

8) Steganography

Christian Cachin

Privacy and Data Security, 2021

(originally presented at CMS 2005)

Steganography ≠ Cryptography

Cryptography hides the content of communication

Adversary knows that communication exists

Conceals data



Steganography hides the existence of communication

Adversary should not discover existence of communication

Conceals also metadata



Information-hiding concepts

Embed information in a "cover" (carrier signal)

Steganography

Hide presence of hidden information

To avoid censorship and surveillance

Presence not known; can easily be removed

Watermarking

Hide information in a robust way

To authenticate (multimedia) data, control ownership

Presence is known; should be difficult to remove

Fingerprinting

Like watermarking, but hides user-specific
information for tracing and litigation

Prisoner's problem



Alice and Bob want to coordinate their escape
Communication through passive **observer** (**war-**
den, censor ...)

innocent communication is allowed
talking about escape plans is forbidden

Formulated by Simmons (1983)

Steganography = hidden communication

How to be convinced that it is "really" hidden?

Many attacks

Correlation, histogram, transforms ...

Analogy to cryptology

Ages of broken cryptosystems until ~1980

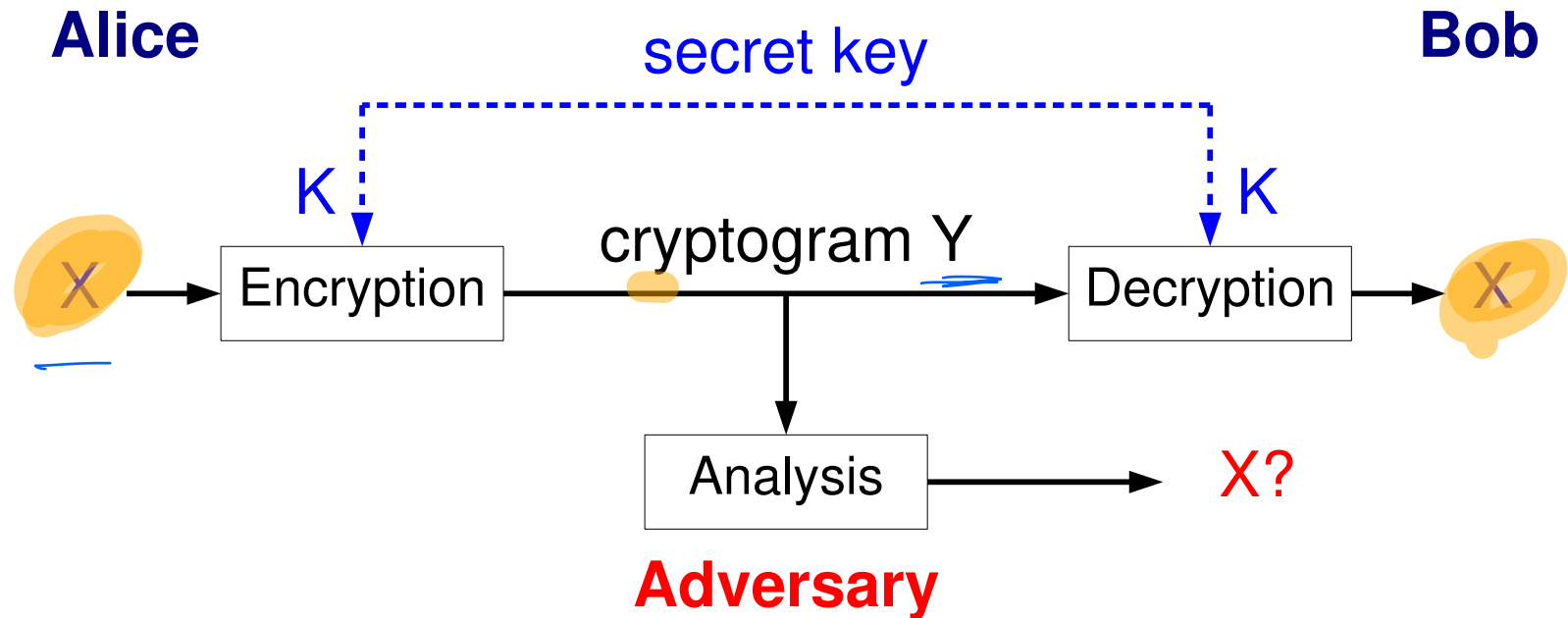
Theory of cryptology with provably secure
cryptosystems since ~1990

→ Formal model for steganography

Models for cryptosystems



Model of a cryptosystem



Shannon (1949)

Adversary is *passive*

Security of cryptosystems

Perfect security (Shannon 1949)

$$\underline{I(X;Y)} = H(X) - H(X|Y) = 0$$

Information-theoretic

Unbounded adversary obtains no information

Implies also that $\underline{H(K)} \geq H(X)$ ☹️

Ex. one-time pad

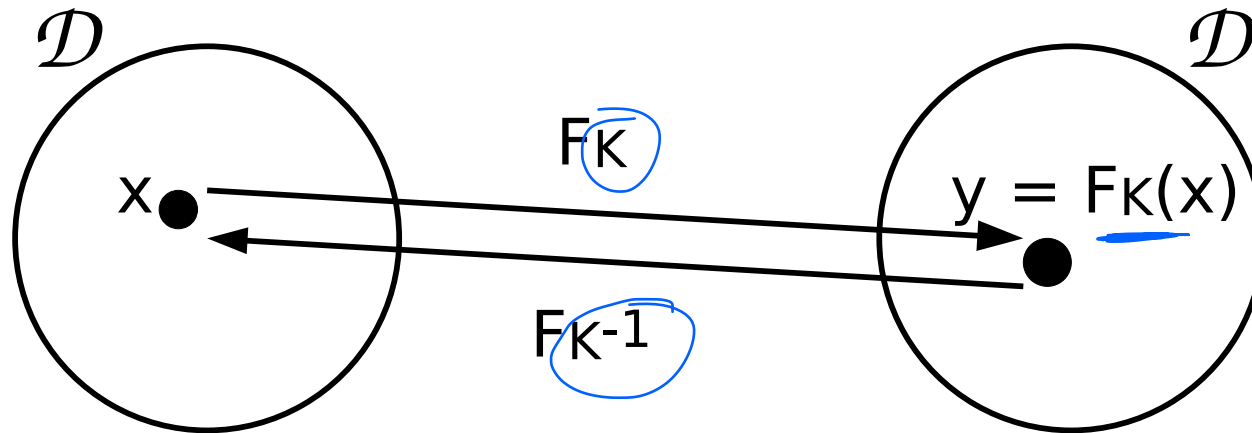
But ... practical cryptography uses block- and stream-ciphers with short keys.

Practical cryptosystems

Family of one-way permutations

$$\mathcal{F} = \{F_K\}, F_K: \mathcal{D} \rightarrow \mathcal{D}$$

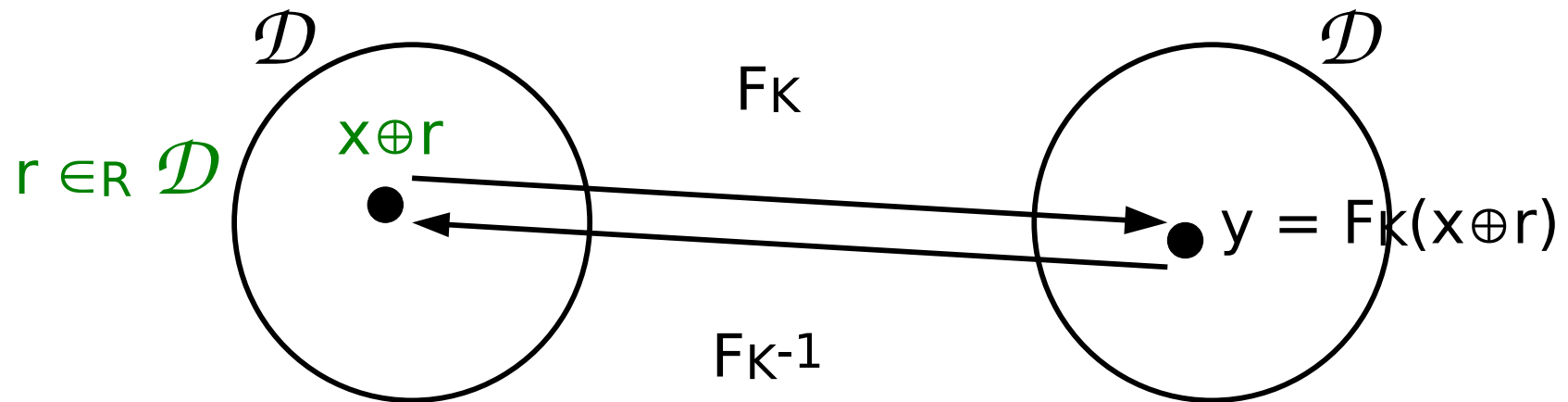
indexed by key K



Encryption and decryption are deterministic
leaks information

Probabilistic encryption

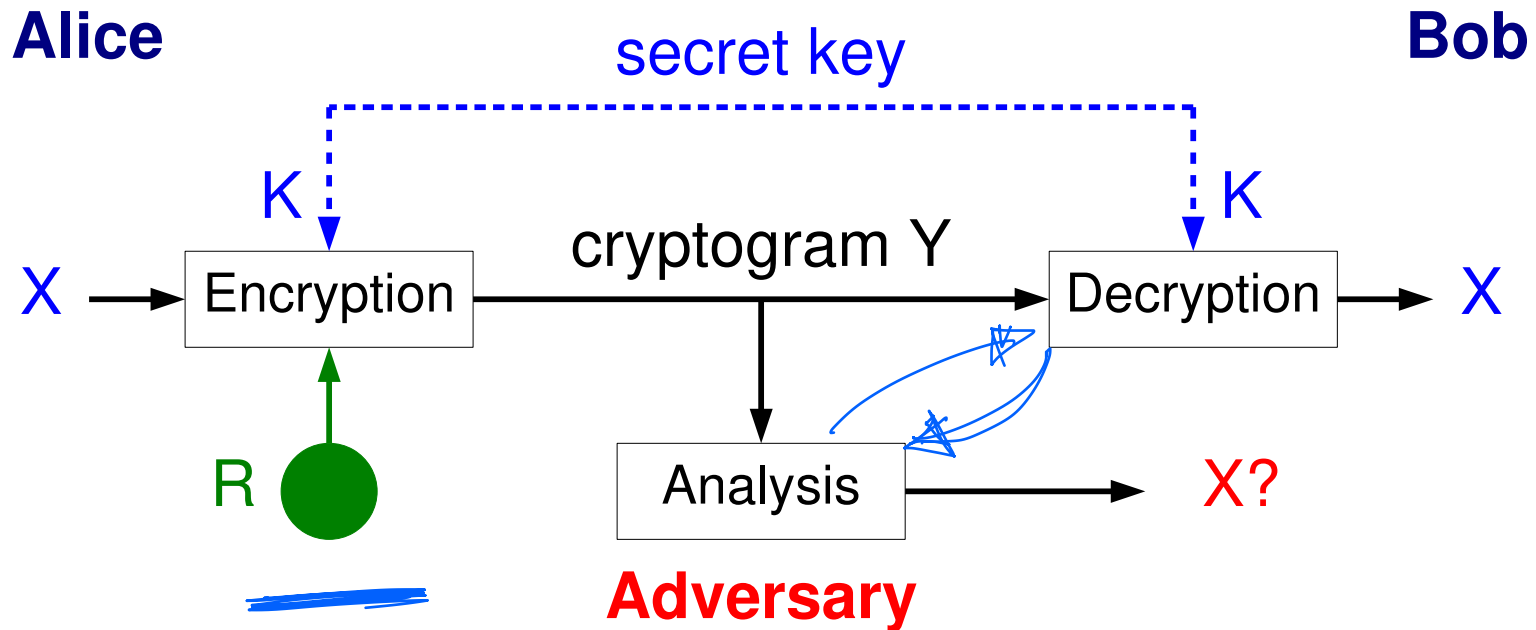
Randomize input to encryption function



Ciphertext is (r, y)

does not repeat for same plaintext x
better protection

Probabilistic cryptosystems



Private random source R

Computational security for cryptosystems

Formal security model (Goldwasser-Micali, 1985)

Chosen-plaintext attacks (passive adversary)

Semantic security \Leftrightarrow indistinguishability of ciphertexts

Defined by experiment with adversary A , a probabilistic polynomial-time (PPT) algorithm:

```
K      ← KG(1k)
(x0,x1) ← A1(K)
b      ∈R {0,1}
c*     ← E(K,xb)
b*     ← A2(x0,x1,c*)
```

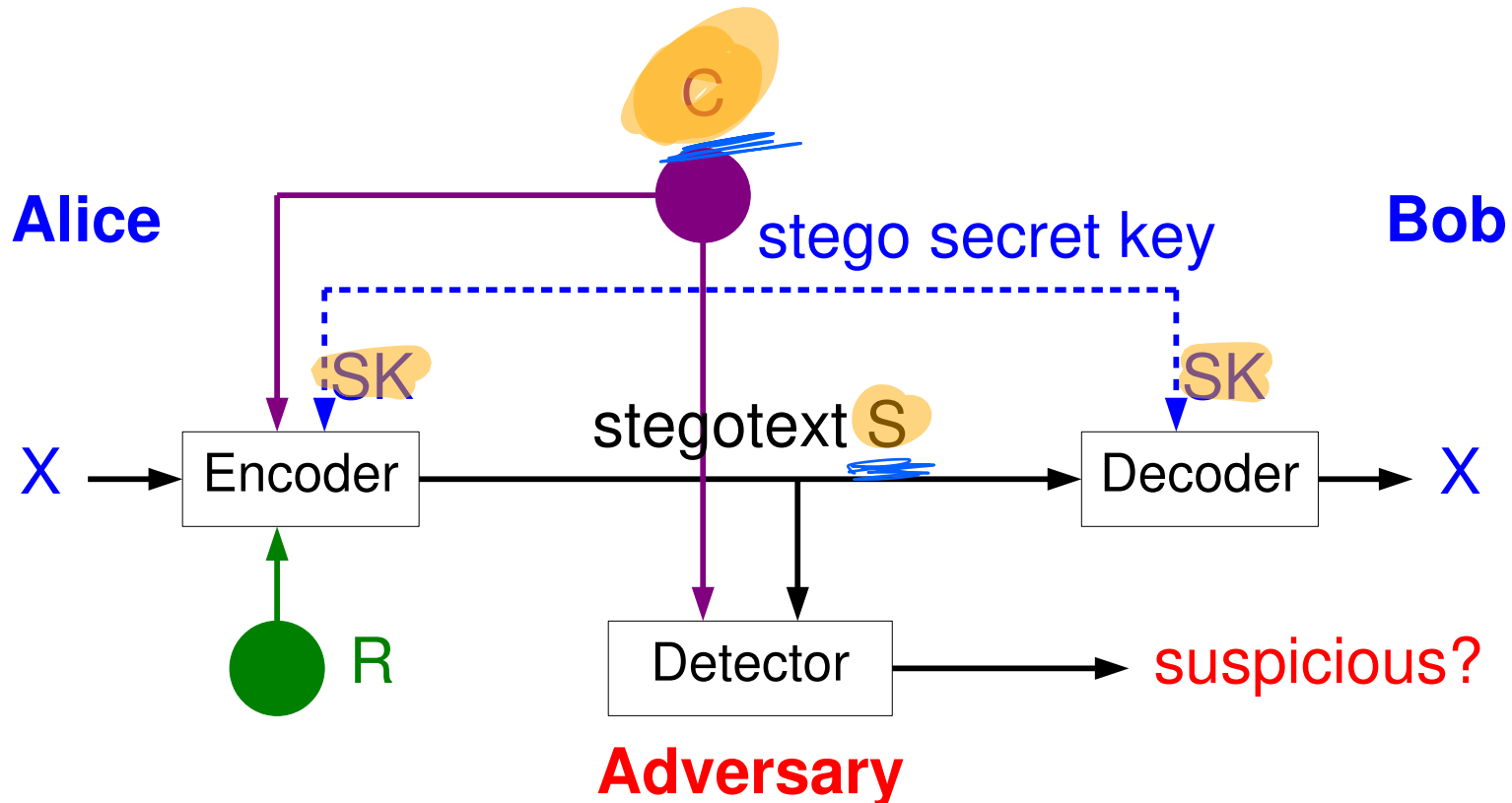
Then

\forall PPT $A=(A_1,A_2) : \Pr[b = b^*] \leq \frac{1}{2} + \text{negl.}$

Models for stegosystems



Model of a stegosystem



Covert text C and stegotext S over domain Y
Adversary is *passive*

How to define security for stegosystems?

Attempt 1

Based on information theory and mutual entropy:

$$I(X;S|C) = 0$$

Fails to differentiate stegosystem from cryptosystem.

Attempt 2

Based on distortion measure $d(C,S) = \|P_C - P_S\|_2$:

$$d(c,s) \rightarrow 0$$

Fails for some examples, e.g., when C is a random n -bit string, but S is a random n -bit string with even parity.

Information-theoretic security

Steganography as hypothesis testing (C. 1998)

Adversary distinguishes stegotext from covertext

Quantified using relative entropy

Perfect security

$$D(P_C \| P_S) = 0$$

Statistical security

$$D(P_C \| P_S) \leq \varepsilon$$

Discrimination or relative entropy

$$D(P_C \| P_S) = \sum_{y \in Y} P_C(y) \log_2(P_C(y)/P_S(y)) \geq 0$$

\mathcal{A}^C \mathcal{A}^S

Bounds on detection

Adversary distinguishes stegotext S (H_1) from
coverttext C (H_0)

Deciding S with signal C is type-I error, prob. α

Deciding C with signal S is type-II error, prob. β

Statistical security $D(P_C || P_S) \leq \epsilon$,

$$d(\alpha, \beta) \leq \epsilon$$

binary discrimination $d(\cdot, \cdot)$

For $\alpha=0$, this implies

$$\beta \geq 2^{-\epsilon}$$

Cover & Thomas
textbook on
information theory

Modeling the covertext

Probabilistic source or channel

$\mathbf{C} = C_0, C_1, C_2 \dots$

Distribution known

\mathbf{C} is a stochastic process

Random variable (i.i.d. \rightarrow ergodic $\rightarrow \dots$)

Unrealistic in practice

Real-world communication does not come with
specification of distribution

Universal stegosystems

What if distribution of **C** is not known?

→ Universal stegosystems

→ No knowledge of cover distribution needed

C is an algorithm, given as oracle

Can be queried on arbitrary history

$$\mathbf{C}(h,n) = \mathbf{C}_{|h|+1}, \mathbf{C}_{|h|+2}, \dots, \mathbf{C}_{|h|+n}$$

Synthetic cover signal **C** from machine learning

Example stegosystem

Known distribution P_C over domain Y

Message $x \in \{0,1\}$

Let

$$\begin{cases} Y_0 = \min_{Y' \subseteq Y} |\Pr[C \in Y'] - \Pr[C \notin Y']| & (\text{prob.} =: \underline{\underline{\varepsilon}}) \\ Y_1 = Y \setminus Y_0 \end{cases}$$

Key $sk \in_R \{0,1\}$

RV C_0 is C restricted to Y_0

RV C_1 is C restricted to Y_1

Stego-encoder: $SE(sk, x) = C_{x \oplus sk}$

Stego-decoder: $SD(sk, y) = 0$ iff $y \in Y_{sk}$

Thm.: This is an $\varepsilon^2/\ln 2$ -statistically secure stegosystem (in terms of relative entropy).

Computationally secure stegosystems

Analogous to computational security model for cryptosystems

Universal stegosystem

coverttext given as oracle **C**

Chosen-plaintext attacks (CPA)

passive adversary

Indistinguishability of coverttext and stegotext

Secret-key stegosystem: SKG, SE, SD

$SKG(1k) \rightarrow sk$

$SE(sk, x) \rightarrow c$

$SD(sk, c) \rightarrow x \text{ or } \perp$

Computational stegosystems (2)

Robustness

$$\underline{SD}(\underline{sk}, \underline{SE}(\underline{sk}, x)) = x$$

SS-CPA security defined by experiment with SA

$K \leftarrow SKG(1^k)$
 $(x^*, s) \leftarrow \underline{SA_1}(K)$
 $\underline{b} \in_R \{0, 1\}$
if $b=0$ then $c^* \leftarrow \underline{SE}(K, x^*)$ else $c^* \leftarrow_R \underline{C}$ fi
 $\underline{b^*} \leftarrow \underline{SA_2}(x^*, c^*, s)$

Then

$$\forall \text{ PPT } SA=(SA_1, SA_2) : \Pr[\underline{b} = \underline{b^*}] \leq \frac{1}{2} + \text{negl.}$$

Stegotext and covertext are indistinguishable

Construction – Assumptions

Given

Sym. cryptosystem (KG, E, D) with pseudo-random
ciphertexts & universal hash function G

SKG(1^k)

Run $KG(1^k) \rightarrow sk$ and pick $g \in_R G \rightarrow (sk, g)$



Construction – Rejection sampler

x

SE(sk, m)

$z \leftarrow E(sk, x)$

for $i = 1 \dots |z|$ **do**

repeat

sample C_i \leftarrow_R C

until $g(C_i) = z_i$

output $C_1 \dots C_{|z|}$

$g : \mathcal{C} \rightarrow \{0, 1\}$

hash function

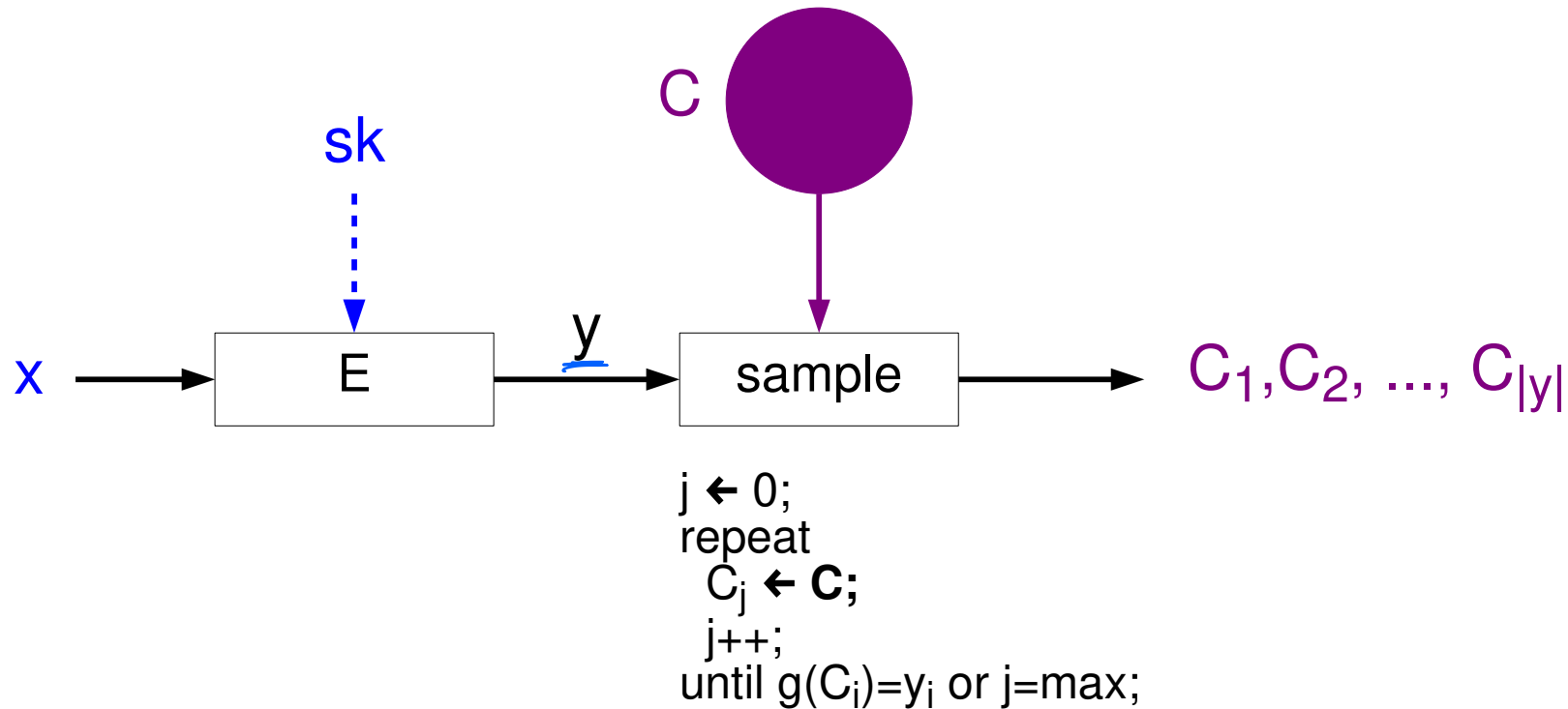
$SD(sk, C_1 \dots C_{|z|})$

for $i = 1 \dots |z|$ **do**

$z_i \leftarrow$ $g(C_i)$

$x \leftarrow D(sk, z)$

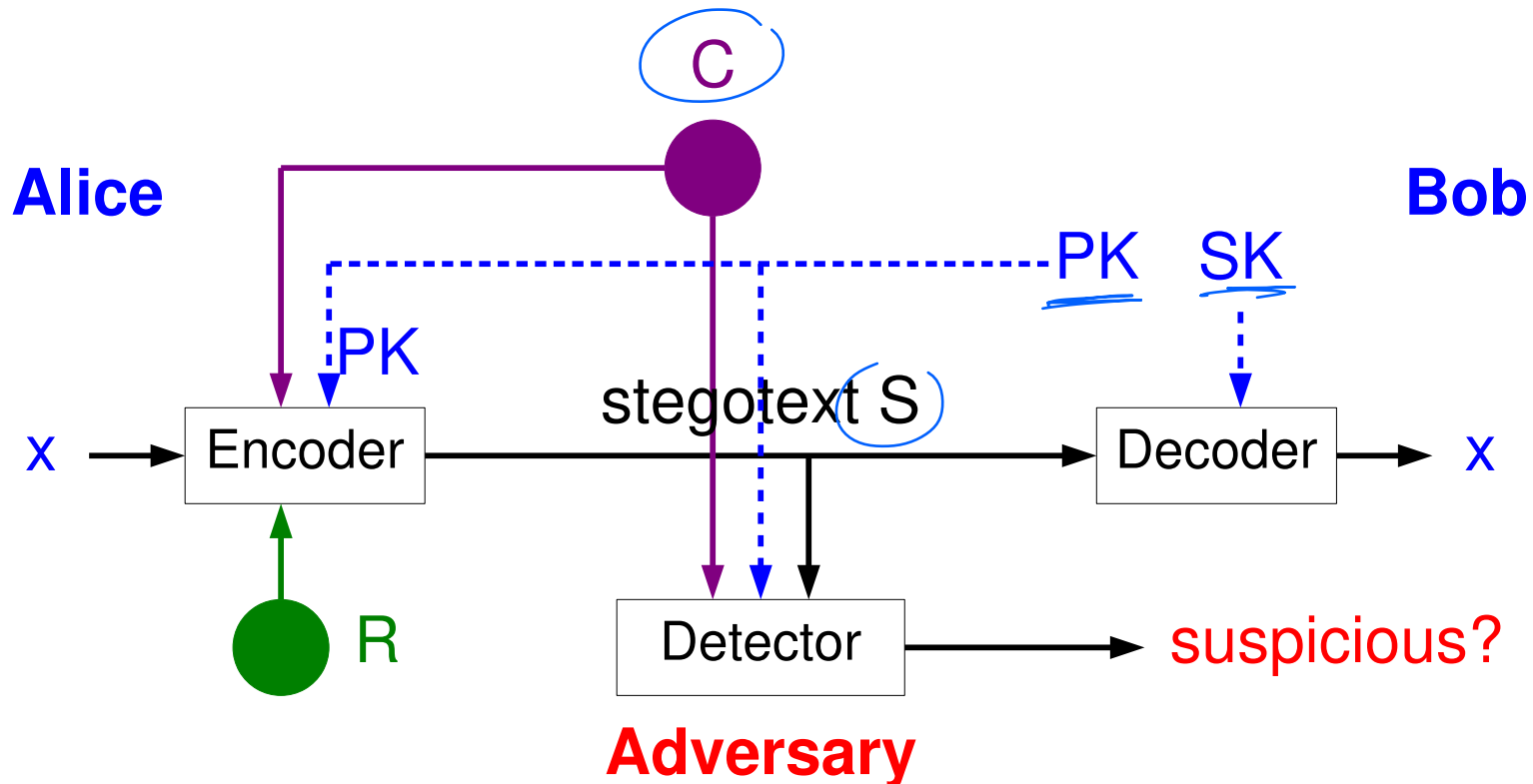
Rejection sampler



If C has large enough min-entropy, then

$$\| \langle g, C_1, C_2 \dots C_{|y|} \rangle - \langle g, C^{|y|} \rangle \| \leq \text{negl.}$$

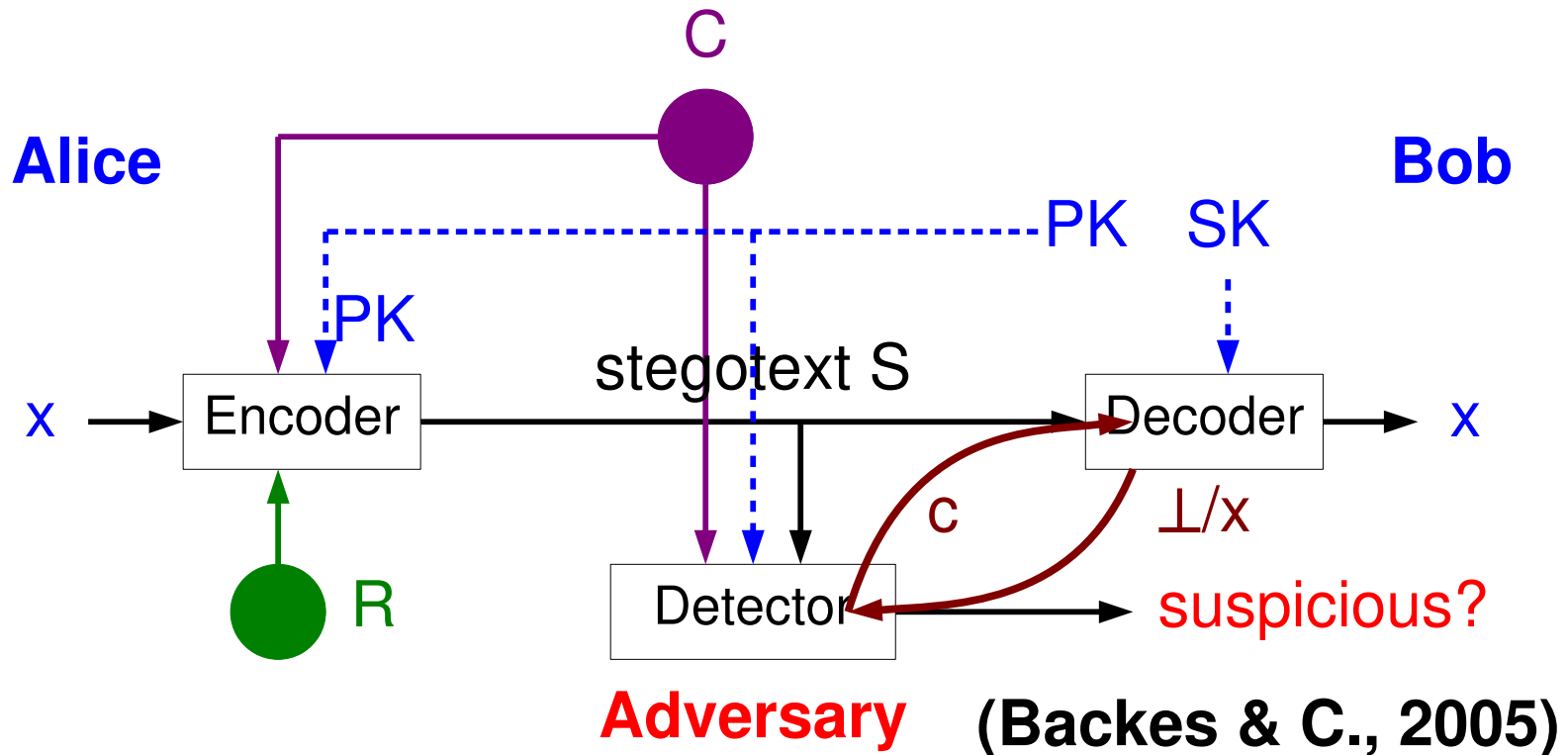
Public-key stegosystems



(v. Ahn & Hopper, 2004)

Bob has public-key/secret-key pair (PK, SK)
Adversary is **passive** (CPA) or **active** (CCA)

Modeling active attacks



Adversary may ask Bob if he considers c to be stegotext
Bob answers 0 (cover) or 1 (stego)

Analogous to adaptive chosen-ciphertext attacks for public-key cryptosystems

Practical steganography

Many tools for image and audio steganography

Ca. 1990-2005

Starting from making imperceptible modifications to the least significant bit (LSB, in pixel or audio data)

Not relevant in practice today

Steganography is often very easy to detect with forensic image analysis tools

Recently – New ideas from machine learning



Stegosystems using ML

Meteor (2021)

<https://meteorfrom.space>

<https://eprint.iacr.org/2021/686>

Sampler uses ML to generate realistic covers

Natural-language text

Uses OpenAI's GPT-2 (Generative Pre-trained Transformer) that produces human-like texts



Meteor

- Examples 1-3
- Long explanation: <https://meteorfrom.space>



Summary

Definition of stegosystems

Security for (secret-key) stegosystems

Perfect → statistical → computational

Always relative to coverttext distribution

Public-key stegosystems

Computational security

Realistic stegosystems using machine learning

Automatically generate coverttext

Permitted coverttext distribution is critical

Executive summary

Stegosystems are cryptosystems with prescribed distribution of ciphertext

