

PERSEUS: A Firefox Plug-in to Fight Botnets ... and Many Other Things

Eddy Deligne & Eric Filiol
{filiol,deligne}@esiea.fr

ESIEA - Laval
Operational Cryptology and Virology Lab $(C + V)^O$

Hack.lu 2009



Introduction

- Many of the computer attacks rely on gathering information on the future targets.
 - Gain technical and/or social intelligence.
 - Steal personal/confidential data.
- This initial step generally relies on listening unprotected data flows (IRC or HTTP flows).
- It is far more simple and far more productive to eavesdrop open (unprotected) traffics than any other more aggressive techniques (intrusion, malware attacks...).
- It is above all passive thus transparent.

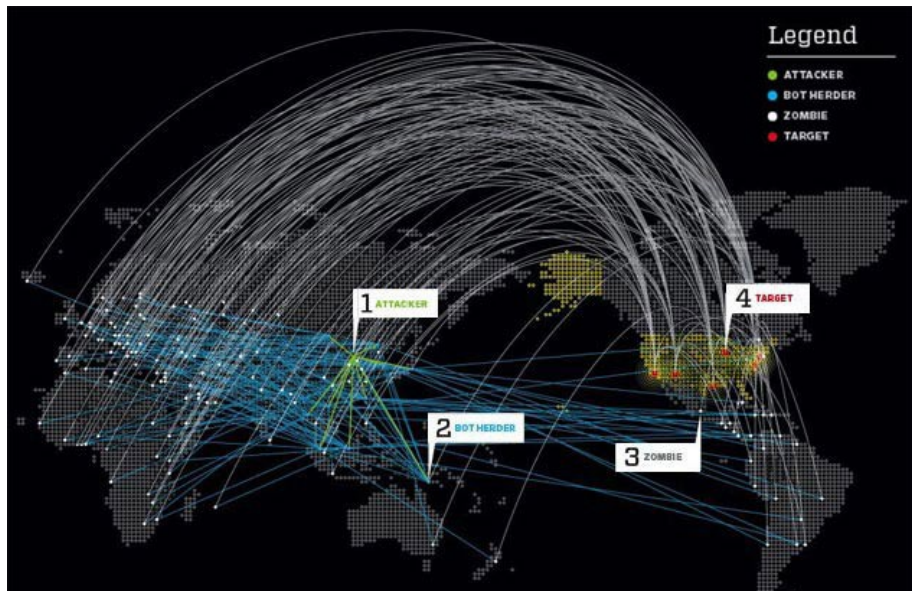


Introduction (2)

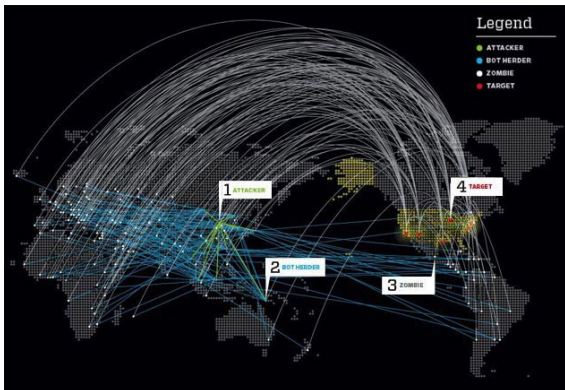
- This issue goes beyond computer attacks and also relates to
 - privacy concerns,
 - misuse of attack techniques by private intelligence companies,
 - abusive gathering of data or citizens' surveillance for commercial purposes (e.g. HADOPI).
- The question is : how to hinder any misuse of eavesdropping techniques while
 - preserving the technical ability of nation states for REAL interior security purposes,
 - be compliant with any existing national laws or regulations (very important issue with respect to transnational data streams).
- The solution is PERSEUS.



Introduction : The Botnet Case



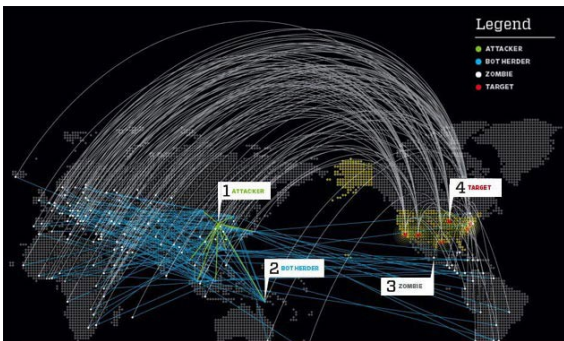
Introduction : The Botnet Case



Obvious solution : use encryption. But can be slow and there might be legal national regulations !



Introduction : The Botnet Case

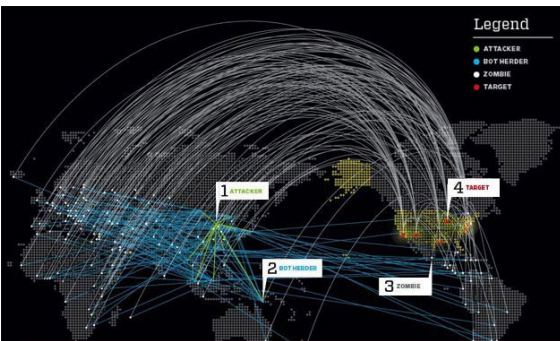


To solve all possible constraints (legal + technical)

How to protect our personal data without infringing any national crypto regulations and preserving the true needs of national security?



Introduction : The Botnet Case



To solve all possible constraints (legal + technical)

How to protect our personal data without infringing any national crypto regulations and preserving the true needs of national security?

Solution

Replace crypto techniques with noisy coding techniques!

Encryption vs Noisy Coding

- “Legal” definition of whether it is cryptography or not, relates in a way or another directly to the following probability

$$P[c_t = m_t + e_t] = P[e_t = 1]$$

where c_t, m_t are the ciphertext and plaintext bits respectively and where e_t can be defined as the noise bit produced by the key and the cryptosystem (at time instant t).

- If $P[e_t = 1] = \frac{1}{2} + \epsilon$ (with ϵ close to 0) then it is cryptography.
 - Otherwise (ϵ significantly different from 0) it is coding theory.
- So the solution is to consider a computationally hard (for the attacker) problem from the coding theory.



Encryption vs Noisy Coding

- Why use noisy encoded data instead of encrypted data ?
 - Encrypted data exhibit a high entropy profile. It is then easy to detect encrypted data.
 - Noisy encoded data exhibit a low entropy profile. The statistical profile is close to that of normal communications (for example cell phone communications).
- This particular statistical profile enabled to bypass any detection by entropy test or any other statistical detection while crypto does not.
 - Some sort of TRANSEC aspect (hide noisy encoded data among other encoded data).
 - Can be applied to bypass any kind of detection based on entropy (malware detection, firewall filtering. . .).



PERSEUS Technology

- Open technology based on punctured convolutional codes with controlled noise.
 - PERSEUS module : Firefox module to protect HTTP streams.
 - First step for practical validation.
- Can be used to protect any kind of protocols (FTP, mail ...).



Plan

- 1 Introduction
- 2 Punctured Convolutional Codes
 - Convolutional Codes
 - Puncturing
 - Convolutional Decoding
 - Convolutional Code Reconstruction
- 3 PERSEUS Description
 - General Description
 - Noise Generation
 - Parameter Management
 - Scenarii
- 4 Experimental Results and Improvements
 - Experimental Results
 - Implementation
- 5 Conclusion



What are Error-Correcting Codes ?

- Mathematical tool introduced by C. E. Shannon (1948) to correct the effect of “natural” noise on a communication channel (Shannon’s 2nd Theorem).

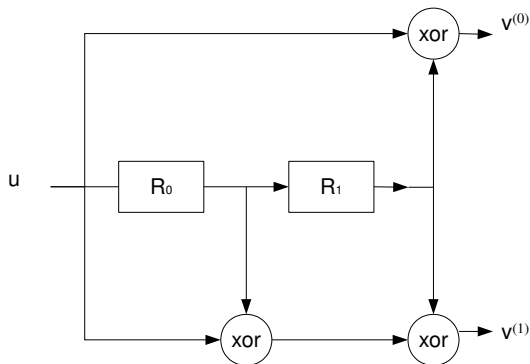
$$m_t \rightarrow m_t \oplus e_t = \hat{m}_t \rightarrow \text{decoding } m'_t$$

- Consists in adding a limited amount of redundancy bits to the message in order to recover from the noise.
 - There always exist codes such that $P[m'_t = m_t]$ tends towards 0.
- Different issues :
 - Efficiency of the decoding (legitimate users).
 - Reconconstruction of unknown encoders (attackers' concern).
- Large variety of codes. We consider one the fastest family for stream coding : convolutional codes.



Convolutional Code

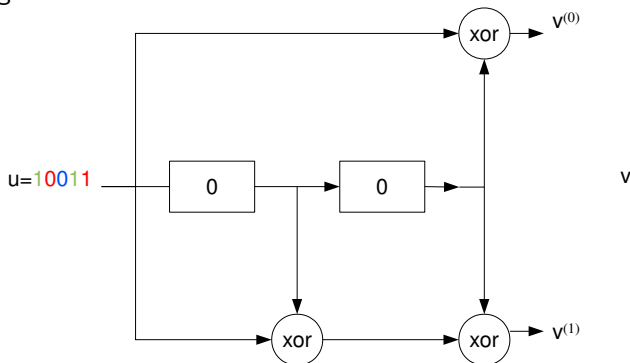
Let us consider a convolutional code \mathcal{C} of rate $\frac{1}{2}$ with a memory size of $M = 2$.



Convolutional Code

Let us consider a convolutional code \mathcal{C} of rate $\frac{1}{2}$ with a memory size of $M = 2$.

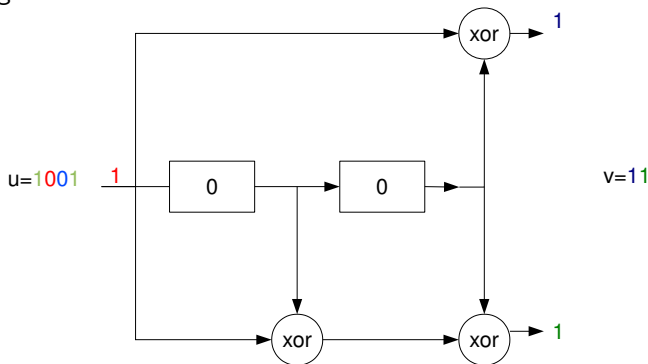
A message $u = 10011$



Convolutional Code

Let us consider a convolutional code \mathcal{C} of rate $\frac{1}{2}$ with a memory size of $M = 2$.

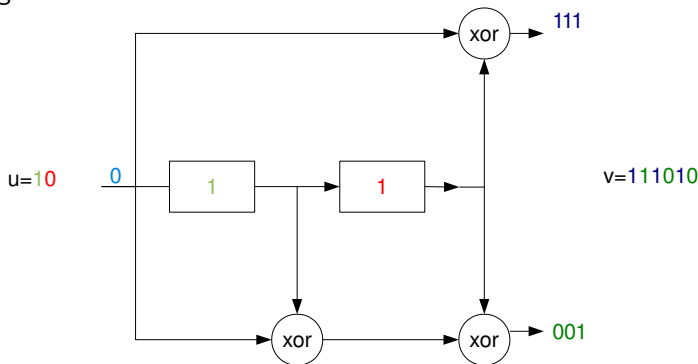
A message $u = 10011$



Convolutional Code

Let us consider a convolutional code \mathcal{C} of rate $\frac{1}{2}$ with a memory size of $M = 2$.

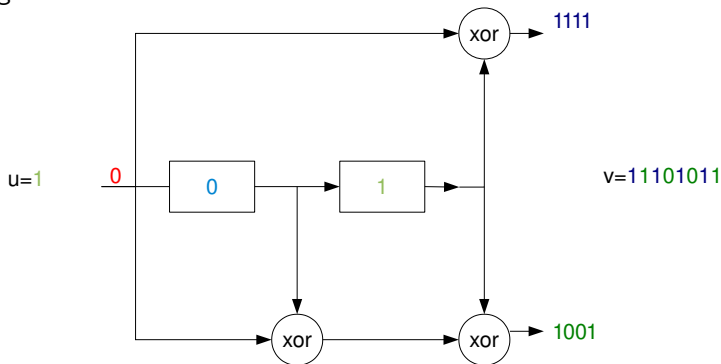
A message $u = 10011$



Convolutional Code

Let us consider a convolutional code \mathcal{C} of rate $\frac{1}{2}$ with a memory size of $M = 2$.

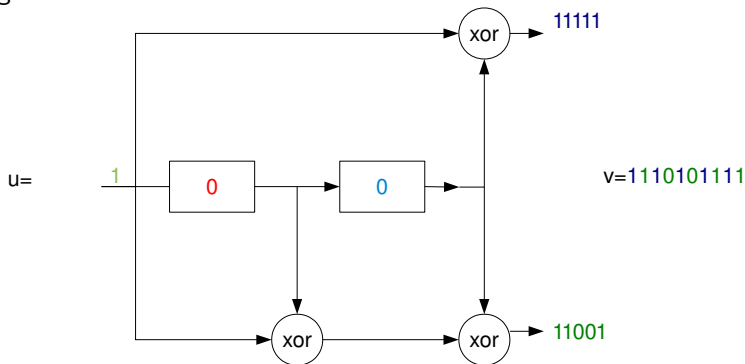
A message $u = 10011$



Convolutional Code

Let us consider a convolutional code \mathcal{C} of rate $\frac{1}{2}$ with a memory size of $M = 2$.

A message $u = 10011$



Presentation

A convolutional code is defined by

- a rate : $\frac{k}{n}$
- a memory size (or constraint length) $K = M + 1$.

Notation

(n, k, K) -convolutional code



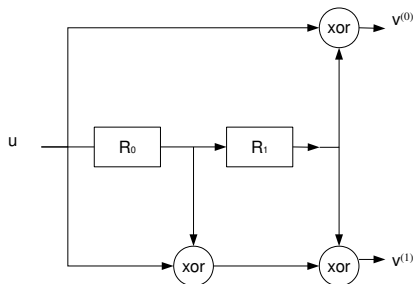
Alternative View

A Convolutional Code

k registers with n polynomials operating on each register.

$n \times k$ polynomials for a (n, k, K) -convolutional code.

The degree of polynomials will be equal to $K - 1$.



$\mathcal{C} : (2, 1, 3)$ -convolutional code

$$v_0 : 1 + x^2$$

$$v_1 : 1 + x + x^2$$



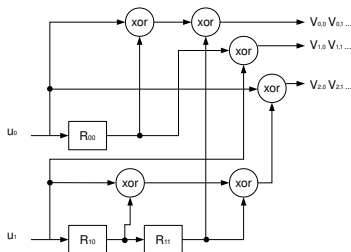
Alternative View

A Convolutional Code

k registers with n polynomials operating on each register.

$n \times k$ polynomials for a (n, k, K) -convolutional code.

The degree of polynomials will be equal to $K - 1$.



$\mathcal{C} : (3, 2, 3)$ -convolutional code

$$v_{0,0} : 1 + x$$

$$v_{0,1} : x^2$$

$$v_{1,0} : x$$

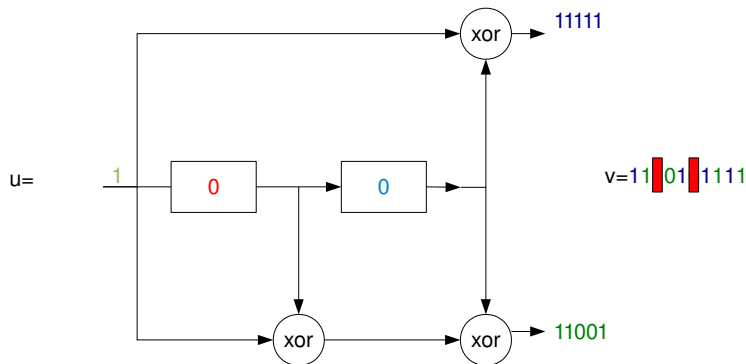
$$v_{1,1} : 1$$

$$v_{2,0} : 1$$

$$v_{2,1} : 1 + x + x^2$$



Puncturing

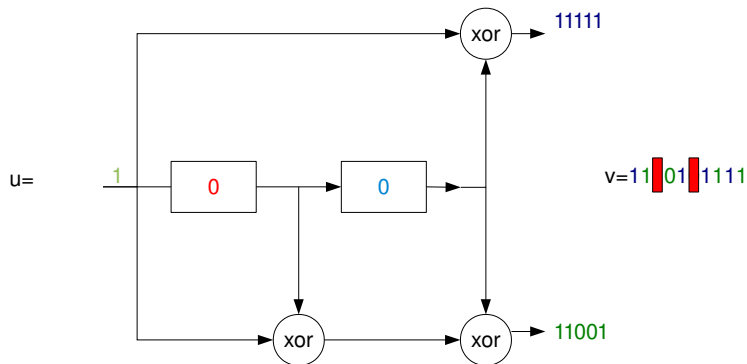


Puncturing pattern

P : a $J \times n$ matrix of weight I .



Puncturing



Puncturing pattern

P : a $J \times n$ matrix of weight I .



Example

Let P be the puncturing pattern given by :

$$P = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$$

and let v be the $(2, 1, 3)$ encoder output sequence :

$$v = \left(\begin{array}{cc|cc|c} 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 \end{array} \right)$$

$$\left(\begin{array}{cc|cc|c} 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 \end{array} \right) \Rightarrow 11010111$$



Why puncturing ?

- 1 Save bandwidth (reduce the redundancy added).
- 2 Produce an equivalent (non punctured) convolutional code which is stronger for our purposes (see further).



Why puncturing?

- 1 Save bandwidth (reduce the redundancy added).
- 2 Produce an equivalent (non punctured) convolutional code which is stronger for our purposes (see further).

Equivalent (non punctured) convolutional code

A (n, k, K) -convolutional code and a $J \times n$ puncturing matrix P of weight I .

$\Rightarrow (I, kJ, K)$ -convolutional code

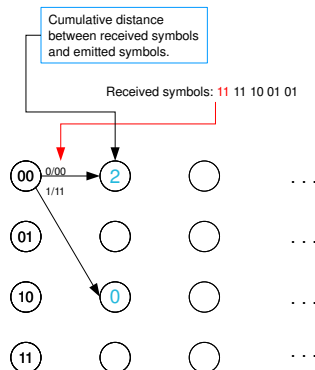


Viterbi Algorithm



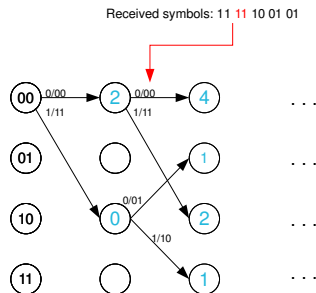
Viterbi Algorithm

Lattice construction



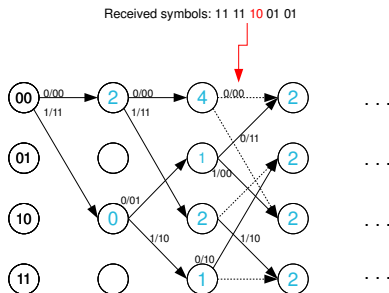
Viterbi Algorithm

Lattice construction



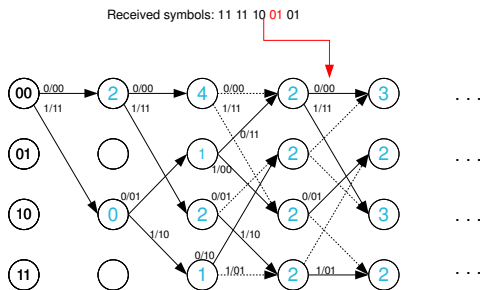
Viterbi Algorithm

Lattice construction



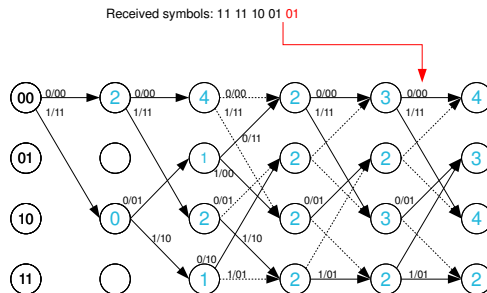
Viterbi Algorithm

Lattice construction



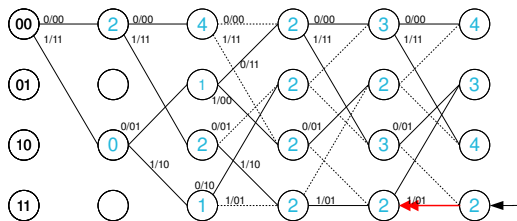
Viterbi Algorithm

Lattice construction



Viterbi Algorithm

Backtracking

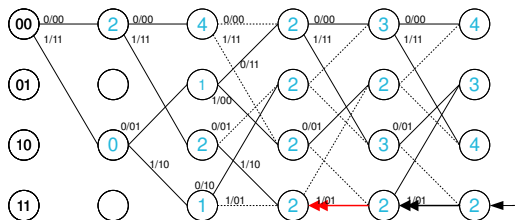


Decoded symbols: **1**



Viterbi Algorithm

Backtracking

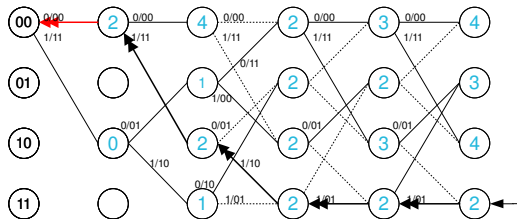


Decoded symbols: 11



Viterbi Algorithm

Backtracking



Decoded symbols: 01111

- Decoding has exponential complexity in K .
- When dealing with puncturing, replace removed bits with zeroes.



Convolutional Code Reconstruction (Filiol 1997 - Barbier 2007)

Aim : recovering all the parameters of an unknown encoder from the encoded data only, to be able to decode data afterwards.

Puncturing effect

Let us consider a (n, k, K) -convolutional code and a $J \times n$ puncturing matrix P of weight I :

$\Rightarrow (I, kJ, K)$ -convolutional code

Reconstruction has the following complexity

$$\mathcal{O}(\alpha \times n^5 \times K^4) \Rightarrow \mathcal{O}(\alpha \times I^5 \times K^4)$$

α : grows exponentially with p , the noise probability



Convolutional Code Reconstruction (Filiol 1997 - Barbier 2007)

Aim : recovering all the parameters of an unknown encoder from the encoded data only, to be able to decode data afterwards.

Puncturing effect

Let us consider a (n, k, K) -convolutional code and a $J \times n$ puncturing matrix P of weight I :

$\Rightarrow (I, kJ, K)$ -convolutional code

Reconstruction has the following complexity

$$\mathcal{O}(\alpha \times n^5 \times K^4) \Rightarrow \mathcal{O}(\alpha \times I^5 \times K^4)$$

α : grows exponentially with p , the noise probability



The Noise Impact

The probability to successfully reconstruct a code exponentially decreases with p . If $p > 10\%$ \Rightarrow online reconstruction is impossible ; offline reconstruction is computationally very hard.

Encoder	Reconstruction Time ($p = 10^{-2}$)	Reconstruction Time ($p = 2 \cdot 10^{-2}$)
(4, 3, 8)	7 min 12 sec	Failure
(4, 3, 9)	6 min 16 sec	Failure

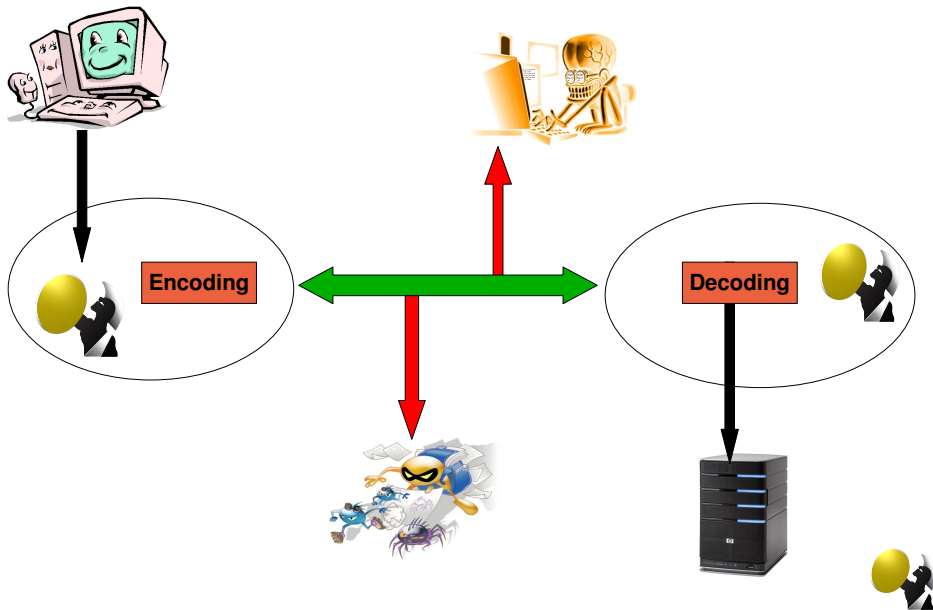
TAB.: Examples of reconstruction times (Pentium IV 2.0 Ghz) for two levels of noise



Plan

- 1 Introduction
- 2 Punctured Convolutional Codes
 - Convolutional Codes
 - Puncturing
 - Convolutional Decoding
 - Convolutional Code Reconstruction
- 3 PERSEUS Description
 - General Description
 - Noise Generation
 - Parameter Management
 - Scenarii
- 4 Experimental Results and Improvements
 - Experimental Results
 - Implementation
- 5 Conclusion





General Principle

- The attacker (botnet client, attacker, eavesdropper...) must face a computationally untractable problem.
- For every different communication, a random encoder is used to encode the HTTP traffic.
- To make the reconstruction computationally untractable we add a secret-based deterministic noise.
- The legitimate knows the exact noise bit indices and can remove the noise to perform a noiseless sequence decoding.



Added Deterministic Noise

To greatly increase the security of data stream as well as the encoder used, noise must be added to the encoded data before transmission.

Problem

Viterbi decoding is easy as long as $p < 5\%$

Encoder reconstruction is impossible as long as $p > 10\%$ (inline mode)



Added Deterministic Noise

To greatly increase the security of data stream as well as the encoder used, noise must be added to the encoded data before transmission.

Problem

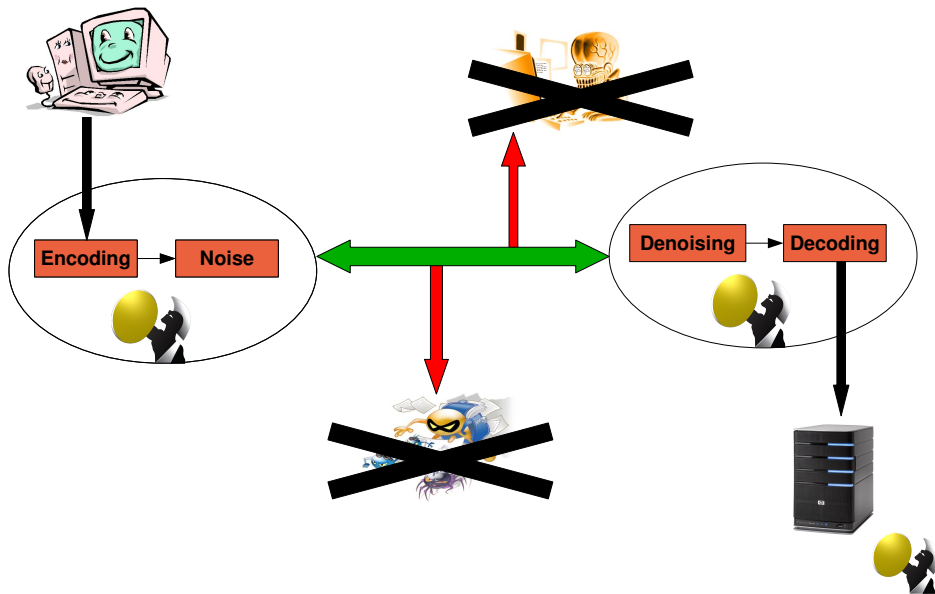
Viterbi decoding is easy as long as $p < 5\%$

Encoder reconstruction is impossible as long as $p > 10\%$ (inline mode)

Solution

We add a deterministic noise with $p \approx 30\%$





What we need

- 1 A noise with probability around 30%.
- 2 A deterministic noise that must be controlled (except for the attackers).

Principle

An array of ten 64-bit integers

Each integer contains approximatively 30% of 1s.

We choose one of the ten integers at random



What we need

- 1 A noise with probability around 30%.
- 2 A deterministic noise that must be controlled (except for the attackers).

Principle

An array of ten 64-bit integers

Each integer contains approximatively 30% of 1s.

We choose one of the ten integers at random



Shared Parameters

- 1 $0 \leq j < 10$ (random) enables to select an integer in the array.
- 2 X_0 a secret random 64-bit integer to initialize the pseudo-random noise generator.

We then have $\approx 2^{67}$ different noise generators.



What are the parameters

For every session :

- ① $2 < n \leq 12$
- ② $1 < k < n$
- ③ $15 < N \leq 30$
- ④ $n \times k$ polynomials of degree $N - 1$
- ⑤ A $J \times n$ -matrix P and of weight $(n \times J) - (J - 1)$
- ⑥ X_0 a 64-bit (secret) integer
- ⑦ j a random integer such that $0 \leq j < 10$



What are the parameters

For every session :

- ① $2 < n \leq 12$
- ② $1 < k < n$
- ③ $15 < N \leq 30$
- ④ $n \times k$ polynomials of degree $N - 1$
- ⑤ A $J \times n$ -matrix P and of weight $(n \times J) - (J - 1)$
- ⑥ X_0 a 64-bit (secret) integer
- ⑦ j a random integer such that $0 \leq j < 10$



What are the parameters

For every session :

- ① $2 < n \leq 12$
- ② $1 < k < n$
- ③ $15 < N \leq 30$
- ④ $n \times k$ polynomials of degree $N - 1$
- ⑤ A $J \times n$ -matrix P and of weight $(n \times J) - (J - 1)$
- ⑥ X_0 a 64-bit (secret) integer
- ⑦ j a random integer such that $0 \leq j < 10$

Sending to the server

Parameters are sent through an initial HTTPS session. The parameter lifetime is limited (either the session has limited length or it is reinitialized with new parameters).

To communicate with the server, PERSEUS uses HTTP headers.

- **check** : enables to ask the server whether the PERSEUS module is present or not.
- **check-ack** : server answer to the check request (means that PERSEUS is present).
- **known** : ask whether the server knows the parameters of the encoder or not. This header is used along with the **Perseus.id :XYZ** where XYZ is a random integer which enables to identify the encoder parameters.
- **known-yes** : the server acknowledges the parameters. Use the **Perseus.id :XYZ** as well.
- **known-no** : the server does not know the parameters.
- **init** : means that the request contains the encoder parameters ("POST" method). Uses the **Perseus.id :XYZ** field as well.
- **init-ack** : answer to the "init" request if everything is OK on the server's side.
- **init-nack** : answer to the "init" request if an error occurred on the server's side.
- **PCC** : specifies that data in the request body are encoded (with a PCC). Uses the **Perseus.id :XYZ** field as well.



Scenarii

PERSEUS

protects GET and POST methods.

Several scenarii have been considered :

- ❶ First (initial) communication between the server and the client.
 - The server does not use the PERSEUS module.
 - The server uses the PERSEUS module.
- ❷ A communication link already exists (or existed).
 - The server already knows the encoder parameters.
 - The server does no longer know the encoder parameters.



Scenarii

PERSEUS

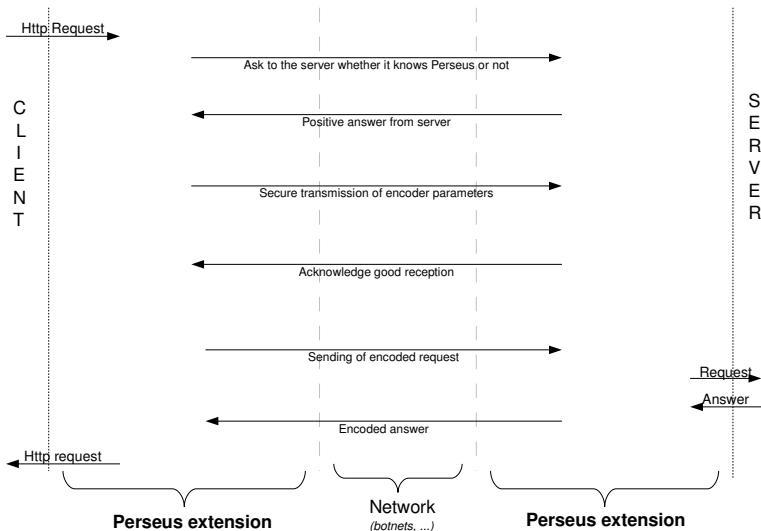
protects GET and POST methods.

Several scenarii have been considered :

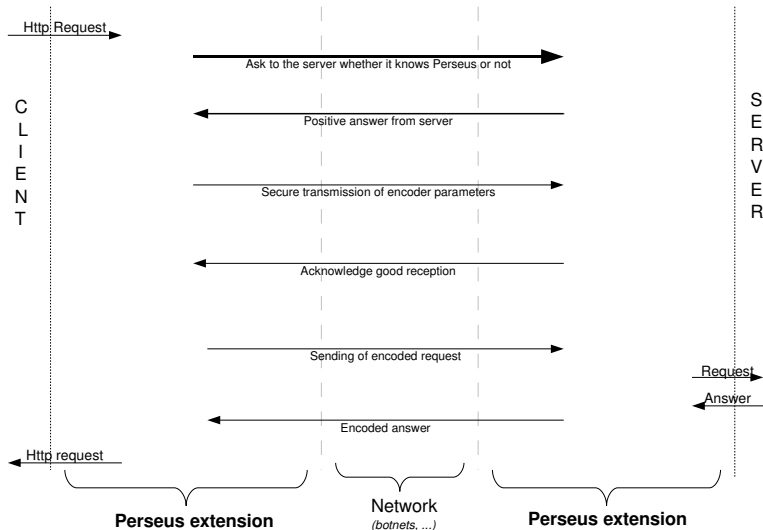
- ① First (initial) communication between the server and the client.
 - The server does not use the PERSEUS module.
 - The server uses the PERSEUS module.
- ② A communication link already exists (or existed).
 - The server already knows the encoder parameters.
 - The server does no longer know the encoder parameters.



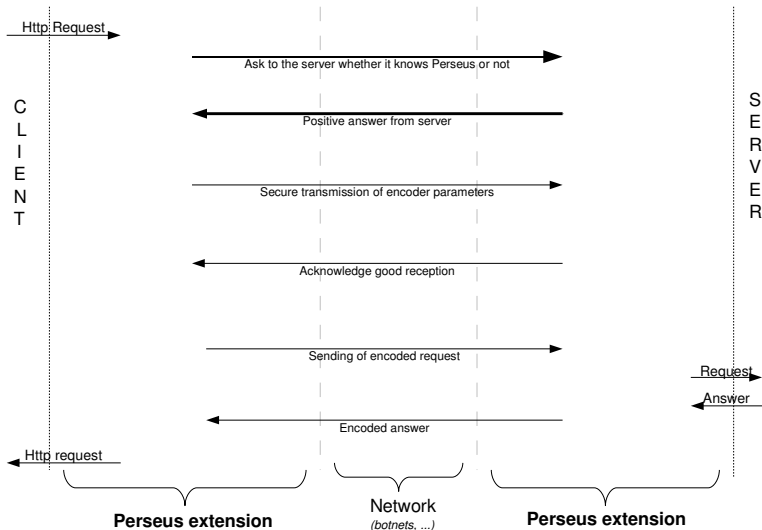
Scenarii (2)



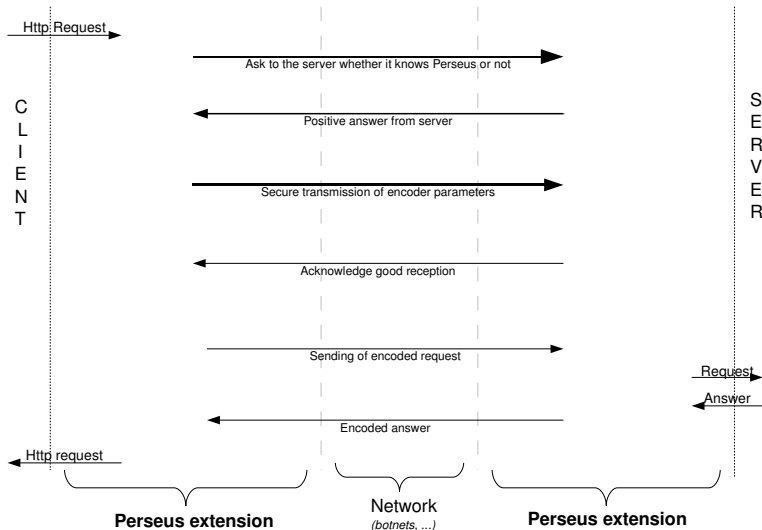
Scenarii (2)



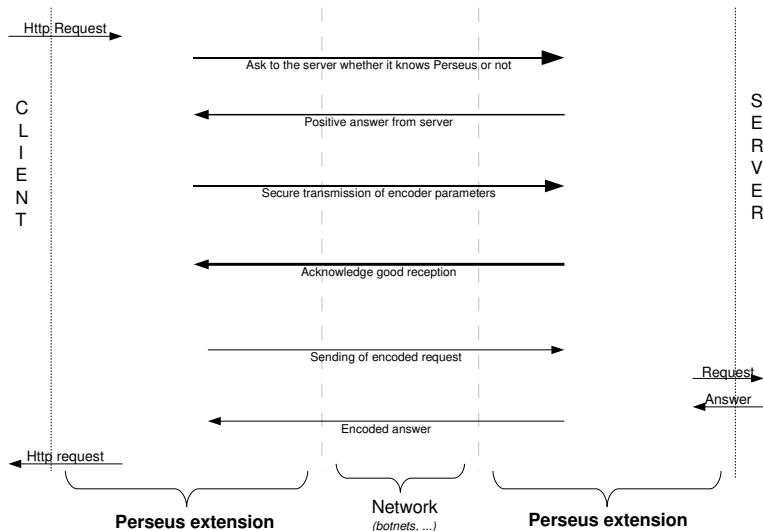
Scenarii (2)



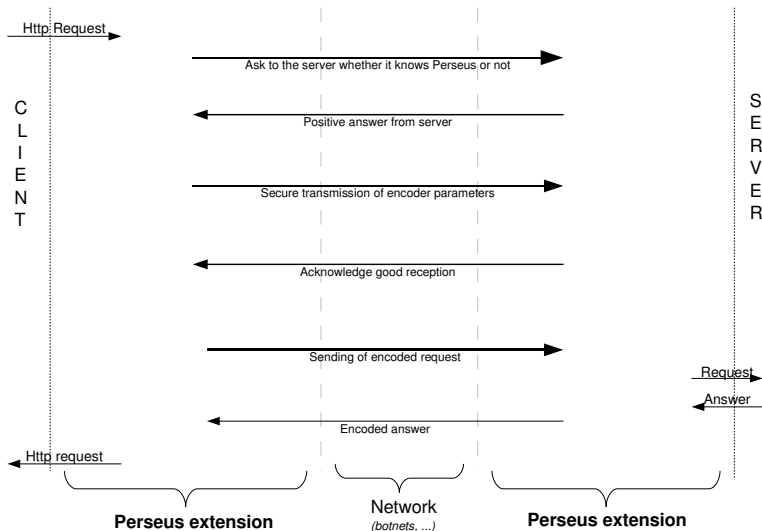
Scenarii (2)



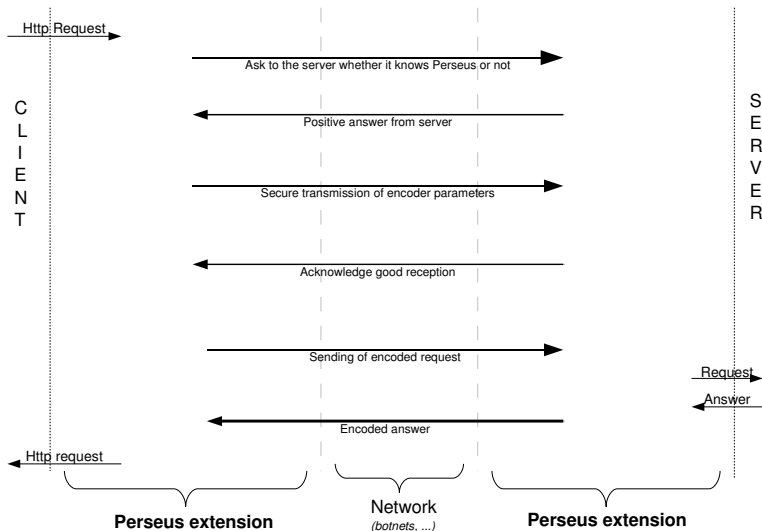
Scenarii (2)



Scenarii (2)



Scenarii (2)



Encoder	Machine	Encoding Time (μs)		Decoding Time (μs)		Data Exchange (s)
		Noise	Total	Noise	Total	
Size of data sent to the server : 231 octets (HTTPS :0, 050s)						
[4, 3, 18, 3, 2]	Client	3460	4463	4370	18362	0,076138
	Server	4429	5591	3167	11520	
[2, 1, 28, 8, 7]	Client	3151	3983	4625	10296	0,075423
	Server	4437	5411	3188	5087	
[4, 1, 22, 9, 8]	Client	8713	10245	12513	20348	0,103056
	Server	12314	13923	8764	11759	
[11, 1, 18, 11, 10]	Client	28892	33535	40169	145110	0,352391
	Server	39911	43751	28674	89893	
[8, 3, 21, 10, 9]	Client	6737	8897	9375	34189	0,117641
	Server	9463	12091	6771	22753	
Size of data sent to the server : 1 Ko (HTTPS :0, 050s)						
[6, 3, 17, 10, 9]	Client	27717	36701	29672	92395	0,289190
	Server	29975	37262	28685	80740	
[3, 2, 23, 7, 6]	Client	17803	23420	13984	87633	0,252633
	Server	17061	21851	17602	88181	

TAB.: Processing and communication times with PERSEUS.

- The noise generator represents about 75% of the processing time.
- Current development of an optimized noise generator is bound to significantly reduce it.



Implementation

PERSEUS

Written in C++.

Under the triple MPL/GPL/LGPL licence.

Compliant with the Mozilla Function programming standards.

Source, documentation and binaries (XPI file) available at

<http://www.mozdev.org/source/browse/perseus>.

Feedback, bug report and contributions requested (contact authors).

Project encouraged by the French DoD.

Future Work and Improvements

Management of UTF-16 encoding.

Management of FTP and FILE protocols.

The PERSEUS module will include a server as well.



Implementation

PERSEUS

Written in C++.

Under the triple MPL/GPL/LGPL licence.

Compliant with the Mozilla Function programming standards.

Source, documentation and binaries (XPI file) available at

<http://www.mozdev.org/source/browse/perseus>.

Feedback, bug report and contributions requested (contact authors).

Project encouraged by the French DoD.

Future Work and Improvements

Management of UTF-16 encoding.

Management of FTP and FILE protocols.

The PERSEUS module will include a server as well.



Demos

PERSEUS Demo 1

Using PERSEUS.

PERSEUS Demo 2

Wireshark analysis without and with PERSEUS.



Demos

PERSEUS Demo 1

Using PERSEUS.

PERSEUS Demo 2

Wireshark analysis without and with PERSEUS.



Conclusion

- PERSEUS gives an elegant answer to a critical issue :
 - How to protect against HTTP traffic eavesdropping by botnet clients...
 - ... and abuses against citizens' privacy fundamental rights...
 - ... without crypto...
 - while preserving TRUE, LEGITIMATE ability for national security enforcement (internal and external) ?
- PERSEUS replaces cryptography by coding theory techniques.
- Only agencies with a strong enough computer power can eavesdrop traffic in an acceptable amount of time... provided that there are not too many communications to deal with at the same time.
- Interesting issue : security is not a technical matter only. Legal aspects are a critical part.





Questions ?

