behavior is unpredictable (looks like a random function) – e.g., AES

Watermarking PRFs [CHNVW16, BLW17, KW17, QWZ18]

# Watermarking Cryptographic Programs

[NSS99, BGIRSVY01, HMW07, YF11, Nis13, CHNVW16, BLW17, KW17, QWZ18, GKMW19]



Program has PRF key  $k$  hard-wired inside it and on input  $x$ , outputs  $\text{PRF}(k, x)$

Watermarking PRFs [CHNVW16, BLW17, KW17, QWZ18]

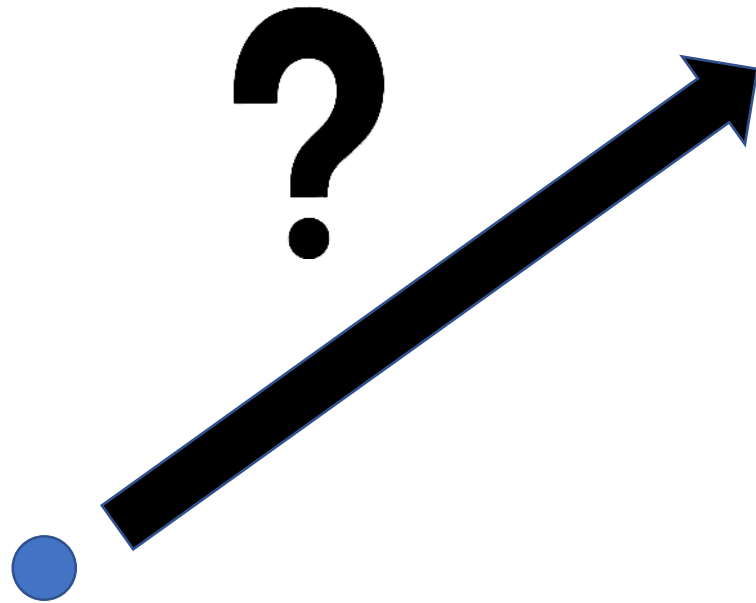
# Watermarkable PRFs

[CHNVW16]: Watermark PRFs from **iO** + **OWFs**  
**Publicly** verifiable

Can we watermark PRFs from standard assumptions?

[KW17]: Watermark PRFs from standard assumptions (LWE)  
**Secretly** verifiable

# Watermarkable PRFs



**Secretly Verifiable**  
Watermarking  
from LWE [KW17]



**Publicly Verifiable**  
Watermarking  
[CHNVW16]

# A Naïve Attempt at Public Verifiability

Just make the  
verification key public!

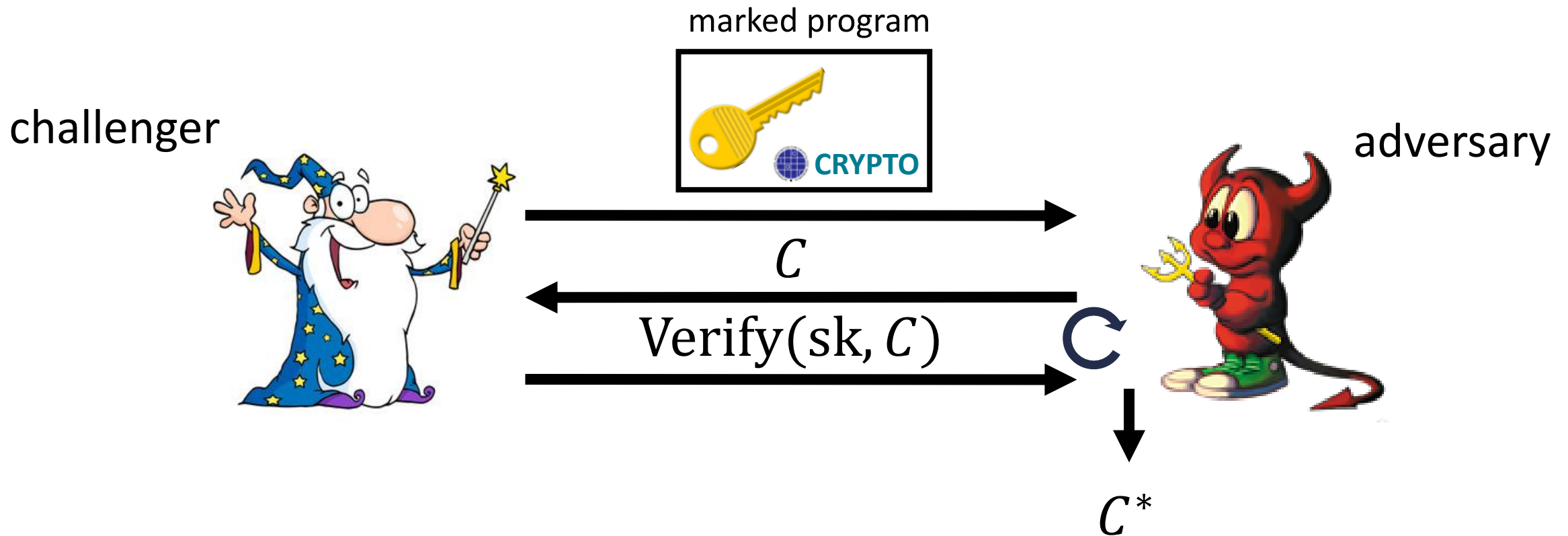
**Problem:** Knowledge of the verification key allows adversary to trivially remove watermark

**In fact:** Even oracle access to the verification key is sufficient to break unremovability (“verifier rejection” problem)



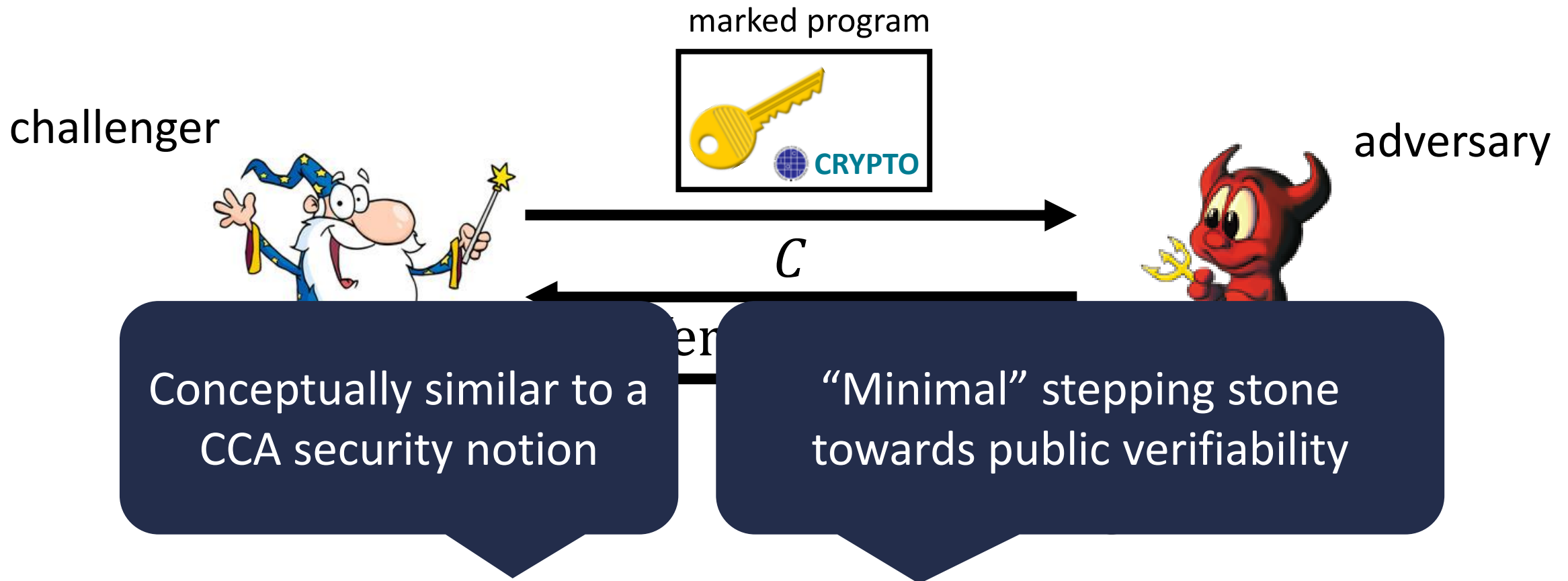
**Secretly** Verifiable  
Watermarking  
from LWE [KW17]

# Between Public and Secret Verification



**Intermediate notion:** Secret verification, but security in the presence of a verification oracle

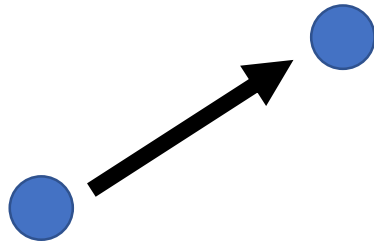
# Between Public and Secret Verification



**Intermediate notion:** Secret verification, but security in the presence of a verification oracle

# Watermarkable PRFs

**Secretly Verifiable**  
Watermarking  
from CCA [QWZ18]



**Secretly Verifiable**  
Watermarking  
from LWE [KW17]



**Publicly Verifiable**  
Watermarking  
[CHNVW16]

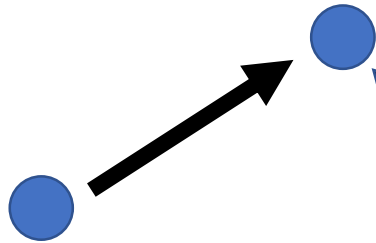


# Watermarkable PRFs

**Secretly Verifiable**  
Watermarking  
from CCA [QWZ18]



**Secretly Verifiable**  
Watermarking  
from LWE [KW17]

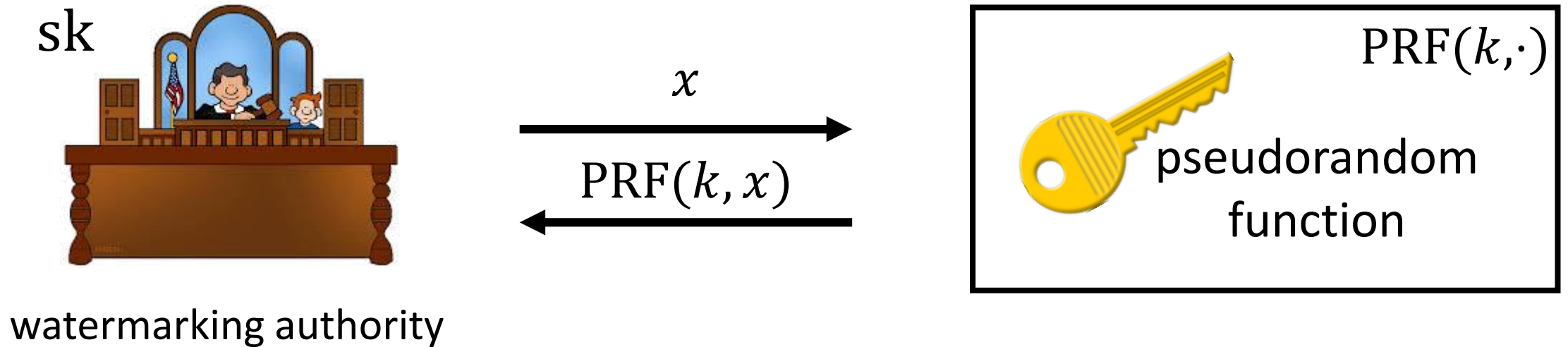


**Good:** Achieves security in the presence of a verification oracle

**Limitation:** Knowledge of the verification key breaks PRF security (even *unmarked* keys)

# Security Against the Authority

[QWZ18]



After seeing single query (on any  $x$ ), authority can distinguish output of PRF from output of random function

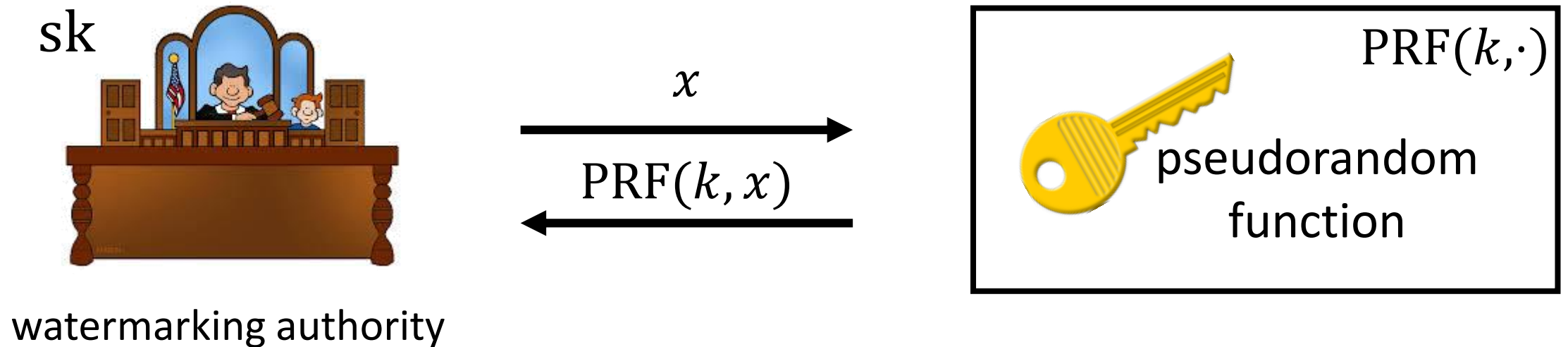
# Security Against the Authority

[QWZ18]

**Implication:** Knowledge of the verification key completely breaks PRF security  
(notion still seems far publicly-verifiable setting)

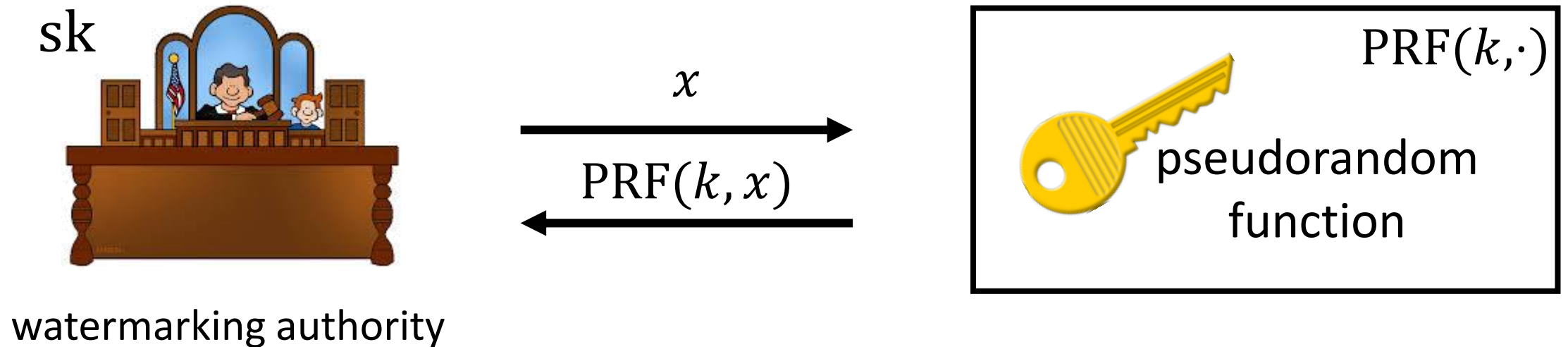
After seeing single query (on any  $x$ ), authority can distinguish output of PRF from output of random function

# Don't We Have to Trust the Authority Anyways?



**Not necessarily:** marking algorithm can be implemented using a two-party computation, so authority never needs to see *any* PRF keys in the clear

# Don't We Have to Trust the Authority Anyways?



**This work:** New watermarkable PRF that provides security even against the watermarking authority

# Our Results

New **secretly verifiable** watermarking for PRF from LWE

- Unremovability holds in the presence of the verification oracle
- **weak pseudorandomness** *even against authority*  
( $T$ -restricted pseudorandomness)
- As secure as any other PRF family from LWE
  - Relies on worst-case lattice problems with **nearly-polynomial** ( $n^{\omega(1)}$ ) approximation factors

# Our Results

New **secretly verifiable** watermarking

- Unremovability holds in the standard model
- **weak pseudorandomness**

Previous constructions (with message-embedding) required private constrained PRFs (which requires quasi-polynomial or sub-exponential approximation factors)

- As secure as any other PRF family from LWE
  - Relies on worst-case lattice problems with **nearly-polynomial** ( $n^{\omega(1)}$ ) approximation factors
- New abstraction: **extractable PRF**

# Starting Point: Puncturable PRF

[BW13, BGI14, KPTZ13]



Punctured key  $k_{x^*}$  can be used to evaluate PRF on all points  $x \neq x^*$   
(value at  $x^*$  is pseudorandom even given  $k_{x^*}$ )

**Private puncturing:** punctured key  $k_{x^*}$  also hides  $x^*$

**Programmability:** program  $F(k_{x^*}, x^*) := y^*$



# From Puncturing to Watermarking

[BLW17, KW17]



## Marking algorithm:

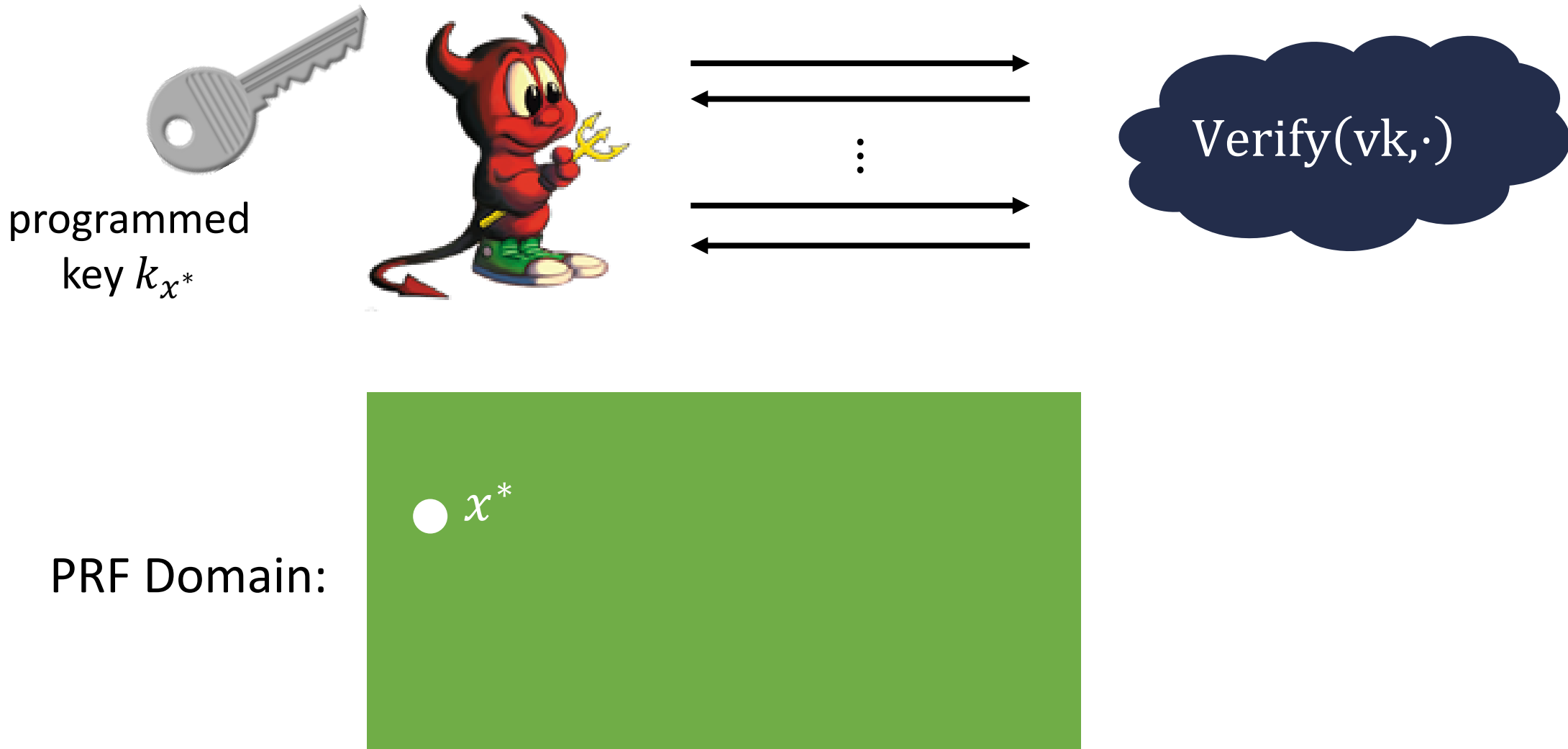
1. Derive a special point  $(x^*, y^*)$  from input/output behavior of PRF
2. Define a marked circuit to be  $F(k_{x^*}, \cdot)$

## Verification algorithm:

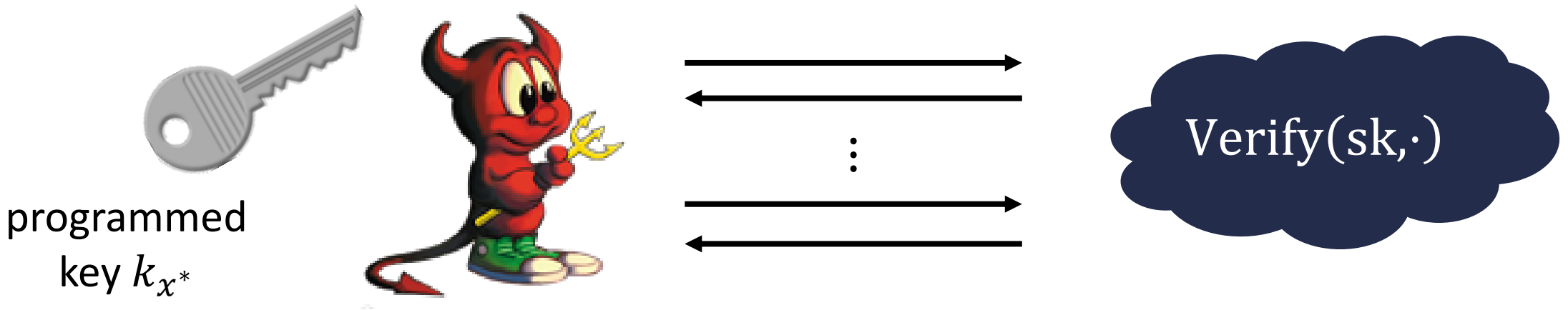
1. Test if  $C(x^*) = y^*$

**Security:** Punctured point  $x^*$  is hidden

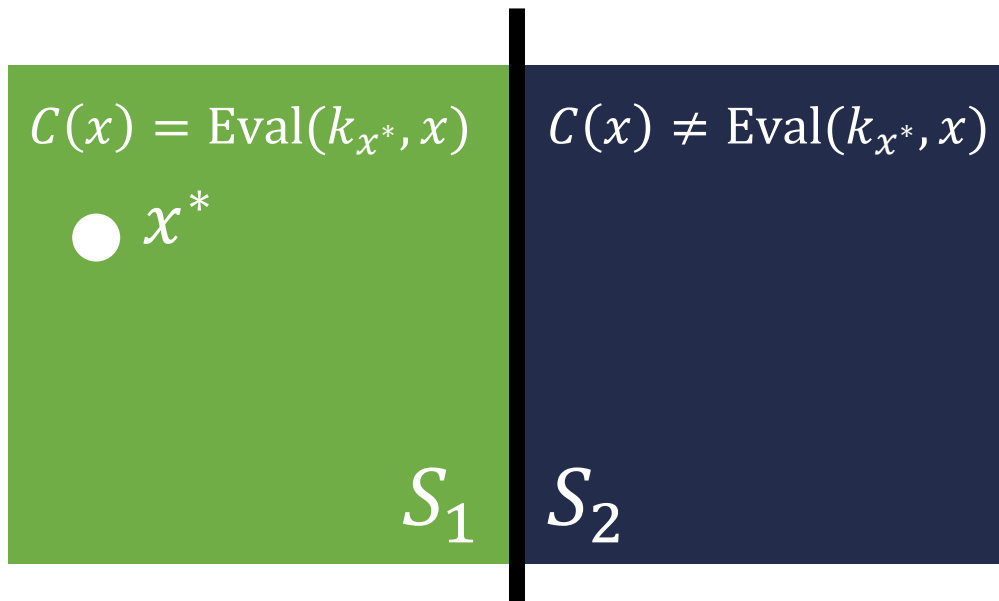
# Intuition: Binary Search Attack



# Intuition: Binary Search Attack



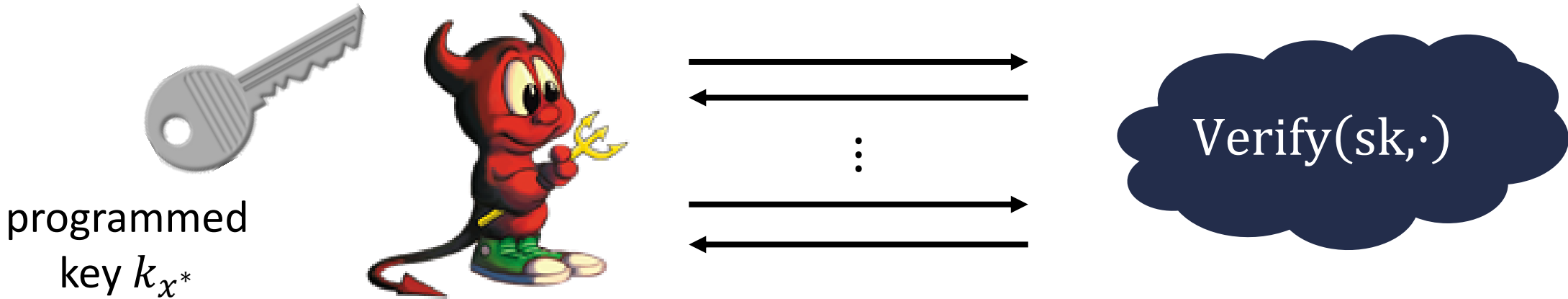
PRF Domain:



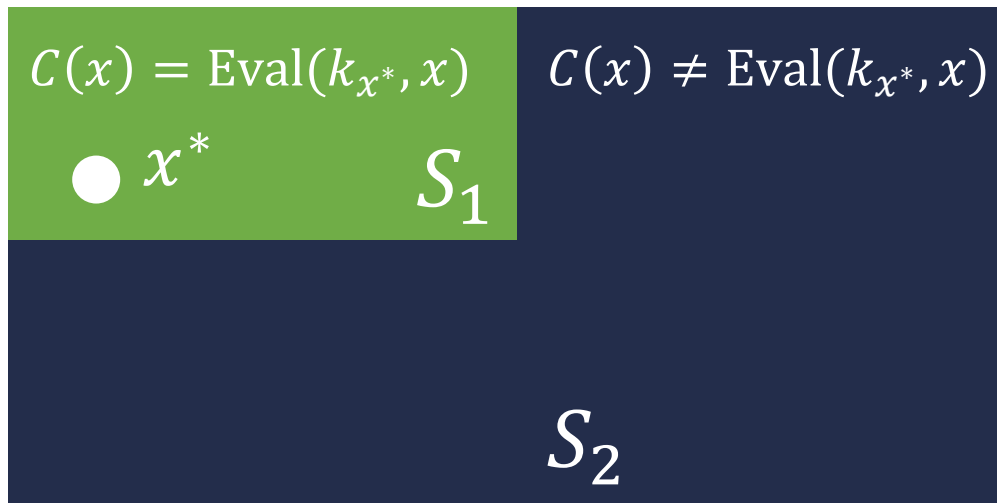
**Intuitively:**

- if  $\text{Verify}(\text{sk}, C) = 1$ , then  $x^* \notin S_2$
- if  $\text{Verify}(\text{sk}, C) = 0$ , then  $x^* \in S_2$

# Intuition: Binary Search Attack



PRF Domain:



**Intuitively:**

if  $\text{Verify}(\text{sk}, C) = 1$ , then  $x^* \notin S_2$

if  $\text{Verify}(\text{sk}, C) = 0$ , then  $x^* \in S_2$

Eventually, adversary recovers special point  $x^*$

# Intuition: Binary Search Attack



programmed  
key  $k_{x^*}$



Very similar to a “verifier rejection” attack encountered in settings like designated-verifier proof systems, CCA-security, etc.

**Solution:** Make the set of “valid” circuits detectable (i.e., cannot change too many points and still preserve mark)

PRF Domain:

$$C(x) = \text{Eval}(k_{x^*}, x)$$

●  $x^*$

$S_1$

$$C(x) \neq \text{Eval}(k_{x^*}, x)$$

$S_2$

**Intuitively:**

if  $\text{Verify}(\text{sk}, C) = 1$ , then  $x^* \notin S_2$

if  $\text{Verify}(\text{sk}, C) = 0$ , then  $x^* \in S_2$

Eventually, adversary recovers special point  $x^*$

# Our Notion: Extractable PRF



Punctured key  $k_{x^*}$  can be used to evaluate PRF on all points  $x \neq x^*$

**Private puncturing:** punctured key  $k_{x^*}$  also hides  $x^*$

**Programmability:** program  $F(k_{x^*}, x^*) := y^*$

**Extractability:** point  $F(k_{x^*}, z) := \text{Encode}(k)$  encode original PRF key  $k$

# Our Notion: Extractable PRF



Punctured key  $k_{x^*}$  can be used to evaluate PRF on all points  $x \neq x^*$

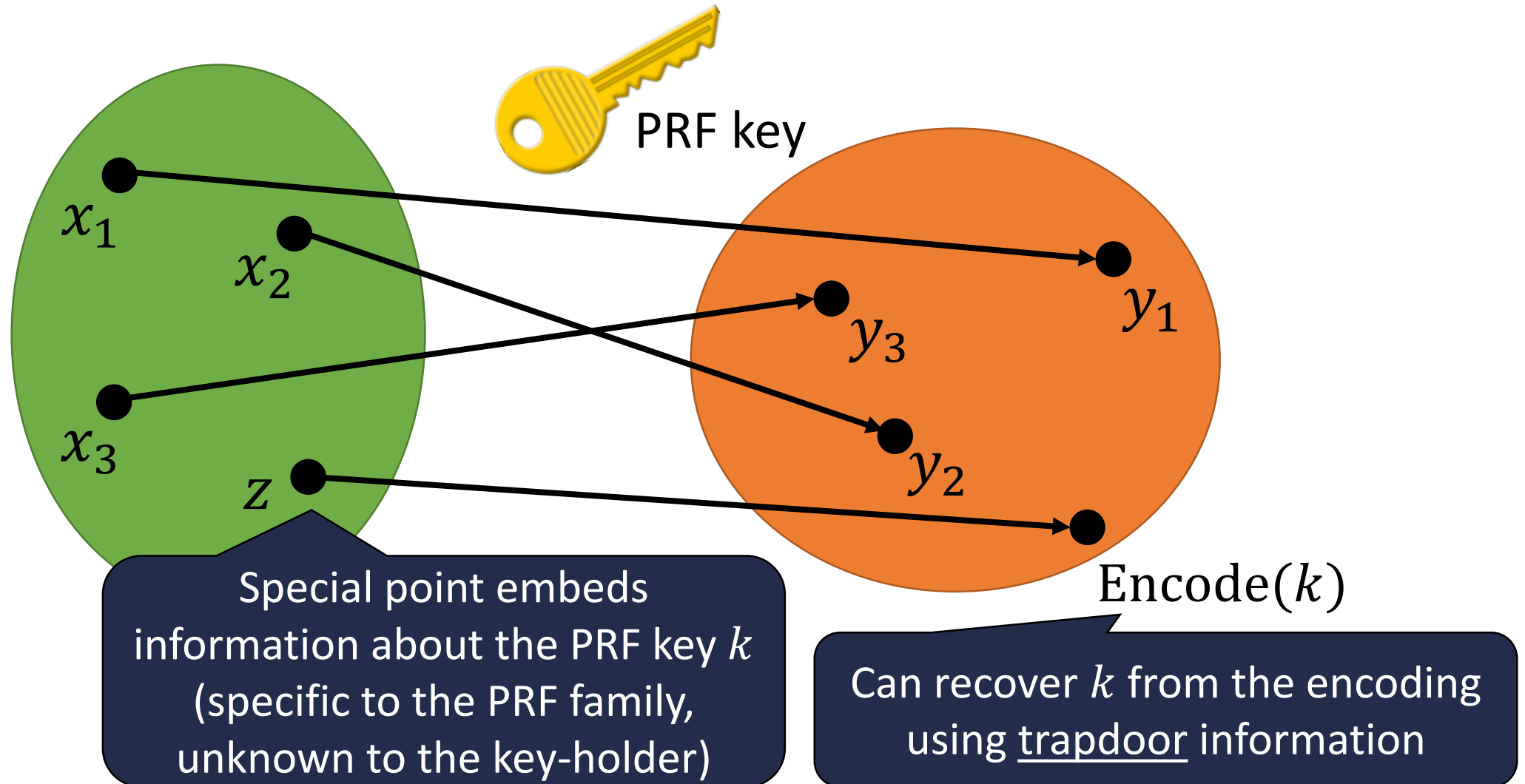
**Private puncturing:** puncture

**Programmability:** program  $F$

**Extractability:** point  $F(k_{x^*}, z) := \text{Encode}(k)$  encode original PRF key  $k$

Decode with trapdoor  $td$   
(part of watermarking  $sk$ )

# Our Notion: Extractable PRF





# Extraction to Watermarking

## Marking algorithm:

1. Derive a special point  $(x^*, y^*)$  from input/output behavior of PRF
2. Define a marked circuit to be  $F(k_{x^*}, \cdot)$

## Verification algorithm:

1. Test if  $C(x^*) = y^*$
2. Extract key  $k$  and test if  $C(\cdot) \approx F(k, \cdot)$   
(output unmarked if key extraction fails)
3. Accept only if both conditions satisfied

Adversary can rule out only a small fraction of domain

**In fact:** extractability enables a simpler marking procedure

# Extraction to Watermarking

## Marking algorithm:

1. Derive a special point  $(x^*, y^*)$  from input/output behavior of PRF
2. Define a marked circuit to be  $F(k_{x^*}, \cdot)$

## Verification algorithm:

1. Test if  $C(x^*) = y^*$

2. Extract key  $k$  : Instead of programming the value at  $x^*$ , puncture the PRF at  $x^*$ : circuit is marked if  $C(x^*) \neq F(k, x^*)$
3. Accept only if  $C(x^*) = y^*$  where  $k$  is the extracted key

**In fact:** extractability enables a simpler marking procedure

# Extraction to Watermarking

## Marking algorithm:

1. Puncture key at  $x^*$  to obtain a key  $k_{x^*}$
2. Define a marked key to be  $F(k_{x^*}, \cdot)$

## Verification algorithm:

1. Extract key  $k$  and test if  $C(\cdot) \approx F(k, \cdot)$   
(output unmarked if key extraction fails)
2. Output marked if  $C(x^*) \neq F(k, x^*)$  and unmarked otherwise

To remove watermark, need to fix the value of the PRF at the punctured point (i.e., guess a pseudorandom value)

**In fact:** extractability enables a simpler marking procedure

# Extraction to Watermarking

## Marking algorithm:

1. Puncture key at  $x^*$  to obtain a key  $k_{x^*}$
2. Define a marked key to be  $F(k_{x^*}, \cdot)$

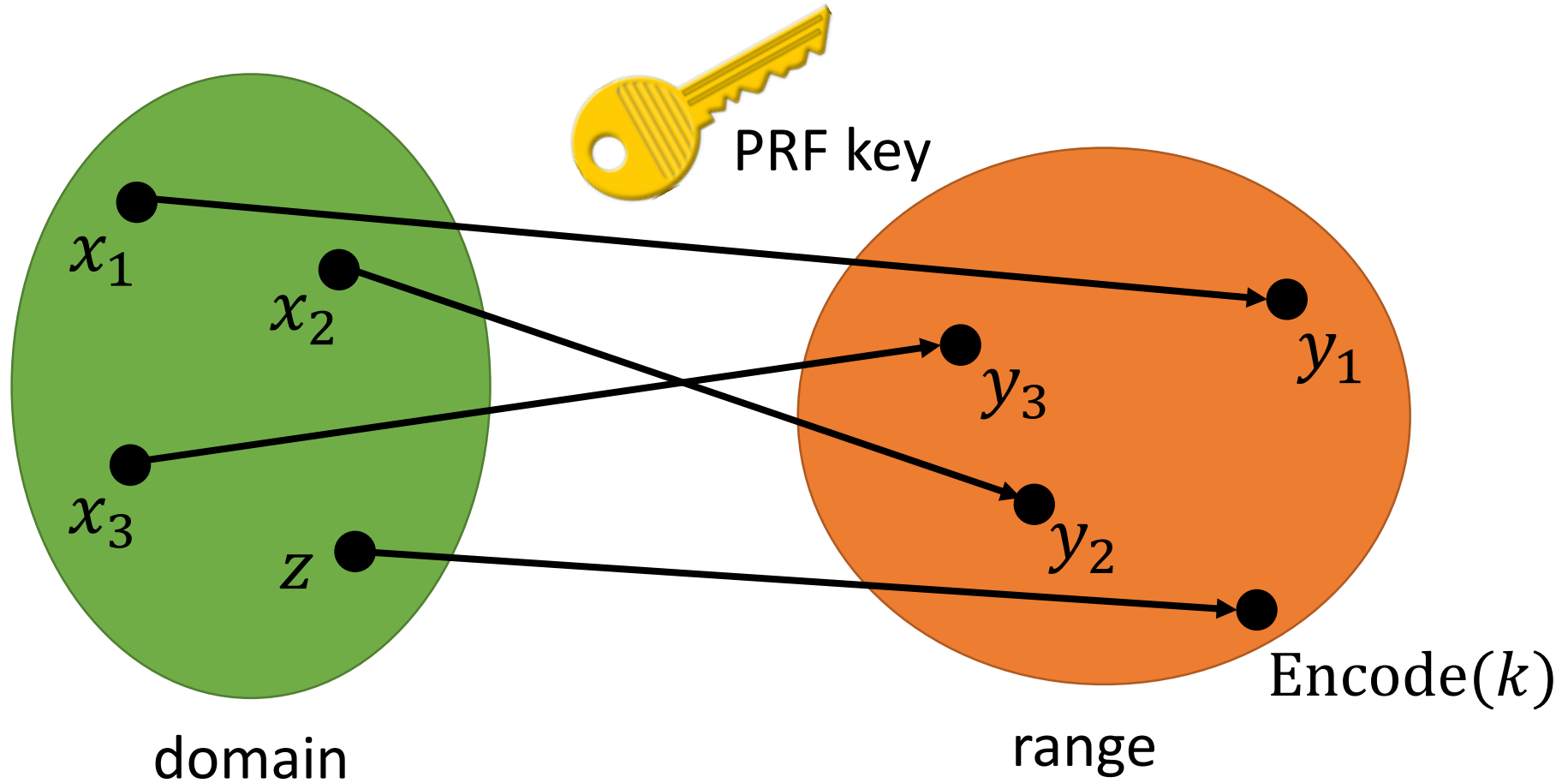
To remove watermark, need to fix the value of the PRF at the punctured point (i.e., guess a pseudorandom value)

## Verification algorithm:

1. Extract key  $k$  and test if  $C(\cdot) \approx F(k, \cdot)$   
(output unmarked if key extraction fails)
2. Output marked if  $C(x^*) \neq F(k, x^*)$  and unmarked otherwise

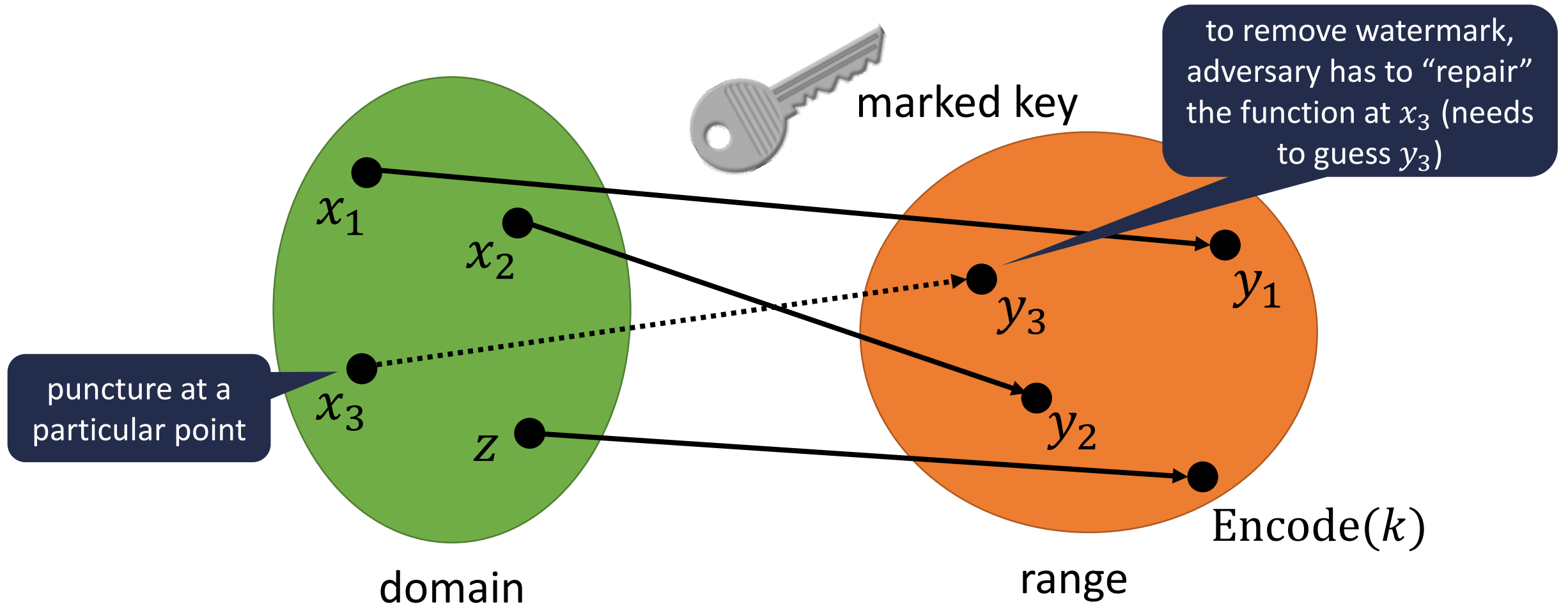
**Advantage:** no longer require private puncturing  
(can base on weaker assumptions)

# Extraction to Watermarking



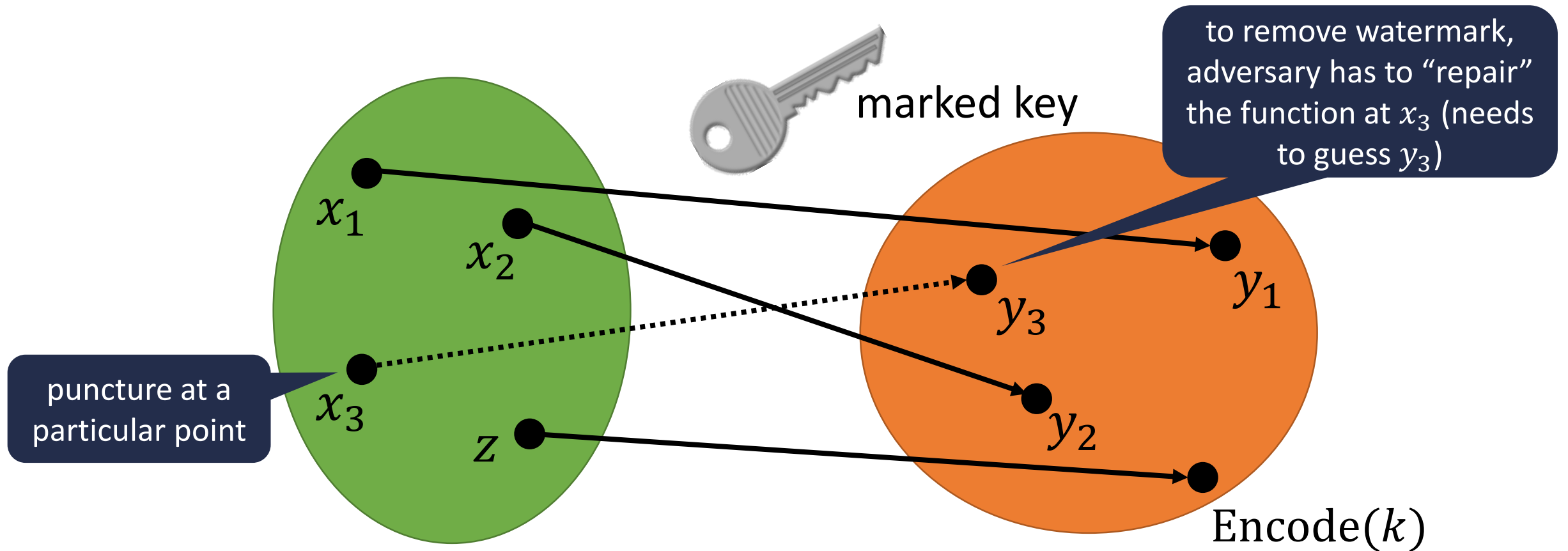
Real PRF key

# Extraction to Watermarking



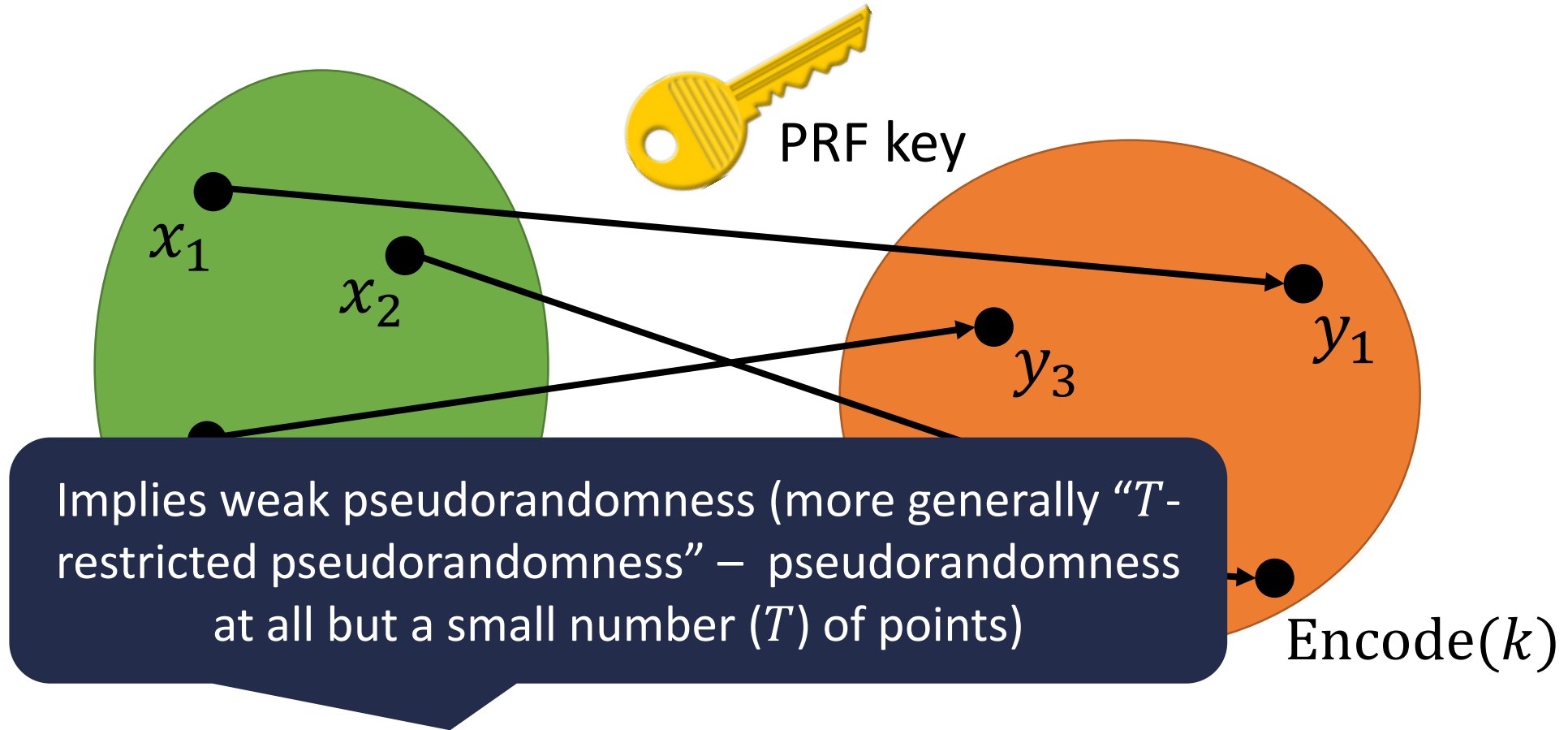
Marked Key

# Extraction to Watermarking



**Preventing verifier rejection:** Queries on circuits that are far away from marked key will always reject, so binary search is no longer effective

# Security Against the Authority



PRF keys are pseudorandom everywhere except at  $z$  (even given the extraction trapdoor)



# Summary

Puncturable  
Extractable PRF



Watermarkable PRF

## High-level overview:

- **Marking:** Puncture PRF key at  $x^*$
- **Verification:** Extract key from circuit, and check correctness of value at  $x^*$

**Unremovability:** Key-extraction succeeds if circuit if adversary's circuit is close to original PRF; removing the mark requires "patching" PRF at punctured point

# Summary

Puncturable  
Extractable PRF



Watermarkable PRF

Property holds even in the presence of the verification oracle

Key at  $x^*$

from circuit, and check correctness of

**Unremovability:** Key-extraction succeeds if circuit if adversary's circuit is close to original PRF; removing the mark requires "patching" PRF at punctured point

# Constructing Extractable PRFs

Structure of lattice PRFs [BV15]:

PRF on  $\ell$ -bit inputs (e.g., domain  $\{0,1\}^\ell$ )

$\mathbf{A}_1, \dots, \mathbf{A}_\ell \in \mathbb{Z}_q^{n \times m}$  public matrices (one for each bit of input)

PRF secret key:  $\mathbf{s} \in \mathbb{Z}_q^n$  (LWE secret)

$$(\mathbf{A}, \mathbf{s}^T \mathbf{A} + \mathbf{e}^T) \approx_c (\mathbf{A}, \mathbf{u}) \text{ where}$$
$$\mathbf{A} \leftarrow \mathbb{Z}_q^{n \times m}, \mathbf{s} \leftarrow \mathbb{Z}_q^n, \mathbf{e} \leftarrow \chi^m, \mathbf{u} \leftarrow \mathbb{Z}_q^m$$

# Constructing Extractable PRFs

Structure of lattice PRFs [BV15]:

PRF on  $\ell$ -bit inputs (e.g., domain  $\{0,1\}^\ell$ )

$\mathbf{A}_1, \dots, \mathbf{A}_\ell \in \mathbb{Z}_q^{n \times m}$  public matrices (one for each bit of input)

PRF secret key:  $\mathbf{s} \in \mathbb{Z}_q^n$  (LWE secret)

PRF evaluation at input  $x$ :  $\text{PRF}(\mathbf{s}, x) := \lfloor \mathbf{s}^T \mathbf{A}_x \rfloor_p$

$\mathbf{A}_x$ : matrix derived from  $\mathbf{A}_1, \dots, \mathbf{A}_\ell, x$

# Constructing Extractable PRFs

**Goal:** embed a trapdoor at  $z$  such that evaluation at  $z$  allows key recovery

Lattice trapdoors [Ajt99, GPV08, AP09, MP12]: can sample

random matrix  $\mathbf{D} \in \mathbb{Z}_q^{n \times m}$

trapdoor  $\text{td}_{\mathbf{D}}$

such that LWE is easy with respect to  $\mathbf{D}$ :

given  $\mathbf{s}^T \mathbf{D} + \mathbf{e}^T$  and  $\text{td}_{\mathbf{D}}$ , can recover LWE secret  $\mathbf{s}$

**Idea:** hide a lattice trapdoor in the public parameters

# Constructing Extractable PRFs

$\mathbf{A}_1, \dots, \mathbf{A}_\ell \in \mathbb{Z}_q^{n \times m}$  public matrices (one for each bit of input)

PRF secret key:  $\mathbf{s} \in \mathbb{Z}_q^n$  (LWE secret)

PRF evaluation at input  $x$ :  $\text{PRF}(\mathbf{s}, x) := \lfloor \mathbf{s}^T \mathbf{A}_x \rfloor_p$

Embed trapdoor at  $z \in \{0,1\}^\ell$ :

Compute  $\mathbf{A}_z$  from  $\mathbf{A}_1, \dots, \mathbf{A}_\ell$

Let  $\mathbf{W} = \mathbf{D} - \mathbf{A}_z$

Include  $\mathbf{W}$  in the public parameters

# Constructing Extractable PRFs

$\mathbf{A}_1, \dots, \mathbf{A}_\ell \in \mathbb{Z}_q^{n \times m}$  public matrices (one for each bit of input)

PRF secret key:  $\mathbf{s} \in \mathbb{Z}_q^n$  (LWE secret)

PRF evaluation at input  $x$ :  $\text{PRF}(\mathbf{s}, x) := \lfloor \mathbf{s}^T \mathbf{A}_x \rfloor_p$

$\mathbf{W}$  hides  $\mathbf{A}_z$  (and thus,  $z$ ) since  $\mathbf{D}$  is statistically close to uniform

Include  $\mathbf{W}$  in the public parameters

PRF evaluation at input  $x$ :  $\text{PRF}(\mathbf{s}, x) := \lfloor \mathbf{s}^T (\mathbf{A}_x + \mathbf{W}) \rfloor_p$

# Constructing Extractable PRFs

$A_1, \dots, A_\ell \in \mathbb{Z}_q^{n \times m}$  public matrices (one for each bit of input)

PRF secret key:  $s \in \mathbb{Z}_q^n$  (LWE secret)

PRF evaluation at input  $x$ :  $\text{PRF}(s, x) := \lfloor s^T A_x \rfloor_p$

Embed trapdoor at  $z \in \{0,1\}^\ell$ :

Co

Le

In

Value everywhere else is still  
pseudorandom

Value at  $z$  is  
 $\lfloor s^T (A_z + W) \rfloor_p = \lfloor s^T D \rfloor_p$ ,  
so can extract  $s$  using trapdoor  $\text{td}_D$

PRF evaluation at input  $x$ :  $\text{PRF}(s, x) := \lfloor s^T (A_x + W) \rfloor_p$



# Summary

Puncturable  
Extractable PRF



Watermarkable PRF

Puncturable extractable PRF can be built from LWE  
(with a nearly polynomial modulus-to-noise ratio)

Yields new watermarking scheme from LWE with  
security in the presence of verification oracle

**Extensions:** Message-embedding, mark-unforgeability [ See paper... ]

# Open Problems

Extractable PRFs from generic techniques?

More applications of extractable PRFs?

Publicly-verifiable watermarking scheme for PRFs?

**Thank you!**

<http://eprint.iacr.org/2018/986>