

How Far Can We Go Beyond Linear Cryptanalysis?

Thomas Baignères, Pascal Junod, and Serge Vaudenay

EPFL

<http://lasecwww.epfl.ch>

Abstract. Several generalizations of linear cryptanalysis have been proposed in the past, as well as very similar attacks in a statistical point of view. In this paper, we define a rigorous general statistical framework which allows to interpret most of these attacks in a simple and unified way. Then, we explicitly construct optimal distinguishers, we evaluate their performance, and we prove that a block cipher immune to classical linear cryptanalysis possesses some resistance to a wide class of generalized versions, but not all. Finally, we derive tools which are necessary to set up more elaborate extensions of linear cryptanalysis, and to generalize the notions of bias, characteristic, and piling-up lemma.

Keywords: Block ciphers, linear cryptanalysis, statistical cryptanalysis.

1 A Decade of Linear Cryptanalysis

Linear cryptanalysis is a known-plaintext attack proposed in 1993 by Matsui [21, 22] to break DES [26], exploiting specific correlations between the input and the output of a block cipher. Namely, the attack traces the statistical correlation between one bit of information about the plaintext and one bit of information about the ciphertext, both obtained linearly with respect to $\text{GF}(2)^L$ (where L is the block size of the cipher), by means of *probabilistic linear expressions*, a concept previously introduced by Tardy-Corfdir and Gilbert [30].

Soon after, several attempts to generalize linear cryptanalysis are published: Kaliski and Robshaw [13] demonstrate how it is possible to combine several independent linear correlations depending on the same key bits. In [31], Vaudenay defines another kind of attack on DES, called χ^2 -attack, and shows that one can obtain an attack slightly less powerful than a linear cryptanalysis, but without the need to know precisely what happens in the block cipher. Harpes, Kramer, and Massey [7] replace the linear expressions with so-called I/O sums, i.e., balanced binary-valued functions; they prove the potential effectiveness of such a generalization by exhibiting a block cipher secure against conventional linear cryptanalysis but vulnerable to their generalization. Practical examples are the attack of Knudsen and Robshaw [15] against LOKI91 and the one of Shimoyama and Kaneko [28] against DES which both use non-linear approximations.

In [8], Harpes and Massey generalize the results of [7] by considering *partitions pairs* of the input and output spaces. Let $\mathcal{X} = \{\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_n\}$ and

$\mathcal{Y} = \{\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_n\}$ be *partitions* of the input and output sets respectively, where \mathcal{X}_i and \mathcal{Y}_i are called *blocks*. The pair $(\mathcal{X}, \mathcal{Y})$ is called a *partition-pair* if all blocks of \mathcal{X} (respectively \mathcal{Y}) contain the same number of plaintexts (respectively ciphertexts). A partitioning cryptanalysis exploits the fact that the probabilities $\Pr [(X, f_k(X)) \in (\mathcal{X}, \mathcal{Y})]$ may not be uniformly distributed for a block cipher f_k when the plaintext X is uniformly distributed. In order to characterize the non-uniformity of a sample distribution, Harpes and Massey consider two “measures” called *peak imbalance* and *squared Euclidean imbalance*. Furthermore, they observe on toy-examples that the latter seems to lead to more successful attacks. These results are completed by Jakobsen and Harpes in [10, 9], where they develop useful bounds to estimate the resistance of block ciphers to partitioning cryptanalysis, with the help of spectral techniques; these bounds are relative to the squared Euclidean imbalance only, but this choice is not motivated in a formal way. To the best of our knowledge, the first practical example of partitioning cryptanalysis breaking a block cipher is the attack known as “stochastic cryptanalysis” [24] proposed by Minier and Gilbert against *Crypton* [17, 18].

In recent papers, Junod and Vaudenay [12, 11] consider linear cryptanalysis in a purely statistical framework, as it was done for the first time by Murphy et al. [25], for deriving optimal key ranking procedures and asymptotic bounds on the success probability of optimal linear distinguishers. A somewhat similar approach is chosen by Coppersmith et al. [1], except that it is adapted to stream ciphers. One can note that tight results about optimal distinguishers allow furthermore to derive useful security criteria.

Finally, the NESSIE effort resulted in a few papers investigating the power of linear (or non-linear) approximations based on different algebraic structures, like \mathbb{Z}_4 . For instance, Parker [27] shows how to approximate constituent functions of an S-box by *any* linear function over *any* weighted alphabet. However, Parker observes that it is not straightforward to piece these generalized linear approximations together. In [29], Standaert et al. take advantage of approximations in \mathbb{Z}_4 by *recombining* the values in order to reduce the problem to the well-known binary case; they obtain more interesting biases comparatively to a classical linear cryptanalysis.

Notation. Throughout this paper, random variables X, Y, \dots are denoted by capital letters, whilst their realizations $x \in \mathcal{X}, y \in \mathcal{Y}, \dots$ are denoted by small letters. The cardinal of a set \mathcal{X} is denoted $|\mathcal{X}|$. The probability function of a random variable X following a distribution D is denoted $\Pr_{\mathsf{D}}[x]$ or abusively $\Pr_X[x]$, when the distribution is clear from the context. For convenience, sequence X_1, X_2, \dots, X_n of n random variables is denoted \mathbf{X}^n . Similarly, a sequence x_1, x_2, \dots, x_n of realizations is denoted \mathbf{x}^n . We call *support* of a distribution D the set of all $x \in \mathcal{X}$ such that $\Pr_{\mathsf{D}}[x] \neq 0$. As usual, “iid” means “independent and identically distributed”. The transpose of a linear function h is denoted ${}^t h$. $\mathbb{1}_A$ is 1 if the predicate A is true, 0 otherwise. Finally, “.” denotes the inner product. The distribution function of the standard normal distribution is denoted

$$\Phi(t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-\frac{1}{2}u^2} du .$$

2 Optimal Distinguisher Between Two Sources

In this section, we shall consider a source generating a sequence of n iid random variables \mathbf{Z}^n following a distribution D and taking values in a set \mathcal{Z} . We wonder whether $D = D_0$ or $D = D_1$ (where D_1 is referred to as an “ideal distribution”), knowing that one of these two hypotheses is true. An algorithm which takes a sequence of n realizations \mathbf{z}^n as input and outputs either 0 or 1 is known as a *distinguisher* limited to n samples. It can be defined by an *acceptance region* $\mathcal{A} \subset \mathcal{Z}^n$ such that the distinguisher outputs 0 (respectively 1) when $\mathbf{z}^n \in \mathcal{A}$ (respectively $\mathbf{z}^n \notin \mathcal{A}$). The ability to distinguish a distribution from another is known as the *advantage* of the distinguisher and is defined by

$$\text{Adv}_{\mathcal{A}}^n = \left| \Pr_{D_0^n} [\mathcal{A}] - \Pr_{D_1^n} [\mathcal{A}] \right| ,$$

which is a quantity an adversary would like to maximize. The distinguisher can make two types of mistakes: it can either output 0 when $D = D_1$ or output 1 when $D = D_0$. We denote α and β the respective error probabilities and $P_e = \frac{1}{2}(\alpha + \beta)$ the *overall probability of error*. We can assume without loss of generality that $P_e \leq \frac{1}{2}$; we easily obtain that $\text{Adv}_{\mathcal{A}}^n = 1 - 2P_e$.

2.1 Deriving an Optimal Distinguisher

We describe here how to derive an optimal distinguisher for the scenario described below [1, 11]. Clearly, $P_e = \frac{1}{2} - \frac{1}{2} \sum_{\mathbf{z}^n \in \mathcal{A}} (\Pr_{D_0^n} [\mathbf{z}^n] - \Pr_{D_1^n} [\mathbf{z}^n])$, and therefore that the set minimizing¹ P_e is

$$\mathcal{A} = \{\mathbf{z}^n \in \mathcal{Z}^n : \text{LR}(\mathbf{z}^n) \geq 1\} \quad \text{where} \quad \text{LR}(\mathbf{z}^n) = \frac{\Pr_{D_0^n} [\mathbf{z}^n]}{\Pr_{D_1^n} [\mathbf{z}^n]} \quad (1)$$

stands for *likelihood ratio*². It defines an optimal distinguisher, i.e., with maximum advantage given a bounded number of samples and with no assumption on the computational power of the adversary.

In order to take a decision, a distinguisher defined by (1) has to keep in memory the results of the n queries, which is not feasible in practice if n grows. Fortunately, it is possible to derive an equivalent distinguisher with $|\mathcal{Z}|$ counter values $N(a|\mathbf{z}^n)$, each one counting the number of occurrence of a certain symbol a of \mathcal{Z} in the sequence \mathbf{z}^n . We summarize this in the following result.

Proposition 1. (*Optimal Distinguisher*) *The optimal acceptance region to test $D = D_0$ against $D = D_1$ is $\mathcal{A}_{\text{opt}} = \{\mathbf{z}^n \in \mathcal{Z}^n : \text{LLR}(\mathbf{z}^n) \geq 0\}$ where*

$$\text{LLR}(\mathbf{z}^n) = \sum_{\substack{a \in \mathcal{Z} \\ \text{s.t. } N(a|\mathbf{z}^n) > 0}} N(a|\mathbf{z}^n) \log \frac{\Pr_{D_0} [a]}{\Pr_{D_1} [a]}$$

¹ Note that we could have equivalently chosen a strict inequality in (1).

² The likelihood ratio builds the core of the Neyman-Pearson lemma [2, Ch. 12].

is the logarithmic likelihood ratio, with the convention that $\log \frac{0}{p} = -\infty$ and $\log \frac{p}{0} = +\infty$ (the $\log \frac{0}{0}$ case can be ignored), and where $N(a|z^n)$ is the number of times the symbol a occurs in the sequence $z^n \in \mathcal{Z}^n$.

Given the number of realizations n , we can compute the exact advantage of the optimal distinguisher. Let $[D_0]^n$ and $[D_1]^n$ be the vectors defined by

$$[D_j]_{(z_1, z_2, \dots, z_n)}^n = \Pr_{D_j} [z_1, z_2, \dots, z_n] \quad \text{with } j \in \{0, 1\} ,$$

which are a specific case of n -wise distribution matrices of the Decorrelation Theory [33] in a simplified case as we have no input here, only outputs z_i . The probability that the distinguisher outputs 0 when $D = D_j$ is $\sum_{z^n \in \mathcal{A}} [D_j]_{z^n}^n$, for $j \in \{0, 1\}$. The advantage is thus $|\sum_{z^n \in \mathcal{A}} ([D_0]_{z^n}^n - [D_1]_{z^n}^n)|$. Since \mathcal{A}_{opt} maximizes the sum, we obtain

$$\text{Adv}_{\mathcal{A}_{\text{opt}}}^n = \frac{1}{2} \| [D_0]^n - [D_1]^n \|_1 ,$$

where the norm $\| \cdot \|_1$ of a vector \mathbf{A} is defined by $\| \mathbf{A} \|_1 = \sum_i |\mathbf{A}_i|$. Note that the statistical framework of Coppersmith et al. [1] is based on this norm.

2.2 Complexity Analysis

In this section, we compute the number of queries the optimal distinguisher needs in order to distinguish D_0 from D_1 , given a fixed error probability P_e .

Definition 2. *The relative entropy or Kullback-Leibler distance between two distributions D_0 and D_1 is defined as*

$$D(D_0 \| D_1) = \sum_{z \in \mathcal{Z}} \Pr_{D_0} [z] \log \frac{\Pr_{D_0} [z]}{\Pr_{D_1} [z]} ,$$

with the convention that $0 \log \frac{0}{p} = 0$ and $p \log \frac{p}{0} = +\infty$ for $p > 0$.

We will refer to this notion using the term relative entropy as, being non-symmetric, it is not exactly a distance. Nevertheless, it is always positive since $-\log$ is convex. Using this notation, the following proposition can be proved.

Proposition 3. *Considering that Z_1, Z_2, \dots is a sequence of iid random variables of distribution D and that D_0 and D_1 share the same support,*

$$\Pr \left[\frac{\text{LLR}(\mathbf{Z}^n) - n\mu}{\sigma\sqrt{n}} < t \right] \xrightarrow{n \rightarrow \infty} \Phi(t) , \quad (2)$$

assuming that $\mu = \mu_j$ with $\mu_0 = D(D_0 \| D_1) \geq 0$ and $\mu_1 = -D(D_1 \| D_0) \leq 0$, and that σ^2 is

$$\sigma_j^2 = \sum_{z \in \mathcal{Z}} \Pr_{D_j} [z] \left(\log \frac{\Pr_{D_0} [z]}{\Pr_{D_1} [z]} \right)^2 - \mu_j^2 , \quad (3)$$

when $D = D_j$ for $j \in \{0, 1\}$.

Proof. We first note that the logarithmic likelihood ratio can be expressed as a sum $\text{LLR}(\mathbf{Z}^n) = R_1 + \dots + R_n$ where

$$R_i = \sum_{z \in \mathcal{Z}} \mathbb{1}_{Z_i=z} \log \frac{\Pr_{D_0}[z]}{\Pr_{D_1}[z]},$$

where every Z_i follows distribution D_j (so that the R_i 's are iid). The Central Limit Theorem then states that $\Pr [(\text{LLR}(\mathbf{Z}^n) - n\mu_j)/(\sigma_j\sqrt{n}) < t]$ converges in distribution towards $\Phi(t)$, where $\mu_j = \mathbb{E}_{D_j}[R_i]$ and $\sigma_j^2 = \text{Var}_{D_j}[R_i]$. Some straightforward computations lead to the announced result. Note that the assumption that both distributions share the same support is necessary for μ_j and σ_j to be well defined. \square

We now assume that the distributions D_0 and D_1 are close to each other, since it is the usual encountered case in practice.

Assumption 4. *Considering that D_0 is close to D_1 , we can write*

$$\forall z \in \mathcal{Z} : \Pr_{D_0}[z] = p_z + \epsilon_z \quad \text{and} \quad \Pr_{D_1}[z] = p_z \quad \text{with} \quad |\epsilon_z| \ll p_z .$$

Note that in such a case we can approximate $\text{LLR}(\mathbf{z}^n)$ by $\sum_a N(a|\mathbf{z}^n)\epsilon_a/p_a$. Proposition 3 can now be simplified using Taylor series.

Proposition 5. *Under the hypothesis of Proposition 3 and of Assumption 4 we have, at order two:*

$$\mu_0 \approx -\mu_1 \approx \frac{1}{2} \sum_{z \in \mathcal{Z}} \frac{\epsilon_z^2}{p_z} \quad \text{and} \quad \sigma_0^2 \approx \sigma_1^2 \approx \sum_{z \in \mathcal{Z}} \frac{\epsilon_z^2}{p_z} .$$

We can finally derive a heuristic theorem giving the number of samples the distinguisher needs, together with the implied probability of error, in order to distinguish close distributions with same support.

Theorem 6. *Let Z_1, \dots, Z_n be iid random variables over the set \mathcal{Z} of distribution D , D_0 and D_1 be two distributions of same support which are close to each other, and n be the number of samples of the best distinguisher between $D = D_0$ or $D = D_1$. Let d be a real number such that*

$$n = \frac{d}{\sum_{z \in \mathcal{Z}} \frac{\epsilon_z^2}{p_z}} \approx \frac{d}{2D(D_0 \| D_1)} \quad (4)$$

(where $p_z = \Pr_{D_1}[z]$ and $p_z + \epsilon_z = \Pr_{D_0}[z]$). Then, the overall probability of error is $P_e \approx \Phi(-\sqrt{d}/2)$.

Proof. If d is such that $\tilde{\mu} = \frac{1}{2}\sqrt{d/n} \tilde{\sigma}$, where $\tilde{\mu}$ and $\tilde{\sigma}$ respectively denote the approximation of μ_0 and σ_0 at order 2, we obtain (4). By definition $P_e = \frac{1}{2}(1 - \Pr_{D_1}[\text{LLR} < 0] + \Pr_{D_0}[\text{LLR} < 0])$. However

$$\Pr_{D_j}[\text{LLR} < 0] = \Pr_{D_j} \left[\frac{\text{LLR} - n\mu_j}{\sigma_j\sqrt{n}} < -\frac{\sqrt{n}\mu_j}{\sigma_j} \right] \approx \Phi \left(-\frac{\sqrt{n}\mu_j}{\sigma_j} \right) ,$$

where we make the usual approximation that the left hand side of (2) can be approximated by $\Phi(t)$. Therefore, as Proposition 5 states that $\mu_0 \approx -\mu_1 \approx \tilde{\mu}$ and that $\sigma_0 \approx \sigma_1 \approx \tilde{\sigma}$, we have $P_e \approx \frac{1}{2} \left(1 - \Phi(\sqrt{d}/2) + \Phi(-\sqrt{d}/2) \right) = \Phi(-\sqrt{d}/2)$. \square

Note that it may be possible to obtain strict tight bounds instead of an approximation for P_e using, for instance, Chernoff bounds.

2.3 Case where the Ideal Source is Uniform

From now on, we assume that D_1 is the uniform distribution. When D_0 is a distribution whose support is \mathcal{X} itself and which is close to D_1 , Theorem 6 can be rewritten with

$$n = \frac{d}{|\mathcal{Z}| \sum_{z \in \mathcal{Z}} \epsilon_z^2} .$$

This shows that the distinguishability can be measured by means of the Euclidean distance between D_0 and D_1 . In the very specific case where $\mathcal{Z} = \{0, 1\}$, we have $\epsilon_0 = -\epsilon_1 = \epsilon$ and one can see that n is proportional to ϵ^{-2} . It is a well accepted fact that the complexity of linear cryptanalysis is linked to the inverse of the square of the bias [21] which is, as we can see, a consequence of Theorem 6. We now recall what appears to be the natural measure of the bias of a distribution, considering the needed number of samples and Assumption 4.

Definition 7. Let $\epsilon_z = \Pr_{D_0} [z] - \frac{1}{|\mathcal{Z}|}$. The Squared Euclidean Imbalance³ (SEI) $\Delta(D_0)$ of a distribution D_0 of support \mathcal{Z} from the uniform distribution is defined by

$$\Delta(D_0) = |\mathcal{Z}| \sum_{z \in \mathcal{Z}} \epsilon_z^2 .$$

It is well-known (see [6, 14]) that a χ^2 cryptanalysis needs $O(1/\Delta(D_0))$ queries to succeed, which is by no means worse, up to a *constant term*, than an optimal distinguisher. Junod observed [11] that a χ^2 statistical test is asymptotically equivalent to a generalized likelihood-ratio test developed for a multinomial distribution; although such tests are not optimal in general, they usually perform reasonably well. Our results confirm this fact: a cryptanalyst will not loose any *essential* information in the case she can describe *only one* of the two distributions, but the precise knowledge of *both* distributions allows to derive an optimal attack. In other words, when it is impossible to derive both probability distributions, or when an attack involves many different distributions and only one is known, the best practical alternative to an optimal distinguisher seems to be a χ^2 attack, as proposed in [31]. This fact corroborates the intuition stipulating that χ^2 attacks are useful when one does not know precisely what happens in the attacked block cipher.

³ Although this appellation coincide with the one of [7], note that the definitions slightly differ.

2.4 Case where the Source Generates Boolean Vectors

We assume here that random variables are bitstrings⁴, so that $\mathcal{Z} = \{0, 1\}^\ell$.

Definition 8. *Following the notations of Assumption 4, let D_0 be the distribution defined by the set $\{\epsilon_z\}_{z \in \mathcal{Z}}$, D_1 being the uniform distribution on \mathcal{Z} . We define the Fourier transform of D_0 at point $u \in \mathcal{Z}$ as*

$$\hat{\epsilon}_u = \sum_{z \in \mathcal{Z}} (-1)^{u \cdot z} \epsilon_z . \quad (5)$$

The involution property of the Fourier transform leads to

$$\epsilon_z = \frac{1}{2^\ell} \sum_{u \in \mathcal{Z}} (-1)^{u \cdot z} \hat{\epsilon}_u . \quad (6)$$

The next property can be compared to Parseval's Theorem.

Proposition 9. *In the case where D_1 is the uniform distribution over $\mathcal{Z} = \{0, 1\}^\ell$, the SEI and the Fourier coefficients are related by:*

$$\Delta(\mathsf{D}_0) = \sum_{u \in \mathcal{Z}} \hat{\epsilon}_u^2 .$$

We now recall the definition of the linear probability [23], which plays a central role in the context of linear cryptanalysis.

Definition 10. *The linear probability of a boolean random variable B is*

$$\text{LP}(B) = (\Pr [B = 0] - \Pr [B = 1])^2 = (2 \Pr [B = 0] - 1)^2 = \left(\mathbb{E} [(-1)^B] \right)^2 .$$

Proposition 11. *Let $\mathcal{Z} = \{0, 1\}^\ell$. If $Z \in \mathcal{Z}$ is a random variable of distribution D_0 , the SEI and the linear probability are related by:*

$$\Delta(\mathsf{D}_0) = \sum_{w \in \mathcal{Z} \setminus \{0\}} \text{LP}(w \cdot Z) .$$

Proof. By using (5) we have $\hat{\epsilon}_u = \mathbb{E}_{\mathsf{D}_0} [(-1)^{u \cdot Z}] - \mathbb{1}_{u=0}$. Proposition 9 gives $\Delta(\mathsf{D}_0) = \sum_{u \in \mathcal{Z} \setminus \{0\}} \left(\mathbb{E}_{\mathsf{D}_0} [(-1)^{u \cdot Z}] \right)^2 = \sum_{w \in \mathcal{Z} \setminus \{0\}} \text{LP}(w \cdot Z)$. \square

Corollary 12. *Let Z be a random variable over $\mathcal{Z} = \{0, 1\}^\ell$ of distribution D_0 and let⁵ LP_{\max}^Z be the maximum of $\text{LP}(w \cdot Z)$ over $w \in \mathcal{Z} \setminus \{0\}$. We have*

$$\Delta(\mathsf{D}_0) \leq (2^\ell - 1) \text{LP}_{\max}^Z .$$

⁴ Note that all the study below extends in a straightforward way to $\mathcal{Z} = \text{GF}(p)^\ell$ for a prime p by replacing (-1) by $e^{\frac{2i\pi}{p}}$ and by using the conjugates of ϵ_z and $\hat{\epsilon}_z$ in (5) and (6) respectively. For simplicity we restrict ourselves to $\text{GF}(2)$.

⁵ We make a slight abuse of notation since LP_{\max}^Z is not a random variable depending on Z , but a real value depending on the *distribution* of Z .

Theorem 6 and Corollary 12 together mean that the complexity of the best distinguisher between two distributions of random bit strings can decrease with a factor up to 2^ℓ when compared to the best linear distinguisher. It is interesting to note that there are cases where this bound is tight. For example if D_0 is such that $\Pr_{D_0}[z]$ is $\frac{1}{2^\ell} + (1 - \frac{1}{2^\ell})\gamma$ if $z = 0$, and $\frac{1}{2^\ell} - \frac{1}{2^\ell}\gamma$ otherwise (where γ is a positive constant), it can be shown that $\text{LP}(w \cdot Z) = \gamma^2$ for all $w \neq 0$. Hence $\Delta(D_0) = (2^\ell - 1)\gamma^2$ and $\text{LP}_{\max} = \gamma^2$.

2.5 Statistical Distinguishers

In the last section, we have been trying to distinguish two random variables following two distinct distributions in a set $\mathcal{Z} = \{0, 1\}^\ell$ where ℓ should not be too large from an implementation point of view. If we try to distinguish two random variables distributed in some set $\{0, 1\}^L$ of large cardinality (e.g. where $L = 128$), we won't be able to implement the best distinguisher of Proposition 1 as the memory requirement would be too high. Instead, we can reduce the source space to a smaller space $\mathcal{Z} = \{0, 1\}^\ell$ by means of a *projection*⁶ $h : \{0, 1\}^L \rightarrow \mathcal{Z}$ defining, for a random variable $S \in \{0, 1\}^L$ of distribution \tilde{D} , a random variable $Z = h(S)$ of distribution D . Here we consider that h is a balanced function and that \tilde{D}_1 is a uniform distribution, so that D_1 is a uniform distribution as well. This is a typical construction in a real-life block cipher cryptanalysis, where the block length is quite large. Now, even though we know which distinguisher is the best to use in order to distinguish D_0 from D_1 , it is still not clear how the projection h has to be chosen. Probably the most classical example arises when $\ell = 1$ and $h(S) = a \cdot S$ for some non-zero $a \in \{0, 1\}^\ell$. We then talk about a *linear distinguisher*. In this case, we note that $\Delta(D_0) = \text{LP}(a \cdot S) \leq \text{LP}_{\max}^S$. Modern ciphers protect themselves against that type of distinguisher by bounding the value of LP_{\max}^S . A natural extension of the previous scheme would be to consider any linear projection onto wider spaces, e.g. to consider $h(S) \in \mathcal{Z} = \{0, 1\}^\ell$ (where $\ell > 1$ is still small) such that h is GF(2)-linear. We then talk about an *extended linear distinguisher*. It seems natural to wonder about the complexity gap between linear cryptanalysis and this extension. The following theorem proves that if a cipher provably resists classical linear cryptanalysis, it is (to some extent) protected against extended linear cryptanalysis.

Theorem 13. *Let S be a random variable over $\{0, 1\}^L$. Whenever the source space is reduced by a projection $h : \{0, 1\}^L \rightarrow \{0, 1\}^\ell$ in a GF(2)-linear way, we have $\Delta(h(S)) \leq (2^\ell - 1)\text{LP}_{\max}^S$.*

Proof. We use Proposition 11 and the fact that $w \cdot h(S) = {}^t h(w) \cdot S$. □

A classical example of a linear space reduction arises when considering *concatenation* of several projections. For example, denoting $D_0^{(i)} = h^{(i)}(\tilde{D}_0)$ for $i \in \{1, \dots, \ell\}$ where $h^{(i)} : \{0, 1\}^L \rightarrow \{0, 1\}$ is linear, we consider $h(S) =$

⁶ We borrow this appellation from Vaudenay [31]; the same expression is used within Wagner's unified view of block cipher cryptanalysis [34] as well.

$(h^{(1)}(S), \dots, h^{(n)}(S))$. This corresponds to the works of Kaliski and Robshaw [13] (where different linear characteristics involving *identical* key bits are merged) and of Junod and Vaudenay [12] (where different linear characteristics involving *different* key bits are merged). In the latter situation, if no assumption is made about the dependency among the $\Delta(\mathsf{D}_0^{(i)})$'s, Theorem 13 tells us $\Delta(\mathsf{D}_0^{(1)} \times \dots \times \mathsf{D}_0^{(\ell)}) \leq (2^\ell - 1)\text{LP}_{\max}^S$. The following proposition tells us what happens in general when the $\mathsf{D}_0^{(i)}$'s are independent but do not necessarily come from a linear projection nor a Boolean projection.

Proposition 14. *Consider the case where $\mathsf{D}_0 = \mathsf{D}_0^{(1)} \times \dots \times \mathsf{D}_0^{(\ell)}$. If $\mathsf{D}_0^{(1)}, \dots, \mathsf{D}_0^{(\ell)}$ are independent distributions, then $\Delta(\mathsf{D}_0) + 1 = \prod_{i=1}^{\ell} (\Delta(\mathsf{D}_0^{(i)}) + 1)$. Therefore, $\Delta(\mathsf{D}_0)$ can be approximated by the sum of the $\Delta(\mathsf{D}_0^{(i)})$'s.*

Proof. For the sake of simplicity, we restrict this proof to the case where $\mathsf{D}_0 = \mathsf{D}_0^{(a)} \times \mathsf{D}_0^{(b)}$. Let $Z = (A, B)$ where A and B are two independent random variable following distributions $\mathsf{D}_0^{(a)}$ and $\mathsf{D}_0^{(b)}$ respectively. As in Proposition 11, we have

$$\begin{aligned} \Delta(\mathsf{D}_0^{(a)} \times \mathsf{D}_0^{(b)}) &= \sum_{(v,w) \in \mathbb{Z}^2 \setminus \{\mathbf{0}\}} (\mathbb{E} [(-1)^{v \cdot A \oplus w \cdot B}])^2 \\ &= \sum_{(v,w) \in \mathbb{Z}^2 \setminus \{\mathbf{0}\}} (\mathbb{E} [(-1)^{v \cdot A}])^2 (\mathbb{E} [(-1)^{w \cdot B}])^2 \\ &= (\Delta(\mathsf{D}_0^{(a)}) + 1) (\Delta(\mathsf{D}_0^{(b)}) + 1) - 1 . \end{aligned}$$

□

This result tells us that merging ℓ independent biases should only be considered when their respective amplitudes are within the same order of magnitude.

In the light of the preceding discussion, the cryptanalyst may wonder if it is possible to find a distinguisher with a high advantage even though the value of LP_{\max}^S is very small. We provide an example for which it is indeed the case.

Example. Consider a source generating a random variable $S = (X_1, \dots, X_{n+1}) \in \mathbb{Z}_4^{n+1}$, where n is some odd large integer, and we represent \mathbb{Z}_4 by $\{0, 1\}^2$ in binary. Here we have $L = 2n+2$. If the source follows distribution D_0 , then $X_1, \dots, X_n \in \mathbb{Z}_4$ are uniform iid random variables and $X_{n+1} = (Y + \sum_{i=1}^n X_i) \bmod 4$, where $Y \in \{0, 1\}$ is a uniformly distributed random variable independent of X_1, \dots, X_n . If the source follows distribution D_1 , $S \in \mathbb{Z}_4^{n+1}$ is uniformly distributed. It can be shown (see Appendix A) that $\text{LP}_{\max}^S = 2^{-(n+1)}$. On the other hand, if we let $h : \mathbb{Z}_4^{n+1} \rightarrow \mathbb{Z}_2$ be such that $h(S) = \text{msb}((X_{n+1} - \sum_{i=0}^n X_i) \bmod 4)$ (where msb stands for most significant bit), we have $\ell = 1$ and a SEI equal to 1, so that $\Delta(\mathsf{D}_0) \gg \text{LP}_{\max}^S$: D_0 can be distinguished from D_1 despite LP_{\max}^S is small.

This example shows that Theorem 13 tells us nothing about the SEI whenever the plaintext space is reduced by a non-linear projection. Therefore, even though

LP_{\max}^S is very low, there may exist some tricky non-linear projections which lead to significant breakdown of the complexity of distinguishers, i.e., there may be something beyond linear cryptanalysis.

3 Optimal Distinguisher Between Two Oracles

So far we discussed how to distinguish random sources. Now we investigate applications for distinguishing random oracles, such as block ciphers, and how to transform this into the previous problem.

We consider the random variable Z taking values in \mathcal{Z} to be a couple of random variables (X, Y) taking values in $\mathcal{X} \times \mathcal{Y}$. As discussed in Sect. 2.5, the couple (X, Y) can be seen like the image of a plaintext/ciphertext couple (P, C) by some balanced projections ϕ and ψ (which actually define the statistical cryptanalysis in use); in other words, the adversary queries the oracle for known-plaintext pairs and compute the projections ϕ and ψ to sample (X, Y) . For simplicity reasons, we focus our study on known-plaintext attacks (such as linear cryptanalysis) and thus, we consider that X is uniformly distributed. The distribution of Y is defined by a transition matrix \mathbf{T} such that

$$[\mathbf{T}]_{x,y} = \Pr [Y = y | X = x] = \Pr [\psi(C) = y | \phi(P) = x] .$$

The transition matrix \mathbf{T} can either be \mathbf{T}_0 or \mathbf{T}_1 , where \mathbf{T}_1 is the uniform transition matrix (i.e., $[\mathbf{T}_1]_{x,y} = \frac{1}{|\mathcal{Y}|}$). The distribution D of Z depends on the transition matrix \mathbf{T} . We will denote it D_0 (respectively D_1) when $\mathbf{T} = \mathbf{T}_0$ (respectively $\mathbf{T} = \mathbf{T}_1$). We can see that if $\mathbf{T} = \mathbf{T}_1$, as X is uniformly distributed, the distribution D_1 of Z is also uniform. Therefore, all the results presented so far can be applied to the particular case we study here. Indeed, if we note that

$$\Pr_D [z] = \Pr [X = x, Y = y] = [\mathbf{T}]_{x,y} \Pr [X = x] .$$

we can express Proposition 1 in terms of the transition matrices.

Proposition 15. (Optimal Binary Hypothesis Test with Transition Matrices) *The optimal acceptance region to test $D = D_0$ against $D = D_1$ (where D_1 is the uniform distribution), that is to test $\mathbf{T} = \mathbf{T}_0$ against $\mathbf{T} = \mathbf{T}_1$, is*

$$\mathcal{A}_{\text{opt}} = \{(\mathbf{x}^n, \mathbf{y}^n) \in \mathcal{X}^n \times \mathcal{Y}^n : \text{LLR}(\mathbf{x}^n, \mathbf{y}^n) \geq 0\}$$

where

$$\text{LLR}(\mathbf{x}^n, \mathbf{y}^n) = \sum_{\substack{(x,y) \in \mathcal{X} \times \mathcal{Y} \\ \text{s.t. } N((x,y)|\mathbf{z}^n) > 0}} N((x,y)|\mathbf{z}^n) \log \frac{[\mathbf{T}_0]_{x,y}}{[\mathbf{T}_1]_{x,y}}$$

with the conventions used in Proposition 1.

In the next sections, we derive the complexity of this distinguisher, discuss the relationship between our model and previous work, and study how Matsui's *Piling-up Lemma* [21] extends to our model.

3.1 Cryptanalysis Complexity

We introduce the notion of *bias matrix* $\mathbf{B} = \mathbf{T}_0 - \mathbf{T}_1$. Note that $\sum_{x \in \mathcal{X}} [\mathbf{B}]_{x,y} = 0$ when X is uniformly distributed and that $\sum_{y \in \mathcal{Y}} [\mathbf{B}]_{x,y} = 0$ in any case. Similarly to Definition 8, the Fourier transform $\widehat{\mathbf{B}}$ of the bias matrix \mathbf{B} is such that

$$[\widehat{\mathbf{B}}]_{u,v} = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (-1)^{u \cdot x \oplus v \cdot y} [\mathbf{B}]_{x,y} .$$

Furthermore, we define LPM, the *linear probability matrix*, by $[\text{LPM}]_{u,v} = 0$ if $u = v = 0$ and by $[\text{LPM}]_{u,v} = \text{LP}(u \cdot X \oplus v \cdot Y)$ otherwise. It can be noted that $[\widehat{\mathbf{B}}]_{u,v}^2 = |\mathcal{X}|^2 [\text{LPM}]_{u,v}$. With the notations we just introduced, it is possible to derive the complexity of the best distinguisher between two oracles as a simple consequence of Theorem 6 and of Proposition 11.

Proposition 16. *Let n be the number of queries of the best distinguisher between \mathbf{T}_0 and \mathbf{T}_1 , which are supposed to be close to each other and of same support. Then the overall probability of error is $P_e \approx 1 - \Phi(\sqrt{d}/2)$, where d is a real number such that $n = d/\Delta(\mathbf{D}_0)$. Furthermore, as*

$$\Delta(\mathbf{D}_0) = \frac{|\mathcal{Y}|}{|\mathcal{X}|} \|\mathbf{B}\|_2^2 = \frac{1}{|\mathcal{X}|^2} \|\widehat{\mathbf{B}}\|_2^2 = \sum_{(u,v) \in \mathcal{X} \times \mathcal{Y}} [\text{LPM}]_{u,v} ,$$

n can be equivalently expressed in terms of the bias matrix, of its Fourier transform, or of the linear probability matrix (and thus, of the linear probabilities).

Matsui's linear expressions are a very particular case of the transition matrices we have defined at the beginning of Sect. 3. Indeed, choosing balanced *linear* projections $\phi, \psi : \{0,1\}^L \rightarrow \{0,1\}$ is equivalent to choose input/output masks on the plaintext/ciphertext bits. The respective shapes of the corresponding bias matrix, of its Fourier transform, and of the LPM matrix are

$$\mathbf{B} = \begin{pmatrix} \epsilon & -\epsilon \\ -\epsilon & \epsilon \end{pmatrix} , \quad \widehat{\mathbf{B}} = \begin{pmatrix} 0 & 0 \\ 0 & 4\epsilon \end{pmatrix} , \quad \text{and} \quad \text{LPM} = \begin{pmatrix} 0 & 0 \\ 0 & 4\epsilon^2 \end{pmatrix} ,$$

where ϵ is nothing but the bias of Matsui's linear expressions. According to Proposition 16, we see that the complexity of the distinguishing attack is proportional to $\|\mathbf{B}\|_2^{-2}$, which is a well known result in linear cryptanalysis, for which $\|\mathbf{B}\|_2^2 = 4\epsilon^2$.

There is an intuitive link between linear probability matrices and *correlation matrices* [3]. Recall that the correlation matrix of a Boolean function $f : \{0,1\}^n \rightarrow \{0,1\}^m$ is the $2^m \times 2^n$ matrix $\mathbf{C}^{(f)}$ such that $[\mathbf{C}^{(f)}]_{u,v} = 2 \Pr [u \cdot f(P) \oplus v \cdot P] - 1$, where the probability holds over the uniform distribution of P , so that $[\mathbf{C}^{(f)}]_{u,v}^2 = \text{LP}(u \cdot f(P) \oplus v \cdot P)$. We see that correlation matrices are strongly related to the linear probability matrices in the specific case where ϕ and ψ are identity functions (i.e., no reduction is performed on the plaintext space).

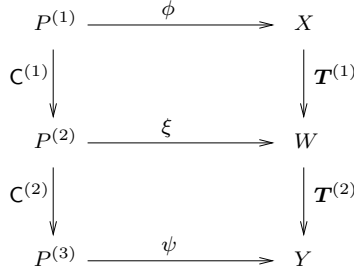


Fig. 1. Two rounds of an iterated block cipher

3.2 Piling-up Transition Matrices

A distinguishing attack on an iterated cipher is practical on condition that the cryptanalyst knows a transition matrix spanning several rounds. In practice, the cryptanalyst will derive a transition matrix on each round and, provided that the projections were chosen carefully, pile them in order to obtain a transition matrix on several rounds of the cipher.

We consider the scenario where a block cipher is represented by a random permutation C over $\{0, 1\}^L$ (L denotes the block size of the cipher), where the randomness comes from the key. Moreover we suppose that the block cipher is made of two rounds corresponding to the succession of two random permutations $C^{(1)}$ and $C^{(2)}$. In other words $C = C^{(2)} \circ C^{(1)}$. We denote $P^{(1)}, P^{(2)} \in \{0, 1\}^L$ the respective inputs of $C^{(1)}$ and $C^{(2)}$, whereas $P^{(3)}$ denote the output of $C^{(2)}$. The random variables X, W , and Y respectively denote $\phi(P^{(1)})$, $\xi(P^{(2)})$, and $\psi(P^{(3)})$, where ϕ, ξ , and ψ are projections onto \mathcal{X}, \mathcal{W} , and \mathcal{Y} , respectively. With these notations, the respective transition matrices of $C^{(1)}, C^{(2)}$, and C are

$$\begin{aligned}
[T^{(1)}]_{x,w} &= \Pr_{W|X} [w | x] \quad , \quad [T^{(2)}]_{w,y} = \Pr_{Y|W} [y | w] \quad , \\
\text{and } [T]_{x,y} &= \Pr_{Y|X} [y | x] \quad .
\end{aligned}$$

This situation is represented on Fig. 1. Note that we use a representation which is very similar to Wagner's commutative diagrams [34]. Under the assumption that $X \leftrightarrow W \leftrightarrow Y$ is a Markov chain (as in [34]), it can easily be shown that successive transition matrices are multiplicative, i.e., $T = T^{(1)} \times T^{(2)}$. Note that this situation is idealistic as, even under the classical assumption that $P^{(1)} \leftrightarrow P^{(2)} \leftrightarrow P^{(3)}$ is a Markov chain [16, 32], $X \leftrightarrow W \leftrightarrow Y$ may not be a Markov chain unless the projection are chosen with care. Nevertheless, under the suitable suppositions, the following lemma shows how the Piling-Up Lemma extends to our model.

Lemma 17. *Let $B^{(1)}, B^{(2)}$, and B be the bias matrices associated with $T^{(1)}, T^{(2)}$, and T respectively, such that $T = T^{(1)} \times T^{(2)}$. In the case of a known-plaintext attack, $B = B^{(1)} \times B^{(2)}$ and $\widehat{B} = \frac{1}{|\mathcal{W}|} \widehat{B}^{(1)} \times \widehat{B}^{(2)}$. Therefore,*

$\| \mathbf{B} \|_2^2 \leq \| \mathbf{B}^{(1)} \|_2^2 \| \mathbf{B}^{(2)} \|_2^2$, with equality if, and only if we can write $[\mathbf{B}^{(1)}]_{x,w} = \alpha_x \gamma_w$ and $[\mathbf{B}^{(2)}]_{w,y} = \gamma_w \beta_y$, for some $\alpha \in \mathbb{R}^{|\mathcal{X}|}$, $\beta \in \mathbb{R}^{|\mathcal{Y}|}$ and $\gamma \in \mathbb{R}^{|\mathcal{W}|}$.

Proof. As $\mathbf{T} = \mathbf{T}^{(1)} \times \mathbf{T}^{(2)}$, we have

$$[\mathbf{B}]_{x,y} = [\mathbf{T}]_{x,y} - \frac{1}{|\mathcal{Y}|} = \sum_{w \in \mathcal{W}} \left([\mathbf{B}^{(1)}]_{x,w} + \frac{1}{|\mathcal{W}|} \right) \left([\mathbf{B}^{(2)}]_{w,y} + \frac{1}{|\mathcal{Y}|} \right) - \frac{1}{|\mathcal{Y}|} .$$

As $\sum_w [\mathbf{B}^{(1)}]_{x,w} = 0$, we obtain $[\mathbf{B}]_{x,y} = [\mathbf{B}^{(1)} \times \mathbf{B}^{(2)}]_{x,y} + \frac{1}{|\mathcal{W}|} \sum_w [\mathbf{B}^{(2)}]_{w,y}$. The fact that $P^{(1)}$ is uniformly distributed implies that $P^{(2)}$ and $P^{(3)}$ are uniformly distributed and thus, as ϕ , ξ , and ψ are balanced, that X , Z , and Y are also uniformly distributed. In that case, we know that $\sum_{w \in \mathcal{W}} [\mathbf{B}^{(2)}]_{w,y} = 0$, which proves that $\mathbf{B} = \mathbf{B}^{(1)} \times \mathbf{B}^{(2)}$. We also have

$$\begin{aligned} \left[\widehat{\mathbf{B}}^{(1)} \times \widehat{\mathbf{B}}^{(2)} \right]_{u,v} &= \sum_{a \in \mathcal{W}} \left[\widehat{\mathbf{B}}^{(1)} \right]_{u,a} \left[\widehat{\mathbf{B}}^{(2)} \right]_{a,v} \\ &= \sum_{\substack{(x,w) \in \mathcal{X} \times \mathcal{W} \\ (w',y) \in \mathcal{W} \times \mathcal{Y}}} (-1)^{u \cdot x \oplus v \cdot y} [\mathbf{B}^{(1)}]_{x,w} [\mathbf{B}^{(2)}]_{w',y} \sum_{a \in \mathcal{W}} (-1)^{a \cdot (w \oplus w')} \\ &= |\mathcal{W}| \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (-1)^{u \cdot x \oplus v \cdot y} \sum_{w \in \mathcal{W}} [\mathbf{B}^{(1)}]_{x,w} [\mathbf{B}^{(2)}]_{w,y} \\ &= |\mathcal{W}| \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} (-1)^{u \cdot x \oplus v \cdot y} [\mathbf{B}]_{x,y} \\ &= |\mathcal{W}| \left[\widehat{\mathbf{B}} \right]_{u,v} , \end{aligned}$$

which proves that $\widehat{\mathbf{B}} = \frac{1}{|\mathcal{W}|} \widehat{\mathbf{B}}^{(1)} \times \widehat{\mathbf{B}}^{(2)}$. Finally, from Cauchy-Schwarz inequality:

$$\begin{aligned} \| \mathbf{B}^{(1)} \times \mathbf{B}^{(2)} \|_2^2 &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left(\sum_{w \in \mathcal{W}} [\mathbf{B}^{(1)}]_{x,w} [\mathbf{B}^{(2)}]_{w,y} \right)^2 \\ &\leq \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} \left(\sum_{w \in \mathcal{W}} [\mathbf{B}^{(1)}]_{x,w}^2 \right) \left(\sum_{w' \in \mathcal{W}} [\mathbf{B}^{(2)}]_{w',y}^2 \right) \\ &= \| \mathbf{B}^{(1)} \|_2^2 \| \mathbf{B}^{(2)} \|_2^2 , \end{aligned}$$

with equality if, and only if, for all $x, y \in \mathcal{X} \times \mathcal{Y}$ there exists some $\lambda_{x,y}$ such that $[\mathbf{B}^{(1)}]_{x,w} = \lambda_{x,y} [\mathbf{B}^{(2)}]_{w,y}$, so that $[\mathbf{B}^{(1)}]_{x,w} = \lambda_{x,0} [\mathbf{B}^{(2)}]_{w,0} = \alpha_x \gamma_w$. Taking β_y equal to $\alpha_0 / \lambda_{0,y}$ when $\lambda_{0,y} \neq 0$ and to zero otherwise leads to the announced result. \square

How to find projections ϕ , ψ and ξ on larger spaces exhibiting such a Markovian property in a given block cipher remains however an open question to us. We may hope to *approximate* such a Markovian process.

4 Distinguishers Versus Key Recovery

In this section we show that our framework can adapt to key recovery instead of distinguishers. Let us consider a process which generates *independent* random variables $Z_{1,K}, \dots, Z_{n,K}$ depending on some key $K \in \{0, 1\}^k$. We assume that for one unknown value $K = K_0$ all $Z_{i,K}$'s follow distribution D_0 , whereas when $K \neq K_0$ all $Z_{i,K}$'s follow distribution D_1 . We consider the simple key ranking procedure which, for all possible $K \in \{0, 1\}^k$, instantiates \mathbf{z}_K^n and ranks K according to the grade $G_K = \text{LLR}(\mathbf{z}_K^n)$. For any $K \neq K_0$ we obtain (similarly to what we had in Theorem 6) that $G_{K_0} - G_K$ is approximately normally distributed with expected value $n\Delta(D_0)$ and standard deviation $\sqrt{2n\Delta(D_0)}$. Hence we obtain $G_{K_0} < G_K$ (i.e., a wrong key K has a better rank than the right key K_0) with probability approximately $\Phi\left(-\sqrt{n\Delta(D_0)/2}\right)$. Let d be such that $n = d/\Delta(D_0)$. This probability becomes $\Phi\left(-\sqrt{d/2}\right)$ which is approximately $e^{-d/4}/\sqrt{2\pi}$ when d is large. So K_0 gets the highest grade with probability approximately equal to $(1 - e^{-d/4}/\sqrt{2\pi})^{2^k-1} \approx \exp(-2^k \cdot e^{-d/4}/\sqrt{2\pi})$, which is high provided that $d \geq 4k \log 2$. Hence we need

$$n \geq \frac{4k \log 2}{\Delta(D_0)}.$$

This formula is quite useful to estimate the complexity of many attacks, e.g. [19, 20]⁷. We can finally note that the expected rank of K_0 (from 1 up to 2^k) is $1 + (2^k - 1)\Phi\left(-\sqrt{n\Delta(D_0)/2}\right)$.

5 Conclusion

Most modern block ciphers are proven to be resistant to linear cryptanalysis in some sense. In this paper, we wonder how this resistance extends to (both known and unknown) generalizations of linear cryptanalysis. For this, we define a sound and rigorous statistical framework which allows us to interpret most of these attacks in a simple and unified way; secondly, we develop a set of useful statistical tools to describe such attacks and to analyze their performance. Recently, our results on GF(2)-linear projections were exploited by [20] to obtain a small improvement factor in an attack on E0, and by [19] in another attack against two-level E0 [19]. In the sequel of this paper, we observe that resistance to linear cryptanalysis implies (a somewhat weaker) resistance to generalizations based on GF(2)-*linear projections*; however this resistance does not extend to *all* statistical cryptanalysis, as demonstrated by our example exploiting correlations in \mathbb{Z}_4 , which confirms observations of Parker and Standaert et al. [27, 29]. The next natural step, which we hope to have rendered easier, will be to exhibit such a practical statistical cryptanalysis against a block cipher immune to linear cryptanalysis, like AES [4].

⁷ Note that [19, 20] use slightly different notations: $\Delta(D_0)$ denotes the Euclidean Imbalance instead of the *Squared* Euclidean Imbalance.

References

- [1] D. Coppersmith, S. Halevi, and C. Jutla. Cryptanalysis of stream ciphers with linear masking. In M. Yung, editor, *Advances in Cryptology - CRYPTO '02*, volume 2442 of *LNCS*, pages 515–532. Springer-Verlag, 2002.
- [2] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley Series in Telecommunications. John Wiley & Sons, 1991.
- [3] J. Daemen, R. Govaerts, and J. Vandewalle. Correlation matrices. In B. Preneel, editor, *Fast Software Encryption - FSE '94*, volume 1008 of *LNCS*, pages 275–285. Springer-Verlag, 1995.
- [4] J. Daemen and V. Rijmen. *The Design of Rijndael*. Information Security and Cryptography. Springer-Verlag, 2002.
- [5] W. Feller. *An Introduction to Probability Theory and Its Applications*. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons, third edition, 1968.
- [6] H. Handschuh and H. Gilbert. χ^2 cryptanalysis of the SEAL encryption algorithm. In E. Biham, editor, *Fast Software Encryption - FSE '97*, volume 1267 of *LNCS*, pages 1–12. Springer-Verlag, 1997.
- [7] C. Harpes, G. Kramer, and J. Massey. A generalization of linear cryptanalysis and the applicability of Matsui's piling-up lemma. In L.C. Guillou and J.-J. Quisquater, editors, *Advances in Cryptology - EUROCRYPT '95*, volume 921 of *LNCS*, pages 24–38. Springer-Verlag, 1995.
- [8] C. Harpes and J. Massey. Partitioning cryptanalysis. In E. Biham, editor, *Fast Software Encryption - FSE '97*, volume 1267 of *LNCS*, pages 13–27. Springer-Verlag, 1997.
- [9] T. Jakobsen. *Higher-order cryptanalysis of block ciphers*. PhD thesis, Department of Mathematics, Technical University of Denmark, 1999.
- [10] T. Jakobsen and C. Harpes. Non-uniformity measures for generalized linear cryptanalysis and partitioning cryptanalysis. In J. Pribyl, editor, *PRAGOCRYPT '96*. CTU Publishing House, 1996.
- [11] P. Junod. On the optimality of linear, differential and sequential distinguishers. In E. Biham, editor, *Advances in Cryptology - EUROCRYPT '03*, volume 2656 of *LNCS*, pages 17–32. Springer-Verlag, 2003.
- [12] P. Junod and S. Vaudenay. Optimal key ranking procedures in a statistical cryptanalysis. In T. Johansson, editor, *Fast Software Encryption - FSE '03*, volume 2887 of *LNCS*, pages 235–246. Springer-Verlag, 2003.
- [13] B. Kaliski and M. Robshaw. Linear cryptanalysis using multiple approximations. In Y.G. Desmedt, editor, *Advances in Cryptology - CRYPTO '94*, volume 839 of *LNCS*, pages 26–39. Springer-Verlag, 1994.
- [14] J. Kelsey, B. Schneier, and D. Wagner. *modn* cryptanalysis, with applications against RC5P and M6. In L. Knudsen, editor, *Fast Software Encryption - FSE '99*, volume 1636 of *LNCS*, pages 139–155. Springer-Verlag, 1999.
- [15] L. Knudsen and M. Robshaw. Non-linear approximations in linear cryptanalysis. In U. Maurer, editor, *Advances in Cryptology - EUROCRYPT '96*, volume 1070 of *LNCS*, pages 224–236. Springer-Verlag, 1996.
- [16] X. Lai, J. Massey, and S. Murphy. Markov ciphers and differential cryptanalysis. In D.W. Davies, editor, *Advances in Cryptology - EUROCRYPT '91*, volume 547 of *LNCS*, pages 17–38. Springer-Verlag, 1991.
- [17] C.H. Lim. CRYPTON: A new 128-bit block cipher. In *The First AES Candidate Conference*. National Institute for Standards and Technology, 1998.

- [18] C.H. Lim. A revised version of CRYPTON: CRYPTON V1.0. In L. Knudsen, editor, *Fast Software Encryption - FSE '99*, volume 1636 of *LNCS*, pages 31–45. Springer-Verlag, 1999.
- [19] Y. Lu and S. Vaudenay. Cryptanalysis of Bluetooth Keystream Generator Two-level E0. In *Advances in Cryptology - ASIACRYPT '04*, LNCS. Springer-Verlag, 2004.
- [20] Y. Lu and S. Vaudenay. Faster correlation attack on Bluetooth keystream generator E0. In M. Franklin, editor, *Advances in Cryptology - CRYPTO '04*, volume 3152 of *LNCS*, pages 407–425. Springer-Verlag, 2004.
- [21] M. Matsui. Linear cryptanalysis method for DES cipher. In T. Helleseth, editor, *Advances in Cryptology - EUROCRYPT '93*, volume 765 of *LNCS*, pages 386–397. Springer-Verlag, 1993.
- [22] M. Matsui. The first experimental cryptanalysis of the Data Encryption Standard. In Y.G. Desmedt, editor, *Advances in Cryptology - CRYPTO '94*, volume 839 of *LNCS*, pages 1–11. Springer-Verlag, 1994.
- [23] M. Matsui. New structure of block ciphers with provable security against differential and linear cryptanalysis. In D. Gollman, editor, *Fast Software Encryption - FSE '96*, volume 1039 of *LNCS*, pages 205–218. Springer-Verlag, 1996.
- [24] M. Minier and H. Gilbert. Stochastic cryptanalysis of Crypton. In B. Schneier, editor, *Fast Software Encryption - FSE '00*, volume 1978 of *LNCS*, pages 121–133. Springer-Verlag, 2000.
- [25] S. Murphy, F. Piper, M. Walker, and P. Wild. Likelihood estimation for block cipher keys. Technical report, Information Security Group, University of London, England, 1995.
- [26] National Institute of Standards and Technology, U. S. Department of Commerce. *Data Encryption Standard*, NIST FIPS PUB 46-2, 1993.
- [27] M. Parker. Generalized S-Box linearity. Technical report [nes/doc/uib/wp5/020/a](https://www.cryptoneessie.org), NESSIE Project, 2003. Available on <https://www.cryptoneessie.org>.
- [28] T. Shimoyama and T. Kaneko. Quadratic relation of S-Box and its application to the linear attack of full round DES. In H. Krawczyk, editor, *Advances in Cryptology - CRYPTO '98*, volume 1462 of *LNCS*, pages 200–211. Springer-Verlag, 1998.
- [29] F.-X. Standaert, G. Rouvroy, G. Piret, J.-J. Quisquater, and J.-D. Legat. Key-dependent approximations in cryptanalysis: an application of multiple Z4 and non-linear approximations. In *24th Symposium on Information Theory in the Benelux*, 2003.
- [30] A. Tardy-Corffdir and H. Gilbert. A known plaintext attack of FEAL-4 and FEAL-6. In J. Feigenbaum, editor, *Advances in Cryptology - CRYPTO '91*, volume 576 of *LNCS*, pages 172–182. Springer-Verlag, 1992.
- [31] S. Vaudenay. An experiment on DES statistical cryptanalysis. In *3rd ACM Conference on Computer and Communications Security*, pages 139–147. ACM Press, 1996.
- [32] S. Vaudenay. On the security of CS-cipher. In L. Knudsen, editor, *Fast Software Encryption - FSE '99*, volume 1636 of *LNCS*, pages 260–274. Springer-Verlag, 1999.
- [33] S. Vaudenay. Decorrelation: a theory for block cipher security. *Journal of Cryptology*, 16(4):249–286, 2003.
- [34] D. Wagner. Towards a unifying view of block cipher cryptanalysis. In B. Roy and W. Meier, editors, *Fast Software Encryption - FSE '04*, volume 3017 of *LNCS*, pages 16–33. Springer-Verlag, 2004.

A A Strange Distribution

We consider a source generating a random variable $S = (X_1, \dots, X_{n+1}) \in \mathbb{Z}_4^{n+1}$, where n is some large integer, which follows either D_0 or D_1 (the uniform distribution). The distribution D_0 is such that the X_1, \dots, X_n are uniformly distributed iid random variables and $X_{n+1} = (Y + \sum_{i=1}^n X_i) \bmod 4$, where $Y \in \{0, 1\}$ is uniformly distributed and independent of X_1, \dots, X_n . We claim that the linear probability of the best linear distinguisher with one query is very small (equal to $2^{-(n+1)}$) whereas it is still possible to find a projection h such that $Z = h(S)$ has a high SEI. In order to simplify the proof, we will suppose that $n + 1$ is a multiple of 4.

Proposition 18. *Let $h : \mathbb{Z}_4^{n+1} \rightarrow \mathbb{Z}_2$ be defined by $h(S) = \text{msb}((X_{n+1} - \sum_{i=1}^n X_i) \bmod 4)$. Then the SEI of $Z = h(S)$ is 1.*

The following lemmas will be used to prove that the best linear distinguisher is drastically less powerful than the distinguisher of Proposition 18.

Lemma 19. *Let $\mathbf{u} = u_1 u_2 \dots u_n$ be a string of n bits. If we denote w the Hamming weight of \mathbf{u} then we have*

$$\sum_{1 \leq j < k \leq n} u_j u_k = \frac{w(w-1)}{2} .$$

Lemma 20. *For any positive integer N , we have:*

$$\begin{aligned} \sum_{j=0}^{\lfloor N/4 \rfloor} \binom{N}{4j} &= \frac{1}{4} (2^N + (1+i)^N + (1-i)^N) \quad \text{and} \\ \sum_{j=0}^{\lfloor (N-1)/4 \rfloor} \binom{N}{4j+1} &= \frac{1}{4} (2^N - i(1+i)^N + i(1-i)^N) , \end{aligned}$$

where i is the imaginary unit equal to $\sqrt{-1}$.

Proposition 21. *When S follows D_0 we have $\text{LP}_{\max}^S = 2^{-(n+1)}$.*

Proof. Each X_i is in \mathbb{Z}_4 so that it can be described by two bits, denoted $X_i^H X_i^L$. If S is considered like a bit string, a linear distinguisher will be defined by a hash function h such that

$$h(S) = \left(\bigoplus_{j=1}^{n+1} a_j X_j^L \right) \oplus \left(\bigoplus_{j=1}^{n+1} b_j X_j^H \right) ,$$

where $a_1, \dots, a_{n+1}, b_1, \dots, b_{n+1} \in \{0, 1\}$ with at least one non-zero value. We easily prove that

$$X_{n+1}^L \oplus Y = \bigoplus_{j=1}^n X_j^L \quad \text{and} \quad X_{n+1}^H = \left(\bigoplus_{j=1}^n X_j^H \right) \oplus \left(\bigoplus_{j < k \leq n} X_j^L X_k^L \right) \oplus \left(\bigoplus_{j=1}^n X_j^L Y \right) .$$

Thus, if B denotes the value of the bit $h(S)$, we have

$$B = \left(\bigoplus_{j=1}^n (a_j \oplus a_{n+1}) X_j^L \right) \oplus \left(\bigoplus_{j=1}^n (b_j \oplus b_{n+1}) X_j^H \right) \oplus a_{n+1} Y \\ \oplus \left(b_{n+1} \bigoplus_{1 \leq j < k \leq n} X_j^L X_k^L \right) \oplus \left(b_{n+1} \bigoplus_{j=1}^n X_j^L Y \right).$$

If $b_{n+1} = 0$ we can see that (as at least one of the $a_1, \dots, a_{n+1}, b_1, \dots, b_n$ is strictly positive) $\Pr_{D_0} [B = 0] = \frac{1}{2}$, hence $\text{LP}(B) = 0$. If $b_{n+1} = 1$, we have

$$B = \left(\bigoplus_{j=1}^n (a_j \oplus a_{n+1}) X_j^L \right) \oplus \left(\bigoplus_{j=1}^n \bar{b}_j X_j^H \right) \oplus a_{n+1} Y \\ \oplus \left(\bigoplus_{1 \leq j < k \leq n} X_j^L X_k^L \right) \oplus \left(\bigoplus_{j=1}^n X_j^L Y \right).$$

If one of the \bar{b}_j 's is non-zero, then B is uniformly distributed so $\text{LP}(B) = 0$. We now assume that $b_j = 1$ for all $j = 1, \dots, n$. We have

$$B = \left(\bigoplus_{j=1}^n (a_j \oplus a_{n+1}) X_j^L \right) \oplus a_{n+1} Y \oplus \left(\bigoplus_{1 \leq j < k \leq n} X_j^L X_k^L \right) \oplus \left(\bigoplus_{j=1}^n X_j^L Y \right).$$

Let us define $U_j = X_j^L \oplus a \oplus a_j$ for $j \in \{0, \dots, n\}$ and $U_j = Y \oplus a$ for $j = n+1$, with $a = \bigoplus_{j=1}^{n+1} a_j$. We can show that

$$B = \left(\bigoplus_{1 \leq j < k \leq n+1} U_j U_k \right) \oplus c,$$

where $c \in \{0, 1\}$ is a constant. Using Lemma 19 and denoting W the Hamming weight of the random string of bits U_1, \dots, U_{n+1} we obtain

$$\Pr [B = c] = \Pr \left[\frac{W(W-1)}{2} \equiv 0 \pmod{2} \right] = \Pr [W \bmod 4 = 0 \text{ or } 1] \\ = \frac{1}{2^{n+1}} \sum_{j=0}^{\frac{n+1}{4}} \binom{n+1}{4j} + \frac{1}{2^{n+1}} \sum_{j=0}^{\lfloor \frac{n}{4} \rfloor} \binom{n+1}{4j+1}.$$

Using Lemma 20 we deduce

$$\Pr [B = c] = \frac{1}{2} + \frac{(1+i)^n + (1-i)^n}{4 \times 2^n} = \frac{1}{2} + \frac{\cos\left(\frac{n\pi}{4}\right)}{2^{\frac{n}{2}+1}} = \frac{1}{2} + \frac{(-1)^{\frac{n+1}{4}}}{2^{\frac{n+3}{2}}},$$

where we used the fact that $n+1$ is a multiple of 4. Finally, $\text{LP}_{\max}^S = 2^{-(n+1)}$. \square