

Bias-based traitor tracing codes

Boris Škorić

Eindhoven University of Technology

Guest lecture

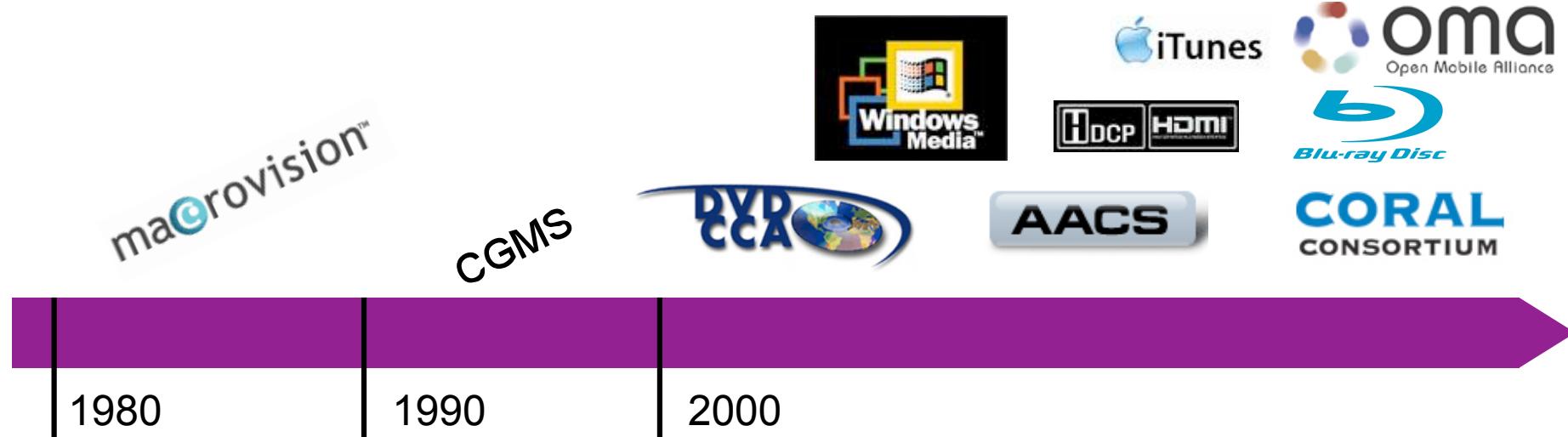
3 December 2013



Outline

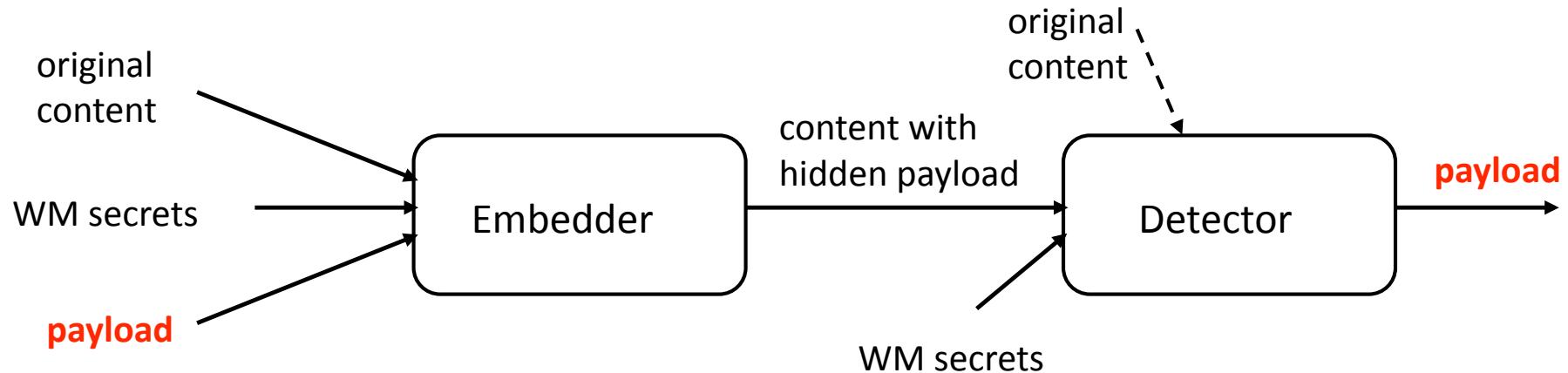
- Collusion attacks on watermarks
- Bias-based codes (Tardos codes)
- Information-theoretic maxmin game
 - Saddlepoint
 - Capacity
- Capacity-achieving score function

Trends in content protection (old slide)



- Consumers increasingly dislike DRM
- Vista content protection spec
"longest suicide note in history" (Gutmann 2006)
- "Disembodied" distribution ⇒ Hard to DRM-protect; Easy to watermark
- April 2007: EMI announces DRM-free music
- Gradual shift from copy prevention to distribution tracking

Watermarking (a.k.a. Fingerprinting)



Forensic tracing

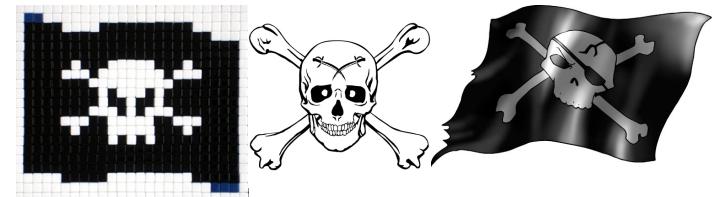
- Payload = unique identifier of recipient
- Redistribution traced back to source

Examples

- Jan.2004: Man arrested for distributing oscar screeners.
- Digital cinema

Collusion attacks

"Coalition of pirates"



- Users pool their content
- Differences point to watermark
- Attackers remove watermark

■ = "detectable positions"

<i>pirate #1</i>	1	1	1	0	1	0	1	0	0	0	0	1
------------------	---	---	---	---	---	---	---	---	---	---	---	---

#2	1	0	1	0	1	0	1	0	1	0	1	1
----	---	---	---	---	---	---	---	---	---	---	---	---

#3	1	0	1	0	1	0	1	0	0	0	1	1
----	---	---	---	---	---	---	---	---	---	---	---	---

#4	1	1	1	0	0	0	1	1	0	0	0	1
----	---	---	---	---	---	---	---	---	---	---	---	---

Attacked Content	1	0/1	1	0	0/1	0	1	0/1	0/1	0	0/1	1
------------------	---	-----	---	---	-----	---	---	-----	-----	---	-----	---

Collusion-resistant watermarking

Requirements

- Resistance against $c \leq c_0$ attackers
- Low False Positive error rate
- Low False Negative error rate
- ... and all that with small watermark payload! (7bits/min. video)

Attack model

- "**Marking assumption**": Modification only at detectable positions
- Several options
 - **Restricted digit model**: Choice from available symbols only
 - Unreadable digit model: Erasure allowed
 - Arbitrary digit model: Arbitrary symbol (but not erasure)
 - General digit model

} equivalent
for
binary
symbols

History of collusion resistance: Code length

Construction



Boneh and Shaw 1998: $m = \mathcal{O}(c_0^4 \ln[n/\eta] \ln[1/\eta]), q = 2$

Tardos 2003: $m = 100c_0^2 \ln(1/\varepsilon_1), q = 2, \varepsilon_2 = \varepsilon_1^{c_0/4}$

Chor et al 2000: $m = 4c_0^2 \log n, q = 2c_0^2$

Staddon et al 2001: $m = c_0^2 \log_q(n), q > m - 2$

Huang + Moulin; Amiri + Tardos 2009: $m = 2\ln 2 \cdot c_0^2 \ln[1/\varepsilon_1], q = 2$

Tardos 2003:

$m = \Omega(c_0^2 \ln[1/\varepsilon_1]), q$ arbitrary

Boneh and Shaw 1998:

$m = \Omega(c_0 \ln[1/c_0\eta]), q = 2$

n = #users
 m = code length in symbols
 q = alphabet size
 ε_1 = Prob[accuse specific innocent]
 η = Prob[not all accused are guilty]
 ε_2 = False Negative prob.

Lower bound

Bias-based code [Tardos 2003, ŠKC 2007]

Alphabet Q

Step 1:

For each position, generate bias vector $\mathbf{p}=(p_\alpha)_{\alpha \in Q}$. $|\mathbf{p}|=1$ $\mathbf{p} \sim F$

Step 2:

For each position and user, draw watermark symbol: $\Pr[\text{symbol } \alpha] = p_\alpha$.

									p_A															
									p_B															
									p_C															
									p_D															
										A														
										C														
										A														
										B														
										B														
										A														
										D														

pirated copy carries watermark y

Step 3:

Find attackers based on X and y

Asymptotically optimal scaling:
code length $\propto c_0^2$

De Tardos code

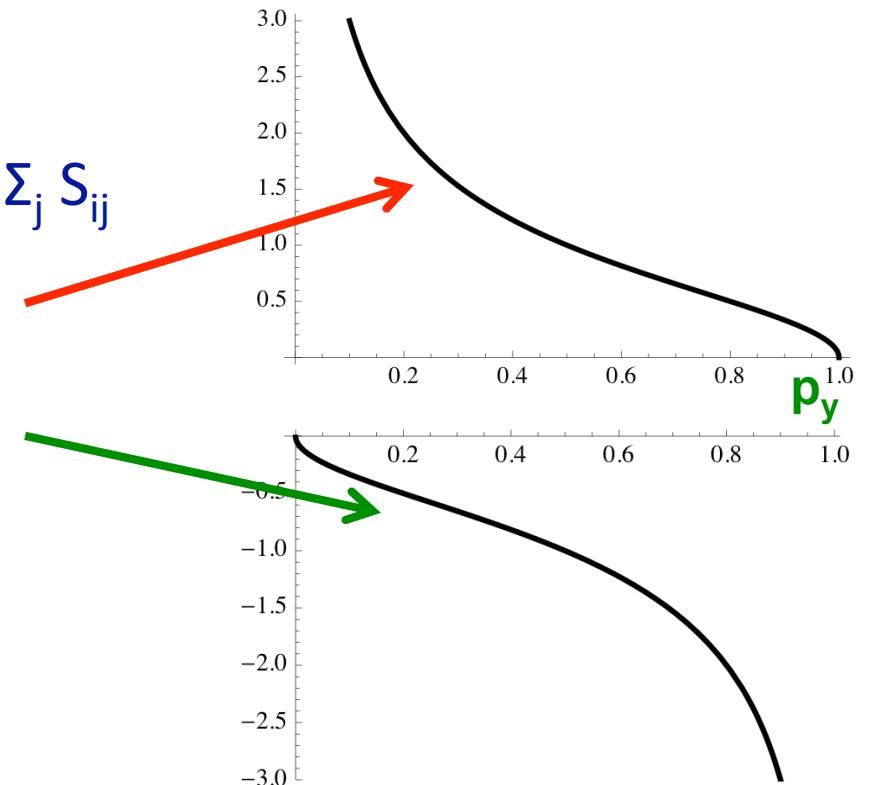
- Allereerste code met $m \propto c^2$ macht en klein alfabet
- In twee opzichten probabilistisch
 - staat kleine kans toe op valse beschuldiging en compleet missen van de aanvallers
 - constructie van de code is gerandomiseerd
- In 2003 verzonnen voor $q=2$,
in 2007 uitgebreid naar algemene q

Tardos code: Het traceren

- Er wordt een "ongeauthoriseerde" kopie gevonden
- Watermerk-detector ziet symbool y_j in segment j
- Reken voor elke klant i de "score" S_i uit
 - som van losse scores per segment: $S_i = \sum_j S_{ij}$

$$S_{ij} = \begin{cases} X_{ij} = y_j : & \sqrt{(1 - p_y) / p_y} \\ X_{ij} \neq y_j : & -\sqrt{p_y / (1 - p_y)} \end{cases}$$

- Klant i is verdacht als score S_i boven een bepaalde grens uitkomt



Tardos code: speciale eigenschappen

Scores van onschuldigen gedragen zich eenvoudig:

- Gemiddelde is nul in elk segment
- Variantie 1 in elk segment

$$E[S_{ij}] = p_y \sqrt{\frac{1 - p_y}{p_y}} - (1 - p_y) \sqrt{\frac{p_y}{1 - p_y}} = 0$$

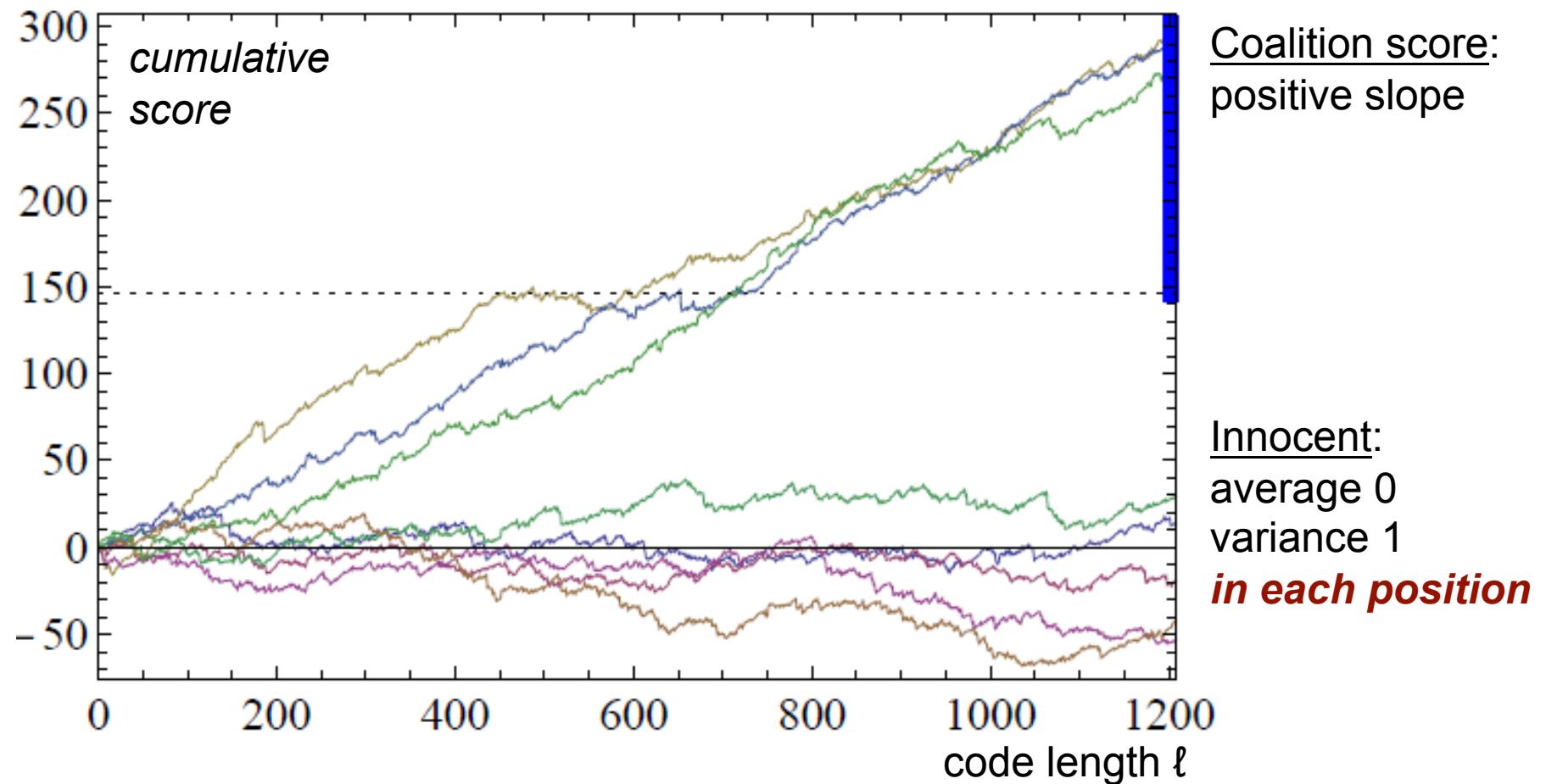
$$E[S_{ij}^2] = p_y \frac{1 - p_y}{p_y} + (1 - p_y) \frac{p_y}{1 - p_y} = (1 - p_y) + p_y = 1$$

- De Tardos score-functie is de enige met deze eigenschap
- Aanvallers hebben geen invloed op gemiddelde en variantie van onschuldigen

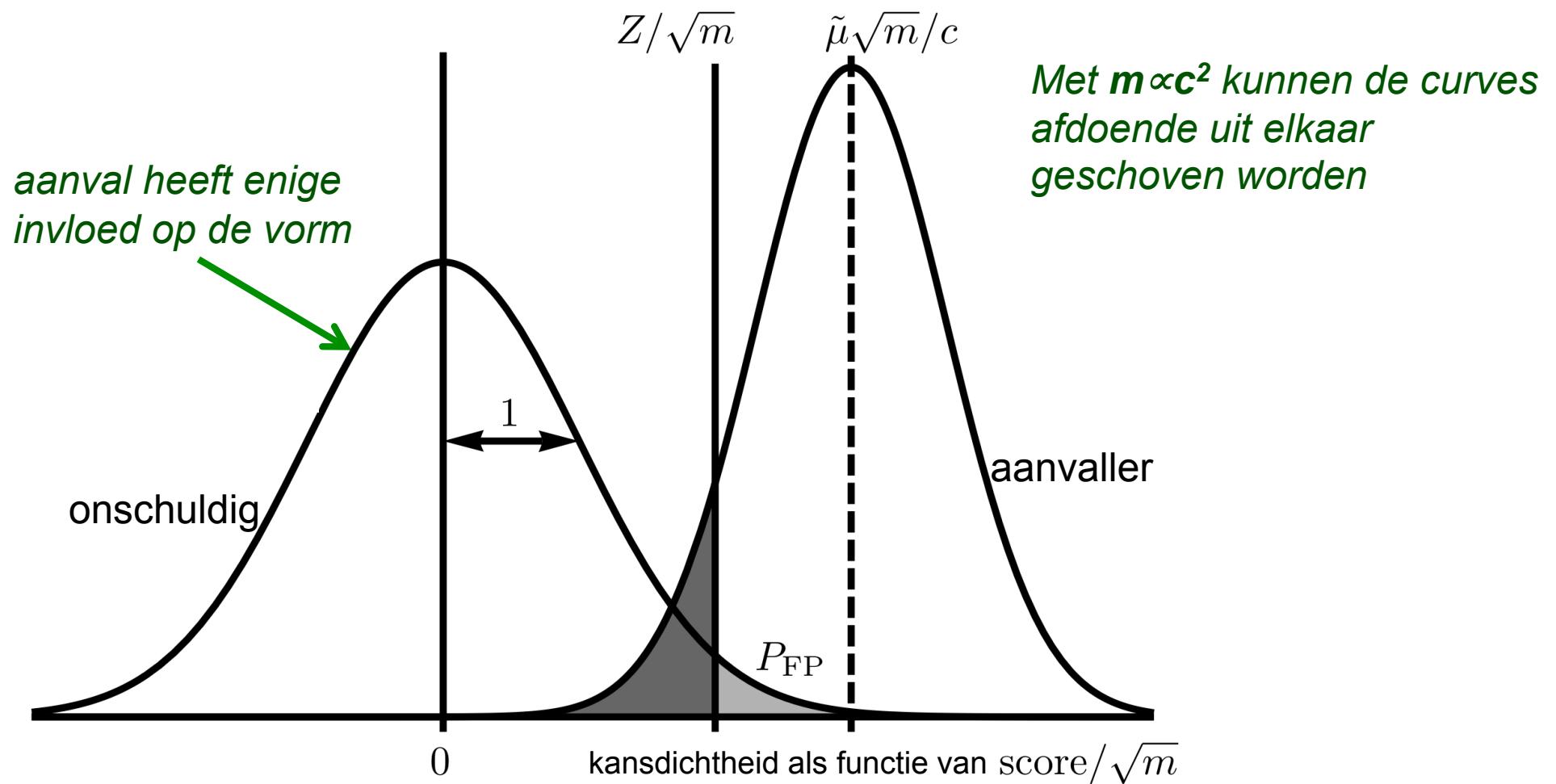
Separating the attackers from the innocents

Binary alphabet; Tardos score function

$$g(x,y,p) = \begin{cases} \sqrt{(1-p_y)/p_y} & \text{if } x = y \\ -\sqrt{p_y/(1-p_y)} & \text{if } x \neq y \end{cases}$$



Gevolgen van de speciale eigenschappen



P_{FP} = kans dat (onschuldige) klant i onterecht wordt beschuldigd

m = #segmenten

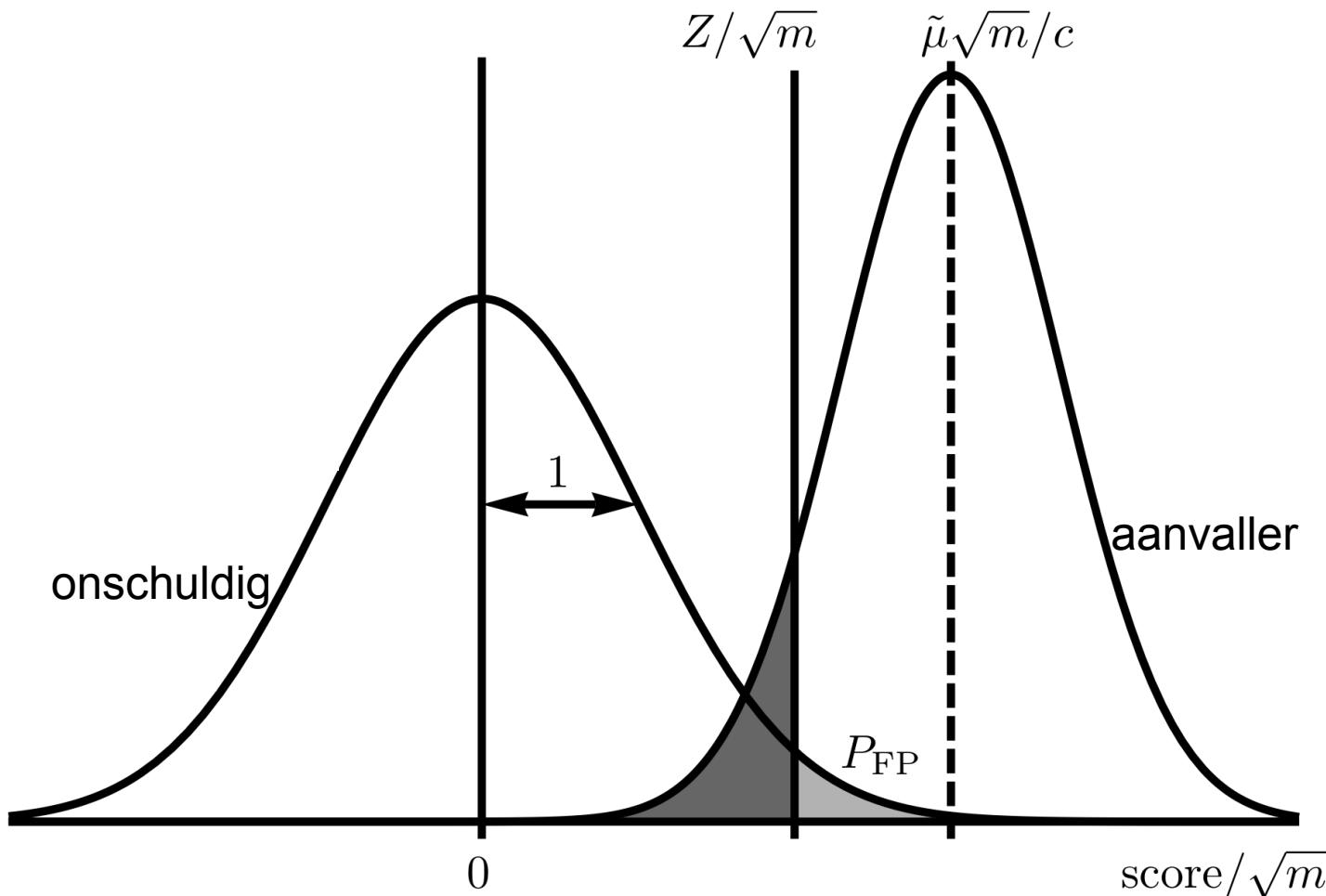
c = #aanvallers

Z = grens

"No framing"

Wat als er meer aanvallers zijn dan geanticipeerd?

- Rechtercurve schuift links van de grens $Z \Rightarrow$ **aanvallers niet gepakt**
- Linkercurve verandert nauwelijks \Rightarrow **geen onschuldigen gepakt**



Collusion channel (Restricted Digit Model)

"Tally" vector \mathbf{m} :

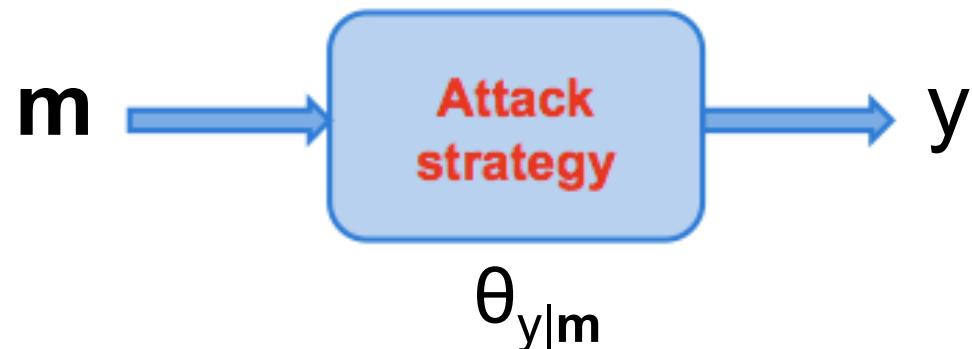
- #colluders = c
- $m_\alpha = \#\alpha$ received by colluders
- $|\mathbf{m}|=c$

Attack:

- same strategy in each position (asymptotically strongest)
- Choose y as a function of \mathbf{m} :
 $\theta_{y|\mathbf{m}} = \text{Prob}[\text{output } y \text{ given } \mathbf{m}]$

pirate code words	A	B		A		C
	C	A		A		A
	A	B		A		B
allowed attack symbols	A	(A)		A		A B C

$\mathbf{m}=(1,2,0)$

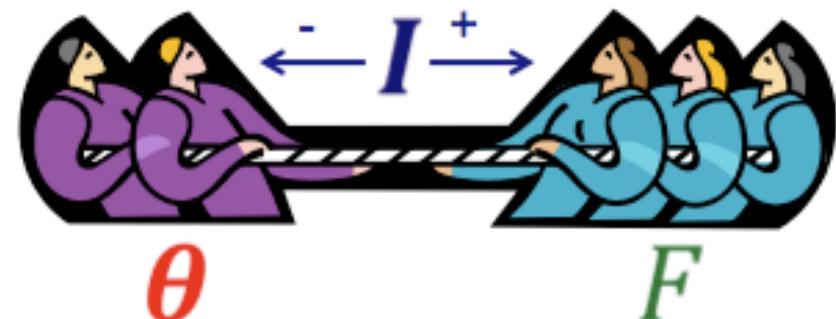


Information theory: capacity

- Collusion attack can be seen as "malicious noise".
- Use techniques from channel coding!
 - How much does Y reveal about \mathbf{M} ?
(\mathbf{M} equivalent to colluder identities)
 - *Mutual information* $I(\mathbf{M}; Y)$

Fingerprinting game:

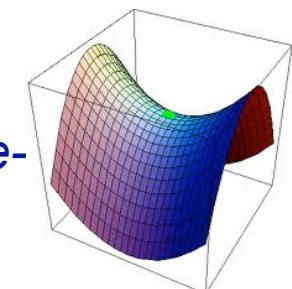
- Pay-off function $I(\mathbf{M}; Y | \mathbf{P}) / c$
- Tracer chooses bias distribution $F(\mathbf{p})$
- Colluders choose strategy θ



Fingerprinting capacity

$$C = \frac{1}{c} \max_F \min_{\theta} I(M; Y | P)$$

saddle-point



Fingerprinting Capacity

Meaning of capacity C:

- Max. achievable code rate
- Asymptotic error rate follows from C and R

$$\text{DEF: } R = \frac{\log_q n}{\ell}$$

n = #users
q = alphabet size
 ℓ = code length

$$P_{\text{err}} \leq q^{-(C-R)\ell}$$

$$\ell_{\text{sufficient}} = \frac{1}{C \ln q} \ln \frac{n}{P_{\text{err}}}$$

Asymptotic capacity and saddlepoint

Asymptotic: $c \rightarrow \infty$

- [Huang & Moulin 2010]. Solution for binary alphabet.

$$F(p_0, p_1) = \frac{1}{\pi} \frac{1}{\sqrt{p_0 p_1}}$$

"arcsine distribution"

$$\theta_{y|m} = \frac{m_y}{c}$$

interleaving attack (pick random attacker)

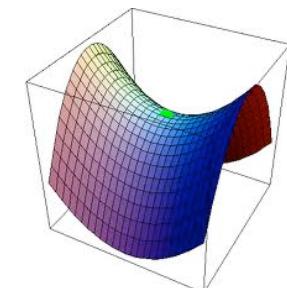
$$C = \frac{1}{c^2 2 \ln 2}$$

- [Boesten & Skoric 2011]. **Capacity** for q-ary alphabet.
Increases with q .
- [Huang & Moulin 2012]. **Saddlepoint** for q-ary case.

$$C = \frac{q-1}{c^2 2 \ln q}$$

$$F(\mathbf{p}) \propto \prod_{\alpha \in Q} p_\alpha^{-1/2} \quad \text{vs. interleaving attack}$$

Dirichlet distribution



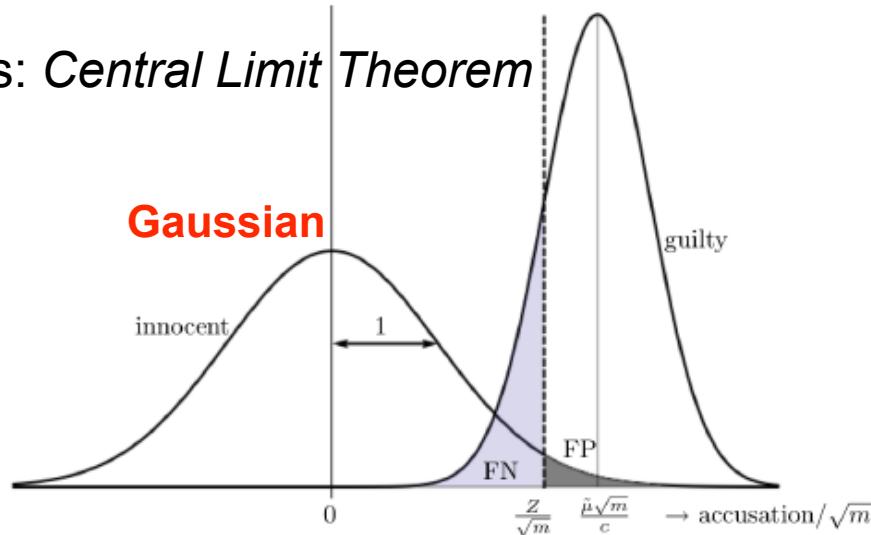
Decoder

- In theory one can achieve sufficient code length
 - if it is known how to trace users!
- "Decoder" algorithm for finding colluders based on X, y, p .
 - **Simple** decoder: each user gets a score
 - **Joint** decoder: triplets of users get a score
- [Tardos 2003, Skoric et al. 2007]: simple decoder
 - far away from capacity (at least factor 2.5)
- [Amiri & Tardos 2009]: Joint decoder for $q=2$
 - capacity-achieving but impractical
- [Huang & Moulin 2012]
 - "simple capacity = joint capacity" (asymptotically)
 - **no recipe**
- [Oosterwijk et al. 2013]
 - **simple decoder that achieves capacity**
 - **q-ary**

$$\ell_{\text{suff}} = \frac{2}{q-1} c^2 \ln \frac{n}{P_{\text{err}}}$$

Finding the optimal score function

Asymptotics: Central Limit Theorem

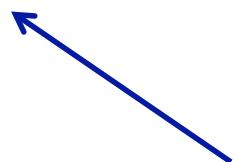


Optimize μ_{guilty} with constraints

- $\mu_{\text{innocent}} = 0$
- $\sigma_{\text{innocent}} = 1$

Euler-Lagrange optimization, with F and θ in saddlepoint

$$L[h, \lambda_1, \lambda_2] = \mu_{\text{guilty}}[h] - \lambda_1 \mu_{\text{inn}}[h] - \frac{1}{2} \lambda_2 (\sigma_{\text{inn}}^2[h] - 1)$$



score function $h(x,y,p)$

x = symbol of user under scrutiny

Optimal score function

$$h(x,y,p) = \frac{1}{\sqrt{q-1}} \left(\frac{\delta_{xy}}{p_y} - 1 \right)$$

- "strongly centered": $E_x[h(x,y,p)] = \sum_x p_x h(x,y,p) = 0$
- If attack=interleaving, then for all F:

$$E_{xy|p}[h_{\text{inn}}^2] = \frac{1}{q-1} E_{y|p}[1/p_y - 1] = 1 \quad (\text{normalized})$$

$$\mu_{\text{coalition}} = \sqrt{q-1}$$

- ✓ Consistency check: indeed a code-rate saddlepoint
 - "ridge" at $\theta=\text{interleaving}$
- σ_{inn} depends on attack strategy

Optimal score function (2)

"Tardos score":
$$g(x, y, p) = \begin{cases} \sqrt{\frac{1 - p_y}{p_y}} & \text{if } x = y \\ -\sqrt{\frac{p_y}{1 - p_y}} & \text{if } x \neq y \end{cases} = \sqrt{\frac{p_y}{1 - p_y}} \left[\frac{\delta_{xy}}{p_y} - 1 \right]$$

Tardos' score is "strongly normalized" version of optimal score!

- Guaranteed to have $E[g^2]=1$ (innocent) for *any* attack and *any* p .
- Demanding strong normalization over-constrains the problem.

Game over?

Did we kill the field?
NO!

Still to be done:

- Validation
 - simulations, provable bounds, etc.
- *Dynamic* traitor tracing
 - is the capacity the same?
 - different conditions, different solutions?
- Finite c ; not just asymptotics
 - find the saddlepoint
 - joint decoder required?
- More realistic attack models
 - Combined Digit Model

