

**THOMSON**

# Copy protection & Statistics

[teddy.furon@thomson.net](mailto:teddy.furon@thomson.net)

[teddy.furon@inria.fr](mailto:teddy.furon@inria.fr)



**THOMSON**

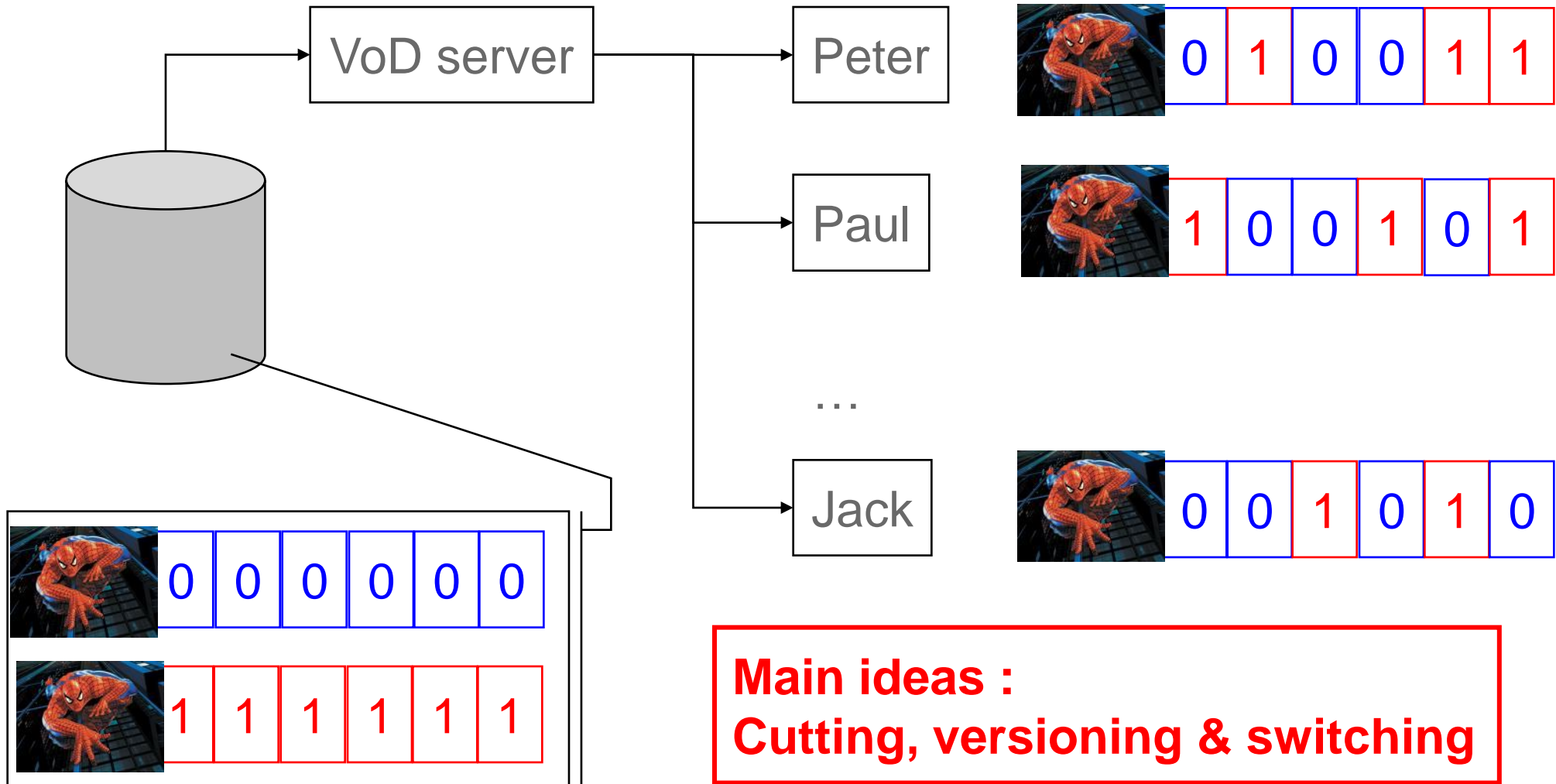
# Introduction: problem

---

- **New trend in copy protection**

- Fight against illegal redistribution of content.
  - content = Hollywood movies
- Find the identity of the hackers amongst  $n$  users.
- Dissuasive weapon.
- a.k.a. : *fingerprinting, content serialization, user forensics, transactional watermarking...*





# Introduction: typical scenario



# Introduction: the collusion

- **Block exchange :**

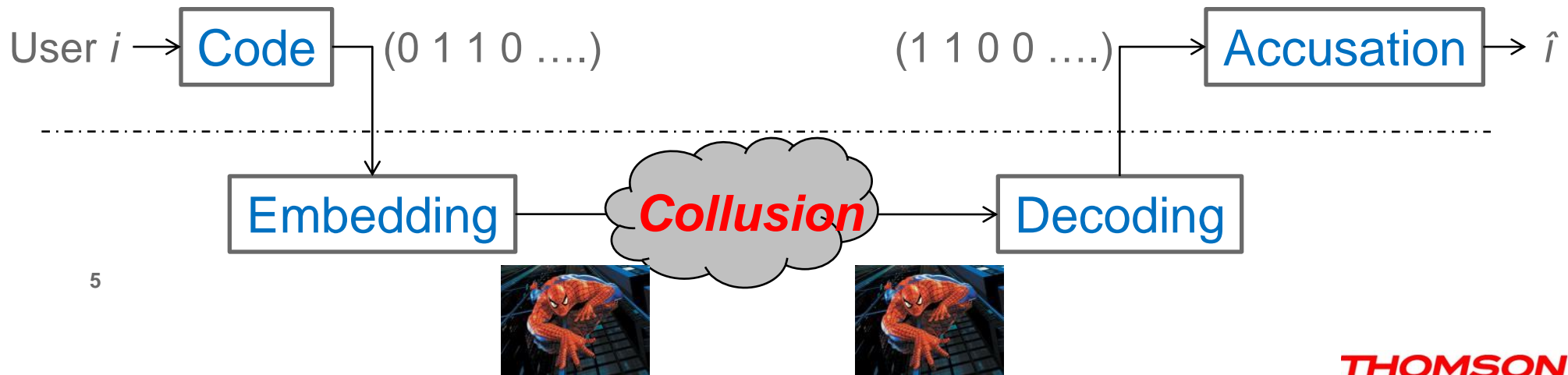
- Colluders cannot create version they don't have.
- The  $i$ -th block in the pirated copy is one of the  $i$ -th blocks from the colluders (« marking assumption »)

Peter		0	1	0	0	1	1
Paul		1	0	0	1	0	1
Jack		0	0	1	0	1	1
<hr/>							
Pirated copy		1	1	1	1	0	1

# Outlines

---

- Introduction
- **Traitor tracing**
  - How to design the code?
  - How to accuse guilty people?
- Digital watermarking
  - How to create two versions of a block?
- False positive probability estimation



# Traitor tracing

---

- **Requirements**

- $n$  users,  $c$  colluders,  $m$  binary code length

- **Code construction**

- $\mathbf{X}$  binary matrix  $n \times m$  (set  $\mathcal{X} \subset \mathcal{B}^m$ )
- $\mathbf{x}_j$  codeword given to user # $j$
- $x_{ji}$   $i$ -th bit of the  $j$ -th codeword

- **Collusion**

- Input:  $C = \{\mathbf{x}_1, \dots, \mathbf{x}_c\} \subset \mathcal{X}$
- Output:  $\mathbf{y} \in \mathcal{B}^m$  pirate sequence

- **Accusation**

- Input:  $\mathbf{y}$  pirate sequence
- Output:  $\mathcal{G}$  set of guilty users

# Cryptographic contributions

---

- **Problem statement** [Fiat&Naor]
- **Terminology (old)** [Pfitzmann]
  - Ex.: Frameproof code  $\nexists C \mid \mathbf{y} \in \mathcal{X} \setminus C$
- **Relationship with Error Correcting Codes** [Stinson]
  - Pirate sequence = codeword + noise
  - Accusation = correcting errors
  - Not efficient: very very very long code
- **Relaxation of the constraint** [Boneh&Shaw]
  - $P_{fa}$  Probability of accusing at least one innocent
  - $P_{mi}$  Probability of missing all colluders
- **Non constructive theorem** [Peikert03]
  - $m \gtrsim O( c^2 \cdot \log( n \cdot P_{fa}^{-1} ) )$

# A revolution coming from the Statistics

---

## ● Probabilistic codes

[Gabor Tardos]

- The first exhibition of a code achieving the Peikert bound
$$m = 100 c^2 \cdot \log(n \cdot P_{fa}^{-1})$$
- An unknown genius... work rediscovered 2 years after.
- No rationale, no clue, no intuition except
  - « *the full power of randomization* »
- Extremely simple : 10 lines of matlab
- Extremely flexible :  $n$ ,  $m$  loosely tightened
- Constant '100' raised suspicion



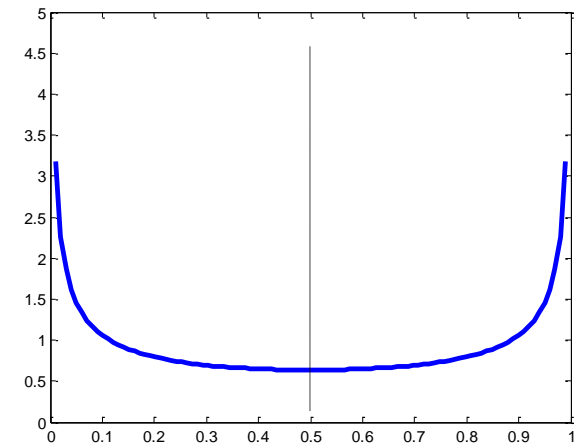
# Tardos codes

## ● Initialization

— Draw randomly:  $p_i \in [0,1]$  ,  $i=1:m$  , i.i.d.,  $p \sim f(p)$

—  $\mathbf{p} = (p_1, \dots, p_m)$  is the secret of the code

—  $f(p) = 1 / \pi(p(1-p))^{1/2}$  (???)



## ● Code construction

— Draw randomly:  $x_{ji} \in \{0,1\}$  ,  $\text{Prob}[x_{ji} = 1] = p_i$

— If  $p_i \sim 0^+$  , then almost all user have a ‘0’ sparse code

— If  $p_i \sim 1/2$  , then as many ‘0’ as ‘1’ dense code

# Tardos accusation

## ● Accusation

— Accuse user #j if

$$S_j > T$$

$$S_j = \sum_i g(y_i, x_{ji}, p_i)$$

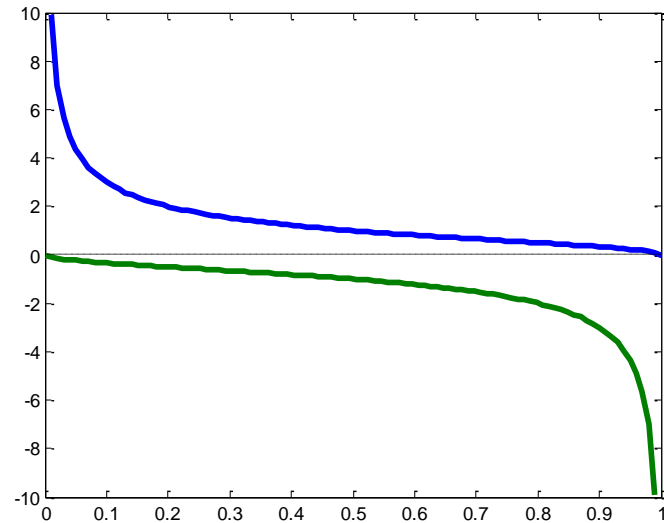
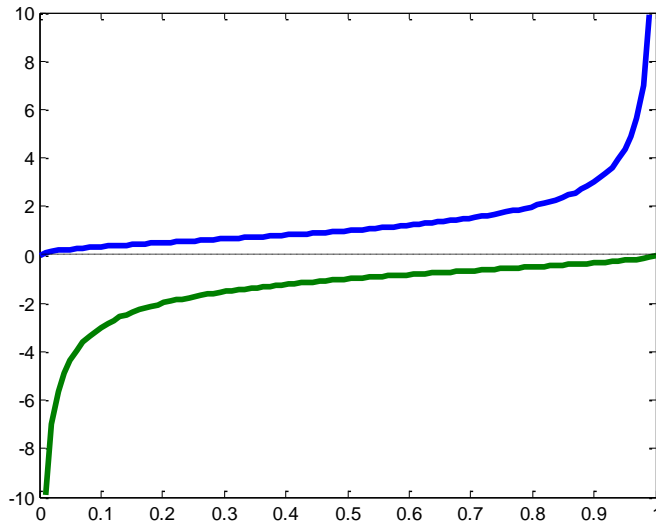
—  $g(0,0,p) = + ( p/(1-p) )^{1/2}$

—  $g(0,1,p) = - ( (1-p)/p )^{1/2}$

$g(1,1,p) = + ( (1-p)/p )^{1/2}$

$g(1,0,p) = - ( p/(1-p) )^{1/2}$

(?!?)



# Tardos code

---

**Why does it work?**

# Mathematical model of the collusion

---

- Assumptions about the collusion

- Memoryless

$$y_i = F_i(x_{1i}, \dots, x_{ci})$$

- Stationary

$$y_i = F(x_{1i}, \dots, x_{ci})$$

- Permutation Invariant

$$y_i = F(s_i)$$

$$s_i = \sum_j x_{ji}$$

- Probabilistic

$$\theta_s = \text{Prob}(y = 1 | s)$$

- Model

$$\theta = (\theta_0, \theta_1, \dots, \theta_c)$$

Marking assumption

$$\theta_0 = 0 \quad \text{and} \quad \theta_c = 1$$

Therefore, the collusion indeed lies in  $[0, 1]^{c-1}$

# Mathematical model of the collusion

---

- Using the model

$$\text{Prob}[y=1 \mid p] = \sum_s \text{Prob}[y=1, s \mid p] = \sum_s \theta_s \cdot \binom{c}{s} p^s (1-p)^{c-s}$$

$$\text{Prob}[y=1 \mid x=1, p] = \sum_s \theta_s \cdot \binom{c-1}{s-1} p^{s-1} (1-p)^{c-s}$$

- 1<sup>st</sup> and 2<sup>nd</sup> order statistics

- Innocent:  $\mathbb{E} [ S_j ] = 0$   $\mathbb{E} [ S_j^2 ] = m$

- Colluder:  $\mathbb{E} [ S_j ] = 2m/\pi c$   $\mathbb{E} [ S_j^2 ] = m$

( here:  $\mathbb{E} [ . ] = \mathbb{E}_P [ \mathbb{E}_X [ \mathbb{E}_Y [ . ] ] ]$  )

- Miracle: independent from  $\theta$

- Markov bound:  $m = 100 c^2 \cdot \log( n \cdot P_{fa}^{-1} )$

- Asymptotics: CLT – Scores are Gaussian distributed

# Statistical interpretations (I)

---

- The collusion process is a nuisance parameter
- The hypothesis test is based on a pivotal quantity...
- ... at least up to 1<sup>st</sup> and 2<sup>nd</sup> order statistics.

$$\begin{aligned} f(p) &= 1 / \pi(p(1-p))^{1/2} \\ g(0,0,p) &= + (p/(1-p))^{1/2} \\ g(1,1,p) &= \dots \end{aligned}$$

$\Rightarrow$

Innocent  $\mathbb{E}[S_j] = 0; \mathbb{E}[S_j^2] = m$

$\Leftarrow$

Colluder  $\mathbb{E}[S_j] = 2m/\pi c; \mathbb{E}[S_j^2] = m$

# Statistical interpretation (II)

---

- **One collusion process – Many code densities**
  - Dense ( $p=1/2$ ): worst attack: minority vote 3 colluders
  - Sparse ( $p=0^+$  or  $1^-$ ): worst attack: majority vote
- **The sequence  $p$  is THE secret!**
- **« How secret is this secret? »**
  - The colluders could estimate it:  $p_i = s_i / c$
  - The colluders know  $f(p)$ : a priori distribution.
- **Wrong idea: already captured by our model.**
- **Shed more light on Tardos choice:**
  - $f(p)$  is the Jeffreys prior, the less informative prior distribution
  - The less useful for the colluders.

# Statistical interpretation (III)

---

- **Hypothesis test**

- $H_0$ : Innocent  $\text{Prob}(\mathbf{y}, \mathbf{x} | \mathbf{p}) = \text{Prob}(\mathbf{y} | \theta, \mathbf{p}) \cdot \text{Prob}(\mathbf{x} | \mathbf{p})$
- $H_1$ : Colluder  $\text{Prob}(\mathbf{y}, \mathbf{x} | \mathbf{p}) = \text{Prob}(\mathbf{y} | \mathbf{x}, \theta, \mathbf{p}) \cdot \text{Prob}(\mathbf{x} | \mathbf{p})$
- Performance criterion

$$R(f; \theta) = \mathbb{E}_P [ D_{\text{KL}}(H_1; H_0 | P, \theta) ]$$

- **Game theory**

- Between designer and colluders
- MaxMin game

$$R(f^*; \theta^*) = \max_f \min_{\theta} R(f; \theta)$$

- **Asymptotically,  $c \rightarrow \infty$ :**

- Equilibrium:  $f^*(p) = (\pi^2 p(1-p))^{-1/2}$  and  $\theta^* = (0, c^{-1}, 2c^{-1}, \dots, 1)$ .



# Trends

---

- **New accusation strategy**

- « Learn and Match »
- Estimate  $\theta$  and use Likelihood ratio test to accuse (E.-M.)

- **K-uplets scores**

- Inf. theory:

$$I(\mathbf{y} ; \mathbf{x} | \theta) \leq I(\mathbf{y} ; \{\mathbf{x}_1, \mathbf{x}_2\} | \theta) \leq \dots \leq I(\mathbf{y} ; \{\mathbf{x}_1, \dots, \mathbf{x}_c\} | \theta)$$

- Practical? Complexity in  $\sim O(n^c)$

- **Q-ary alphabet**

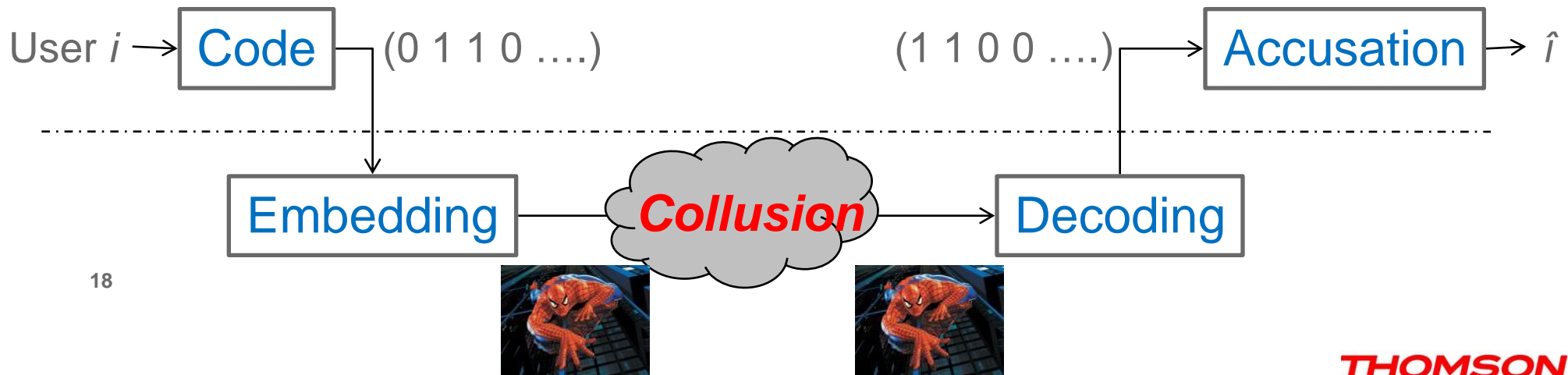
- **Other collusion models**

- Erasures (cut movie scenes)

# Outlines

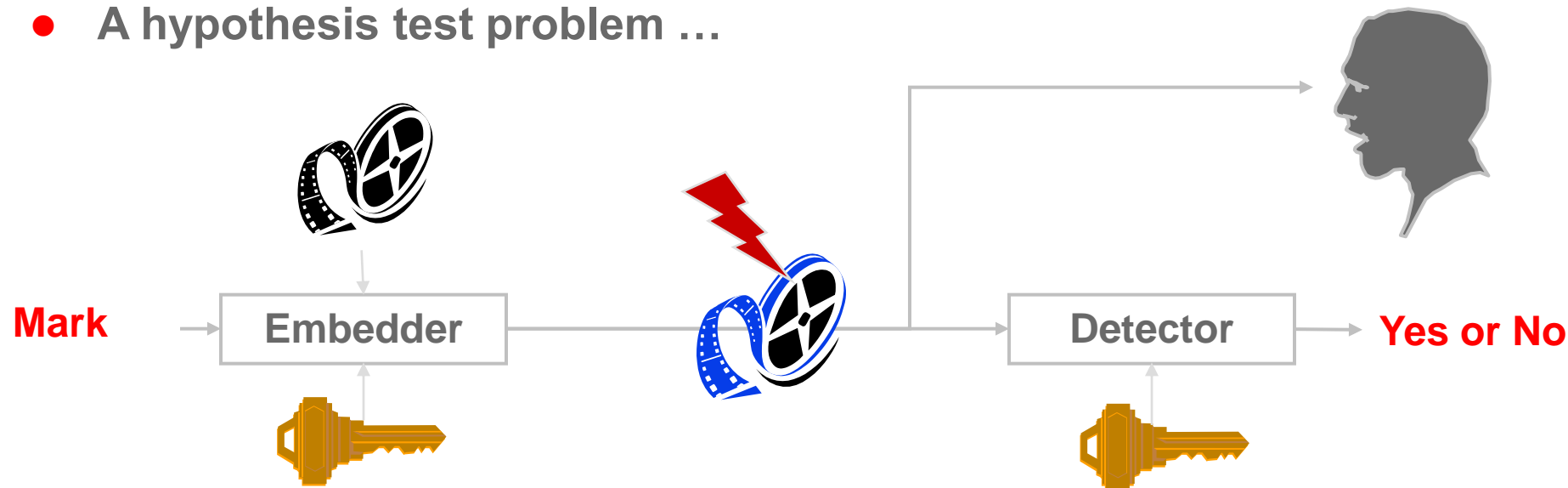
---

- Introduction
- Traitor tracing
  - How to design the code?
  - How to accuse people?
- **Digital watermarking**
  - How to create two versions of a block?
- False positive probability estimation



# Definition of digital watermarking

- Data hiding: art and science of hiding data in multimedia digital contents.
- A hypothesis test problem ...



- ... under some special constraints:
  - non perceptibility (watermark, *filigrane* in French)
  - robustness (*tatouage* in French)
  - security

# Assumptions: feature extraction

---

- **Watermark Embedding**

- From a content block  $B$ , extract some meaningful features

$$\mathbf{s} = \text{Ext} ( B ), \quad \text{with } \mathbf{s} \in \mathbb{R}^L$$

- Modify  $\mathbf{s}$  into  $\mathbf{x} = \mathbf{s} + \mathbf{w}$  s.t.  $\|\mathbf{w}\|^2 \leq L \cdot P_w$

- The vector  $\mathbf{w}$  is the secret of the watermarking scheme
- Put back the features into the content

$$B_w = \text{Ext}^{-1}( \mathbf{x}, B )$$

- **Attack on the watermarked image**

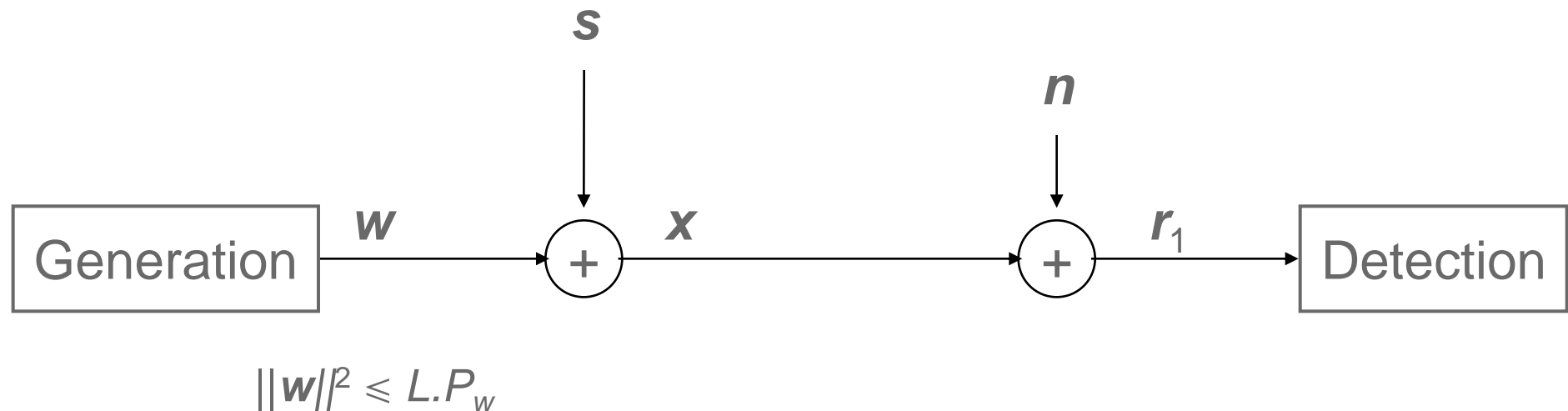
- Distort  $\mathbf{x}$  into  $\mathbf{z} = \mathbf{x} + \mathbf{n}$  s.t.  $\mathbb{E} \|\mathbf{n}\|^2 \leq L \cdot P_n$

- **Detection**

- $H_0: \mathbf{r}_0 = \mathbf{s} + \mathbf{n}$  (given by Nature)
- $H_1: \mathbf{r}_1 = \mathbf{s} + \mathbf{w} + \mathbf{n}$

# Naïve idea

---



- **Gaussian setup**

$$r_0 = s + n \sim \mathcal{N}(\mathbf{0}, P_s + P_n) \quad \text{vs.} \quad r_1 = s + w + n \sim \mathcal{N}(w, P_s + P_n)$$

- Performances limited by the Kullback-Leibler distance

$$D_{\text{KL}}(r_0 \| r_1) = L.P_w / 2 (P_n + P_s)$$

- Data processing theorem, Stein Lemma.

# An Information theoretic revolution [Costa]



$$\mathbb{E}_{\mathbf{s}} \|\mathbf{w}(\mathbf{s})\|^2 \leq L \cdot P_w$$

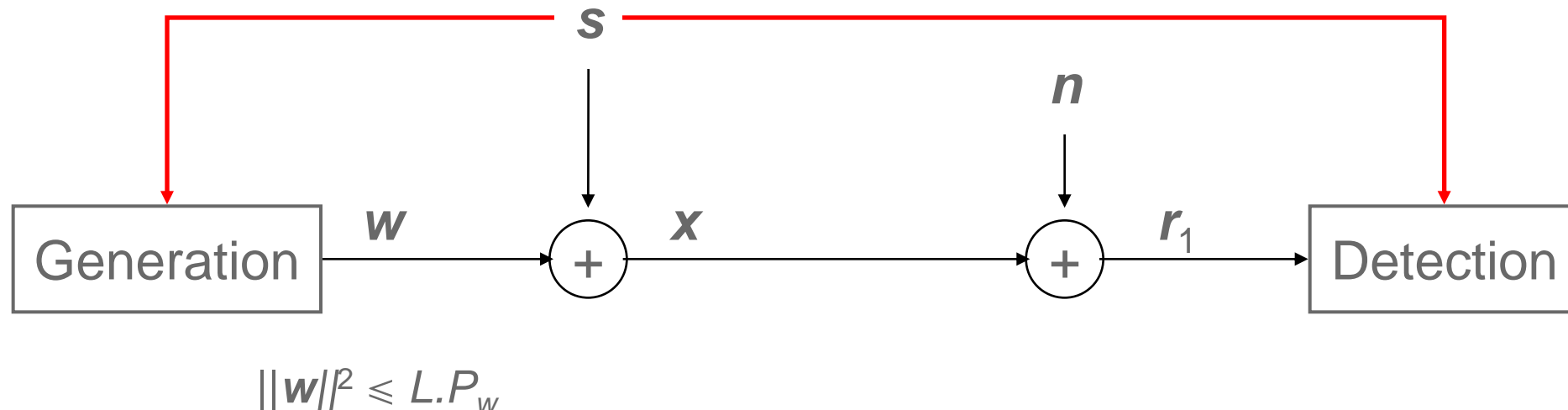
- **Gaussian setup**

$$\mathbf{r}_0 = \mathbf{s} + \mathbf{n} \sim \mathcal{N}(\mathbf{0}, P_s + P_n) \quad \text{vs.} \quad \mathbf{r}_1 = \mathbf{s} + \mathbf{w}(\mathbf{s}) + \mathbf{n}$$

— Performances limited by the Kullback-Leibler distance

$$\max_{\mathbf{w}(\cdot)} D_{\text{KL}}(\mathbf{r}_0 \parallel \mathbf{r}_1) = ???$$

# The informed setup



- **Gaussian setup**

- $r_0 - s = n \sim \mathcal{N}(\mathbf{0}, P_n)$  vs.  $r_1 - s = w + n \sim \mathcal{N}(w, P_n)$
- Performances limited by the Kullback-Leibler distance

$$D_{\text{KL}}(r_0 \| r_1) = P_w / 2.P_n$$

# The side informed setup

[Costa]



$$\mathbb{E}_{\mathbf{s}} \|\mathbf{w}(\mathbf{s})\|^2 \leq L \cdot P_w$$

## • Gaussian setup

- Performances limited by the Kullback-Leibler distance

$$P_w / 2 \cdot (P_n + P_s) \leq \max D_{\text{KL}}(r_0 \| r_1) \leq P_w / 2 \cdot P_n$$

## • Philosophical question

- Is  $\mathbf{s}$  a channel noise *Detection* or a channel state? *Generation*



# Asymptotical Gaussian case

---

- Suppose the following watermark embedding

$$\mathbf{x} = \mathbf{s} + \mathbf{w}(\mathbf{s}) = \mathbf{s} + (\alpha - \lambda \cdot \mathbf{s}^T \mathbf{u}) \cdot \mathbf{u} \quad \text{with } \|\mathbf{u}\|^2 = 1$$

- $\lambda \in [0, 1], \alpha > 0$
- “Push and cancel” mixed strategy
- Power constraint:

$$\alpha^2 + \lambda^2 \cdot P_s = L \cdot P_w$$

- KL- distance

- $D_{\text{KL}}(\mathbf{r}_0 \| \mathbf{r}_1) = F(\lambda, \alpha) = F(\lambda)$ .
- Optimize :  $\lambda^* = \arg \max_{\lambda} F(\lambda)$
- Take to the limit

$$\lim_{L \rightarrow \infty} F(\lambda^*) = P_w / 2 \cdot P_n$$

# How to put this into practice?

---

- **Gaussian setup with fixed length  $L$**

- How to maximize  $D_{\text{KL}}$ ?
- How to design the detector?

[Merhav]

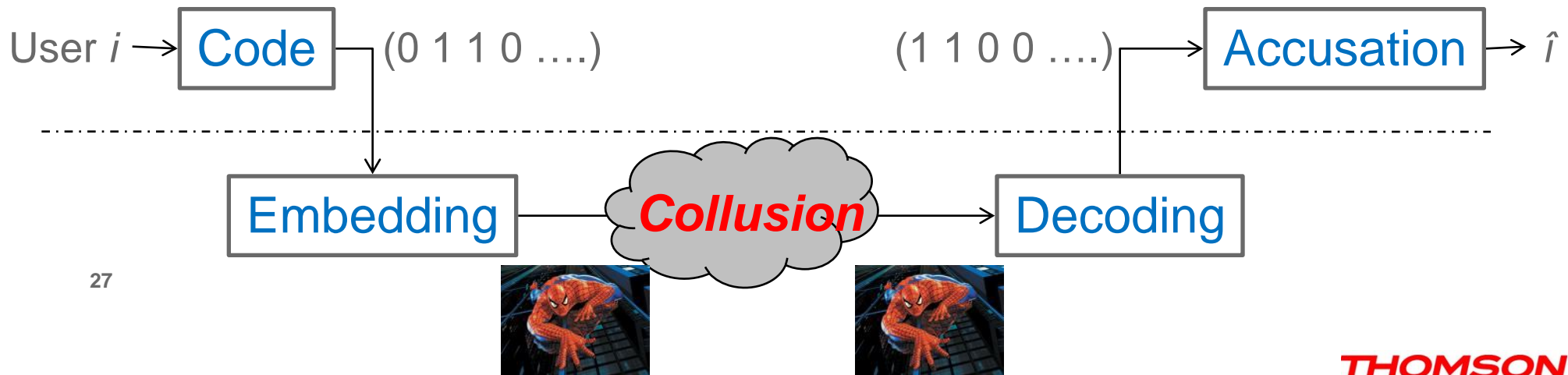
- **Real world**

- Nuisance parameters:  $P_s$ ,  $P_w$ ,  $P_n$ , type of attack
- No longer Gaussian r.v.
- No longer Euclidean distance, but perceptual distance

# Outlines

---

- Introduction
- Traitor tracing
  - How to design the code?
  - How to accuse people?
- Digital watermarking
  - How to create two versions of a block?
- **False positive probability estimation**



# False positive estimation

---

- **2 techniques have a common issue: very very small  $P_{fa}$**

- Traitor tracing: probability to accuse an innocent.

$$P_{fa} = \text{Prob} [ d(\mathbf{x}) > T ]$$

- with  $d(.)$  Tardos scoring
- $\mathbf{x}$  the sequence of a user,  $\text{Prob}(\mathbf{x}_j) = \prod_i p_i^{x_{ji}} \cdot (1-p_i)^{(1-x_{ji})}$

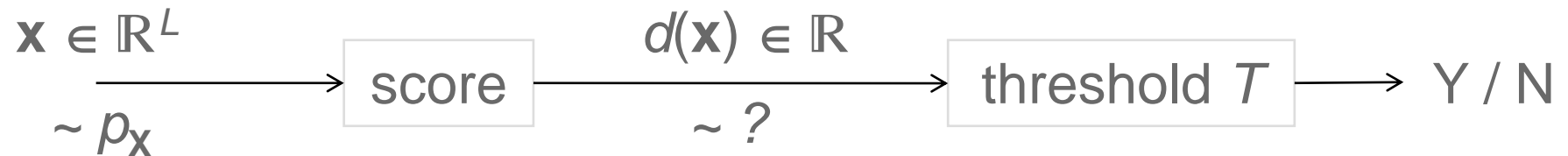
- Watermarking: probability of detecting the mark in a non watermarked content

$$P_{fa} = \text{Prob} [ d(\mathbf{x}) > T ]$$

- with  $d(.)$  likelihood of being watermarked
- $\mathbf{x}$  features extracted from a content block,  $\mathbf{x} \sim p_{\mathbf{x}}$

# Monte Carlo Method estimation

- Structure of the detector



- MCM estimation

- Run  $n$  experiments
- Increment  $k$  when  $d(\mathbf{x}) > T$
- Estimate  $\hat{P}_{fa} = k/n$

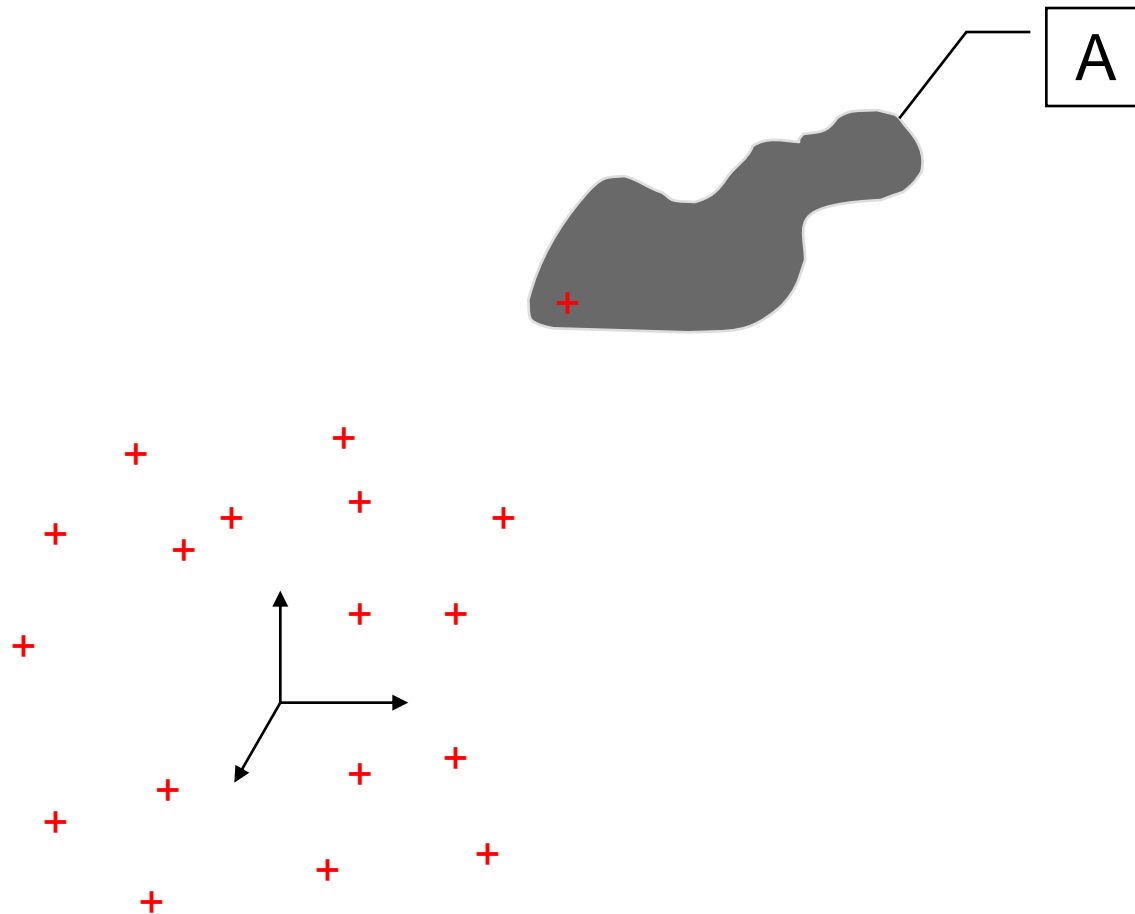
- Issues

- $k \neq 0$   $n = O(1/P_{fa})$
- Relative std of  $\hat{P}_{fa} = (P_{fa} \cdot n)^{-1/2}$   $n \sim 100/P_{fa}$

**⇒ The smaller the probability, the harder its estimation**

# Geometric interpretation

---



# Key idea of our algorithm

---

## ■ Divide and Conquer

$$\begin{aligned} P_{fa} &= \Pr(A) = \Pr(A, A_{N-1}) && \text{if } A \text{ implies } A_{N-1} \\ &= \Pr(A | A_{N-1}) \cdot \Pr(A_{N-1}) && \text{2 estimations, less difficult} \\ &= \Pr(A | A_{N-1}) \cdot \Pr(A_{N-1} | A_{N-2}) \cdot \Pr(A_{N-2} | A_{N-3}) \dots \Pr(A_1) \\ &&& A \Rightarrow A_{N-1} \Rightarrow A_{N-2} \Rightarrow \dots \Rightarrow A_1 \end{aligned}$$

$$\begin{aligned} P_{fa} &= \Pr(d(\mathbf{x}) \geq T) \\ &= \Pr(d(\mathbf{x}) \geq T | d(\mathbf{x}) \geq T_{N-1}) \cdot \Pr(d(\mathbf{x}) \geq T_{N-1} | d(\mathbf{x}) \geq T_{N-2}) \dots \Pr(d(\mathbf{x}) \geq T_1) \\ &&& T > T_{N-1} > T_{N-2} > \dots > T_1 \end{aligned}$$

$$\hat{P}_{fa} = \hat{a}_N \cdot \hat{a}_{N-1} \cdot \dots \cdot \hat{a}_1$$

# Our estimator

---

- **Inputs**

- Distribution of input data  $p_{\mathbf{x}}$ , score  $d(\cdot)$ , threshold  $T$

- **Outputs**

- Estimation  $\hat{P}_{fa}$  of  $\Pr(d(\mathbf{x}) > T)$  for  $\mathbf{x} \sim p_{\mathbf{x}}$

- **Ingredients: 3 subroutines**

- SCORE

- Function  $d(\cdot): \mathbb{R}^L \rightarrow \mathbb{R}$
- ‘Smooth’

- GENERATE

- Generate input data  $\mathbf{x}$  distributed  $\sim p_{\mathbf{x}}$

- MODIFY

- $\mathbf{y} = f(\mathbf{x})$ , random function
- such that  $\mathbf{y} \sim p_{\mathbf{x}}$  and  $\mathbf{y}$  is ‘close’ to  $\mathbf{x}$



# Our estimator

---

- **Divide and conquer**

$$P_{fa} = \Pr(d(\mathbf{x}) \geq T | d(\mathbf{x}) \geq T_{N-1}) \dots \Pr(d(\mathbf{x}) \geq T_{j+1} | d(\mathbf{x}) \geq T_j) \dots \Pr(d(\mathbf{x}) \geq T_1)$$

- **Initialization**

- ‘small’ Monte Carlo Method Estimator
- Generate  $n$  input data (particles)  $\mathbf{x}^{(1)} \sim p_{\mathbf{x}}$
- Count the number of times the score is above the 1<sup>st</sup> threshold

$$k_1 = |\{\mathbf{x}_i^{(1)} \mid d(\mathbf{x}_i^{(1)}) > T_1\}|$$

$$\hat{a}_1 = k_1/n$$

# Our estimator

---

- **Divide and conquer**

$$P_{fa} = \Pr(d(\mathbf{x}) \geq T | d(\mathbf{x}) \geq T_{N-1}) \dots \Pr(d(\mathbf{x}) \geq T_{j+1} | d(\mathbf{x}) \geq T_j) \dots \Pr(d(\mathbf{x}) \geq T_1)$$

- **Iteration  $j \rightarrow j+1$**

- We start with  $k_j = |\{\mathbf{x}_i^{(j)} \mid d(\mathbf{x}_i^{(j)}) > T_j\}|$  particles in region  $A_j$ .
- DUPLICATE: We randomly select  $n$  particles in this set
- MODIFY:  $\mathbf{z} = f(\mathbf{x}_i^{(j)})$
- SELECTION:
  - If  $d(\mathbf{z}) > T_j$  then  $\mathbf{x}_i^{(j+1)} = \mathbf{z}$
  - Else  $\mathbf{x}_i^{(j+1)} = \mathbf{x}_i^{(j)}$
- We now have  $n$  particles  $\sim p_{\mathbf{x}}$  in  $A_j$ .

# Our estimator

---

- Iteration  $j \rightarrow j+1$



- THRESHOLD (or MCM estimator)

- Count the number of times the score is above the  $(j+1)^{\text{th}}$  threshold

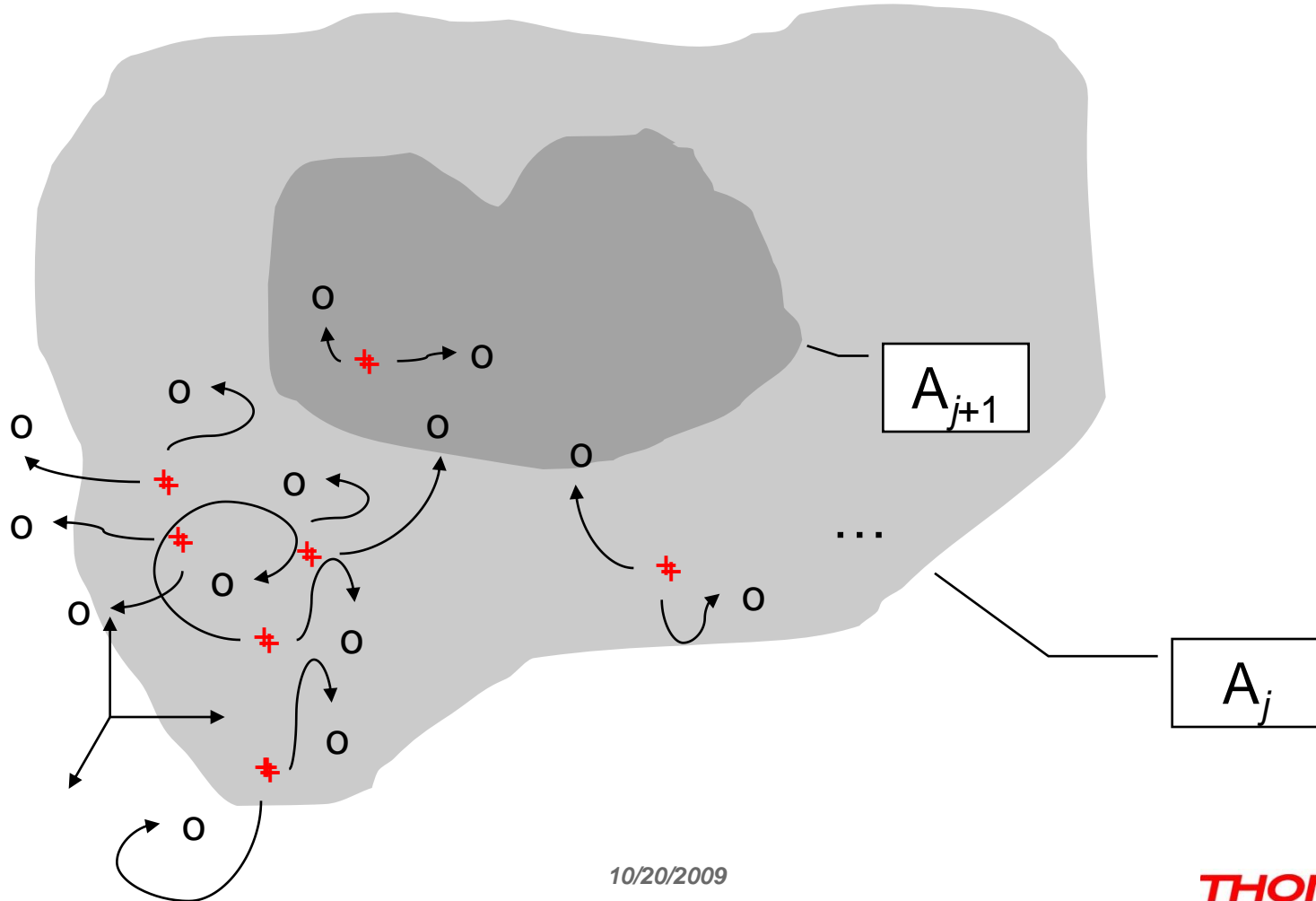
$$k_{j+1} = |\{\mathbf{x}_i^{(j+1)} \mid d(\mathbf{x}_i^{(j+1)}) > T_{j+1}\}|$$

$$\hat{a}_{j+1} = k_{j+1}/n$$

# Our estimator

---

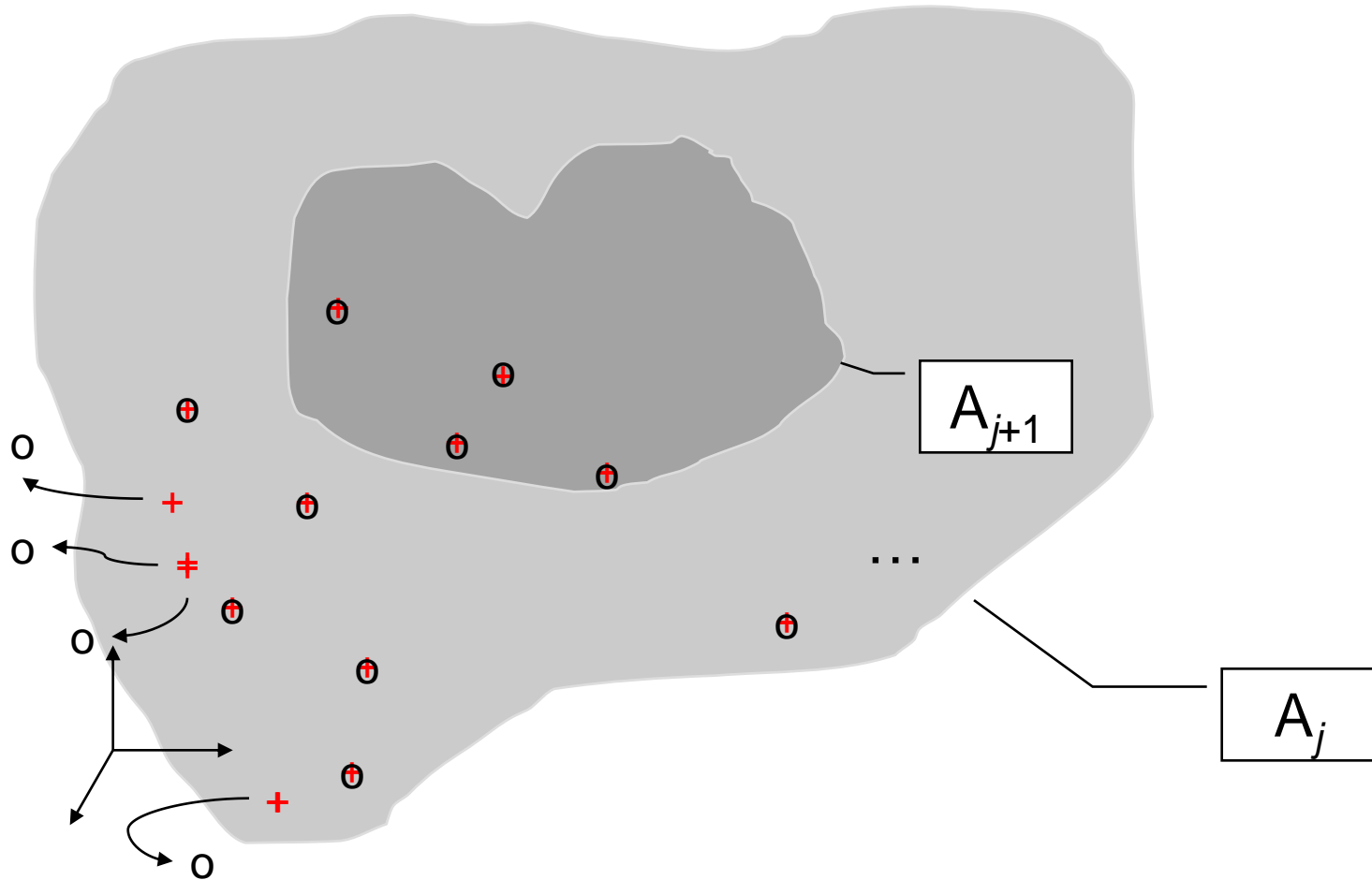
- Geometric interpretation



# Our estimator

---

- Geometric interpretation



# Our estimator

---

- **Last trick**

- How to define the thresholds  $T_j$ ?
- Variance of the estimation  $P_{fa}$  is minimized if  $a_j = \text{cte}$
- Inverse:  $T_j$  is the  $k$ -th biggest score out of  $n$ :  $\hat{a}_j = k/n$

- **Divide and conquer**

$$P_{fa} = \Pr(d(\mathbf{x}) \geq T \mid d(\mathbf{x}) \geq T_{N-1}) \dots \Pr(d(\mathbf{x}) \geq T_{j+1} \mid d(\mathbf{x}) \geq T_j) \dots \Pr(d(\mathbf{x}) \geq T_1)$$

- **Stopping condition**

- When  $T_N > T$ .
- Count the number  $k'$  of scores really above  $T$
- $\hat{P}_{fa} = (k/n)^{N-1} \cdot k' / n$

# Properties

---

- **Small number of trials**

$$nN = n[\log P_{fa}^{-1} / \log (k/n)^{-1}]$$

- **Asymptotic consistency**

$$\hat{P}_{fa} \rightarrow P_{fa} \text{ as } n \rightarrow \infty$$

(almost surely)

- **Asymptotic Normality**

$$n^{1/2} \cdot (\hat{P}_{fa} - P_{fa}) \rightarrow \mathcal{N}(b, \sigma^2) \text{ as } n \rightarrow \infty$$

(in law)

- **Asymptotic statistics**

- Relative std in  $O( (\log(P_{fa}^{-1}) / n)^{1/2} )$
- Fast vanishing relative bias in  $O(1/n)$

# Experiment #1: Digital Watermarking

---

- **Toy example**

- $\mathbf{x} \sim p_{\mathbf{x}}$  isotropic (for instance, white Gaussian noise  $\mathcal{N}(\mathbf{0}, \mathbf{I}_L)$ )
- $d(\mathbf{x}) = \mathbf{x}^T \mathbf{u} / \|\mathbf{x}\|$
- Ground truth:  $P_{fa} = 1 - \text{IncBeta}(T^2, 1/2, (L-1)/2)$

- **Ingredients**

- GENERATE: Matlab randn
- MODIFY:

$$\mathbf{y} = (\mathbf{x} + q \cdot \mathbf{n}) / (1 + q^2)^{1/2}$$

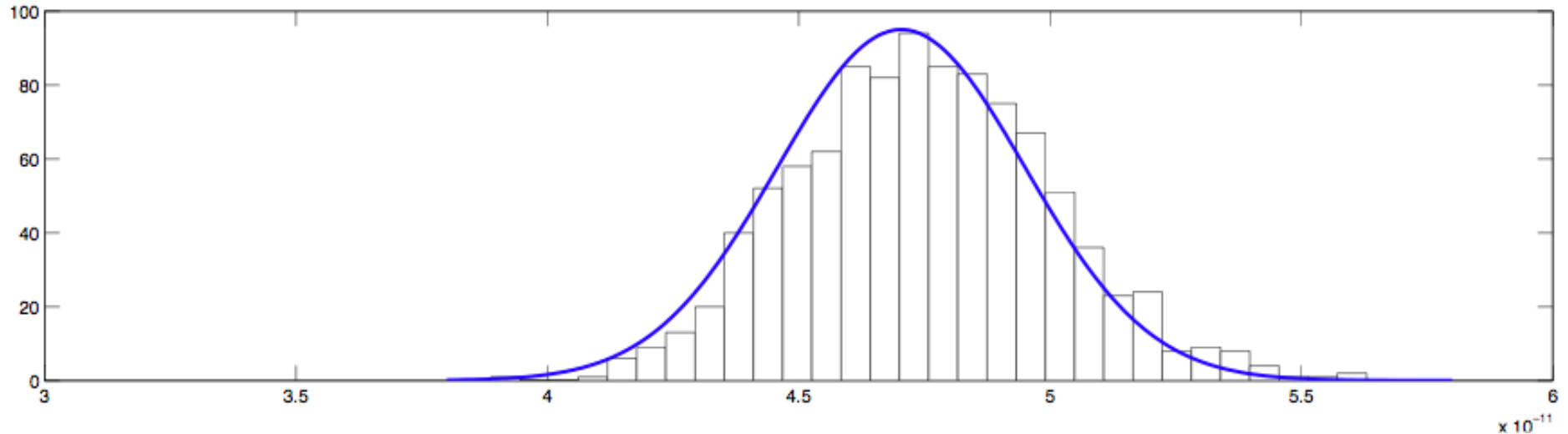
$$\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_L)$$

- $q$  fixes the strength of the modification.



# Experiment #1: Asymptotic normality

---



$$L = 20$$

$$T = 0.95$$

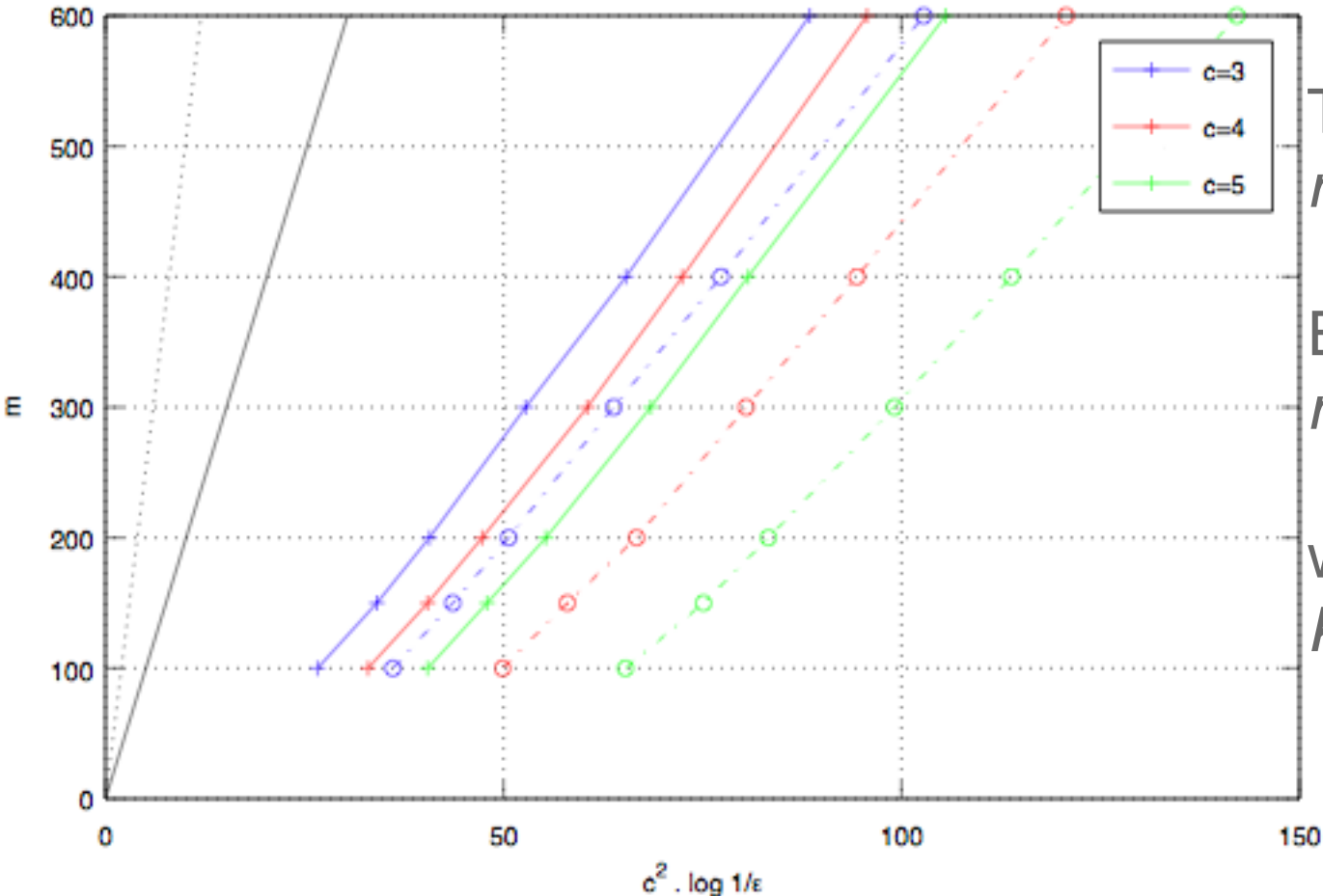
$$P_{fa} = 4.7 \cdot 10^{-11}$$

$$k/n = 3/4$$

$$n = 50,000$$

200 runs

# Experiment #2: Estimating minimum length



Tardos

$$m = 100 \cdot c^2 \log(n/P_{fa})$$

Experimental

$$m \approx K_1 \cdot c^2 \log(n/P_{fa}) + K_0(c)$$

with

$$K_1 \approx 7.6$$

# Conclusion

---

# ***THOMSON***

# Thank you for your attention

This document is for background informational purposes only. Some points may, for example, be simplified. No guarantees, implied or otherwise, are intended



***THOMSON***