# How to Encrypt with the LPN Problem

Henri Gilbert, Matt Robshaw, and Yannick Seurin

Orange Labs

ICALP 2008 – July 9, 2008

&ftgroup

orange™

# the context

- the authentication protocol $HB^+$ by Juels and Weis [JW05] recently renewed interest in cryptographic protocols based on the LPN (*Learning Parity with Noise*) problem, the problem of learning an unknown vector $x$ given noisy versions of its scalar product $a \cdot x$ with random vectors $a$

- this problem seems promising to obtain efficient protocols since it implies only basic operations on GF(2)

- in this work, we present a probabilistic symmetric encryption scheme, named LPN-C, whose security against chosen-plaintext attacks can be proved assuming the hardness of the LPN problem

# outline

- the LPN problem: a brief survey

- description and analysis of the encryption scheme LPN-C

- concrete parameters, practical optimizations

- conclusion & open problems

# the LPN problem

> Given $q$ noisy samples $(\mathbf{a}_i, \mathbf{a}_i \cdot \mathbf{x} \oplus \nu_i)$, where $\mathbf{x}$ is a secret $k$-bit vector, the $\mathbf{a}_i$'s are random, and $\Pr[\nu_i = 1] = \eta$, find $\mathbf{x}$.

- similar to the problem of decoding a random linear code (NP-complete)

- best solving algorithms require $T, q = 2^{\Theta(\frac{k}{\log k})}$: Blum, Kalai, Wasserman [BKW03], Levieil, Fouque [LF06]

- a variant by Lyubashevsky [L05] requires $q = \mathcal{O}(k^{1+\epsilon})$ but $T = 2^{\mathcal{O}(\frac{k}{\log \log k})}$

- numerical examples:

  ▸ for $k = 512$ and $\eta = 0.25$, LF requires $T, q \simeq 2^{89}$
  ▸ for $k = 768$ and $\eta = 0.01$, LF requires $T, q \simeq 2^{74}$

# previous schemes based on LPN

- PRNG by Blum et al. [BFKL93]

- public-key encryption scheme by Regev [R05] based on the LWE problem, the generalization of LPN to GF($p$), $p > 2$

- the HB family of authentication protocols:

  - HB [HB01]
  - HB$^{+}$ [JW05]
  - HB$^{++}$ [BCD06]
  - HB$^{*}$ [DK07]
  - HB$^{\#}$ [GRS08]
  - Trusted-HB [BC07]
  - PUF-HB [HS08]

# description of LPN-C

- **public components:** a (linear) error-correcting code $C : \{0,1\}^r \to \{0,1\}^m$ of parameters $[m, r, d]$ and the corresponding decoding algorithm $C^{-1}$

- **secret key:** a $k \times m$ binary matrix $M$

- **encryption:**

  - $r$-bit plaintext $x$, encode it to $C(x)$

  - draw a random $k$-bit vector $a$ and a noise vector $\nu$ where $\Pr[\nu[i] = 1] = \eta$

  - ciphertext $(a, y)$, where $y = C(x) \oplus a \cdot M \oplus \nu$

- **decryption:** on input $(a, y)$, compute $y \oplus a \cdot M$ and decode the resulting value, or output $\perp$ if unable to decode

# security intuition

- $y = C(x) \oplus a \cdot M \oplus \nu$

- in a chosen-plaintext attack, the adversary only learns $a_i \cdot M \oplus \nu_i$ for random vectors $a_i$

- hardness of the LPN problem implies that the adversary cannot guess $a \cdot M$ for a new random $a$ better than with *a priori* probability ("MHB puzzle" [GRS08]), hence will have no information on a challenge ciphertext $(a, C(x) \oplus a \cdot M \oplus \nu)$

# decryption failures

- decryption failures happen when $\mathsf{Hwt}(\boldsymbol{v}) > t$, where $t = \left\lfloor \frac{d-1}{2} \right\rfloor$ is the correction capacity of the code

- when the noise is randomly drawn,

$$P_{\mathsf{DF}} = \sum_{i=t+1}^{m} \binom{m}{i} \eta^i (1-\eta)^{m-i}$$

  is negligible for $\eta m < t$

- for eliminating decryption failures, the Hamming weight of the noise vector can be tested before being used and regenerated when $\mathsf{Hwt}(\boldsymbol{v}) > t$, but this may impact the security proof

# quasi-homomorphic encryption

- the scheme enjoys some kind of "homomorphism" property

- given two plaintexts

$$(\mathbf{a}, \mathbf{y}) = (\mathbf{a}, C(\mathbf{x}) \oplus \mathbf{a} \cdot M \oplus \boldsymbol{\nu})$$
$$(\mathbf{a}', \mathbf{y}') = (\mathbf{a}', C(\mathbf{x}') \oplus \mathbf{a}' \cdot M \oplus \boldsymbol{\nu}'),$$

one has:

$$\mathbf{y} \oplus \mathbf{y}' = C(\mathbf{x} \oplus \mathbf{x}') \oplus (\mathbf{a} \oplus \mathbf{a}') \cdot M \oplus (\boldsymbol{\nu} \oplus \boldsymbol{\nu}')$$

so that $(\mathbf{a} \oplus \mathbf{a}', \mathbf{y} \oplus \mathbf{y}')$ is a valid ciphertext for $\mathbf{x} \oplus \mathbf{x}'$ if $\mathsf{Hwt}(\boldsymbol{\nu} \oplus \boldsymbol{\nu}') \leqslant t$

- $\boldsymbol{\nu} \oplus \boldsymbol{\nu}'$ is a noise vector with noise parameter $\eta' = 2\eta(1-\eta)$; if $\eta'm < t$, the homomorphism property holds with overwhelming probability

# security notions

- security goals: indistinguishability (IND) and non-malleability (NM)

- adversaries run in two phases; at the end of the first phase they output a distribution on the plaintexts and receive a ciphertext challenge

- they are denoted $P X$-$C Y$ according to the oracles (P for encryption, C for decryption) they can access

  ▸ $X, Y = 0$: the adversary can never access the oracle

  ▸ $X, Y = 1$: the adversary can only access the oracle during phase 1 (non-adaptive)

  ▸ $X, Y = 2$: the adversary can access the oracle during phases 1 and 2, *i.e.* after having seen the challenge ciphertext (adaptive)

# security notions

- relations between different types of attacks have been studied by Katz and Yung [KY06]:

- IND-P1-C Y $\Leftrightarrow$ IND-P2-C Y and NM-P1-C Y $\Leftrightarrow$ NM-P2-C Y

- IND-P2-C2 $\Leftrightarrow$ NM-P2-C2

# security proof: a useful lemma

- notations:

  ▸ $U_{k+1}$ will be the oracle returning uniformly random $(k+1)$-bit strings

  ▸ $\Pi_{s,\eta}$ will be the oracle returning the $(k+1)$-bit string $(a, a \cdot s \oplus v)$, where $a$ is uniformly random and $\Pr[v = 1] = \eta$

- we have the following decision-to-search lemma (Regev [R05], Katz and Shin [KS06]):

  **lemma:** if there is an efficient oracle adversary distinguishing between the two oracles $U_{k+1}$ and $\Pi_{s,\eta}$, then there is an efficient adversary solving the LPN problem

# IND-P2-C0 security proof

- P2-C0 adversary $\mathcal{A}$ breaking the indistinguishability of the scheme

- we use it to distinguish between $U_{k+1}$ and $\Pi_{s,\eta}$ as follows:

  - ▸ draw a random $j \in [1..m]$ and a random $k \times (m - j)$ binary matrix $M'$

  - ▸ use the following method to encrypt:

    - • get a sample $(\mathbf{a}, z)$ from the oracle $\mathcal{O}$

    - • form the $m$-bit masking vector $\mathbf{b} = \mathbf{r}\|z\|(\mathbf{a} \cdot M' \oplus \boldsymbol{\nu})$ where $\mathbf{r}$ is a random $(j - 1)$-bit string and $\boldsymbol{\nu}$ an $(m - j)$-bit noise vector

    - • return the ciphertext $(\mathbf{a}, C(\mathbf{x}) \oplus \mathbf{b})$

  - ▸ play the indistinguishability game with $\mathcal{A}$; if $\mathcal{A}$ distinguishes, return 1, otherwise return 0

# IND-P2-C0 security proof

- masking vector $\mathbf{b} = \mathbf{r}\|z\|(\mathbf{a} \cdot M' \oplus \mathbf{v})$

- when $\mathcal{O} = U_{k+1}$, the $j$ first bits of $\mathbf{b}$ are random and the $m - j$ last ones are distributed according to an LPN distribution; for $j = m$ the ciphertexts are completely random

- when $\mathcal{O} = \Pi_{\mathbf{s},\eta}$, the $j - 1$ first bits of $\mathbf{b}$ are random and the $m - j + 1$ last ones are distributed according to an LPN distribution; for $j = 1$ the encryption is perfectly simulated

- when expressing the advantage of this distinguisher, the terms for $j = 2$ to $(m - 1)$ cancel and we obtain advantage $\delta/m$ if the advantage of the original distinguisher $\mathcal{A}$ was $\delta$

# malleability

- as is, the scheme is clearly malleable (P0-C0 attack):

- given a ciphertext $(a, y)$ corresponding to some plaintext $x$, the adversary can simply modify it to $(a, y \oplus C(x'))$, which will correspond to the plaintext $x \oplus x'$

- since IND-P2-C2 $\Leftrightarrow$ NM-P2-C2, the scheme cannot be IND-P2-C2 or even IND-P0-C2 either

- what about non-adaptive ciphertext attacks?

# an IND-P0-C1 attack

- idea: query the decryption oracle on $(\mathbf{a}, \mathbf{y_i})$ many times with the same $\mathbf{a}$ and random $\mathbf{y_i}$'s to get approximate equations on $\mathbf{a} \cdot M$

- when $\mathbf{y_i} \oplus \mathbf{a} \cdot M$ is at Hamming distance less than $t$ from a codeword, the decryption oracle will return $\mathbf{x_i}$ such that $\mathrm{Hwt}(C(\mathbf{x_i}) \oplus \mathbf{y_i} \oplus \mathbf{a} \cdot M) \leqslant t$

- this will give an approximation of each bit of $\mathbf{a} \cdot M$ with noise parameter less than $t/m$; repeating the experiment sufficiently many times with the same $\mathbf{a}$ enables to retrieve $\mathbf{a} \cdot M$ with high probability, hence to retrieve the secret key $M$

- this attack works only if the probability that a random $m$-bit string is decodable is sufficiently high, *i.e.* if the code is good enough

# P2-C2 security

- one can obtain an IND/NM-P2-C2 scheme by appending a MAC to the ciphertext (*Encrypt-then-MAC* paradigm studied by Bellare et al. [BN00])

- we propose the following MAC based on the LPN problem:

  ▸ let $M$ be a $l \times l'$ secret binary matrix and $H$ be a one-way function

  ▸ for $X \in \{0,1\}^*$ define $\mathrm{MAC}_M(X) = H(X) \cdot M \oplus \boldsymbol{\nu}$, where $\boldsymbol{\nu}$ is a noise vector of parameter $\eta$

- one can prove the security of this MAC in the random oracle model for $H$, using the hardness of the "MHB puzzle" [GRS08]

Given $q$ noisy samples $(\boldsymbol{a_i}, \boldsymbol{a_i} \cdot M \oplus \boldsymbol{\nu_i})$, where $M$ is a secret $k \times m$ matrix and $\Pr[\boldsymbol{\nu_i}[j] = 1] = \eta$, and a random challenge $\boldsymbol{a}$, find $\boldsymbol{a} \cdot M$.

# example parameters

- expansion factor $\sigma = \dfrac{|\text{ciphertext}|}{|\text{plaintext}|} = \dfrac{m+k}{r}$

| $k$ | $\eta$ | $m$ | $r$ | $d$ | expansion factor | key size | key size (Toeplitz) | $P_{DF}$ |
|------|--------|------|------|------|------------------|-----------|---------------------|----------|
| 512 | 0.125 | 80 | 27 | 21 | 21.9 | $40,960$ | 591 | 0.42 |
| 512 | 0.125 | 160 | 42 | 42 | 16 | $81,920$ | 671 | 0.44 |
| 768 | 0.05 | 80 | 53 | 9 | 16 | $61,440$ | 847 | 0.37 |
| 768 | 0.05 | 160 | 99 | 17 | 9.4 | $122,880$ | 927 | 0.41 |
| 768 | 0.05 | 160 | 75 | 25 | 12.4 | $122,880$ | 927 | 0.06 |

# possible variants and optimizations

- use of Toeplitz matrices to reduce the key size

$$\begin{pmatrix} & & & t_3 & t_2 & t_1 \\ & & & & t_3 & t_2 \\ & & \ddots & & & t_3 \\ t_{k+m-1} & & & & & \end{pmatrix}$$

- Toeplitz matrices have good randomization properties: $(x \to x \cdot T)_T$ is a $1/2^m$-balanced function family (for any non-zero vector $a$, $a \cdot T$ is uniformly distributed)

- possibility to pre-share the random vectors $a$ used to encrypt, or to re-generate them from a PRNG and a small seed; then $\sigma = \frac{m}{r}$, the expansion factor of the error-correcting code

# conclusion & open problems

- we presented LPN-C, a probabilistic symmetric encryption scheme whose security relies on the LPN problem

- it extends the range of cryptographic protocols based on the LPN problem

- implementation would be quite efficient but practical problems remain: expansion of the ciphertext, high key size

- open problems include:

  ▸ understand the impact of the use of Toeplitz matrices on the security of the scheme

  ▸ devise an efficient MAC whose security relies only on the LPN problem to obtain an IND/NM-P2-C2 secure encryption scheme

# thanks for your attention!

comments ∨ questions?