

# TEN YEARS OF CUBE ATTACKS

MARCO CIANFRIGLIA, ELIA ONOFRI, SILVIA ONOFRI, AND MARCO PEDIKINI

**ABSTRACT.** In 2009, Dinur and Shamir proposed the cube attack, an algebraic cryptanalysis technique that only requires black box access to a target cipher. Since then, this attack has received both many criticisms and endorsements from crypto community; this work aims at revising and collecting the many attacks that have been proposed starting from it. We categorise all of these attacks in five classes; for each class, we provide a brief summary description along with the state-of-the-art references and the most recent cryptanalysis results. Furthermore, we extend and refine the new notation we proposed in 2021 and we use it to provide a consistent definition for each attack family. Finally, in the appendix, we provide an in-depth description of the *kite attack framework*, a cipher independent tool we firstly proposed in 2018 that implements the kite attack on GPUs. To prove its effectiveness, we use Mickey2.0 as a use case, showing how to embed it in the framework.

## 1. INTRODUCTION

Modern cryptographic algorithms are usually based on problems whose difficulty is provable. A golden problem, extensively used in many cryptographic applications, consists in solving large systems of multivariate polynomial equations.

The problem is in fact NP-complete also under basic assumptions, like quadratic equations in  $\mathbb{F}_2$  only.

On the other hand, it is also true that any algorithm may be seen, extensionally, as a black-box computing a boolean function. This is a-fortiori true for cryptographic algorithms where outputs can be sketched as the bits generated by inputs evaluation regardless of the intrinsic structure of algorithm itself. Moreover, under coercion of domain and codomain having the algebraic structure of finite field, such a representation of the algorithm can be built from the evaluation of enough points in the domain. This analysis paves the way for a wide variety of cryptanalysis techniques based on the reformulation of a crypto-system as a polynomial function over  $\mathbb{F}_2$ .

Such a cardinal problem in cryptography has seen the development of many different techniques in recent years that aimed to solve it efficiently. We annoverate Gröbner bases (see [34]) and linearisation techniques (see XL in [19] and XSL in [20]) as some of the most promising techniques actually fading after the same fate: Gröbner bases, in particular, despite being a very general and versatile solution, are unfeasible in many practical cases due to its computational cost.

Such failures led to the objective of finding useful algebraic relations between cryptographic schemes' input and output as a research topic meagre in results for a number of years.

---

1991 *Mathematics Subject Classification.* Primary: 94A60; Secondary: 11T71.

*Key words and phrases.* Cryptanalysis, cube attacks, kite attack, algebraic attacks, GPU implementation, tweakable black box polynomials, division property, Mickey2.0.

The first author has been partially supported by Institute for Applied Computing 'Mauro Picone' IAC - CNR.

In 2009, however, everything changed as a novel approach introduced by Dinur and Shamir at Eurocrypt’09 brought nourishment to algebraic cryptanalysis research branch. In [27], the first paper of a long run, authors introduced, in fact, the *Cube Attack*, a fresh technique that in a few years became a big family of attacks, consisting of many variants including property testers [6] and differential trails [94] amongst the others.

The newly born family was characterized by the just introduced concept of analysing a cypher as a black-box tweakable polynomial, a polynomial where some variables could be set at will during the attack. We feel comfortable to admit that the valuable original idea here exploited is not the effort from the sole Dinur and Shamir. They surely have the merit of clearly pointing out crucial steps and making organised this technique in all its aspects; higher-order differential cryptanalysis mentioned by Lai [51] and Knudsen [50] in the 1990s probably contributed to its successful development, as well as Vielhaber’s AIDA (algebraic IV differential attacks) [86].

The scope of this paper is to give a wide view on the complex and entangled development of the cube attack family, by unravelling the various contributions and presenting them in a more cohesive view. The work is organised as a survey of the many concepts developed around the cube attack. Nevertheless, we work all the examples of application of the cube attack, to uniform to a common notation we introduced in [65].

In Section 2, we revise and extend our novel notation, already introduced in [65], and we revise the original cube attack applied on the binary field  $\mathbb{F}_2$ . Then we extend it to a general finite field  $\mathbb{F}_q$ .

Section 3 provides a summary of various techniques that fall in the macro-family of cube attacks.

In Section 4, we revise the few attempts of implementing cube attack techniques, while we correlative it with Section 5, where we provide a list of the best known attacks on real-world ciphers due to the various families techniques.

Finally, we conclude this work with some considerations in Section 6 and we provide in Appendix A a spot “dive-in” into one of the most recent frameworks for cube attacks, implementing the *kite attack*.

## 2. THE CUBE ATTACK

**2.1. Notation.** Any field of research has its own nomenclature and cryptography makes no exception. However, cube attacks in particular never saw a common agreement in how to refer to its various key concepts. For this reason, each research line developed its own language and notation, often incoherent with the others. Do mind, as a simple example, that the name “Conditional Cube Attacks” is overloaded in literature: usually, it refers to an extension of “Cube Testers”, however it is used also in relation to cube attacks where conditions are imposed a priori.

Diversity in notations makes it difficult to agree upon the specific novel contributions brought by each research; it is not unusual that concepts are claimed as novel and revolutionary, while simply being reformulations of already known results.

We found this as a valid motivation to propose in [65] a novel nomenclature that encloses the various approach of this field. In the following, we revise such notation while enriching it to make it even more inclusive.

Usually, a cipher is a function  $f$  defined over  $\mathbb{F}_q$  as:

$$(1) \quad f : \begin{array}{ccc} \mathbb{F}_q^n \times \mathbb{F}_q^m & \rightarrow & \mathbb{F}_q \\ (\underline{x}, \underline{v}) & \mapsto & c \end{array} .$$

where  $\underline{x} = (x_1, \dots, x_n) \in \mathbb{F}_q^n$  represents a private key,  $\underline{v} = (v_1, \dots, v_m) \in \mathbb{F}_q^m$  represents a public vector and  $c$  is the output value. It is possible to have vectors or stream of values as well (*e.g.*, a keystream); in such cases we consider a function  $f$  per output component.

In a natural way we are able to reformulate the function  $f$  as its Algebraic Normal Form (ANF) representation: a polynomial  $\mathbf{p}$  defined in the equivalence classes of the polynomial ring with variables in  $x_1, \dots, x_n, v_1, \dots, v_m$  and coefficients in  $\mathbb{F}_q$ :

$$\mathbb{F}_q[x_1, \dots, x_n, v_1, \dots, v_m] \text{ modulo } x_1^q - x_1, \dots, x_n^q - x_n, v_1^q - v_1, \dots, v_m^q - v_m .$$

We omit the modulus for the sake of readability by considering exponents of variables  $y_i^{e_i}$  in integers  $0 \leq e_i \leq q - 1$ . We cast  $\mathbf{p} \in \mathbb{F}_q[\underline{x}, \underline{v}]$  when it is important to distinguish between private and public part, and  $\mathbf{p} \in \mathbb{F}_q[\underline{y}]$ , when it is not (with  $y = (y_1, \dots, y_N), N = n + m$ ).

In the following, we use capitalised  $I$  and  $J$  to identify sets of variable indices in  $\{1, \dots, n\}$ ,  $\{1, \dots, m\}$ , or  $\{1, \dots, N\}$  (as it is clear from the context). Such sets are particularly useful when referring to monomials. This is particularly straightforward in the binary setting  $q = 2$  (see later in Section 2.3 for the same notation in fields of higher-order  $q = p^k > 2$ ); in fact, monomials  $\mathbf{m} \in \mathbb{F}_2[\underline{y}]$  are in bijective correspondence with subsets  $I \subset \{1, \dots, N\}$ ,  $|I| = d$  so that any index set  $I$  corresponds to a monomial:

$$(2) \quad \mathbf{m}_I := \prod_{i \in I} y_i \in \mathbb{F}_2[\underline{y}] ,$$

where  $d := \deg(\mathbf{m}_I) = |I|$ .

We then introduce the following notations:

**Zero vector ( $\underline{0}$ ):** is a generic-length 0 vector, meaning that

$$\underline{y} = \underline{0} \quad \text{represents} \quad \underline{y} = (0, \dots, 0)$$

**Unit vector ( $\underline{i}$ ):** is a unitary basis vector, meaning that

$$\underline{y} = \underline{i} \quad \text{represents} \quad \underline{y} = (0, \dots, 0, 1, 0, \dots, 0) ,$$

where the  $i$ -th variable only is set to 1.

**Unit vector set ( $\underline{I}$ ):** is the set of unit vectors obtained from the indices  $i \in I$ ; namely, the underlined notation is mapped through  $I$ , *i.e.*

$$\underline{I} := \{\underline{i} | i \in I\} .$$

**Explicit concat ( $::$ ):** is the concatenation notation (omitted when it is not necessary), meaning that:

$$(1, 1, 0) :: (1, 0, 1) = (1, 1, 0, 1, 0, 1) .$$

**Vector copy ( $\underline{y}^l$ ):** is the notation to  $l$ -times repeat  $y$ , meaning that

$$\underline{y}^3 = \underline{y} :: \underline{y} :: \underline{y} .$$

Do also note that  $\underline{x} = \underline{0}$  is equivalent to  $\underline{x} = 0^n$ .

**Addition ( $+$ ):** is the component wise sum meaning that

$$(1, 1, 0) + (1, 0, 1) = (2, 1, 1) .$$

The sum behaviour depends on the underlying finite field; *e.g.* in  $\mathbb{F}_2$ , it operates as the **xor**, meaning that

$$1\ 1\ 0 + 1\ 0\ 1 = 0\ 1\ 1 ,$$

where we write binary vectors as bit sequences.

Do also note that the following equality always holds:

$$\underline{x} :: \underline{v} = \underline{x} :: 0^m + 0^n :: \underline{v}$$

**Polynomial mapping  $\mathfrak{p}(S)$ :** is a compact notation for applying a polynomial or a function through a set, as it occurs to unit vector sets, meaning that

$$(3) \quad \mathfrak{p}(S) = \{\mathfrak{p}(s) \mid s \in S\} .$$

**Partial assignment  $(\underline{y}[I \rightarrow \underline{a}])$ :** is a compact notation to perform the assignment of specific variables, meaning that

$$(4) \quad \underline{y}[I \rightarrow \underline{a}] \quad \text{represents} \quad y_i = a_i, \quad i \in I .$$

The partial assignment is particularly useful when it is combined with (3) to perform *partial polynomial evaluations*<sup>1</sup>:

$$\mathfrak{p}(\underline{y}[I \rightarrow \underline{a}]) = \mathfrak{p}(\alpha_1, \dots, \alpha_N), \quad \text{where} \quad \alpha_i = \begin{cases} a_{j_i} & \text{if } i \in I \\ y_i & \text{otherwise} \end{cases} .$$

The same notation also applies when considering public and private variables separately; consider the following polynomial:

$$\mathfrak{p}(x_1, x_2, x_3, v_1, v_2, v_3) = x_1x_2 + v_1v_2 + x_1v_3 + x_2v_2 \in \mathbb{F}_q[\underline{x}, \underline{v}] ,$$

then, given two index sets  $J = \{1\}$  and  $I = \{1, 3\}$ , and two vectors  $\underline{a} = (1)$  and  $\underline{b} = (1, 0)$  we have that  $\mathfrak{p}(\underline{x}[J \rightarrow \underline{a}] :: \underline{v}[I \rightarrow \underline{b}])$  evaluates to:

$$\mathfrak{p}(\underline{x}[\{1\} \rightarrow (1)] :: \underline{v}[\{1, 3\} \rightarrow (1, 0)]) = x_2 + v_2 + x_2v_2 .$$

**Cube notation  $(\underline{y}[I \rightarrow A])$ :** defines a set of copies of the vector  $\underline{y}$  where components specified by  $I$  are assigned to every possible value in a given set  $A$ , meaning that:

$$\underline{y}[I \rightarrow A] = \{\underline{y}[I \rightarrow \underline{a}], \text{ for all } \underline{a} \in A\} .$$

The size  $d = |I|$  is called *dimension of the cube* and each vector in  $A$  has exactly  $d$  components. Usually, the set  $A$  is built as a cartesian product of binary assignments per each variable:

$$A = A_1 \times A_2 \times \dots \times A_d, \quad A_i = \{0, a_i\}, \quad a_i \in \mathbb{F}_q .$$

In particular, since it is a common case, if  $A = \mathbb{F}_2^d = \{0, 1\}^d$ , we omit the  $A$ , and we write:

$$(5) \quad \underline{y}[I] := \underline{y}[I \rightarrow \mathbb{F}_2^d] = \{\underline{y}[I \rightarrow \underline{a}], \text{ for all } \underline{a} \in \{0, 1\}^d\} .$$

To further clarify this cardinal notion, we provide the following two examples:

$$\begin{aligned} \underline{y}[\{2, 4, 5\}] &= \{y_1 :: z_2 :: y_3 :: z_4 :: z_5 :: y_6 :: \dots :: y_N, \\ &\quad \text{for all } \underline{z} = \{z_2, z_4, z_5\} \in \{0, 1\}^3\} . \end{aligned}$$

and

$$\begin{aligned} \underline{y}[\{2, 4\} \rightarrow (\{0, 1\} \times \{0, 3\})] &= \{y_1 :: z_2 :: y_3 :: z_4 :: y_5 :: \dots :: y_N, \\ &\quad \text{for all } z_2 \in \{0, 1\}, z_4 \in \{0, 3\}\} . \end{aligned}$$

---

<sup>1</sup>We should say  $\alpha_i = a_{j_i}$  if  $i \in I$ , where  $j_i$  is the index of the element of  $\underline{a}$  which corresponds to the element  $i$  in  $I$ . However, we trust in readers' adaptability.

**Sum reduce** ( $\sum S$ ): is a compact notation for the sum of all the elements within a set  $S$ , namely

$$\sum S = \sum_{s \in S} s .$$

The notation fits particularly well with cubes notation:

$$\sum \underline{y}[I] = \sum_{\underline{z} \in \underline{y}[I]} \underline{z} = \sum_{\underline{a} \in \{0,1\}^d} \underline{y}[I \rightarrow \underline{a}].$$

Shorten summation notation can be used to derive the characteristic vector of an index set too, meaning that

$$\underline{y} = \sum \underline{I} \quad \text{represents} \quad \underline{y} = (y_1, \dots, y_N) \quad y_i = \begin{cases} 1 & \text{if } i \in I \\ 0 & \text{otherwise} \end{cases} .$$

**Cube & partial assignment** ( $\underline{y}[J \rightarrow \underline{a}, I]$ ): is a compact notation to perform both partial assignment and cube evaluation  $I \cap J = \emptyset$ , namely:

$$\underline{y}[J \rightarrow \underline{a}, I] \quad \text{represents} \quad \underline{y}[J \rightarrow \underline{a}, I \rightarrow \mathbb{F}_2^d] .$$

Do mind that, as we certify at the end of this section, it is common to have the variables with indices in  $I^c = \{1, \dots, N\} \setminus I$  set to zero. In such a case, do note that the following equality holds:

$$(6) \quad \underline{y}[I^c \rightarrow 0^{N-d}, I] = \underline{0}[I] .$$

**Cube & partial assignment in  $\mathbb{F}_2$**  ( $\underline{y}[I_0, I_1, I]$ ): is an even more compact notation applicable in  $\mathbb{F}_2$  to perform both partial assignment and cube evaluation. In  $\mathbb{F}_2$ , in fact, each variable can only be set to either 0 or 1, therefore it makes sense to distinguish three sets: the set ( $I_0$ ) of indices of variables substituted by 0, the set ( $I_1$ ) of indices variables substituted by 1 and the set of indices  $I$  of cube variables. We can then shorten the previous notation as follows:

$$\underline{y}[I_0, I_1, I] = \underline{y}[I_0 \rightarrow 0^{|I_0|}, I_1 \rightarrow 1^{|I_1|}, I] .$$

In general,  $I_0 \cup I_1 \cup I \neq \{1, \dots, N\}$  therefore the result can still depend of some variables. If it is not the case, however, we remove the set  $I_0$  from the list, since it can be derived from the context, therefore writing:

$$\underline{y}[I_1, I] = \underline{y}[(I_1 \cup I)^c \rightarrow \underline{0}, I_1 \rightarrow 1^{|I_1|}, I] = \underline{0}[I_1 \rightarrow 1^{|I_1|}, I] .$$

Finally, if  $I_1 = \emptyset$  ( $I_0 = \{1, \dots, N\} \setminus I = I^c$ ), we adopt the strategy of (6):

$$\underline{y}[I^c, \emptyset, I] = \underline{y}[\emptyset, I] = \underline{0}[I]$$

**Monomial generation** ( $\underline{y}^{\underline{s}}$ ): is a notation to generate a monomial from a characteristic vector  $\underline{s} = (s_1, \dots, s_N)$ ,  $s_i \in \{0, 1\}$ , in other words:

$$\underline{y}^{\underline{s}} = \prod_{i \in \{1, \dots, N\}} y_i^{s_i} .$$

Note that (2) can be written this way as well:

$$\underline{m}_I = \underline{y}^{\sum \underline{I}} .$$

**Monomial set generation** ( $\underline{y}^S$ ): is the natural mapping of the previous notation to a set  $S$  of vectors in  $\mathbb{Z}^N$ , in other words:

$$\underline{y}^S = \{\underline{y}^{\underline{s}} \mid \underline{s} \in S\} .$$

Note that  $\underline{y}^I$  by monomial set generation is the set of variables  $y_i$  such that  $i \in I$ . This notation is particularly useful in the Division Property setting

since, fusing the notation with (5), we get the set of all monomials which divide the monomial  $\mathbf{m}_I$ :

$$\underline{y}^{0[I]} = \{\underline{y}^{\underline{s}} \mid \underline{s} = (s_1, \dots, s_N) \text{ where } s_i \in \{0, 1\} \text{ if } i \in I \text{ and } s_i = 0 \text{ if } i \notin I\}.$$

**Monomials assignments** ( $\mathbf{m}[\circ]$ ): all the assignment notations given for sets of variables are still valid also for monomials, like the partial assignment  $\mathbf{m}[I \rightarrow \underline{a}]$ , the cube operator  $\mathbf{m}[I \rightarrow A]$  and their combinations.

**2.2. Cube attack forefather.** Shamir and Dinur introduced in [27] the cube attack by considering any encryption function represented as a polynomial  $\mathbf{p}$  on the binary field  $\mathbb{F}_2$ . This is a very convenient setting since  $z^2 = z$  and  $z + z = 0$  for any  $z \in \mathbb{F}_2$  and, as in (2), each monomial can be represented by the set of indices of its variables. The attack splits into a key-independent (offline) phase and a key-specific (online) phase.

During the offline phase, the attacker has access to an oracle ciphering machine for  $\mathbf{p}$  and can set  $\underline{x}$  and  $\underline{v}$  at will; the goal of this phase is to find appropriate values of  $\underline{v}$  to get at least  $n$  independent linear equations on the  $\underline{x}$  unknowns. The online phase takes place when the defender sets a specific key vector  $x$ . Once again, we suppose the attacker as able to set the public vector  $\underline{v}$  at will: this assumption requires either a chosen plaintext setting (as it could be for Message Authentication Code generation) or enough spoofing time on randomly generated  $\underline{v}$  (as it could be for authentication challenges, *e.g.*, in Wi-Fi handshaking). The goal of this phase is to reconstruct the  $n$  equations found during the off-line phase and solve the corresponding linear system to retrieve (a portion of)  $\underline{x}$ .

We now describe the two phases in detail:

*Offline phase.* Let  $\mathbf{m}_I$  be the monomial generated by a set of variable indices  $I \subseteq \{1, \dots, N\}$ ,  $|I| = d$ . In actual applications,  $I$  should address public variables only ( $I \subseteq \{1, \dots, m\}$ ); however, since the methodology we are now describing is general, we prefer (also for ease of notation) to consider  $I$  as referring to both private and public variables.

Given the cipher  $\mathbf{p}$  and the monomial  $\mathbf{m}_I$ , according the Division Algorithm there exist  $\mathbf{q}_I$  and  $\mathbf{r}_I$  such that:

$$(7) \quad \mathbf{p}(\underline{y}) = \mathbf{m}_I \cdot \mathbf{q}_I(\underline{y}) + \mathbf{r}_I(\underline{y}) ,$$

where none of the monomials in the reminder  $\mathbf{r}_I$  is divisible by  $\mathbf{m}_I$  (all of them miss at least a variable from  $\underline{y}^I$ ) and none of the variables  $\underline{y}^I$  can be found in the quotient  $\mathbf{q}_I$  (since all of the variables in  $\mathbf{p}$  are of degree 1).

We call  $\mathbf{q}_I$  the *superpoly of  $I$  in  $\mathbf{p}$*  and, if  $\mathbf{q}_I$  is linear, we refer to  $\mathbf{m}_I$  as a  *$d$ -degree maxterm of  $\mathbf{p}$* .

In particular, applying  $\mathbf{p}$  to the cube  $\underline{y}[I]$  exterminates the reminder  $\mathbf{r}_I$  while keeping (when  $\underline{y}[I \rightarrow 1^d]$ ) a single instance of the superpoly  $\mathbf{q}_I$ . We then claim:

**Proposition 1** (cfr.[27]). *The superpoly  $\mathbf{q}_I$ , defined in (7), can be retrieved as:*

$$(8) \quad \mathbf{q}_I(\underline{y}) = \sum \mathbf{p}(\underline{y}[I]) .$$

Surprisingly, provided that  $\mathbf{m}_I$  is a maxterm, (8) gives us a method to numerically determine the ANF of  $\mathbf{q}_I$ , even when  $\mathbf{p}$  is given as a black-box. In fact, since  $\mathbf{q}_I$  is linear it does not contain any of the variables  $\underline{y}^I$  and its ANF is given by:

$$(9) \quad \mathbf{q}_I = a_0 + \sum_{j \notin I} a_j y_j, \quad a_j \in \{0, 1\} .$$

We then claim:

**Proposition 2** (cfr.[27]). *Let  $\mathbf{q}_I$  be the superpoly defined in (7). Its coefficients, as defined in (9) are given by*

$$(10) \quad a_0 = \sum \mathbf{p}(\underline{0}[I]) \quad \text{and} \quad a_j = \sum \mathbf{p}(\underline{j}[I]) - a_0 .$$

As stated above, the linear relations we are building are exploited later in the online phase, where the attacker's aim is to retrieve the private key  $\underline{x}$ . For this reason, cube variables  $I$  are in general chosen amongst the public ones ( $I \subset \{1, \dots, m\}$ ), while the remaining ones ( $I^c = \{1, \dots, m\} \setminus I$ ) are usually tweaked to zero or to any other value  $\underline{s}$  to lower the complexity of the resulting system [30]. In this case, the ANF of  $\mathbf{p}$  assumes the following form:

$$\mathbf{p}(\underline{x} :: \underline{v}[I^c \rightarrow \underline{s}]) = \mathbf{m}_I \cdot \mathbf{q}_{I,\underline{s}}(\underline{x}) + \mathbf{r}_{I,\underline{s}}(\underline{x} :: \underline{v}[I^c \rightarrow \underline{s}]) ,$$

and, therefore, the superpoly only depends on the private key vector  $\underline{x}$ :

$$\mathbf{q}_{I,\underline{s}}(\underline{x}) = \sum \mathbf{p}(\underline{x} :: \underline{v}[I^c \rightarrow \underline{s}, I]) .$$

where (10) assume the following form:

$$(11) \quad a_0 = \sum \mathbf{p}(\underline{0} :: \underline{v}[I^c \rightarrow \underline{s}, I]) \quad \text{and} \quad a_j = \sum \mathbf{p}(\underline{j} :: \underline{v}[I^c \rightarrow \underline{s}, I]) - a_0 .$$

In particular, when  $\underline{s} = \underline{0}$ , we omit to report the subscript  $\underline{s}$ , obtaining:

$$\mathbf{q}_I(\underline{x}) = \sum \mathbf{p}(\underline{x} :: \underline{0}[I]) , \quad a_0 = \sum \mathbf{p}(\underline{0} :: \underline{0}[I]) , \quad a_j = \sum \mathbf{p}(\underline{j} :: \underline{0}[I]) - a_0 .$$

*Online phase.* In online phase we suppose the unknown key  $\underline{x} = \underline{k}$  to be set and secret while the public vector  $\underline{v}$  to be settable at will. Online phase is made of two parts: (i) for each maxterm  $\mathbf{m}_I$  found in the offline phase, let us obtain an evaluation of the superpoly  $\mathbf{q}_I$  via

$$b_{I,\underline{s}} = \mathbf{q}_{I,\underline{s}}(\underline{k}) = \sum \mathbf{p}(\underline{k} :: \underline{v}[I^c \rightarrow \underline{s}, I])$$

and (ii) solving the corresponding linear equations system with  $\underline{x}$  as unknown variables, yielded by

$$a_0 + \sum a_i \cdot k_i = b_{I,\underline{s}} .$$

The complexity of the first part depends on the size and the number of the cubes, as well as on the practical difficulty to set  $\underline{v}$ . The second part can be tackled by means of basic linear algebra algorithms in order to retrieve the full key or a portion of it (depending on the number of independent equations found); gaussian elimination algorithm requires *e.g.*,  $\mathcal{O}(n^3)$  steps.

Both parts of the online phase, despite being computationally intense, are computationally bounded by the complexity of the offline one. Therefore, being able to carry out the offline phase actually breaks (or weakens) a specific cipher in all of its instances.

**2.3. Cube attack in higher order fields.** As we highlight in the previous section, the standard cube attack works when polynomials are given over the binary field  $\mathbb{F}_2$ . This restriction is required to prove (8) which is crucial to the cube attack and derives from fundamental equations in the binary field, namely  $y^2 = y$  and  $y + y = 0$ .

When we place the encryption function in the finite field  $\mathbb{F}_q$  those equations assume a different form that depends on the order  $q$  and on the characteristic  $p$  *i.e.*, if  $q = p^k$ , then  $y^q = y$  and  $p \cdot y = 0$ . Consequently, monomials are no longer one-to-one with the indices sets since each variable can have exponent up to  $q - 1$ . A monomial is therefore defined by a vector  $\underline{s} = (s_1, \dots, s_N)$  of exponents where  $s_i \in \mathbb{Z}_q$ , obtaining the following equation equivalent to (2):

$$(12) \quad \mathbf{m}_{\underline{s}} = \underline{y}^{\underline{s}} ,$$

where the set of variables involved are, as always, denoted by  $I$ , namely

$$I = \{i \in \{1, \dots, N\} \mid s_i \neq 0\} .$$

Given a cipher  $\mathfrak{p}$  and a monomial  $\mathfrak{m}_{\underline{s}}$ , we can derive via Division Algorithm an equation analogous to (7), namely

$$(13) \quad \mathfrak{p}(\underline{y}) = \mathfrak{m}_{\underline{s}} \cdot \mathfrak{q}_{\underline{s}}(\underline{y}) + \mathfrak{r}_{\underline{s}}(\underline{y}) + \mathfrak{r}_I(\underline{y}) ,$$

where none of the monomials in the reminder is divisible by  $\mathfrak{m}_{\underline{s}}$ . Though, such monomials can contain some of the variables  $\underline{y}^I$ , therefore, in order to resemble an analogy with  $\mathbb{F}_2$ , it makes sense to divide the reminder into two parts: monomials that do contain all the variables from  $I$  ( $\mathfrak{r}_{\underline{s}}$ ) and monomials that do not ( $\mathfrak{r}_I$ ).

Dinur and Shamir claimed in [27] the possibility of extending the attack to a generic field, however, the first proof of this approach can be found in [3] due to Agnesse and Pedicini. Their main contribution consists in reworking Proposition 1 to extend it by considering the relation given by (13):

**Proposition 3** (cfr.[3]). *Given a set  $\underline{s}$  of exponents working on the variables defined by  $I$ , the superpoly defined in (7) can be retrieved as:*

$$(14) \quad \mathfrak{q}_I(\underline{y}) = \mathfrak{q}_{\underline{s}}(\underline{y}[I \rightarrow 1^d]) + \mathfrak{r}_{\underline{s}}^I(\underline{y}[I \rightarrow 1^d]) = \sum \mathfrak{p}(\underline{y}[I]^0) - \sum \mathfrak{p}(\underline{y}[I]^1) ,$$

where the set  $\underline{y}[I]$  is partitioned as  $\underline{y}[I] = \underline{y}[I]^0 \cup \underline{y}[I]^1$  and  $\underline{y}[I]^0$  contains those vectors with the same parity as  $\underline{y}[I \rightarrow 1^d]$ .

The parity of  $\underline{y}[I \rightarrow 1^d]$  is a symbolic parity since it depends on the values assigned to  $I^c$  variables, namely:

$$\underline{y}[I]^0 = \{a \in \underline{y}[I] \mid a[I^c \rightarrow 0] = d \pmod{2}\} .$$

This key concept is resumed in 2012 Vargiu Master Thesis [85] and later expanded in [64] where Onofri presented many proofs and computational bounds when  $\mathfrak{m}_{\underline{s}}$  is chosen as in standard cube attack, *i.e.*  $s_i = 1, i \in I$ . If this condition holds, in fact, the cube attack straightforward extends from  $\mathbb{F}_2$  to  $\mathbb{F}_q$  up to a factor  $-1$  (which depends on the parity of the specific element within the cube evaluation); in fact  $\mathfrak{r}^I = 0$ , then (14) shortens to

$$(15) \quad \mathfrak{q}_I(\underline{y}) = \mathfrak{q}_{\underline{s}}(\underline{y}) = \sum \mathfrak{p}(\underline{y}[I]^0) - \sum \mathfrak{p}(\underline{y}[I]^1)$$

and, in particular, we can state, analogously to Proposition 2, that

**Proposition 4** (cfr.[64]). *For any polynomial  $\mathfrak{p}$  in  $\mathbb{F}_q[x]$  and cube  $I$  yielding a maxterm  $\mathfrak{m}_I$ , the superpoly has ANF*

$$\mathfrak{q}_I(\underline{y}) = a_0 + \sum_{j \notin I} a_j y_j, \quad a_j \in \mathbb{F}_q .$$

*Coefficients can be numerically evaluated by*

$$a_0 = \sum \mathfrak{p}(0[I]^0) - \sum \mathfrak{p}(0[I]^1) \quad \text{and} \quad a_j = \sum \mathfrak{p}(j[I]^0) - \sum \mathfrak{p}(j[I]^1) - a_0 .$$

An analogous approach to [3] to extend cube attack in higher order fields can be found in [69], where authors fuse standard cube attack with higher order differentiation technique introduced by Lai in [51]. The key concepts are, to the best of our knowledge, totally comparable; however, the processing is performed under the point of view of differentiation techniques. Here, the main contribution is given by the following observation:



**Proposition 5** (cfr.[69]).

$$(16) \quad \mathbf{q}_I = \Delta_{\underline{m} \times \underline{I}}^{(\sum \underline{m})} \mathbf{p}$$

where we are denoting with  $\underline{m} \times \underline{I}$  the multiset of single-variable single-step differentiation:

$$(17) \quad \underline{m} \times \underline{I} = \{\underline{i} \text{ taken } m_i \text{ times}\}_{i \in I}$$

and the  $\Delta^{(k)}$  notation is the standard definition of multi-differentiation:

$$(18) \quad \Delta_{\underline{a}_1, \dots, \underline{a}_k}^{(k)} \mathbf{p} = \Delta_{\underline{a}_1} \dots \Delta_{\underline{a}_k} \mathbf{p}$$

and the standard differentiation is given by

$$(19) \quad \Delta_{\underline{a}} \mathbf{p}(\underline{y}) = \mathbf{p}(\underline{y} + \underline{a}) - \mathbf{p}(\underline{y}).$$

**2.4. Searching for cubes.** A cardinal point of the cube attack is to efficiently determine if (9) holds, or, in other words, whether the superpoly  $\mathbf{q}_I$  is linear or not. Two main approaches are used in this context: (i) retrieve the maximum degree  $\delta$  of  $\mathbf{p}$  and then consider cubes of dimension  $d = \delta - 1$  or (ii) employ stochastic tests to guess the linearity of  $\mathbf{q}_I$ .

The first approach was firstly introduced and used in [26], where Dinur *et al.* applied it to attack Keccak sponge functions; however, it has limited applications since correctly determining the degree is hard when the  $\mathbf{p}$  is complex. As we see later in the next section, this approach is often reversed instead, employing the cube attack itself to probabilistically determine  $\delta$  (see [6]).

The latter approach is the widely used instead. In the original paper [27], Dinur and Shamir employed the Bloom-Luby-Rubinfeld Test from [13], a linearity test originally developed as a Self-Testing/Correcting with Applications algorithm. Later, a novel test optimised by reusing computations is proposed in [30]. In this sense, however, a notable contribution is by Winter, Salagean, and Phan who proposed in [93] an improved linearity test based on higher order differentiation, enhanced by Moebius transform. Srinivasan *et al.* propose instead a three step algorithm in [73], where filters are applied one after the other to “prove” the linearity of a given black box polynomial at a computational cost of  $\mathcal{O}(2^{d+1}(n^2 + n))$ .

Testing the linearity of randomly chosen superpolies  $\mathbf{q}_I$  proves however, to be inefficient. For this reason, in [27], variables with indices in  $I$  were originally picked accordingly to a random walk on the monomial lattice. “Moving” aleatory, however, does not guarantee a success, therefore efforts were devoted to find a pattern to efficiently select monomials  $\mathbf{m}_I$  while looking for the maxterms. Aumasson *et al.* proposed in [5] an evolutionary algorithm to search for cubes that maximise the number of rounds after which the superpoly is still unbalanced. Also Wang *et al.* in [92] propose a new methodology to find more linear equations from the same cube set. Following this trend, Cianfriglia *et al.* developed in [16] a CUDA framework to parallel control all the cubes within a specific *kite shaped* region of the lattice (see Appendix A).

However, sophisticated algorithms were also developed to avoid manipulating such large cubes directly. Stankowski in [74], for example, introduced a greedy bit set algorithm with  $\mathcal{O}(2^{n+c})$  complexity, later expanded in [45].

By talking about heuristics, cryptographers also tried to reduce the density of the ANFs empirically: two examples can be found in [35] and [59].

Following a different research path, more recently, Ye and Tian developed in [100] a novel algebraic criterion to recover the exact superpoly of useful cubes.

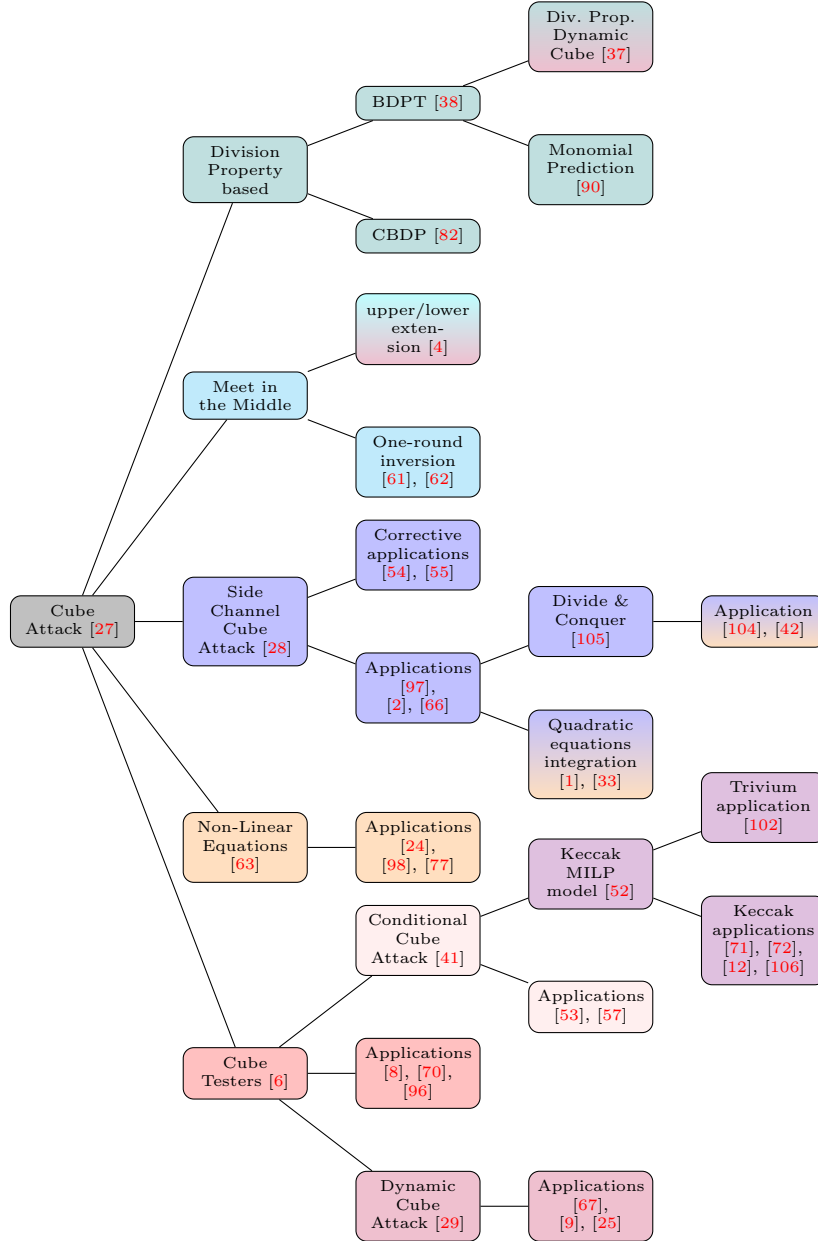


FIGURE 1. Cube Attacks family branchings.

### 3. CUBE ATTACKS FAMILY

The cube attack is a powerful but expensive approach to break ciphers. The idea behind is, however, very solid and flexible and can be combined with many different other approaches to enhance their efficiency. Figure 1 presents various branches traversed by cube attacks.

**3.1. Dynamic cube attack and cube testers.** The first approach in this sense can be found in [6] where cube building is mixed with efficient property-testers in order to detect non-randomness in cryptographic primitives (or either mount cipher distinguishers). The cube framework can, in fact, be exploited in order

to test global properties of a black-box polynomial without retrieving the formal expression of the polynomial.

Such strategy is under the name of *Cube Testers* and combines a property tester on the superpoly (for some property  $\mathcal{P}$ ) with a statistical decision rule that probabilistically recognises whenever the superpoly is  $\delta$ -far from  $\mathcal{P}$ . Namely, the linearity test exploited in the canonical cube attack is itself a cube tester; other examples of properties realisable as cube testers are polynomial randomness (*i.e.* the superpoly coefficients are balanced) and the test of presence of neutral variables (*i.e.* the superpoly does not depend of such a variable).

Cube testers are the basis to create flexible distinguishers as can be seen in [8] first and in [70] later, where authors develop, relying on [74], distinguishers for Trivium based ciphers; however, cube testers main contribution to cryptanalysis can be found in [29] where they are exploited to create *Dynamic Cube Attacks*. The main observation here is that the resistance of many ciphers to cube testers depends on a few number of non-linear operations that usually take place in the latest stages of the encryption process; this is especially true if inputs variables are not mixed enough during the encryption process. Such a behaviour reflects in a very few high order monomials in the ANF of  $\mathfrak{p}$  that, if identified at early stages of the encryption process, can be efficiently killed by vanishing specific input bits – often called *Dynamic variables*. Such dynamic variables, forming a disjoint set from cube variables, usually belong both to public and private ones: the public ones can be set at will during online phases; private ones must instead be guessed and, in particular, these guesses can be confirmed or refuted by cube testers themselves. This, therefore, allows the cryptographer to eventually retrieve key bits without solving any algebraic system at the cost of a more complex offline phase.

Effective usage of this approach can be found, for example, in [67] where authors attack a reduced version of Simon lightweight cipher by using deterministic distinguishers based on cube testers, or in [9] where the author proposes a bi-dimensional dynamic cube attack against 105 round Grain v1 that retrieves nine secret-key-bits of the cipher.

A similar approach is also developed in [41] on Keccak sponge functions, where authors combine cube testers with bit-tracing method (see [91]) to create *Conditional Cube Testers*. Here authors impose a further classification of cube variables, dividing those that mix together after the second encryption round (conditional cube variables) from those that are not multiplied with each other after the first round and are not multiplied with any conditional cube variable after the second round (ordinary cube variables).

The same approach is further improved firstly in [53], where the limitation that no mutual multiplication between cube variables occur in the first round is removed, and then in [57], where more constraints on the number of conditions involving the secret bits are added.

Conditional cube testers were also fused with Mixed Integer Linear Programming by Li *et al.* in [52] and [12], by Song *et al.* in [71] and [72] and, more recently, by Zhao *et al.* in [106] where ciphers of the Keccak family were attacked.

Many efforts focussed recently on novel methods for finding cube testers. A possible strategy is to esteem the probability for the superpoly in selected rounds, as in [23]. Another interesting approach can be found in [96] where Liu *et al.* extended its numeric mapping method for estimating the algebraic degree of NFSR-based cryptosystems (presented in [58]) with the works [74] and [45] by Stankowsky *et al.*

Liu’s numeric mapping method was also employed in a more recent work by Kesarwani *et al.* where, in [46], the authors propose a new algorithm for cube generation following the research branch of [70].

**3.2. Exploiting of non linear equations.** As we pointed out in the previous section, linearity is not the only property we can require on the degree of the superpoly. In particular, we can consider maxterms up to a certain small degree and still recover a polynomial system whose resolution is feasible.

As an example, Mroczkowski and Szmidt propose in [63] an improvement to the cube attack concerning both linear and quadratic equations. They employed a Quadracity Test to retain discarded non-linear equations and use key bits obtained via linear equations to solve “by hand” the quadratic ones. This solution highly enhances the number of key bits the attacker can recover while still limiting the cube search phase: in their application to Trivium-709, for example, they claim no brute-force is needed to recover the whole key.

Combining the ideas of [63] with their previous work [92], Wang *et al.* proposed in [24] a new methodology which makes use of those common variables in two different dimensional cubes to induce maxterms of higher-order from those of lower-order, thus recovering more key bits and reducing the search complexity.

It is also worth of mention the work by Ye and Tian [98], where an experimental approach is employed against Trivium-like ciphers. The authors focus on improving nonlinear superpolys recovery by means of linearisation techniques. Under this setting, several linear and quadratic superpolys are claimed for the 802-round Trivium as well as the possibility of finding a quadratic superpoly for Kreyvium is shown. Relying on specific features discovered on Trivium also an enhanced method to attack Trivium-like ciphers is presented, claiming a generic method of choosing useful nonlinear key expressions.

Clearly, the higher degree equation found this way can be used in many different ways; two more interesting approaches on this side are given by Sun and Guan in [77] where cube attacks are exploited to find new linear relations for linear cryptanalysis purposes and by Eskandari and Ghaemi Bafghi in [32] where non-linear equations are treated as linear equations with noise to attack KATAN lightweight cipher.

**3.3. Cube attacks on side channel attacks.** The original version of cube attack has no free quarters for uncertainty or measurements errors. However, cube attacks have a natural error correction mechanism (see [28]): by considering a cube  $K$  large enough during the offline phase and by evaluating all of its sub-cubes  $I \subset K$  yielding linear relations it is in fact possible to gather redundant linear equations. In the online phase, assuming a per-round leakage with uncertainty (as it happens when Hamming weight only is available), the summation of all the leaked bits from a specific sub-cube assignment yields a new linear equation in the  $\underline{x}$  with a known term depending on the assignment of the known leaked bits. These new relations can be equated to the corresponding linear combination of key variables  $\underline{k}$  pre-evaluated during the offline phase, obtaining a linear system of equations in the  $\underline{x}$  and  $\underline{k}$  variables that the attacker can exploit.

This approach was applied on many block ciphers by exploiting their specific structure starting from [28], where many linear relations were found for Serpent and AES. Later the same year also Yang *et al.* used this approach to analyse Present Lightweight cipher in [97].

The same approach led Abdul-Latip *et al.* to produce two works: in [2] they halved the complexity of NOEKEON block cipher by considering a single bit information leakage from the internal state after the second round; in [1], the authors

modified the cube attack used in [97] by employing some low-degree non-linear equation (*e.g.* quadratic equations) to exploit leakages on PRESENT.

The theory from the previous section was also combined with side channel cube attacks by Fan and Gong in [33] where the security of the Hummingbird-2 cipher (an ultra-lightweight cryptographic algorithm) is discussed. In particular, they describe an efficient term-by-term quadraticity test for extracting simple quadratic equations besides linear ones to be exploited along with the bit-leakage model in a fast GPU model.

Concurrently, also Zhao *et al.* produced an attack to PRESENT in [105] relying on [97] where a two-layer “divide and conquer” strategy is used concurrently with a sliding window approach and an iterated version of the attack is proposed. Their iterative method was further refined in [104] where the authors also propose a model based on non-linear equations.

During the same year, the full version of LBlock was also attacked by means of these techniques in [42].

Li *et al.* also approached side channel cube attacks on PRESENT the same year in [54] after their preliminary work on LBlock of the year before [56]. However, their work focussed on data refinement employing the maximum likelihood decoding algorithm in order to correct the side channel outputs by considering it as a linear code transmitted through a binary symmetric channel with crossover probability depending on the accuracy of the measurements. A 50% success rate is achieved in [55] even when data are more than 40% dirty.

**3.4. Meet in the middle techniques.** One more interesting approach to cube attack is the possibility to fuse it with *meet-in-the-middle* techniques. Firstly suggested in the original paper [27], the first implementation of this approach is to the iconic 120-bit Courtois Toy cipher (CTC) due to Mroczkowski and Szmidi in [61]. Here the offline phase is performed against four rounds of encryption, by recovering many linear equations as usual (here more than 600 linear equations were found). On the online phase, however, the defender encrypts the messages with a five-round encryption. The explicit inversion is therefore performed by obtaining the ciphertext bits after four rounds of encryption by means of equations in the key bits as unknowns and ciphertext bits as known variables. Due to the simplicity of the encryption round rule, these equations are linear in the key bits so, by equating these polynomials to the one gathered in the offline phase, the result is still a linear system that can be solved as in usual cube attack approaches.

Later on, just mimicking what they did earlier in [61], the authors extended the technique to 255-bit Courtois Toy cipher 2 in [62].

A different approach about dynamic cube attack on stream ciphers that is somehow related to MitM techniques can also be found in [4]. Here Ahmadian *et al.* proceed in the opposite direction to usual MitM, by explicitly splitting the cipher into three sections computed independently: an upper extension part, an intermediate section where cube variables are chosen, and a lower extension part.

**3.5. Cube attacks based on division property.** Finally, one of the most recent and promising extensions of cube attack family consists in fusing it with the *division property*, a tool originally introduced by Todo in [80] (later formalised in [81]) as an improvement over Integral Cryptanalysis (see [49]). A multiset  $A$  with elements in  $\mathbb{F}_2^N$  is said to have the *division property*  $\mathcal{D}_k^N$ , with  $0 \leq k \leq N$ , if:

$$\deg(\underline{y}^{\underline{s}}) = \sum \underline{s} < k \quad \text{implies} \quad \sum \underline{y}^{\underline{s}}[\{1, \dots, N\} \rightarrow A] = 0,$$

at the varying of  $\underline{s} \in \mathbb{Z}_2^N$ .

Before pairing it with cube attacks, the concept was firstly extended at FSE 2016, where Todo and Morii in [84] applied it to SIMONS family, introducing the *conventional bit-based division property* (CBDP) and the *bit-based division property using three subsets* (BDPT) and exploiting for the first time the zero-sum property: this solution was more robust than the classical division property, even though it was not efficient enough to carry out a feasible attack.

In order to overcome the efficiency issue, Xiang *et al.* introduced in [94] the *division trails* *i.e.* the propagation of the division property through the rounds of the cipher, along with an approach to evaluate them through MILP (Mixed Integer Linear Programming) models, hence enabling faster computations. More formally, given  $R$  rounds of a cipher, an input  $\underline{y}$  generates, for each round  $r$ , an internal state  $\underline{y}^{(r)}$ . Analogously, a set  $A$  with elements in  $\mathbb{F}_2^N$ , generates  $R$  sets  $(A^{(0)} = A, A^{(1)}, \dots, A^{(R-1)})$ , where  $A^{(r)} = \{\underline{y}^{(r)} \mid \underline{y} \in A\}$ . A division trail for the set  $A$  is a vector  $\underline{k} = (k_0, \dots, k_{R-1})$ , with  $0 \leq k_r \leq N$ , such that the division property  $\mathcal{D}_{k_r}^N$  holds for the set  $A^{(r)}$ , for all  $0 \leq r < R$ . Analysing the round function of the cipher, we can build relations between elements of a trail (*i.e.*, study the propagation of the division property): by writing such relations in a MILP way (relying on the three basic operations of `and`, `xor` and `copy`) we obtain a system of linear inequalities which solutions correspond to valid trails.

The introduction of MILP models allowed Todo *et Al.* in [82] to efficiently apply Division Property along with cube attack, exploiting the CBDP: the non-blackbox representation of the cipher allowed, jointly with the efficient MILP interpretation of the division trail, to obtain unexpected results, hence enabling the authors to break 832-round Trivium. Further results were obtained the following year at Crypto'18, where Wang *et al.* in [87] improved the attack up to 839 rounds of the same cipher.

In [39] Wang *et al.* introduced a new algorithm to find better cubes: to do so, they used a particular MILP model to find division trails based on SAT (see [75]) and on the flag technique (see [88]).

Always thanks to MILP, also BDPT become exploitable in feasible time, as it is shown more recently in [38] where up to 841-round of Trivium were successfully broken.

Later, Wang *et al.* introduced in [90] a novel algebraic version of the division property under the name of *monomial prediction*, also showing its strict similarities with BDPT itself: here, the state variables  $\underline{y}^{(r)} = (y_0^{(r)}, y_1^{(r)}, \dots)$  of round  $r$  are considered as polynomial components  $y_i^{(r)} = \mathbf{p}_{r,i}(\underline{y}^{(r-1)})$  representing the update function of the  $i$ -th component of the state at round  $r$  depending on state components at round  $r - 1$  (hence,  $\mathbf{p}$  can be obtained iterating composition of  $\mathbf{p}_{r,i}$ , round-by-round); these formal relations are then exploited to determine whether specific input state variables  $y_i^{(0)}$  (or, possibly monomials  $\mathfrak{s}$  in the  $\underline{y}^{(0)}$  variables) do or do not propagate to the upcoming rounds. This task can be achieved by analysing round-by-round whether  $\mathbf{p}_{r,i}$  does contain first-degree monomials  $y_j^{(r-1)}$  for some  $j$  or does not, that is, considering the set  $P_{r,i}$  of the monomials  $\mathbf{p}_{r,i}$  is made of (namely, such that  $\mathbf{p}_{r,i} = \sum P_{r,i}$ ), we say that  $y_j^{(r-1)}$  is *monomial predicted* at round  $r$  if:

$$(20) \quad y_j^{(r-1)} \in P_{r,i} \quad \text{for some } j, \text{ component index of the state at round } r - 1 .$$

A set of variables  $(y_{i_0}^{(0)}, y_{i_1}^{(1)}, \dots, y_{i_{R-1}}^{(R-1)})$  such that (20) pairwise holds (*i.e.*,  $y_{i_{r-1}}^{(r-1)}$  is monomial predicted in  $y_{i_r}^{(r)}$ ) is said a *monomial trail*. We then claim:

**Proposition 6** (cfr.[90]). *A given first-round state variable  $y_{i_0}^{(0)}$  can be found in  $\mathbf{p}_{r,i_r}$  if and only if the number of monomial trails connecting them is odd.*

The same proposition holds if considering monomials  $\mathfrak{s}^{(0)}$  in  $\underline{y}^{(0)}$  variables too.

Given a cube  $I$ , Proposition 6 gives us a method to evaluate the superpoly  $q_I$  of the cube attack by exploiting the monomial trails and, hence, by adopting efficient MILP models. In fact, if we consider the set  $P$  of all monomials  $\mathfrak{p}$  is made of (namely,  $\mathfrak{p} = \sum P$ ), we can then reformulate (7) as follows:

$$q_I = \sum M, \quad M = (\underline{y}^{0[I \rightarrow 1^d, I^c]} \cap P) / \mathfrak{m}_I = \{\mathfrak{m} \in \underline{y}^{0[I^c]} \mid \mathfrak{m} \cdot \mathfrak{m}_I \in P\} .$$

The speed-up obtained via MILP modelling, allowed the authors to break Trivium reduced up to 842-rounds [90].

Division property is also linked with other kinds of cube attacks, such as dynamic cube attacks: in [37], Hao *et al.* introduced on one hand a heuristic algorithm using flag technique division property that permits to find superpolies with low bias, on the other hand, a new MILP model method for division property using nullification strategies. With this approach, it was possible to define a new dynamic cube attack on Grain-128 with a success probability of 99.83% and to use the new MILP modelling to attack 892 rounds of Kreyvium.

#### 4. FRAMEWORKS & IMPLEMENTATIONS

Since from its very first introduction, cube attack was presented not only as a theoretical attack, but also as a practical methodology to break real-world ciphers.

For this reason, Aumasson *et al.* built in [5] a first cube tester framework on field-programmable gate array (FPGA) capable of attack 237 rounds in Grain-128 (out of 256) in  $2^{54}$  cipher runs. The idea behind this implementation is hereafter to split the computation in an input generator, an output collector and a controller unit that employs an evolutionary algorithm for the cube searching.

Later FPGA implementation of dynamic cube attacks can also be found in [25] (later revised in [36]) where RIVYERA computing system is adopted.

The main contribution of the previous approaches was however given by the possibility of simultaneously evaluate multiple instances of the cipher in order to fraction the execution times. Following this trend, GPUs were for example employed to test SHA-3 candidates against unbalances, as reported in [44].

Cipher evaluations occurring in cube construction are highly related one to the other and often repeated. In [15], [16], and [14] the Cranic Computing group<sup>2</sup> worked out a complete refactoring of the computation on GPUs in view of repurposing of values already computed. Main contributions are in the organisation of the cube attack as a Time Memory Data Trade-off algorithm, named *kite attack*, to optimise the computation in accord with the structure of GPU memory layers. The development of a CUDA framework for the cube attack resulted in an open source framework enabling the finding of an 800-rounds superpoly in Trivium [17].

As highlighted by Zhu *et al.* in [107], the framework development is a key point not only to check attacks feasibility, but also to show the correctness of many unfitting assumptions cryptographers may claim. In particular, their contribution is under a python-based web application (unfortunately no longer accessible by now) to test cube attacks-like (in particular linearity of given superpoly) on different ciphers (Trivium only was implemented, however, simple extension could be made to integrate other ciphers).

Other notable cube attacks implementations are introduced in [4] as we discuss earlier in Section 3.4 and in [42] where Islam *et al.* develop a GUI toolkit which can load stream or block cipher and can check its resistance against the cube attack.

<sup>2</sup><https://www.cranic.it>

Attack Family	Attack Type	Rounds out of 1152	Maxterms/ Key bits	Time	Bibliography term
cube attack	Key Recovery	735	53 M	$2^{30}$	[27]
cube attack	Key Recovery	767	35 M	$2^{45}$	[27]
cube tester	Distinguisher	790	–	$2^{30}$	[6]
cube tester	Non-randomness	885	–	$2^{27}$	[6]
cube tester	Distinguisher	806	–	$2^{44}$	[74]
cube tester	Non-randomness	1078	–	$2^{54}$	[74]
cube tester	Distinguisher	806	–	–	[48]
non-linear eqs.	Key Recovery	709	full	–	[63]
non-linear eqs.	Key Recovery	799	full	$2^{39}$	[35]
linear extension	Direct Key Rec.	576	26 M	–	[24]
cube like	Distinguisher	839	–	$2^{37}$	[59]
cube like	Key Recovery	576	69 M	–	[73]
cube like	Key Recovery	703	–	–	[93]
bias cube tester	Distinguisher	823	–	$2^{42.74}$	[8],[70]
cube attack	Key Recovery	576	69K	$2^{12.63}$	[43]
Kite Attack	Key Recovery	799	15 M	$2^{45.3}$	[15]
Kite Attack	Key Recovery	800	1 M	$2^{46.3}$	[16]
MILP CDBF cube	Key Recover	832	–	–	[82]
non-linear eqs.	Key Recovery	802	7 M	–	[98]
Division property	Distinguisher	838	–	–	[83]
MILP CDBP cube	Key Recovery	839	–	–	[87]
Div. prop. framework	Key Recovery	805	full	$2^{41.4}$	[101]
Algebraic recovery	Key Recovery	838	5 M	$2^{37}$	[100]
cube tester	Distinguisher	850	–	–	[46]
MILP BDPT cube	Key Recovery	839	full	$2^{78.6}$	[40]
MILP BDPT cube	Key Recovery	841	–	–	[38]
MILP BDPT cube	Key Recovery	978	1 K	$2^{28.5}$	[102]
MILP BDPT cube	Non-randomness	1108	–	$2^{28.5}$	[102]
MILP monomial pred.	Key Recovery	842	–	–	[90]
MILP monomial pred.	Key Recovery	843	2 M	$2^{79}$	[78]

TABLE 1. Results on Trivium cipher.

Grain-128					
Attack Family	Attack Type	Rounds out of 256	Maxterms/ Key bits	Time	Bibliography term
FPGA tester	Distinguisher	237	–	$2^{54}$	[5]
cube tester	Distinguisher	246	–	$2^{42}$	[74]
cube tester	Non-randomness	full	–	–	[74]
dynamic cube	Key Recovery	207	80 K	$2^{31}$	[29]
dynamic cube	Key Recovery	250	theo	$2^{101}$	[29]
dynamic cube	Key Recovery	full	theo	$2^{113}$	[29]
dynamic cube	Key Recovery	full	full	$2^{90}$	[25],[36]
Kite Attack	Key Recovery	160	70000 M	–	[17]
DP dynamic cube	Key Recovery	full	3	$2^{97.86}$	[37]
cube tester	Distinguisher	191	–	$2^{33.86}$	[23]
Grain-v1					
Attack Family	Attack Type	Rounds out of 160	Maxterms/ Key bits	Time	Bibliography term
cube tester	Distinguisher	90	–	$2^{39}$	[74]
cube tester	Non-randomness	96	–	$2^7$	[74]
cube attack	Key Recovery	75	19 M	–	[92]
dynamic cube	Key Recovery	105	9 K	$2^{34}$	[9]
dynamic cube	Key Recovery	100	full	$2^{47}$	[68]

TABLE 2. Results on Grain-128 and Grain-v1 ciphers.

Finally, Ye and Tian introduced in [101] a framework for Trivium efficient key-recovery where Stankovski’s Greedy bit set algorithm fuses with division property and Improved Moebius Transformation to construct potentially good cubes.

## 5. APPLICATIONS

The cube attack family focussed since its beginning on stream ciphers like Trivium (Kreyvium, Quavium, ...) and Grain (Grain-v1, Grain-128, ...). We report the respective main results in Table 1 and Table 2.

Also PRESENT cipher is entitled of an honourable mention, as many developments in side-channel cube attacks were performed on this cipher. Table 3 reports the principal contributions.



Attack Family	Leakage round	Leaked data	Error toll.	Key bits	Time bound	Data required	Biblio. term
PRESENT-80							
cube attack	3rd	0,1,2,3	0%	48	$2^{32}$	$2^{15}$	[97]
non-linear eqs.	after 1	Hamming	0%	64	$2^{16}$	$2^{13}$	[1]
cube attack	3rd	4,8,12	0%	48	–	$2^{11.92}$	[105]
iterated cube	4rd	0	0%	72	–	$2^{15.154}$	[105]
non-linear iterated	after 3	Hamming	0%	72	–	$2^{8.95}$	[104] <sup>†</sup>
max likelihood	after 1	LSB	0.6%	64	$2^{21.6}$	$2^{18.9}$	[54]
max likelihood	after 2	2nd LSB	0.4%	64	$2^{20.6}$	$2^{23.1}$	[54]
max likelihood	after 1	LSB	19.4%	64	$2^{21.6}$	$2^{10.2}$	[54]
max likelihood	after 1	LSB	23.2%	64	$2^{31.6}$	$2^{10.1}$	[55]
max likelihood	after 1	LSB	29.5%	64	$2^{27.6}$	$2^{16.2}$	[55]
max likelihood	after 1	LSB	40.5%	64	$2^{27.6}$	$2^{21.2}$	[55]
PRESENT-128							
non-linear eqs.	after 1	Hamming	0%	64	$2^{64}$	$2^{13}$	[1]
iterated cube	4rd	0	0%	85	–	$2^{15.156}$	[105]
non-linear iterated	after 3	Hamming	0%	121	–	$2^{9.78}$	[104] <sup>†</sup>

LSB states the Least significant bit in the hamming weight of the internal state bytes.

Error tolerant methods all have success probability above 50%.

<sup>†</sup>: tested on real devices with SC countermeasures like random delay and masking.

TABLE 3. Key recovery results via side channel attack on PRESENT cipher with key length of 80 and 128 bit.

Attack Family	Attack Type	Rounds out of 24	Time	Memory	Bibliography term
Keccak-MAC-128					
Cube like	Key Recovery	6	$2^{66}$	$2^{32}$	[26]
Conditional Cube	Key Recovery	6	$2^{40}$	–	[41]
Divide-and-Conquer	Key Recovery	6	$2^{45}$	$2^{13}$	[98]
MILP-aided Cube-like	Key Recovery	6	$2^{42}$	$2^9$	[12]
Cube-like	Key Recovery	7	$2^{97}$	$2^{32}$	[26]
Conditional Cube	Key Recovery	7	$2^{72}$	–	[41]
Divide-and-Conquer	Key Recovery	7	$2^{84}$	$2^{64}$	[98]
MILP-aided Cube-like	Key Recovery	7	$2^{80}$	$2^{15}$	[12]
Cube-like	Forgery	7	$2^{65}$	–	[26]
Keccak-MAC-256					
Cube-like	Forgery	8	$2^{129}$	–	[26]
Keccak-MAC-512					
Conditional Cube	Key Recovery	6	$2^{58.3}$	–	[52]
Conditional Cube	Key Recovery	6	$2^{40}$	–	[72]
Conditional Cube	Key Recovery	7	$2^{111}$	–	[71]
Conditional Cube	Key Recovery	7	$2^{112.6}$	$2^{47}$	[12]
Conditional Cube	Key Recovery	7	$2^{72}$	–	[53]
MILP-aided Cube-like	Key Recovery	7	$2^{108}$	$2^{108}$	[106]

TABLE 4. Results on Keccak sponge function.

Even if cube attacks work on ciphers by considering them as black-box polynomial and therefore are suitable to attack nearly any cryptosystem, they can also exploit specific cipher vulnerabilities. It is the case, for example, of the work performed by Dinur and Shamir first, and by many other cryptographers later, on Keccak family (Ketje, Keyak, ...). We report in Table 4 the principal results obtained against Keccak sponge function.

Many other cipher has been attacked via cube family. It is the case of lightweight and ultra-lightweight ciphers like SIMONs ([84], [67], ...), Simeck ([103]) KATAN ([48], [90], ...), Subterranean 2.0 ([57]), Hitag2 ([76]), LBlock ([95], [42], ...), Hummingbird-2 ([33]), TinyJAMBU ([79]), Ascon ([31]), MORUS ([39]) and many others.

## 6. CONCLUSIONS

In this paper, we revise and improve a novel notation for cube attacks family. We employ this notation to analyse and provide a cohesive review of the state-of-the-art for this wide family of cryptanalysis techniques.

We discuss the original Dinur and Shamir’s attack in  $\mathbb{F}_2$  and we extend it in a generic finite field  $\mathbb{F}_q$ , also providing a description of recent methodologies employed to find cubes. We summarise the family of attacks in five principal research branches: (i) Cube Testers and its extensions (Dynamic and Conditional Cube Attacks), (ii) Cube Attacks with non linear equations, (iii) Cube Attacks with information leakages, (iv) Meet in the Middle cube attacks, and (v) Cube Attacks based on the Division Property and its extensions (based on Division Trails and Monomial Prediction). For what concerns the latter, we also focus on formalising the contributions with the introduced notation, lightening the wordiness of the original one. Later we provide an overview of the few frameworks and implementations currently available. We devote a single appendix to describe in detail our framework implementation of the Kite Attack, where we also present Mickey2.0 as a test case. Finally, we resume in convenient tables all of (to the best of our knowledge) the most significant results obtained through the various approaches applied to the principal attacked ciphers, namely: Trivium, Grain, Present, and Keccak. We believe that cube attacks, in particular combined with DPs, still have a long road to run across.

#### REFERENCES

- [1] S. F. Abdul-Latip, M. Reyhanitabar, W. Susilo, and J. Seberry. Extended cubes: Enhancing the cube attack by extracting low-degree non-linear equations. In *Proceedings of the 6th ACM Symposium on Information, Computer and Communications Security*, pages 296–305, 03 2011.
- [2] S. F. Abdul-Latip, M. R. Reyhanitabar, W. Susilo, and J. Seberry. On the security of NOEKEON against side channel cube attacks. *Information Sec. Practice and Exp.*, 2010.
- [3] A. Agnesse and M. Pedicini. Cube attack in finite fields of higher order. *CRPIT*, 116:9–14, 2011.
- [4] Z. Ahmadian, S. Rasoolzadeh, M. Salmasizadeh, and M. R. Aref. Automated dynamic cube attack on block ciphers: Cryptanalysis of SIMON and KATAN. *IACR C. ePrint A.*, 2015.
- [5] J.-P. Aumasson, I. Dinur, L. Henzen, W. Meier, and A. Shamir. Efficient FPGA implementations of high-dimensional cube testers on the stream cipher Grain-128. *SHARCS09*, 2009.
- [6] J.-P. Aumasson, I. Dinur, W. Meier, and A. Shamir. Cube Testers and key recovery attacks on reduced-round MD6 and Trivium. *Lecture Notes in Computer Science*, pages 1–22, 2009.
- [7] S. Baggage and M. Dodd. The stream cipher mickey 2.0, revised ecrypt stream cipher submission. <http://www.ecrypt.eu.org/stream/p3ciphers/mickey/mickeyp3.pdf>.
- [8] A. Baksi, S. Maitra, and S. Sarkar. New distinguishers for reduced round Trivium and trivium-sc using cube testers. In P. Charpin, N. Sendrier, and J.-P. Tillich, editors, *WCC2015 - 9th International Workshop on Coding and Cryptography 2015*, Proceedings of the 9th International Workshop on Coding and Cryptography 2015, pages 1–10. Anne Canteaut, Gaëtan Leurent, Maria Naya-Plasencia, 04 2015.
- [9] S. Banik. A dynamic cube attack on 105 round Grain v1. *Appl. Stat*, 34(2):49–50, 2014.
- [10] M. Belmonte. Twiddle code. Accessed: 2020-11-12.
- [11] T. Beyne, A. Canteaut, I. Dinur, M. Eichlseder, G. Leander, G. Leurent, M. Naya-Plasencia, L. Perrin, Y. Sasaki, Y. Todo, and F. Wiemer. Out of oddity – new cryptanalytic techniques against symmetric primitives optimized for integrity proof systems, 2020. <https://eprint.iacr.org/2020/188>.
- [12] W. Bi, X. Dong, Z. Li, R. Zong, and X. Wang. MILP-aided cube-attack-like cryptanalysis on keccak keyed modes. *Designs, Codes and Cryptography*, 87(6):1271–1296, 2019.
- [13] M. Blum, M. Luby, and R. Rubinfeld. *Linearity Testing/Testing Hadamard Codes*, pages 1107–1110. Springer, Berlin, Heidelberg, 2016.
- [14] M. Cianfriglia. *Exploiting GPUs to speed up cryptanalysis and machine learning*. PhD thesis, Roma Tre University, 2017/18.
- [15] M. Cianfriglia and S. Guarino. Cryptanalysis on gpus with the cube attack: Design, optimization and performances gains. In *2017 International Conference on High Performance Computing & Simulation (HPCS)*, pages 753–760. IEEE, 07 2017.
- [16] M. Cianfriglia, S. Guarino, M. Bernaschi, F. Lombardi, and M. Pedicini. *A Novel GPU-Based Implementation of the Cube Attack*, pages 184–207. Springer, 2017.

- [17] M. Cianfriglia, S. Guarino, M. Bernaschi, F. Lombardi, and M. Pedicini. Kite attack: re-shaping the cube attack for a flexible gpu-based maxterm search. *J. Crypt. Eng.*, 2019.
- [18] M. Cianfriglia and M. Pedicini. Unboxing the kite attack. In R. La Scala, M. Pedicini, and A. Visconti, editors, *De Cifris Cryptanalysis Selected papers from the ITASEC2020 Workshop De Cifris Cryptanalysis: Cryptanalysis a Key Tool in Securing and Breaking Ciphers*, volume 1 of *Collectio CiphRARum*, pages 31–38. Aracne editrice, 2022.
- [19] N. Courtois, A. Klimov, J. Patarin, and A. Shamir. Efficient algorithms for solving overdefined systems of multivariate polynomial equations. In B. Preneel, editor, *Advances in Cryptology — EUROCRYPT 2000*, pages 392–407. Springer Berlin Heidelberg, 2000.
- [20] N. Courtois and J. Pieprzyk. Cryptanalysis of block cyphers with overdefined systems of equations. In Y. Zheng, editor, *ASIACRYPT 2002*, pages 267–287, 2002.
- [21] Using shared memory in CUDA C/C++. <https://devblogs.nvidia.com/using-shared-memory-cuda-cc/>. Accessed: 2020-11-12.
- [22] Nvidia cuda gpu capability. <https://developer.nvidia.com/cuda-gpus>. Accessed: 2020-11-12.
- [23] D. K. Dalai, S. Pal, and S. Sarkar. Some conditional cube testers for grain-128a of reduced rounds. *IEEE Transactions on Computers*, 2021.
- [24] L. Ding, Y. Wang, and Z. Li. Linear extension cube attack on stream ciphers. *Malaysian J. Math. S.*, 9:139–156, 2015.
- [25] I. Dinur, T. Güneysu, C. Paar, A. Shamir, and R. Zimmermann. An experimentally verified attack on full Grain-128 using dedicated reconfigurable hardware. In D. H. Lee and X. Wang, editors, *ASIACRYPT 2011*, pages 327–343, 2011.
- [26] I. Dinur, P. Morawiecki, J. Pieprzyk, M. Srebrny, and M. Straus. Practical complexity cube attacks on round-reduced keccak sponge function. *IACR C. ePrint A.*, 2014.
- [27] I. Dinur and A. Shamir. Cube attacks on tweakable black box polynomials. *EUROCRYPT 2009*, pages 278–299, 2009.
- [28] I. Dinur and A. Shamir. Side channel cube attacks on block ciphers. *IACR C. ePrint A.*, 2009:127, 2009.
- [29] I. Dinur and A. Shamir. Breaking Grain-128 with dynamic cube attacks. In A. Joux, editor, *Fast Software Encryption*, pages 167–187, 2011.
- [30] I. Dinur and A. Shamir. Applying cube attacks to stream ciphers in realistic scenarios. *Crypt. Comm.*, 4:217–232, 2012.
- [31] J. E. Duarte-Sanchez and B. Halak. A cube attack on a trojan-compromised hardware implementation of ascon. In *Hardware Supply Chain Security*, pages 69–88. Springer, 2021.
- [32] Z. Eskandari and A. Ghaemi Bafghi. Extension of cube attack with probabilistic equations and its application on cryptanalysis of katan cipher. *The ISC International Journal of Information Security*, 12(1):1–12, 2020.
- [33] X. Fan and G. Gong. On the security of hummingbird-2 against side channel cube attacks. In *Western European Workshop on Research in Cryptology*, pages 18–29. Springer, 2011.
- [34] J.-C. Faugere. A new efficient algorithm for computing Gröbner bases (F4). *Journal of pure and applied algebra*, 139(1-3):61–88, 1999.
- [35] P.-A. Fouque and T. Vannet. Improving key recovery to 784 and 799 rounds of Trivium using optimized cube attacks. In *Fast Software Encryption*, pages 502–517. Springer, 2013.
- [36] T. Güneysu, T. Kasper, M. Novotný, C. Paar, L. Wienbrandt, and R. Zimmermann. High-performance cryptanalysis on RIVYERA and COPACOBANA computing systems. In *HPC Using FPGAs*, pages 335–366. Springer, 2013.
- [37] Y. Hao, L. Jiao, C. Li, W. Meier, Y. Todo, and Q. Wang. Links between division property and other cube attack variants. *IACR Transactions on Symmetric Cryptology*, pages 363–395, 2020.
- [38] Y. Hao, G. Leander, W. Meier, Y. Todo, and Q. Wang. Modeling for three-subset division property without unknown subset: Improved cube attacks against Trivium and Grain-128aead. In *Lect. N. Computer S.*, volume 12105 LNCS, pages 466–495. Springer, 2020.
- [39] Y. He, G. Wang, W. Li, and Y. Ren. Improved cube attacks on some authenticated encryption ciphers and stream ciphers in the internet of things. *IEEE Access*, 8:20920–20930, 2020.
- [40] K. Hu, S. Sun, M. Wang, and Q. Wang. An algebraic formulation of the division property: Revisiting degree evaluations, cube attacks, and key-independent sums (full version), 2020.
- [41] S. Huang, X. Wang, G. Xu, M. Wang, and J. Zhao. Conditional cube attack on reduced-round keccak sponge function, 2017.
- [42] S. Islam, M. Afzal, and A. Rashdi. On the security of lblock against the cube attack and side channel cube attack. In *International Conference on Availability, Reliability, and Security*, pages 105–121. Springer, 2013.

- [43] S. Islam and I. U. Haq. Cube attack on Trivium and A5/1 stream ciphers. In *13th IBCAST*, pages 409–415, 2016.
- [44] A. Kaminsky. Gpu parallel statistical and cube test analysis of the sha-3 finalist candidate hash functions. In *15th SIAM (PP12)*, pages 1–15, 2012.
- [45] L. Karlsson, M. Hell, and P. Stankovski. Improved greedy nonrandomness detectors for stream ciphers. *ICISSP*, 2017.
- [46] A. Kesarwani, D. Roy, S. Sarkar, and W. Meier. New cube distinguishers on nfsr-based stream ciphers. *Designs, Codes and Cryptography*, 88(1):173–199, 2020.
- [47] The official kite-attack github repository. <https://github.com/iac-cranic/kite-attack>. Accessed: 2020-11-12.
- [48] S. Knellwolf, W. Meier, and M. Naya-Plasencia. Conditional differential cryptanalysis of Trivium and KATAN. In *International Workshop on Selected Areas in Cryptography*, pages 200–212. Springer, 2011.
- [49] L. Knudsen and D. Wagner. Integral cryptanalysis. In *Fast Software Encryption*, pages 112–127. Springer, 2002.
- [50] L. R. Knudsen. Truncated and higher order differentials. In *Fast Software Encryption*, pages 196–211. Springer Berlin Heidelberg, 1995.
- [51] X. Lai. *Higher Order Derivatives and Differential Cryptanalysis*, pages 227–233. Springer US, Boston, MA, 1994.
- [52] Z. Li, W. Bi, X. Dong, and X. Wang. Improved conditional cube attacks on keccak keyed modes with MILP method. In *Int. C. Th. Application of Crypt. Information Security*, pages 99–127. Springer, 2017.
- [53] Z. Li, X. Dong, W. Bi, K. Jia, X. Wang, and W. Meier. New conditional cube attack on keccak keyed modes. *IACR Transactions on Symmetric Cryptology*, pages 94–124, 2019.
- [54] Z. Li, B. Zhang, J. Fan, and I. Verbauwhede. A new model for error-tolerant side-channel cube attacks. In *International Conference on Cryptographic Hardware and Embedded Systems*, pages 453–470. Springer, 2013.
- [55] Z. Li, B. Zhang, A. Roy, and J. Fan. Error-tolerant side-channel cube attack revisited. In *International Conference on Selected Areas in Cryptography*, pages 261–277. Springer, 2014.
- [56] Z. Li, B. Zhang, Y. Yao, and D. Lin. Cube cryptanalysis of lblock with noisy leakage. In T. Kwon, M.-K. Lee, and D. Kwon, editors, *ICISC 2012*, pages 141–155, 2013.
- [57] F. Liu, T. Isobe, and W. Meier. Cube-based cryptanalysis of subterranean-sae. *IACR Transactions on Symmetric Cryptology*, pages 192–222, 2019.
- [58] M. Liu. Degree evaluation of nfsr-based cryptosystems. In *Annual Int. Crypt. C.*, pages 227–249. Springer, 2017.
- [59] M. Liu, D. Lin, and W. Wang. Searching cubes for testing boolean functions and its application to Trivium. In *2015 IEEE ISIT*, pages 496–500. IEEE, 2015.
- [60] The Mickey2.0 eSTREAM source code. [http://www.ecrypt.eu.org/stream/p3ciphers/mickey/mickey\\_p3source.zip](http://www.ecrypt.eu.org/stream/p3ciphers/mickey/mickey_p3source.zip). Accessed: 2020-11-12.
- [61] P. Mroczkowski and J. Szmids. Cube attack on courtois toy cipher. *IACR C. ePrint A.*, 2009:497, 01 2009.
- [62] P. Mroczkowski and J. Szmids. The cube attack in the algebraic cryptanalysis of CTC2, 2011.
- [63] P. Mroczkowski and J. Szmids. The cube attack on stream cipher Trivium and quadracity tests. *Fundamenta Informaticae*, 114(3-4):309–318, 2012. Republic of MroczkowskiSzmids10.
- [64] E. Onofri. A computational investigation of the cube attack in general finite fields. Master’s thesis, Roma Tre Univ., 3 2020. available at [bit.ly/3FMXPan](http://bit.ly/3FMXPan).
- [65] E. Onofri and M. Pedicini. Novel notation on cube attacks. *Collectio CiphRARum, De Cifris Cryptanalysis, selected papers from the ITASEC2020 workshop*, 2021.
- [66] K.-A. Pang and S. F. Abdul-Latip. Key-dependent side-channel cube attack on craft. *ETRI Journal*, 43(2):344–356, 2021.
- [67] R. Rabbaninejad, Z. Ahmadian, M. Salmasizadeh, and M. R. Aref. Cube and dynamic cube attacks on SIMON32/64. In *11th ISC*, pages 98–103, 2014.
- [68] M. Rahimi, M. Barmshory, M. H. Mansouri, and M. R. Aref. Dynamic cube attack on Grain-v1. *IET Information Security*, 10(4):165–172, 2016.
- [69] A. Sălăgean, M. Mandache-Sălăgean, R. Winter, and R. Phan. Higher order differentiation over finite fields with applications to generalising the cube attack. *Designs, Codes and Cryptography*, 84, 10 2014.
- [70] S. Sarkar, S. Maitra, and A. Baksi. Observing biases in the state: case studies with Trivium and Trivia-SC. *Designs, Codes and Cryptography*, 82(1-2):351–375, 2017.
- [71] L. Song and J. Guo. Cube-attack-like cryptanalysis of round-reduced Keccak using MILP. *IACR Transactions on Symmetric Cryptology*, 2018(3):182–214, 2018.

- [72] L. Song, J. Guo, D. Shi, and S. Ling. New MILP modeling: Improved conditional cube attacks on keccak-based constructions. In *Int. C. Th. Application of Crypt. Information Security*, pages 65–95. Springer, 2018.
- [73] C. Srinivasan, U. Pillai, K. Lakshmy, and M. Sethumadhavan. Cube attack on stream ciphers using a modified linearity test. *J. of Discrete Mathematical Sciences and Cryptography*, 18:301–311, 2015.
- [74] P. Stankovski. Greedy distinguishers and nonrandomness detectors. In *International Conference on Cryptology in India*, pages 210–226. Springer, 2010.
- [75] L. Sun, W. Wang, and M. Wang. Automatic search of bit-based division property for ARX ciphers and word-based division property. In *ASIACRYPT 2017*. Springer, 2017.
- [76] S. Sun, L. Hu, Y. Xie, and X. Zeng. Cube cryptanalysis of hitag2 stream cipher. In *International Conference on Cryptology and Network Security*, pages 15–25. Springer, 2011.
- [77] W.-L. Sun and J. Guan. Novel technique in linear cryptanalysis. *ETRI J.*, 37:165–174, 02 2015.
- [78] Y. Sun. Cube attack against 843-round trivium. *IACR Cryptol. ePrint Arch.*, 2021:547, 2021.
- [79] W. L. Teng, I. Salam, W.-C. Yau, J. Pieprzyk, and R. C.-W. Phan. Cube attacks on round-reduced tinyjambu. *Cryptology ePrint Archive*, 2021.
- [80] Y. Todo. Structural Evaluation by Generalized Integral Property. *Proceedings of EUROCRYPT Part I*, pages 287–314, 2015.
- [81] Y. Todo. Integral Cryptanalysis on Full MISTY1. *J. of Cryptology*, 30(3):920–959, 2017.
- [82] Y. Todo, T. Isobe, Y. Hao, and W. Meier. Cube attacks on non-blackbox polynomials based on division property. In *CRYPTO 2017*, pages 250–279. Springer, 2017.
- [83] Y. Todo, T. Isobe, Y. Hao, and W. Meier. Cube attacks on non-blackbox polynomials based on division property. *IEEE Transactions on Computers*, 67(12):1720–1736, 2018.
- [84] Y. Todo and M. Morii. Bit-based division property and application to simon family. In *International Conference on Fast Software Encryption*, pages 357–377. Springer, 2016.
- [85] M. Vargiu. Fast algebraic cryptanalysis in finite fields of higher order with the cube attack. In *100 tesi di crittografia e codici in Italia. 2008-2017*, Crittografia book series. Murru, N. and Bartoli, D. and Pavese, F., 2020.
- [86] M. Vielhaber. Breaking ONE.FIVIUM by AIDA an algebraic IV differential attack, 2007.
- [87] Q. Wang, Y. Hao, Y. Todo, C. Li, T. Isobe, and W. Meier. Improved division property based cube attacks exploiting algebraic properties of superpoly. *CRYPTO 2018*, 2018.
- [88] Q. Wang, Y. Hao, Y. Todo, C. Li, T. Isobe, and W. Meier. Improved division property based cube attacks exploiting algebraic properties of superpoly. *Lect. N. Computer S.*, 10991 LNCS, 2018.
- [89] S. Wang, B. Hu, J. Guan, K. Zhang, and T. Shi. A Practical Method to Recover Exact Superpoly in Cube Attack. *IACR C. ePrint A.*, 2019.
- [90] S. Wang, B. Hu, J. Guan, K. Zhang, and T. Shi. Exploring secret keys in searching integral distinguishers based on division property. *IACR Transactions on Symmetric Cryptology*, 2020(3):288–304, 2020.
- [91] X. Wang and H. Yu. How to break MD5 and other hash functions. In *Int. C. Th. applications of Crypt. Tech.*, pages 19–35, 2005.
- [92] Y. Wang, L. Ding, W. Han, and X. Wang. The improved cube attack on Grain-v1. *IACR C. ePrint A.*, 2013:417, 2013.
- [93] R. Winter, A. Salagean, and C.-W. Phan, Raphael. Comparison of cube attacks over different vector spaces. In J. Groth, editor, *Cryptography and Coding*, pages 225–238, 2015.
- [94] Z. Xiang, W. Zhang, Z. Bao, and D. Lin. Applying MILP method to searching integral distinguishers based on division property for 6 lightweight block ciphers. In *ASIACRYPT 2016*, pages 648–678. Springer, 2016.
- [95] Z. Xiang, W. Zhang, Z. Bao, and D. Lin. Applying MILP method to searching integral distinguishers based on division property for 6 lightweight block ciphers. *Lect. N. Computer S.*, 10031 LNCS:648–678, 2016.
- [96] J. Yang, M. Liu, and D. Lin. Cube cryptanalysis of round-reduced acorn. In *International Conference on Information Security*, pages 44–64, 2019.
- [97] L. Yang, M. Wang, and S. Qiao. Side channel cube attack on PRESENT. In J. A. Garay, A. Miyaji, and A. Otsuka, editors, *Cryptology and Network Security*, pages 379–391, Berlin, Heidelberg, 2009. Springer.
- [98] C.-D. Ye and T. Tian. A new framework for finding nonlinear superpolies in cube attacks against Trivium-like ciphers. *IACR C. ePrint A.*, 2018. <https://eprint.iacr.org/2018/174>.
- [99] C. D. Ye and T. Tian. Revisit division property based cube attacks: Key-recovery or distinguishing attacks? *IACR Transactions on Symmetric Cryptology*, 2019(3):81–102, 2019.

- [100] C.-D. Ye and T. Tian. Algebraic method to recover superpolies in cube attacks. *IET Information Security*, 14(4):430–441, 2020.
- [101] C.-D. Ye and T. Tian. A practical key-recovery attack on 805-round trivium. *IACR Cryptol. ePrint Arch.*, 2020:1404, 2020.
- [102] C.-D. Ye, T. Tian, and F.-Y. Zeng. The MILP-aided conditional differential attack and its application to Trivium. *Des. Codes Cryptogr.*, page 89, 2020.
- [103] M. Zaheri and B. Sadeghiyan. Smt-based cube attack on round-reduced simeck32/64. *IET Information Security*, 14(5):604–611, 2020.
- [104] X. Zhao, S. Guo, F. Zhang, T. Wang, Z. Shi, H. Liu, K. Ji, and J. Huang. Efficient hamming weight-based side-channel cube attacks on PRESENT. *J. of Systems and Software*, 86(3):728–743, 2013.
- [105] X.-j. Zhao, T. Wang, and S. Guo. Improved side channel cube attacks on PRESENT. *IACR C. ePrint A.*, 2011:165, 2011.
- [106] Z. Zhao, S. Chen, M. Wang, and W. Wang. Improved cube-attack-like cryptanalysis of reduced-round ketje-jr and keccak-mac. *Information Processing Letters*, 171:106124, 2021.
- [107] B. Zhu, W. Yu, and T. Wang. A practical platform for cube-attack-like cryptanalyses. *IACR C. ePrint A.*, 2010:644, 01 2010.

## APPENDIX A. UNBOXING THE KITE ATTACK

Here we describe the *Kite-Attack* framework focusing on its source code and how to extend it. The framework has been designed to be cipher independent; as shown in [15] the cost of the attack differs only by a constant factor when different ciphers are used. This appendix aims at providing a detailed guideline on how to extend the framework support to new ciphers; we believe this work can be useful to the crypto-community as the framework provides an easy way to test/analyse ciphers strength w.r.t. the cube attack. A brief description of the structure of the code framework appeared in [18], however here we provide a more detailed version along with all the steps to add a new cipher. We organise the appendix as follows. We start with a brief introduction to Nvidia GPUs and CUDA jargon<sup>3</sup>. Then we describe in detail the framework and its code structure, we define and describe all the steps needed to add a new cipher and, finally, we show how to add it, Mickey2.0, to the framework by crossing these steps.

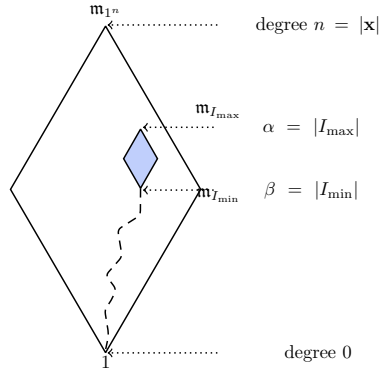
**A.1. CUDA and GPU.** For a better understanding of our work, we report a few, basic, information about the micro-architecture of NVIDIA GPUs as exposed through the CUDA software framework, since this is the solution used in our study. From a hardware standpoint, an NVIDIA GPU is an array of *Streaming Multiprocessors* (SMs); each SM contains a certain number of CUDA *cores*. From a software perspective, a CUDA program is a sequence of computing *grids*; in turn, each grid is split into *blocks*, and each block comprises a certain number of *threads*. Each function executed on the GPU on behalf of the CPU is called *kernel*. To attain a significant fraction of the theoretical peak performance, *occupancy* (*i.e.*, the fraction of active computing elements at a given time) must be consistently kept high, in such a way that thousands of threads must be ready to be scheduled at any time. Threads are executed by an SM in groups of 32 units called *warps*, and performance improves significantly when threads in the same warp execute the same code with no divergence and access memory according to patterns that privilege *threads* locality, *i.e.*, if threads belonging to the same warp access consecutive memory locations (memory *coalescing* in CUDA jargon). Any thread may access data from multiple memory spaces: (*private*) *registers*, (*private*) *local memory*, *shared memory*, *global memory*, and *constant*, *texture* memories that are read-only. *Global memory* is the biggest but slowest memory available and it is persistent across kernel launches by the same application; it can be accessed by all the threads. *Shared memory* is visible to all threads of a block and it has the same lifetime as the block. It is roughly  $100 \times$

---

<sup>3</sup>Readers familiar with the subject may safely skip Section A.1

faster than *global memory* and it can be used for caching or to facilitate memory coalescing in cases where it is not possible otherwise [21]. *Local memory* is actually part of the *global memory* and it is used to provide threads private memory whenever registers are not enough. *Registers* are the fastest memory and they are also used for the warp-level operations called *shuffle* that allow threads belonging to the same warp to exchange data using registers without passing through higher-latency components of the memory hierarchy.

**A.2.  $(\alpha, \beta)$ -Kite attack.** The  $(\alpha, \beta)$ -kite attack, introduced in [17], is based on the choice of a set  $I_{\max}$  of  $\alpha$  public variables and a proper subset  $I_{\min}$  of  $\beta$  variables, with  $\alpha > \beta$ . These two sets represent a *maximal* and a *minimal* cube, respectively  $\underline{x}[I_{\max}]$  and  $\underline{x}[I_{\min}]$ . The name *kite* comes from the observation that the choice of  $I_{\min}$  and  $I_{\max}$  defines a *diamond-shaped* subspace of all possible monomials, with the bottom vertex of the diamond being  $\mathbf{m}_{I_{\min}}$  and the top vertex being  $\mathbf{m}_{I_{\max}}$ . This subspace, schematically depicted in Fig. 2, is made of all monomials  $\mathbf{m}_I$ 's such that  $I_{\min} \subseteq I \subseteq I_{\max}$  and it is exhaustively explored by our attack.



The large diamond is the space of all possible monomials, with the bottom vertex being the constant monomial 1 of degree 0, and the top vertex  $\mathbf{m}_{1^n}$  being the product  $\mathbf{m}_{1^n} = \prod_{i=1}^n x_i$  of all the public variables; the smaller blue diamond is the subspace defined by  $I_{\min}$  and  $I_{\max}$ , whose bottom vertex is  $\mathbf{m}_{I_{\min}}$  of degree  $|I_{\min}| = \beta$  and whose top vertex is  $\mathbf{m}_{I_{\max}}$  of degree  $|I_{\max}| = \alpha$ , which contains all and only the monomials which are divisible by  $\mathbf{m}_I$  and divide  $\mathbf{m}_{I_{\max}}$ .

FIGURE 2. A schematic representation of the *kite attack*.

This definition of the *kite* naturally leads us to a Time Memory Data Trade-Off algorithm where first

- (1) for the given minimal index set  $I_{\min}$  and an initial vector  $\underline{v}$ , we compute many variants of the cube on the index set  $I_{\min}$ : one for each possible combination  $I$  of the indices in  $I_{\max} \setminus I_{\min}$ , we evaluate the encryption function in each cube  $\underline{x} :: \underline{v}[I, I_{\min}]$  for all possible increments of index set  $I \subset \{I_{\max} \setminus I_{\min}\}$  and for any value of  $x \in \{0, 1, \dots, n\}$ ; values are stored in memory to be accessed in a successive moment,
- (2) we iteratively combine previously computed results to evaluate coefficients of the superpoly and test its linearity on larger cubes, namely if we want to step from  $I = \{i_1, \dots, i_d\}$  to  $I' = I \cup \{i_{d+1}\}$ , by keeping the setting of remaining variables as specified by the index set  $I_1$  of variables assigned to 1, we apply the following differentiation formula:

$$\sum \mathfrak{p}(\underline{x} :: \underline{v}[I_1^-, I']) = \sum \mathfrak{p}(\underline{x} :: \underline{v}[I_1^+, I]) + \sum \mathfrak{p}(\underline{x} :: \underline{v}[I_1^-, I])$$

where  $I_1^- := I_1 \setminus \{i_{d+1}\}$  and  $I_1^+ := I_1 \cup \{i_{d+1}\}$  and the increment variable  $i_{d+1} \in I_{\max} \setminus I$  in such a way that  $I'$  always falls in the kite-area ( $I_{\min} \subseteq I' \subseteq I_{\max}$ ).

The implementation following this idea leads to two distinct CUDA kernels: **Kernel1** which is responsible for running (1) and **Kernel2** which runs (2).

A schematic representation of the two kernels, in the case of a minimal example with  $I_{\min} = \{2\}$  and  $I_{\max} = \{1, 2, 3\}$  is reported in Figure 3 and in Figure 4.

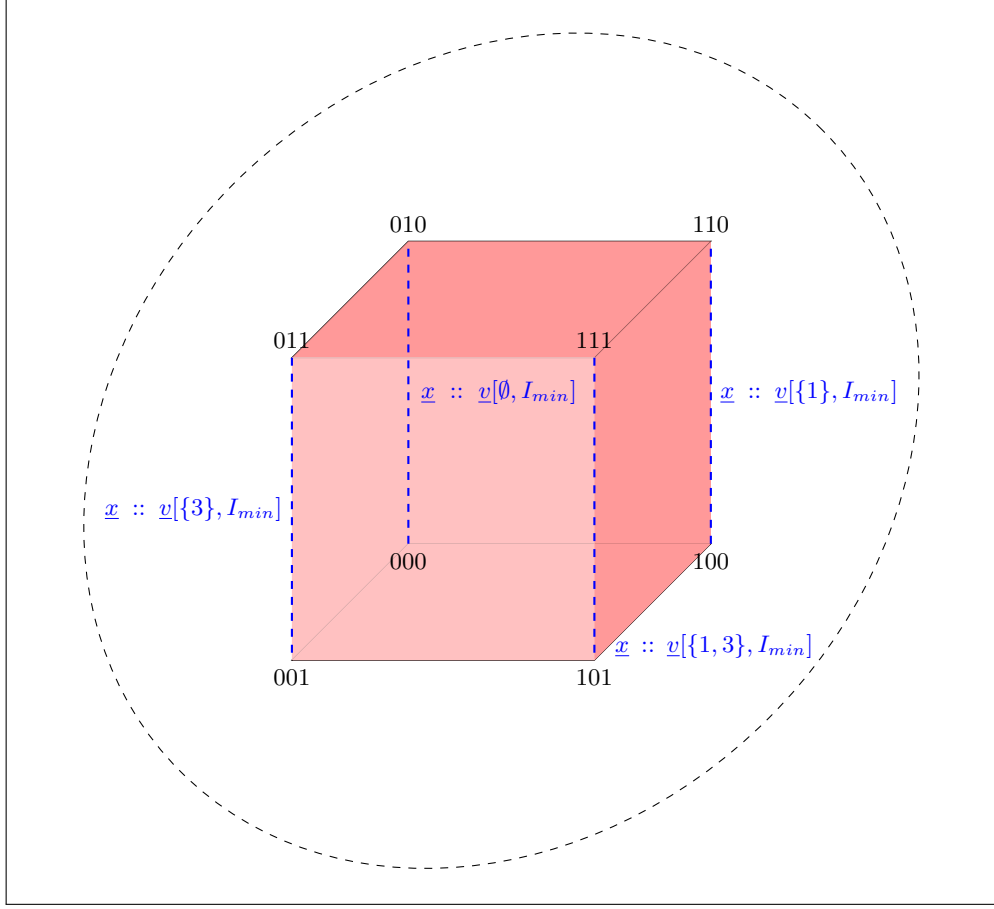


FIGURE 3. A schematic representation of the **Kernel1**. For a given index set  $I_{min}$  computes all the cubes  $\underline{x} :: v[I, I_{min}]$  for all possible increment of index set  $I \subset \{I_{max} \setminus I_{min}\}$ . In the picture  $|I_{min}| = 1$  therefore any cube on  $I_{min}$  contains just two elements and corresponds to blue dashed edges of the three-dimensional cube.

**A.3. Framework code overview.** The framework is composed of three source files (`cubaCUDA.cu`, `twiddle.c` and `auxiliary_functions.c`) along with their corresponding header files. There are two header files: `def.h` which contains all the definitions, macros and includes needed by all the sources, and `key_table.h` must contain two arrays describing how the keys are combined in the linearity tests. Furthermore, each cipher requires the source code for the CUDA implementation and another source file containing auxiliary functions specific to the cipher; for instance, the `setBit` function described below.

The file called `cubeCUDA.cu` includes the two CUDA *kernels*, the main function and other high level functions useful in managing the attack and several steps of computation concerning superpolys. The file called `auxiliary_functions.c` consists of all the auxiliary functions and wrappers used in the framework. `twiddle.c` contains the functions to generate all combinations of  $M$  elements drawn without replacement from a set of  $N$ . This code has been written by M. Belmonte and the original version can be downloaded from [10].



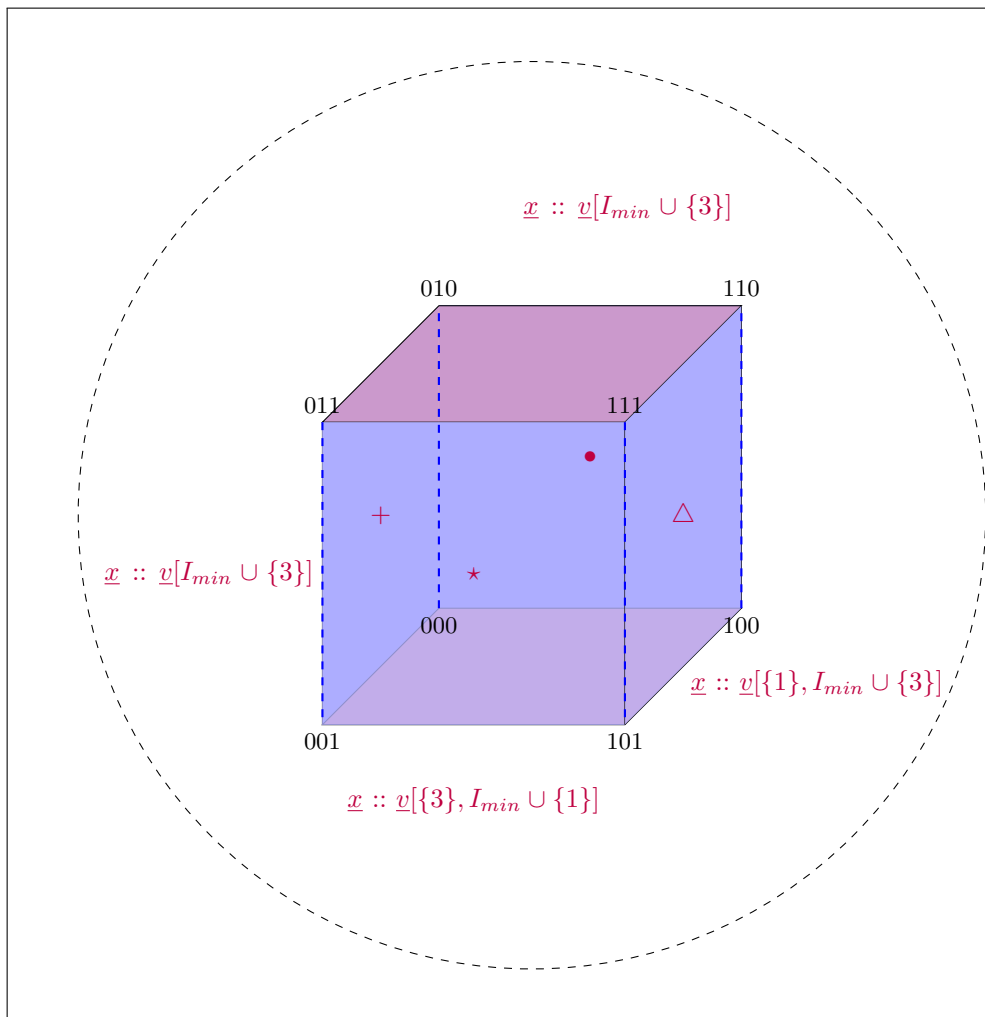


FIGURE 4. A schematic representation of the **Kernel2**. Starting from evaluation results of the first kernel in cubes on  $I_{min}$  stored in memory, it combines the results to obtain values on larger cubes. To obtain evaluation of a face we sum results of evaluations on edges: for instance to obtain evaluations on the face labelled as  $\star$  which is  $\underline{x} :: \underline{v}[\{3\}, I_{min} \cup \{1\}]$ : we have to combine results coming from the two edges  $\underline{x} :: \underline{v}[\{3, 1\}, I_{min}]$  and  $\underline{x} :: \underline{v}[\{3\}, I_{min}]$ . For each cube it also performs the linearity test by exploiting evaluations of the selected cube with different assignments to  $\underline{x}$ .

A GPU run is the sequential call of the `runAttack` function and then of the `computeSuperpoly` one. The first one is in charge of the real attack, it launches the CUDA *kernels* to compute partial sums over  $\underline{x} :: \underline{v}[I_{min}]$  and combines them to test linearity. It also dumps on a binary file all the candidate maxterms found and returns the number of them. The latter function is responsible for reading the binary file containing candidate maxterms, computing the corresponding superpolys, and printing them in human readable format.

Before running the attack, the framework parses the configuration file which contains the following information: the target cipher, the number of initialisation

rounds, the indices belonging to  $I_{\min}$  and those belonging to  $I_{\max} \setminus I_{\min}$ , and an ID string to identify the run.

The framework provides some scripts to interactively generate configuration files. After the setup is complete, it verifies the selected CUDA device is able to run the attack and, if so, it generates all the data needed for the attack and copy them on device memory. For instance, it initialises (i) the vector containing the set of keys used in the attack, (ii) the mask representing  $I_{\min}$  which is composed by setting the  $\beta$  bits with indexes in  $I_{\min}$ , and (iii) the  $2^{\alpha-\beta}$  masks that represent all the possible monomials which are divisible by  $\mathbf{m}_{I_{\min}}$  and divide  $\mathbf{m}_{I_{\max}}$  (see Figure 2). Keys and initial vectors ( $IV$ ) are mapped on contiguous unsigned integer (`u32`); in particular,  $\lceil key\_size/32 \rceil$  and  $\lceil IV\_size/32 \rceil$  unsigned integers are respectively used for each key and  $IV$ . Moreover it allocates the memory to store the output of linearity tests or, in the case of the second kernel, to store coefficients of the superpoly.

Every cipher implementation may adopt its own layout to map keys and  $IV$  bit indexes. For instance, assuming we have a cipher with  $IV\_size$  of `u32`; the bit  $iv_0$  could be mapped to the most significant bit (*msb*) of the most significant byte (*MSB*), to the less significant bit (*lsb*) of *MSB* or to the *lsb* of the less significant byte (*LSB*), and so on. As the framework cannot predict which layout will be used by the cipher, an auxiliary function for each cipher enables the framework to correctly manage any layout. This function basically takes three input parameters: the index of the key/ $IV$  to set, the value to set and the pointer to the first unsigned integer that represents the key/ $IV$ . We are used to name these auxiliary functions `setBit<cipher name>`. For each supported cipher, the following information has to be provided to the framework through the `def.h` file:

- `KEY_SIZE` and `IV_SIZE`: define the number of bits representing respectively the key and the  $IV$ ;
- `KEY_ELEM` and `IV_ELEM`: represent the number of `u32` needed to contain respectively one key and one  $IV$ ;
- `CIPHER_NAME`: is a quoted string containing the cipher name;
- `CIPHER`: is an unquoted string containing the cipher name. This is used to automatically select the `setBit` corresponding to the cipher.
- `KEYS_COEFFICIENT`: is equal to `KEY_SIZE + 1`. It is used for superpolys computations.
- `TOTAL_KEYS`: represents the smallest multiple of 32 greater than `KEYS_COEFFICIENT`
- `RESIDUAL_KEYS`: contains the value `TOTAL_KEYS - KEYS_COEFFICIENT`.

We define some preprocessing macros that automatically select the appropriate `setBit` function once the cipher specific properties are specified in the file `def.h`; of course, the file containing the function implementation should be added to the `Makefile`.

The framework is ready to work with ciphers that support key and  $IV$  of length up to 256 bits. We use other preprocessor macros to setup the framework and kernel functions accordingly to the key and  $IV$  sizes. We adopt this method as it lets us provide optimised code for any size while, at the same time, it keeps the code simple and easy to read. We use the preprocessor macros also to define which cipher function has to be called by the kernels for the above reason.

**A.4. Porting the cipher.** We now describe the hardest step, adapting the cipher function to CUDA. Given a target cipher  $E$ , the first essential step is the definition of the CUDA device function that implements  $E$ . This kernel function should require the key and the  $IV$  as input parameters and should return the corresponding keystream. If the key cannot be stored in just one `u32` word (i.e.  $\geq 32$  bit) it should

be provided as multiple `u32` variables<sup>4</sup> rather than an array of `u32` (i.e. it is better `key1, key2, ..., keyN` than `key[N]`). In this way, elements of the key are placed in registers (if available) otherwise, data are stored on the global memory which has higher latency access time. Of course, `IV` should be treated in the same way.

The function implementing the cipher is called by hundreds of threads simultaneously on different inputs; for this reason, this function must be self-contained, i.e. it should use only (thread) local variables to store partial computations and it should not do anything that can interfere with other computations. In other words, our goal is to implement the target cipher  $E$  efficiently in CUDA such that it can be executed concurrently by thousands of threads. The implementation should be efficient as the cipher function is called  $2^\alpha \times 2^\beta$  times<sup>5</sup> by each thread involved in the computation; so any effort on optimising it will not be vain.

To maximise the attack throughput the cipher function should return 32 bits of keystream, so to fully exploit framework’s capability to test the linearity of 32 polynomials simultaneously; however, this is usually trivial to do.

Finally, as mentioned in Section A.3, the throughput is maximum when all the threads in the same warp execute the same code with no divergence. For this reason, any `if/else` statement should be avoided unless you are sure that the result of the condition is the same at warp level, i.e. all the threads in the warp obtain the same result when tests the condition. In the case this cannot be guaranteed, the `if/else` block of code should be carefully analysed, and, if possible, redesigned with an equivalent block of code that does not contain the branch. An example is reported in Section A.6.

**A.5. Mickey2.0: cipher definition.** MICKEY (Mutual Irregular Clocking KEYstream generator) belongs to eSTREAM portfolio. It is an hardware-efficient stream cipher designed by S. Baggage and M. Dodd [7]. It takes two input parameters, an 80-bit secret key  $K$  and an  $IV$  with variable length between 0 and 80 bits. It is composed of two registers  $R$  and  $S$  of 100-bits each called respectively the *linear* and *non-linear* registers. It defines two functions `clock_R` and `clock_S` to update  $R$  and  $S$  respectively. Differently from other ciphers like Trivium and Grain128, Mickey2.0 does not initialise the registers with key and  $IV$ ; it relies on one specific function instead, called *clock.kg*, that updates both  $R$  and  $S$  by calling `clock_R` and `clock_S`. In the initial steps, the  $IV$  and the key bits are used as input; after these clocks, it runs for 100 more clocks with input 0. Interested readers may find more details in [7].

**A.6. Mickey2: porting to CUDA.** We use as a reference for our porting the *faster* version of Mickey, source code provided in [60]. This version has the advantage that already works with `u32` and efficiently updates register states.

For the other supported ciphers (Trivium and Grain128) we adopted the layout that maps key and  $IV$  of index 0 to the *msb* of the *MSB*; to avoid maintaining multiple layouts, we adopt the same layout also for Mickey2.0 and we define the masks representing the update sequences *COMP0*, *COMP1*, *FB0* and *FB1* and the mask defining the *RTAPS* vector accordingly to the selected layout.

With respect to the original implementation, we do not use auxiliary functions for clocking  $R$  and  $S$  or to initialise the cipher with key and  $IV$ ; we implement all the steps inside the cipher function to avoid the overhead of calling auxiliary functions. Moreover, for the reasons explained in Section A.4, our function do not use arrays for keys and  $IV$  but multiple `u32` words; for instance three `u32` for both

<sup>4</sup>The number of `u32` words is equal to  $\lceil key\_size/32 \rceil$

<sup>5</sup> $2^\alpha$  cubes each of dimension  $2^\beta$ . The cipher function has to be computed for each vertex of the cube.

a key or an IV. All these choices, however, induced us to split the *load IV* and *load key* steps in three loops each. In this way, we duplicate the code but we do not need extra computations to identify which u32 variable is used at every step. We also define a loop that implements *pre-clock* and another one for *keystream generator*. This choice allows us to perform optimisations as described below.

The functions `CLOCK_KG`, `CLOCK_R` and `CLOCK_S` as defined in Mickey2.0 specification, contain some `if` statements. In the follow, we analyse each of them:

- `CLOCK_KG`: the value of the `MIXING` parameter determines how to compute `INPUT_BIT_R`. However the result of this check is known *a-priori* as it is always `TRUE` in the initialisation phase and `FALSE` in keystream generation mode. As we managed the initialisation and keystream generation phases in different loops, we can safely skip the check of `MIXING` parameter and set the correct value of `INPUT_BIT_R` in the loops;
- `CLOCK_R`: there are two `if`; the first one checks the `RTAPS` vector to determine which states have to be `xored` with the value of the `FEEDBACK_BIT`. The second one checks the value of `CONTROL_BIT_R` to determine if the new states need to be bitwise-`xored` with the older ones. We apply the same approach used in [60] in both the checks; we perform a `xor` operation between states and results of the multiplication of the masks representing `RTAPS` and `FEEDBACK_BIT` for the first check, and the result of the older states multiplied by `CONTROL_BIT_R`;
- `CLOCK_S`: here we have an `if-then-else` statement. This is a little bit different w.r.t. the other examples mentioned above as we need to manage also the `else` case. The control statement checks the value of the `CONTROL_BIT_S` variable; if it is 1, the states of the registers are updated by computing the `xor` of  $\hat{s}_i$  with the result of the multiplication of `FB1` and the `FEEDBACK_BIT`; if it is 0, the `xor` is computed between the state and the result of the multiplication of `FB0` and the `FEEDBACK_BIT`. We rewrite this check in the following way

```
S0 ^= (!contr_s & 0x1) * (S_MASK0_0 * feedback);
S0 ^= (contr_s * (S_MASK1_0 * feedback));
```

where

- `S0` contains the states  $s_0 \dots s_{31}$ ,
- `contr_s` is the `CONTROL_BIT_S`,
- `feedback` is the `FEEDBACK_BIT`, and
- `S_MASK0_0 S_MASK0_0` contain respectively  $FB0_0 \dots FB0_{31}$  and  $FB1_0 \dots FB1_{31}$  bits.

These operations are equivalent to the original `if-then-else` statement.

The above rewriting of each `if` statement grants that all the threads of a warp execute the same instruction on different data at the same time. Please notice that if we had left the original `if` statements we could not have the same assurance. This is due to the fact that values in `INPUT_BIT_R`, `INPUT_BIT_S`, `CONTROL_BIT_R` and `CONTROL_BIT_S` are determined from the states of the registers; as each thread executes the cipher on a unique couple of key and IV, each thread may have different values for the variables and consequently yields different results on statement checks.

**A.7. Framework installation and test case.** Here we describe all the steps to install the framework and run a test case. Please notice that you need a Linux computer equipped with an Nvidia GPU of CUDA compute capability  $\geq 3.5$ . Moreover, you need `gcc`  $\geq 4.5.0$  and `CUDA`  $\geq 7$ . To verify the CUDA compute capability please refer to [22].

Before starting please download the latest version of the framework from our repository [47]. You can download it as a zip file or you can clone it from the git repository.

A `Makefile` is provided to install the framework, it instructs the compiler to generate optimised code for most of the compute capabilities. If you have one of the latest GPUs or a Jetson Board, please check if the compute capability of your device is listed on the `Makefile`; if not, please add it to `CUDA_FLAGS` variable with the `-gencode arch=compute_X, code=sm.X`, where `X` is your compute capability. To install the kite-attack framework, simply run `make install`.

Once it is successfully installed, you may test it using one of the test configuration files provided in the `config` directory or you may generate a new configuration file for your customised attack.

An interactive `BASH` script, called `genConfigFile.sh`, is provided inside the `scripts` directory. This script helps users to customise their attacks. It allows to choose the target cipher, the number of initialisation rounds to attack, the  $I_{\max}$  and  $I_{\min}$  indexes, the run identifier and the path-name of the file where the chosen configuration is stored. An example of how to use the script is provided below.

To launch the attack, run the binary file corresponding to the selected target cipher, provide the configuration file, the output directory, and the id of the CUDA device you selected for the attack. If your system has only one device the value to pass is 0, in the case your system has more than one CUDA device provide the id of the chosen device. You may use `nvidia-smi` tool to obtain the list of all devices of your system along with the id and some other details.

In the following, we provide a complete session as list of commands, including instructions to get the framework, install it, generate a custom configuration file to attack the Mickey2.0 cipher, and run the attack.

```
$ git clone https://github.com/iac-cranic/kite-attack
$ cd kite-attack
$ make install
```

The interactive script `genConfigFile.sh` asks the user some questions and generates the configuration file accordingly to the answers. In the following we report the list of questions and corresponding answers along with the generated configuration file:

```
$ scripts/genConfigFile.sh
- [Q]: Where do you want to save the configuration?
(default newKiteAttack.conf):
- [A]: testMickey2.conf
- [Q]: Select the target cipher
- [A]: 3
- [Q]: Insert the number of initialization rounds
for the selected cipher: (default 100):
- [A]: 20
- [Q]: Insert run Identifier:
(automatically generated: KITE_xSI6UZjp8m):
- [A]: KITE_MICKEY_TEST_CASE
- [Q]: Insert the value of I_max :
- [A]: 5
- [Q]: Insert the value of I_min :
- [A]: 2
- [Q]: Insert the 0-th value that belongs to I_min
(please note that the indexes start from 0):
- [A]: 0
- [Q]: Insert the 1-th value that belongs to I_min
(please note that the indexes start from 0):
- [A]: 1
- [Q]: Insert the 0-th value that belongs to I_max
```

```

but not to I_min (i.e. (I_max\ I_min))
(please note that the indexes start from 0):
- [A]: 2
- [Q]: Insert the 1-th value that belongs to I_max
but not to I_min (i.e. (I_max\ I_min))
(please note that the indexes start from 0):
- [A]: 3
- [Q]: Insert the 2-th value that belongs to I_max
but not to I_min (i.e. (I_max\ I_min))
(please note that the indexes start from 0):
- [A]: 4
The configuration file testMickey2.conf has been
successfully generated
=====
TARGET_CIPHER=Mickey2
INIT_ROUNDS=20
RUN_IDENTIFIER=KITE_MICKY_TEST_CASE
I_max=5
I_max_minus_I_min=3
I_min=2
I_MAX_SET={0,1,2,3,4}
I_MAX_minus_I_MIN_SET={2,3,4}
I_MIN_SET={0,1}
=====

```

The last step is to run the attack, this can be done with the following command:

```
$ bin/kite_attack_mickey2 0 testMickey2.conf out_dir
```

In the `config` directory there are also configuration files to test Trivium and Grain128; with these configuration files the framework finds several superpolys of reduced rounds Trivium and Grain128.

*E-mail address:* {mcianfriglia, eonofri, mpedicini}@uniroma3.it

DEPT. OF MATHEMATICS AND PHYSICS, ROMA TRE UNIV., L. S. L. MURIALDO 1, ROME – ITALY

*E-mail address:* silvia.onofri@sns.it

SCUOLA NORMALE SUPERIORE, PIAZZA DEI CAVALIERI 7, PISA – ITALY