

Fast Exhaustive Search for Quadratic Systems in \mathbb{F}_2 on FPGAs

Charles Bouillaguet¹, Chen-Mou Cheng², Tung Chou³,
Ruben Niederhagen⁴(✉), and Bo-Yin Yang⁴

¹ Université de Lille, Lille, France
charles.bouillaguet@lil1.fr

² National Taiwan University, Taipei, Taiwan
doug@crypto.tw

³ Technische Universiteit Eindhoven, Eindhoven, The Netherlands
blueprint@crypto.tw

⁴ Academia Sinica, Taipei, Taiwan
ruben@polycephaly.org, by@crypto.tw

Abstract. In 2010, Bouillaguet *et al.* proposed an efficient solver for polynomial systems over \mathbb{F}_2 that trades memory for speed [BCC+10]. As a result, 48 quadratic equations in 48 variables can be solved on a graphics processing unit (GPU) in 21 min. The research question that we would like to answer in this paper is how specifically designed hardware performs on this task. We approach the answer by solving multivariate quadratic systems on reconfigurable hardware, namely Field-Programmable Gate Arrays (FPGAs). We show that, although the algorithm proposed in [BCC+10] has a better asymptotic *time* complexity than traditional enumeration algorithms, it does not have a better asymptotic complexity in terms of silicon *area*. Nevertheless, our FPGA implementation consumes 20–25 times less energy than its GPU counterpart. This is a significant improvement, not to mention that the monetary cost per unit of computational power for FPGAs is generally much cheaper than that of GPUs.

Keywords: Multivariate quadratic polynomials · Solving systems of equations · Exhaustive search · Parallelization · Field-Programmable Gate Arrays (FPGAs)

1 Introduction

Solving a system of m nonlinear multivariate polynomial equations in n variables over \mathbb{F}_q is called the MP problem. It is known to be NP-hard even if $q = 2$ and if we restrict ourselves to multivariate quadratic equations (in which case we call the problem MQ). These problems are mathematical problems of natural interest to cryptographers since an NP-hard problem whose random instances seem hard

See IACR ePrint Archive, Report 2013/436 [BCC+13], for an extended version.

could be used to design cryptographic primitives. Indeed, a seldom challenged standard conjecture is “any probabilistic Turing machine has negligible chance of successfully solving a random MQ instance with a given sub-exponential (in n) complexity when m/n is a constant” [BGP06].

This led to the development of *multivariate public-key cryptography* over the last decades, using one-way trapdoor functions to build cryptosystems such as HFE [Pat96], SFLASH [CGP02], and QUARTZ [PCG01]. It also led to the study of “provably-secure” stream ciphers like QUAD [BGP06].

In *algebraic cryptanalysis*, on the other hand, one distills from a cryptographic primitive a system of multivariate polynomial equations with the secret among the variables. This does not break AES as first advertised, but does break KeeLoq [CBW08], for a recent example. Fast solving would also be a very useful subroutine in attacks such as [BFJ+09].

Fast Exhaustive Search. When evaluating a quadratic system with n variables over \mathbb{F}_2 , each variable can be chosen as either 0 or 1. Thus, a straight forward approach is to evaluate each equation for all of the 2^n choices of inputs and to return any input that is evaluated to 0 by every single equation. The 2^n inputs can be enumerated by, e.g., using the binary representation of a counter of n bits where bit i is used as value for x_i . Since there are $\frac{n \cdot (n-1)}{2}$ pairs of variables and since in a generic (random) system each coefficient is 1 with probability $\frac{1}{2}$, each generic equation has about $\frac{n \cdot (n-1)}{2} \cdot \frac{1}{2}$ quadratic terms. Therefore, this approach has an asymptotic time complexity of $O(2^n \cdot m \cdot \frac{n \cdot (n-1)}{2} \cdot \frac{1}{2})$. Obviously, the second equation only needs to be evaluated in case the first one evaluates to 0 which happens for about 50% of the inputs. The third one only needs to be evaluated if the second one evaluated to 0 and so forth. The expected number of equations that need to be evaluated per iteration is $\sum_{i=1}^m 2^{1-i} < 2$. Thus, the overall complexity can be reduced to $O(2^n \cdot 2 \cdot \frac{n \cdot (n-1)}{2} \cdot \frac{1}{2}) = O(2^{n-1} (n-1)n)$ or more roughly $O(2^n n^2)$. Observe that the asymptotic time complexity is independent of m , the number of equations in the system, and only depends on n , the number of variables. This straight forward approach will be called *full-evaluation* approach in the remainder of this paper.

The full-evaluation approach requires a small amount of memory. The equation system is known beforehand and can be hard-coded into program code. It requires only n bits to store the current input value plus a small number of registers for the program state and temporary results. Thus, it has an asymptotic memory complexity of $O(n)$.

However, [BCC+10] suggests that we can trade memory for speed. The full-evaluation approach has the disadvantage that computations are repeated since the input of two consecutive computations may be only slightly different. For example, for a counter step from 16 (10000_b) to 17 (10001_b) only the least significant bit and thus the value of x_0 has changed; all the other inputs do not change, the computations not involving x_0 are exactly the same as in the previous step. In other examples, e.g., stepping from 15 (01111_b) to 16 (10000_b) more bits and therefore more variables are affected. Nevertheless, it is not important in which

order the inputs are enumerated. The authors of [BCC+10] point out that, by enumerating the inputs in Gray-code order, we ensure that between two consecutive enumeration steps only exactly one bit and therefore only one variable is changed. Therefore only those parts of an equation need to be recomputed that are affected by the change of that single variable x_i . Being in \mathbb{F}_2 , we only need to add $\frac{\partial f}{\partial x_i}(x)$ to the previous result. This reduces the computational cost from evaluating a *quadratic* multivariate equation in each enumeration step to evaluating a *linear* multivariate equation.

Furthermore, the authors of [BCC+10] prove that between two consecutive evaluations of $\frac{\partial f}{\partial x_i}(x)$ for a particular variable x_i only one other variable x_j of the input has changed. That is, the partial derivative of each variable is also evaluated in Gray-code order, and hence the trick can be applied *recursively*. Thus, by storing the result of the previous evaluation of $\frac{\partial f}{\partial x_i}(x)$, we only need to compute the change in regard to that particular variable x_j , i.e., the second derivative $\frac{\partial^2 f}{\partial x_i \partial x_j}(x)$, which is a *constant* value for quadratic equations.

Therefore, we can trade larger memory for less computation by storing the second derivatives in respect to all pairs of variables in a constant lookup table and by storing the first derivative in respect to each variable in registers. This requires $\frac{n \cdot (n-1)}{2}$ bits for the constant lookup table of the second derivatives and n bits of registers for the first derivatives. The computational cost of each iteration step is reduced to two additions in \mathbb{F}_2 , one for updating the value of a particular first derivative and one for computing the result of the evaluation.

The computational cost for each equation is independent from the values of n and m and thus will be considered constant for asymptotic estimations. However, since a state is updated in every iteration, all equations need to be computed (in parallel, e.g., using the bitslicing technique as suggested in [BCC+10]) in every single iteration. Therefore, the asymptotic time complexity for this approach is $O(2^n \cdot m)$. The asymptotic memory complexity is $O(m \cdot (\frac{n(n-1)}{2} + n)) = O(\frac{mn(n+1)}{2})$ or more roughly $O(n^2m)$.

Note that both the Gray-code approach and the full-evaluation approach can be combined by using only m_g equations for the Gray-code approach, thus producing about 2^{n-m_g} solution *candidates* to be tested by the remaining $m-m_g$ equations using full evaluation.

Lastly, we note that Gröbner-basis methods like XL [CKP+00] and F_5 [Fau02] using sparse linear solvers such as Wiedemann might have better performance than exhaustive search even over \mathbb{F}_2 . For example, they are claimed to asymptotically outperform exhaustive search when $m = n$ with guessing of $\approx 0.45n$ variables [YCC04, BFS+13]. However, as with all asymptotic results, one must carefully check all explicit and implicit assumptions to see how they hold in practice. When taking into account the true cost of Gröbner-basis methods, e.g., communication involved in running large-memory machines, the cross-over point is expected to be much higher than $n = 200$ as predicted in [BFS+13]. However, even systems in 200 variables are out of reach for today's computing capabilities.

The Research Question. The implementation of the Gray-code approach described in [BCC+10] for x86 CPUs and GPUs solves 48 quadratic equations in 48 binary variables using one NVIDIA GTX 295 graphics card in 21 min. The research question that we would like to answer in this paper is *how specifically designed hardware would perform on this task*. We approach the answer by solving multivariate quadratic systems on reconfigurable hardware, namely Field-Programmable Gate Arrays (FPGAs).

While the Gray-code approach has a lower asymptotic time complexity than full evaluation and is — given a sufficient amount of memory — the best choice for a software implementation, we show in Sect. 2.2 that both approaches have the same asymptotic *area* complexity. Therefore, for an FPGA implementation the choice of using either Gray code or full evaluation depends on the specific parameters and the target architecture of the implementation. We motivate our choice and describe our implementation for the Xilinx Spartan-6 FPGA in Sect. 2.

For a massively parallel implementation on FPGAs, it is most efficient to work on a set of input values in a batch. This increases the probability of having collisions of solutions during the computation, i.e. cases in which more than one input value in a batch is a solution candidate for the equation system. The implementation must guarantee that no solution is silently dropped. We discuss this effect in detail in Sect. 3, followed by the discussion of the implementation results and the conclusion of this paper in Sect. 4.

Source Code. The source code of our implementation is available under MIT License at <http://www.polycephaly.org/forcemq/>.

2 Implementation

The target hardware platform of our implementation is a Xilinx FPGA of the Spartan-6 architecture, device xc6slx150, package fgg676, and speed grade -3. The Spartan-6 architecture offers three different types of logic slices: SLICEX, SLICEL, and SLICEM.

The largest amount with about 50% of the slices is of type SLICEX. These slices offer four 6-input lookup tables (LUTs) and eight flip-flops. The LUTs can either be treated as logic or as memory: Seen as logic, each LUT-6 is computing the output value of any logical expression in 6 binary variables; seen as memory, each LUT-6 uses the 6 input wires to address a bit in a 64-bit read-only memory. Alternatively, each LUT-6 can be used as two LUT-5 with five identical input wires and two independent output wires.

About 25% of the slices are of type SLICEL, additionally offering wide multiplexers and carry logic for large adders. Another roughly 25% of the slices are of type SLICEM, which offer all of the above; in addition, the LUTs of these slices can be used as shift registers or as distributed *read-and-write* memory.

Please refer to [UG384] for more details on the Spartan-6 architecture.

2.1 Parallelization Using Accelerators

Exhaustive search for solutions of multivariate systems is embarrassingly parallel — all inputs are independent from each other and can be tested in parallel on as many computing devices as physically available. Furthermore, resources can be shared during the computation of inputs that have the same value for some of the variables.

Assume that we want to compute 2^i instances in parallel. We simply *clamp* the values of i variables such that x_{n-i}, \dots, x_{n-1} are constant for each instance, e.g., in case $i = 4$ for instance $5 = 0101_b$ variable $x_{n-1} = 0$, $x_{n-2} = 1$, $x_{n-3} = 0$, and $x_{n-4} = 1$. Therefore, the 2^n inputs for computations of a system in n variables can be split into 2^i new systems of 2^{n-i} inputs for $n - i$ variables using precomputation. These 2^i independent systems can either be computed in parallel on 2^i computing devices or sequentially on any smaller number of devices. (Obviously there is a limit on the efficiency of this approach; choosing $i = n$ would result in solving the whole original system during precomputation.) The same procedure of fixing variables can be repeated to cut the workload into parallel instances to exploit parallelism on each computing device.

After fixing variables x_{n-i}, \dots, x_{n-1} , all 2^i instances of one polynomial share the same quadratic terms; all terms involving x_{n-i}, \dots, x_{n-1} become either linear terms or constant terms. Therefore, the computations of the quadratic terms can be shared: For the Gray-code approach, the second derivatives can be shared between all instances while one set of first derivatives needs to be stored per instance; for full evaluation, the logic for the quadratic terms can be shared while the logic for the linear terms differs between the instances. Sharing resources requires communication between the instances and therefore is particularly suitable for computations on one single device. Given a sufficient amount of instances, the total area consumption is dominated by the instances doing the linear computations rather than by the shared computations on the quadratic part; therefore, the computations on the linear part require the most attention for an efficient implementation.

In the following, we investigate the optimal choices of n and m and the number of instances to exhaust the resources of *one single* FPGA most efficiently. Larger values of n can easily be achieved by running such a design several times or in parallel on several FPGAs. Larger values of m can be achieved by computing solutions for a subset of equations on the FPGA and forwarding those solution candidates from the FPGA to a host computer to be checked for the remaining equations. The flexibility in choosing n and m allows to cut the total workload into pieces that take a moderate amount of computation time on a single FPGA. This has the benefit of recovering from hardware failures or power outages without loss of too many computations.

2.2 Full Evaluation or Gray Code?

The asymptotic time and memory complexities of the full-evaluation approach and the Gray-code approach are summarized in Table 1. Considering a software

Table 1. Asymptotic complexities of the two approaches for exhaustive search.

	Time	Memory	Comp. logic	Area
Full evaluation	$O(2^n n^2)$	$O(n)$	$O(n^2 m)$	$O(n^2 m)$
Gray code	$O(2^n m)$	$O(n^2 m)$	$O(m)$	$O(n^2 m)$

implementation, for larger systems the Gray-code approach obviously is the more efficient choice, since it has a significantly lower time complexity and it is rather computational than memory bound. Because the memory complexity is much smaller than the time complexity, the memory demand can be handled easily by most modern architectures for such choices of parameters n and m that can be computed in realistic time.

However, a key measure for the complexity of a hardware design is the *area consumption* of the implementation: A smaller area consumption of a single instance of the implementation allows either to reduce cost or to increase the number of parallel instances and thus to reduce the total runtime. The area can be estimated as the sum of the area for computational logic and the area required for memory: The asymptotic complexity for the computational logic of the full-evaluation approach is about $O(n^2)$ for each equation, thus in total $O(n^2 m)$. The memory complexity is $O(n)$, so the area complexity is $O(n + n^2 m) = O(n^2 m)$. We point out that in contrast to the *time* complexity, the *area* complexity *depends* on m . The asymptotic complexity for the computational logic of the Gray-code approach is $O(m)$, the memory complexity is $O(n^2 m)$; the area complexity in total is $O(n^2 m + m) = O(n^2 m)$. Therefore, the asymptotic area complexity of the full-evaluation approach is equal to the area complexity of the Gray-code approach.

In contrast to a software implementation, it is not obvious from the asymptotic complexities, which approach eventually gives the best performance for specific hardware and specific values of n and m . The question is: which approach is using the resources of an FPGA more efficiently.

Choosing the Most Efficient Approach. For fixed parameters n and m we want to run as many parallel instances as possible on the given hardware. Since the quadratic terms are shared by the instances, the optimization goal is to minimize the resource requirements for the computations on the linear terms.

The main disadvantage of the Gray-code approach is that it requires access to *read-and-write memory* to keep track of the first derivatives. The on-chip memory resources, i.e., block memory and distributed memory using slices of type SLICEM, are quite limited. In contrast, the full-evaluation approach “only” requires *logic* that can be implemented using the LUTs of all types of slices.

However, each LUT in a SLICEM stores 64 bits; this is sufficient space for the first derivatives of 64 variables using the Gray-code approach. On the other hand, there are four times more logic-LUTs than RAM-LUTs. Four LUT-6 cover the evaluation of at most 24 variables. Therefore, the Gray-code approach is using the available input ports more efficiently. This is due to the fact that the

inputs for the Gray-code approach are addresses of width $O(\log n)$, whereas full evaluation requires $O(n)$ inputs for the variables. Thus, the Gray-code approach requires smaller bus widths and buffer sizes for pipelining.

Finally, the Gray-code approach allows to easily reuse a placed and routed design for different equation systems by exchanging the data in the lookup tables. An area-optimized implementation of the full-evaluation approach only requires logic for those terms of an equation that have a non-zero coefficient. To be able to use the same design for different equation systems, one would have to provide logic for *all* terms regardless of their coefficients, thus roughly doubling the required logic compared to the optimal solution. The Xilinx tool chain does not include a tool to exchange the LUT data from a fully placed and routed design, so we implemented our own tool for this purpose.

All in all, the Gray-code approach has several benefits compared to the full-evaluation approach which make it more suitable and more efficient for an FPGA implementation on a Spartan-6. The figures might be different, e.g., for an ASIC implementation or for FPGAs with different LUT sizes. We decided to use the Gray-code approach for the main part of our implementation to produce a number of solution candidates from a subset of the equations. These candidates are then checked for the remaining equations using full evaluation, partly on the FPGA, partly on the host computer.

2.3 Implementation of the Gray-Code Approach

As described in Sect. 1, the Gray-code approach trades larger memory for less computation. Algorithm 1 shows the pseudo code of the Gray-code approach (see the extended version of [BCC+10]). In case of the FPGA implementation, the initialization (Algorithm 1, lines 20 to 35) is performed offline and is hard-coded into the program file. Figure 1 shows the structure of the module *solver* for solving a system of m equations in n variables with 2^i instances of m_g equations using the Gray-code approach and $m - m_g$ full evaluations for the remaining equations.

The implementation of the Gray-code approach works as follows: First and second derivatives in respect to each variable are stored in lookup tables d' and d'' (Algorithm 1, lines 27 and 32). The second derivatives are constant and thus only require read-only memory. They require a quadratic amount of bits depending on the number of variables n . The first derivatives are computed in each iteration step based on their previous value (Algorithm 1, line 16). Therefore, the first derivatives are stored in a relatively small random access memory with a size linear to n .

Due to the structure of the Gray code, when looking at two consecutive values v_{i-1}, v_i in Gray-code enumeration, the position k_1 of the least-significant non-zero bit in the binary representation of i is the particular bit that is toggled when stepping from v_{i-1} to v_i . Therefore, the first derivative $\frac{\partial f}{\partial x_{k_1}}$ in respect to variable x_{k_1} needs to be considered for the evaluation. Furthermore, since the last time the bit k_1 had toggled, only the bit at the position k_2 of the second

```

1: function RUN( $f, n$ )
2:    $s \leftarrow \text{INIT}(f, n)$ ;
3:   while  $s.i < 2^n$  do
4:     NEXT( $s$ );
5:     if  $s.y = 0$  then
6:       return  $s.i \oplus \text{SHR}_1(s.i)$ ;
7:     end if
8:   end while
9: end function
10:
11: function NEXT( $s$ )
12:    $s.i \leftarrow s.i + 1$ ;
13:    $k_1 \leftarrow \text{BIT}_1(s.i)$ ;
14:    $k_2 \leftarrow \text{BIT}_2(s.i)$ ;
15:   if  $k_2$  valid then
16:      $s.d'[k_1] \leftarrow s.d'[k_1] \oplus s.d''[k_1, k_2]$ ;
17:   end if
18:    $s.y \leftarrow s.y \oplus s.d'[k_1]$ ;
19: end function
20: function INIT( $f = a_{n-1, n-2}x_{n-1}x_{n-2} +$ 
     $a_{n-1, n-3}x_{n-1}x_{n-3} + \dots + a_{1,0}x_1x_0 +$ 
     $a_{n-1}x_{n-1} + a_{n-2}x_{n-2} + \dots + a_0x_0 + a, n$ )
21:   state  $s$ ;
22:    $s.i \leftarrow 0$ ;
23:    $s.x \leftarrow 0$ ;
24:    $s.y \leftarrow a$ ;
25:   for all  $k, 0 < k < n$ , do
26:     for all  $j, 0 \leq j < k$ , do
27:        $s.d''[k, j] \leftarrow a_{k,j}$ ;
28:     end for
29:   end for
30:    $s.d'[0] \leftarrow a_0$ ;
31:   for all  $k, 1 \leq k < n$ , do
32:      $s.d'[k] \leftarrow s.d''[k, k-1] \oplus a_k$ ;
33:   end for
34:   return  $s$ ;
35: end function

```

Algorithm 1. Pseudo code for the Gray-code approach (see [BCC+10]). The functions BIT_1 and BIT_2 return the positions of the first and second least-significant non-zero bits respectively and SHR_1 is a logical shift right by one position.

least-significant non-zero bit in i has changed. So we need to access $\frac{\partial^2 f}{\partial x_{k_1} \partial x_{k_2}}$ in the static lookup table.

To compute k_1 and k_2 (Algorithm 1, lines 13 and 14), we use a module *counter* (Fig. 1, bottom) that is incrementing a counter by 1 in each cycle (cf. Algorithm 1, line 12). The counter counts from 0 to 2^{n-i} . To determine its first and second least-significant non-zero bits, we feed the counter value to a module called *gray_tree* that derives the index positions of the first and the second non-zero bit based on a divide-and-conquer approach. The output of the *gray_tree* module are buses k_1 and k_2 of width $\lceil \log_2(n) \rceil$ and two wires enable_1 and enable_2 (not shown in the figure) indicating whether k_1 and k_2 contain valid information (e.g., for all counter values $2^j, j \geq 0$, the output k_2 is invalid since the binary representation of 2^j has only one non-zero bit).

Next, we compute the address addr of the second derivative in the lookup table from the values k_1 and k_2 as $\text{addr} = k_2(k_2 - 1)/2 + k_1$ (cf. Algorithm 1, line 16). The computation is implemented fully pipelined to guarantee short data paths and a high frequency at runtime. The modules *counter* and *gray_tree* and the computation of the address for the lookup in the table are the same for all instances of all equations and therefore are required only once.

Now, the address is forwarded to the logic for the first equation eq_0 . Here, the buses addr and k_1 are buffered and in the next cycle forwarded to the lookup table of equation eq_1 and so on. The address is fed to the constant memory that returns the value of the second derivative d''_0 . We implement the constant

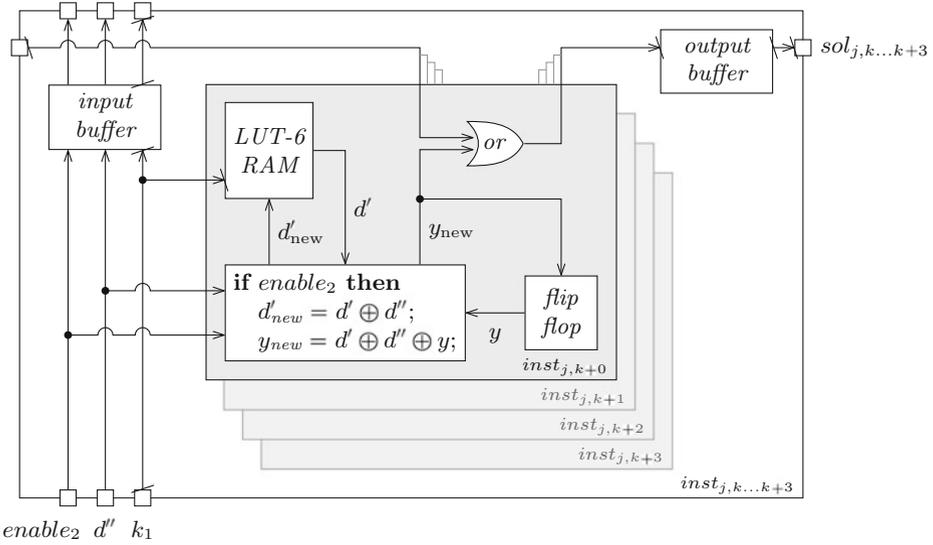


Fig. 2. Schematic of a Gray-code instance group.

two LUT-5, giving four inputs d'' , d' , $enable_2$, and y and receiving two outputs for the new first derivative and for the new y . Furthermore, we compute the or with the solutions of the previous equations as $sol_{j,k} = sol_{j-1,k} \vee y$. The inputs d'' , $enable_2$, and k_1 as well as the output are buffered using flip-flops.

Finally, the buffered result $sol_{j,k}$ is forwarded to $inst_{j+1,k}$ of eq_{j+1} in the next cycle. After the result of $inst_{j+1,k}$ has been computed as described before, the cumulated result $sol_{j+1,k} = sol_{j,k} \vee y$ is computed and forwarded to instance $inst_{j+2,k}$ and so on.

Each SLICEM has four LUTs that can be addressed as memory. However, they can only be written to if they all share the same address wires as input. Therefore, we combine four instances $inst_{j,k...k+3}$ of an equation j in one SLICEM using the same data as input. As a side effect, this reduces the number of buffers that are required for the now four-fold shared inputs. All in all, for up to 64 variables, a group of four instances for one equation requires 4 slices, one of them being a SLICEM.

Eventually, the four results $sol_{m_g-1,k...k+3}$ of an instance group are put on a bus together with a group ID that defines the value of the clamped variables. If more than one instance group finds a solution candidate in the same enumeration step, there might be a collision on the bus. We describe these collisions and our counter measures in detail in Sect. 3.

In each cycle, the solution candidates together with their ID are forwarded from one bus segment to the next, each connected to an instance group, until they eventually are leaving the bus after the last segment has been reached.

We are using the remaining resources of the FPGA to compute the actual solutions of the equation system. The computations on a subset of m_g equations using the Gray-code approach drastically reduce the search space from 2^n to 2^{n-m_g} . Therefore, we only need single instances of the remaining equations to check the solution candidates we receive from the Gray-code part. Since the inputs are quasi-random, we use full evaluation to check the candidates. If the system has more equations than we can fit on the FPGA, the remaining solution candidates are eventually forwarded to the host for final processing.

In order to check a solution candidate on the FPGA, we need to compute the actual Gray code for the input first. Since the design is fully pipelined, the value of each solution candidate from the Gray-code part is uniquely defined by the cycle when it appears. Therefore, we use a second counter (*counter₂*) that runs in sync, delayed by the pipeline length, to the original counter (*counter*). We compute the corresponding Gray code from the value *ctr₂* of this counter as $x = ctr_2 \oplus \text{SHR}_1(ctr_2)$ (see Algorithm 1, line 6). This value is appended to the ID and $x_{m_g-1} = (id, x)$ is forwarded into a FIFO queue.

To give some flexibility when fitting the design to the physical resources on the FPGA, our design allows the instances to be split into several *pillars*, each with their own bus, *gray_code* module and FIFO queue. The data paths are merged by selecting one solution candidate per cycle from the FIFO queues in a round-robin fashion.

Each solution candidate is forwarded to a module *eq_{m_g}* which simply evaluates equation m_g for the given input. To implement the full evaluation, we use a Greedy algorithm to map the terms of each equation to as few LUTs as possible (for more details please refer to the extended version [BCC+13] of this paper). The result of *eq_{m_g}* is or-ed to *sol_{m_g-1}* and forwarded to *eq_{m_g+1}* together with a buffered copy of x_{m_g-1} and so on.

Eventually, a vector x , its solution *sol*, and a warning signal *warn* are returned by the module *solver*. In case *sol* is equal to zero, i.e., all equations evaluated to zero for input x , the vector x is sent to the host.

3 Collisions, or Overabundance of Solution Candidates

Our implementation is akin to a map-reduce process, wherein $V = 2^n$ input vectors (n being the number of variables) are passed to many instances that each screen a portion of the inputs against a subset of m_g equations. Solution candidates which pass this stage move to a stage where they are checked against the remaining equations.

As mentioned in the previous section, the solution candidates that are computed in the first stage are collected by a sequential bus that connects all instances. This bus must provide a sufficient amount of resources to transfer all solution candidates to the second stage. Problems occur, if two or more instances find a solution candidate in the same iteration step. Every input processed by each screening instance may become a solution candidate with probability 2^{-m_g} . This is a highly unlikely event for a specific instance, but with a large number of instances, the probability of finding a solution candidate grows.

3.1 Expected Collisions

Let us assume that each of $V = 2^n$ input vectors is checked against m_g equations by $I = 2^i$ instances. A reasonable setup on a Spartan-6 FPGA might have $(n, m_g, i) = (48, 28, 9)$ or $(n, m_g, i) = (48, 14, 10)$.

A back-of-the-envelope calculation would go as follows: There are approximately $V/2^{m_g} = 2^{n-m_g}$ solution candidates, randomly spread over $V/I = 2^{n-i}$ iteration steps. The birthday paradox says that we may reasonably expect one or more collisions from x balls (solution candidates) in y bins (iteration steps) as soon as $x \gtrsim \sqrt{2y}$. Therefore, we should expect a small but non-zero number of “collisions”, iteration steps that have more than one solution candidate.

To articulate the above differently, each input vector has a probability of 2^{-m_g} to pass screening, and the event for each vector may be considered independent. Thus, the probability to have two or more solutions among I instances in the same iteration step is given by the sum of all coefficients of the quadratic and higher terms in the expansion of $(1 + (x - 1)/2^{m_g})^I$. The quadratic term represents the probability of having a collision of two values, the cubic term the probability of three values, and so on. The quadratic coefficient can be expected to be the largest and contribute to most of the sum. The expected number of collisions among all inputs is V/I times this sum, which is roughly

$$(V/I) [x^2] \left((1 - 2^{-m_g}) + 2^{-m_g} x \right)^I = (1 - 2^{-m_g})^{I-2} \frac{(I-1)}{2^{2m_g-n+1}} \approx 2^{i+n-2m_g-1},$$

where $[x^k] f(x)$ denotes the coefficient of x^k in the expansion of f .

The last approximation holds when $(1 - 2^{-m_g})^{I-2} \approx \exp(2^{-(m_g-i)}) \approx 1$ and $I \gg 1$. We can judge the quality of this approximation by the ratio between the quadratic and cubic term coefficients, which is $(I-2)2^{-m_g}/3 \lesssim 2^{-(m_g+1-i)}$. In other words, if $m_g - i > 3$, the number of expected collisions is roughly $2^{n+i-(2m_g+1)}$ with an error bar of 5% or less. Similarly the expected number of c -collisions (with at least c solutions within the same iteration step) is

$$(V/I) [x^c] \left((1 - 2^{-m_g}) + 2^{-m_g} x \right)^I \approx 2^{n-ck+(c-1)i}/c!$$

3.2 Choosing Parameters

In case of the Spartan-6 xc6slx150-fgg676-3 FPGA, the slices are physically located in a quite regular, rectangular grid of 128 columns and 192 rows. The grid has some large gaps on top and in the bottom as well as several vertical and horizontal gaps. By picking a subset of slices in the center of the FPGA we obtain a regular grid structure of 116 columns and 144 rows. Each row has 29 groups of 4 slices: one SLICEM, one SLICEL and two SLICEX. Such a group has enough resources for four Gray-code instances each. Therefore, the whole region can be used for up to $29 \cdot 4 \cdot 144 = 16,704$ Gray-code instances. The area below the slices for the Gray-code instances is used for the modules *counter* and *gray_tree*, for the computation of the address, and for the second-derivative tables. The area above

the instances contains enough logic for evaluating the remaining equations using full evaluation and for the logic required for FPGA-to-host communication.

Using 128 rows, we could fit $I = 128 \cdot 4 = 512 = 2^9$ instances of $m_g = 28$ equations of the Gray-code approach onto the FPGA — one equation per column, four instances per row — while guaranteeing short signal paths, leaving space of four slice columns for the bus, and giving more space for full evaluation on the top. With 28 equations in 2^9 instances, collisions of two solutions during one cycle are very rare and easy to handle by the host CPU. The obvious optimization to double the performance is to double I , the number of system instances, and to halve m_g , the number of equations evaluated with the Gray-code approach. However, this optimization introduces additional complications: Even if we can fit $I = 2^{10}$ instances of $m_g = 14$ equations using the Gray-code approach into the FPGA, Sect. 3.1 shows that one collision appears every 2^{10} cycles on average. We can no longer use the simple approach of re-checking all blocks with collisions on the CPU, we have to handle collisions on the FPGA. We describe in the following how to achieve 2^{10} instances for up to 14 (actually only 12) equations.

3.3 Handling of Collisions

Due to the physical layout of the FPGA and in order to save space for input-buffers, our implementation groups four instances with the same inputs together into an instance group. Instead of resolving a collision within an instance group right away, we forward a word of four bits, one for each instance, to the bus and cope with those collision later.

Whenever there is a collision at a instance group j , i.e., there is already a solution candidate on the bus in bus segment j , the candidate of group j is postponed giving precedence to the candidate on the bus. However, the actual input giving this solution candidate is not stored in the Gray-code instances but is later derived from the cycle in which the solution was found. Therefore, delaying the solution distorts the computation of the input value. Computing the input value immediately at each bus segment would require a lot of logic and would increase the bus width to n . Instead, we count how many cycles each solution candidate is delayed before being placed on the bus. Since the resources are limited, we can use at most 4 bits for this counter. For a push-back of up to 14 cycles, we can compute the exact corresponding solution candidate from each counter value. In case of 15 or more cycles of push-back, 1111_b is put on the bus to signal an error condition to the follow-up logic.

Since the delay has a very limited maximum number of cycles, we can not use classical bus congestion techniques like exponential backoff; we must ensure that candidates are pushed onto the bus as soon as possible. This leads to high congestion in particular at the end of the bus.

Due to the push-back, our collision pool has become temporal as well as spatial. That is, it might happen that another solution candidate is produced by the same instance group before the previous one is handed to the bus. Therefore, we provide four buffer slots for each instance group to handle the rare cases where

candidates are pushed back for several cycles while further candidates come up. If there are more candidates than there are buffer slots available, a warning signal is fired up and the involved input values are recomputed by the host.

All in all, the bus is transporting $i+7$ signals for 2^i instances; $i-2$ signals for the instance-group ID of the solution candidate, 4 signals for the push-back counter, 4 signals for the four outputs of a group of instances, and 1 warning signal.

Figure 3 shows a schematic of a bus segment. The solutions from an instance group of equation eq_{m_g-1} are sent in from the left using signal sol ; the inputs from the previous bus segment are shown in the bottom. Whenever there is no signal on the bus, i.e., sol_in is all high, the control logic sets the signal $step$ to high and a buffered result is pushed onto the bus; further delayed results are forwarded to the next buffer. If an available result can not be sent to the bus because there is already data on the bus, the step signal is set to low and each cycle counter in the counter buffers is incremented by one.

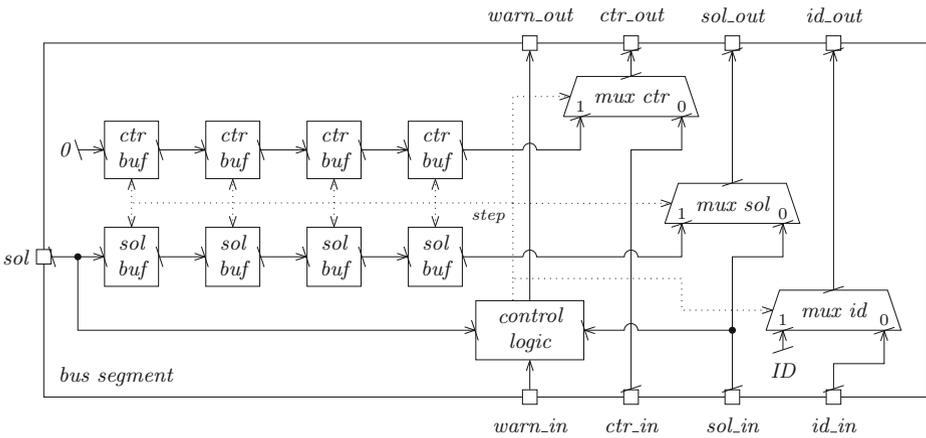


Fig. 3. Schematic of a bus segment.

The logic for each bus segment covering a group of 4 instances requires 10 slices (the area of 2.5 Gray-Code instance groups, i.e. four instances of a single equation) including buffers, counters, and multiplexers. Therefore, even though 29 Gray-Code instance groups would fit into one row on the FPGA, with two buses and two pillars of instances per row, we can only fit instances of 12 equations: $2 \cdot (12 + 2.5) = 29$. However, we achieve the goal of the desired $2^{10} = 1024$ parallel instances.

At the end of the buses, two FIFOs buffer the solution candidates so that the two data streams can be joined safely to forward a single data stream to the following logic for further handling of solution candidates (see Fig. 1). Here also the occasional collisions of solutions are resolved that might occur in an instance group of four instances as described above. Since the Gray-code part is using 2^{10} instances and 12 equations, there is one solution candidate on average every $2^{12-10} = 4$ cycles going into full evaluation.

With each bus averaging $1/8$ new entries and being capable of dispatching 1 entry every cycle, the buses should not suffer from too much congestion (confirmed by simulations and tests). With a push-back of at most 14 cycles, an unhandleable super-collision should only happen if 15 candidates appear within 15 consecutive 4-instance groups each with probability 2^{-10} , *all within 15 cycles*. We can do a back-of-the-envelope calculation for the probability like in Sect. 3.1 to find an upper bound of $\binom{225}{15} (2^{-10})^{15} \approx 6.4 \times 10^{-21}$. Just to be very sure, such super-collisions are still detected and passed to the host CPU, which re-checks the affected inputs. In all our tests and simulations, we did not detect any super-collisions, which confirms that our push-back buffer and counter sizes are sufficient.

We are able to fit logic for full evaluation of at least 42 more equations on the chip, giving 54 equations in the FPGA in total. This reduces the amount of outgoing solution candidates from the FPGA to the host computer to a marginal amount. Therefore, the host computer is able to serve a large amount of FPGAs even for a large total amount of equations in the system.

4 Performance Results and Concluding Remarks

We tested our implementation on a “RIVYERA S6-LX150 FPGA Cluster” from SciEngines. The RIVYERA has a 19-inch chassis of 4U height with an off-the-shelf host PC that controls 16 to 128 Spartan-6 LX150 FPGAs (xc6slx150-fgg676-3) and up to 256 FPGAs in the high density version “RIVYERA S6-LX150 HD”; our RIVYERA S6-LX150 has 16 FPGAs. The FPGAs are mounted on extension cards of 8 FPGAs (16 in the HD version) each with an extra FPGA exclusively for the communication with the host via PCIe.

We are using LOC constraints to explicitly place the Gray-code instances. Therefore, we achieve a tight packing and short data paths which allows us to run our design at up to 200 MHz. Due to the overall architecture we can compute systems of up to 64 variables using 2^{10} instances on a single FPGA. At a clock frequency of 200 MHz, solving a system in 64 variables in a single run requires $2^{64-10}/200 \text{ MHz} \approx 1042$ days; to reduce data loss in case of system failures or power outages, we recommend to divide the workload into smaller pieces with a shorter runtime. Our reference design is using $n = 54$ variables and $m = 54$ equations. In this case, a single run is finished after $2^{54-10}/200 \text{ MHz} \approx 24.5$ h. Preparing and exchanging the LUT data in the program file using our tools takes about 10 s. Therefore a system of 64 variables can be solved in $2^{64-54} = 1024$ separate runs with a negligible overhead.

Area Consumption. The Spartan-6 LX150 FPGA has 23,038 slices. In total, our logic occupies 18,613 slices (80.79 %) of the FPGA. We are using 63.44 % of the LUTs and 44.47 % of the registers.

The logic for the Gray-code evaluation occupies the largest area with 15,281 slices (67.43 %). Only 253 of those slices are used for the second-derivative tables,

Table 2. Comparison of the runtime and cost for systems in 48, 64, and 80 variables.

		Time	Energy	Energy cost	
				Germany	USA
48 variables	Spartan-6	23 min	3.4 Wh	–	–
	GTX 295	21 min	82.3 Wh	–	–
64 variables	Spartan-6	1,042 days	216 kWh	€56	US\$28
	GTX 295	956 days	5,390 kWh	€1,401	US\$701
80 variables	Spartan-6	187,182 years	14.4 GWh	€3.7 mil.	US\$1.9 mil.
	GTX 295	171,603 years	353.3 GWh	€91.8 mil.	US\$45.9 mil.

the counter, and address calculation. The bus occupies 2,740 slices, the remaining 12,288 slices are used for the 1,024 instances of 12 equations.

The logic for full evaluation of the remaining 42 equations, the FIFO queues, and the remaining solver logic requires 1,702 slices (7.39%). Each equation in 54 variables requires 88 LUTs for computational logic, thus about 22 slices. All these slices are located in an area above the Gray-code logic. More than 50% of the slices in this area are still available, leaving space to evaluate more equations using full evaluation if required.

The logic for communication with the host using SciEngine’s API requires 1,377 slices (5.98%).

Performance Evaluation. Our Spartan-6 FPGA design runs at 200 MHz. The design is fully pipelined and evaluates 2^{10} input values in each clock cycle. Thus, we compute all solutions of a system of 48 variables and 48 equations by evaluating all possible 2^{48} input values in $2^{48-10}/200 \text{ MHz} = 23 \text{ min}$ with a single FPGA. The GPU implementation of [BCC+10] computes all solutions on a GTX 295 graphics card in 21 min. Therefore, the Spartan-6 performs about the same as the GTX 295.

However, total runtime is not the only factor that affects the overall cost of the computation; power consumption is another important factor. We measured both the power consumptions of the Spartan-6 FPGA and the GTX 295 during computation: Our RIVYERA requires 305 W on average during the computation using all 16 FPGAs. The host computer with all FPGA cards removed requires 165 W. Therefore, a single FPGA requires $(305 \text{ W} - 165 \text{ W})/16 = 8.8 \text{ W}$ on average, including communication overhead. We measured the power consumption of the GTX 295 in the same way: During computation on the GTX 295, the whole machine required 357 W on average. Without the graphics card, the GPU-host computer requires 122 W. Therefore, the GTX 295 requires 235 W on average during computation. For a system of 48 variables, a single Spartan-6 FPGA requires $8.8 \text{ W} \cdot 23 \text{ min} = 3.4 \text{ Wh}$ for the whole computation. The GPU requires $235 \text{ W} \cdot 21 \text{ min} = 82.3 \text{ Wh}$. Therefore, the Spartan-6 FPGA requires about 25 times less energy than the GTX 295 graphics card.

For a system of 64 variables, the very same FPGA design needs about $2^{64-10}/200 \text{ MHz} = 1042 \text{ days}$ and therefore about 216 kWh. For this system, the

GPU requires about 965 days and roughly 5,390 kWh. A single kWh costs, e.g., about €0.26 in Germany¹ and about US\$0.13 in the USA². Therefore, solving a system of 64 variables with an FPGA costs about $216 \text{ kWh} \cdot e0.26/\text{kWh} = e56$ in Germany and $216 \text{ kWh} \cdot \text{US}\$0.13/\text{kWh} = \text{US}\28 in the US. Solving the same system using a GTX 295 graphics card costs €1,401 or US\$701. Table 2 shows an overview for systems in 48, 64, and 80 variables.

Development of GPU Hardware. The GPU implementation of [BCC+10] from 2010 uses a GTX 295 graphics card. We also measured the performance of their CUDA program on a GTX 780 graphics card which is state-of-the-art in 2013. However, the computations took slightly more time, although the GTX 780 should be more than three times faster than the GTX 295: the GTX 780 has 2304 ALUs running at 863 MHz while the GTX 295 has 480 ALUs running at 1242 MHz.

We suspect that the relative decrease of SRAM compared to the number of ALUs and the new instruction scheduling of the new generation of NVIDIA GPUs is responsible for the tremendous performance gap. To get full performance on the GTX 780 a thorough adaption and hardware-specific optimization of the algorithm would be required; the claim of NVIDIA that CUDA kernels can just be recompiled to profit from new hardware generations does not apply.

Nevertheless, running the code from [BCC+10] on a GTX 780 graphics card requires about 20% less energy than on the GTX 295, about 20 times more than our FPGA implementation.

80-bit Security. We want to point out that it is actually feasible to solve a system in 80 variables in a reasonable time: using $2^{80-64} = 2^{16} = 65,536$ FPGAs in parallel, such a system could be solved in 1042 days. Building such a large system is possible; e.g., the Tianhe-2 supercomputer has 80,000 CPUs.

Each RIVYERA S6-LX150 HD has up to 256 FPGAs; therefore, this computation would require $2^{16}/256 = 256$ RIVYERA-HD computers. The list price for one RIVYERA-HD is €110,000, about US\$145,000; the price for 256 machines is at most $256 \cdot \text{US}\$145,000 \approx \text{US}\37 million. Therefore, solving a system in 80 variables in 1042 days costs US\$39 million, including the electricity bill of US\$2 million for a continuous supply of $((256 \cdot 8.8\text{W}) + 165\text{W}) \cdot 256 = 620\text{kW}$. For comparison, the budget for the Tianhe-2 supercomputer was 2.4 billion Yuan (US\$390 million), *not* including the electricity bill for its peak power consumption of 17.8 MW. Therefore, 80-bit security coming from solving 80-variable systems over \mathbb{F}_2 is, as more cryptographers gradually acknowledge, no longer secure against institutional attackers and today's computing technology.

Acknowledgments. We would like to thank SciEngines for their support to get our design running on the RIVYERA, Ralf Zimmermann for his answers to our FPGA-related questions, and the anonymous reviewers for their valuable feedback. This research was partially sponsored by Academia Sinica under Bo-Yin Yang's Career Award and by the National Science Council project 100-2628-E-001-004-MY3.

¹ Average in 2012 according to the Agentur für Erneuerbare Energien.

² Average in 2012 according to the Bureau of Labor Statistics.

References

- [BCC+10] Bouillaguet, C., Chen, H.-C., Cheng, C.-M., Chou, T., Niederhagen, R., Shamir, A., Yang, B.-Y.: Fast exhaustive search for polynomial systems in \mathbb{F}_2 . In: Mangard, S., Standaert, F.-X. (eds.) CHES 2010. LNCS, vol. 6225, pp. 203–218. Springer, Heidelberg (2010). Extended Version: <http://www.lifl.fr/~bouillag/pub.html>
- [BCC+13] Bouillaguet, C., Cheng, C.-M., Chou, T., Niederhagen, R., Yang, B.-Y.: Fast exhaustive search for quadratic systems in \mathbb{F}_2 on FPGAs. Extended Version. IACR Cryptology ePrint Archive, Report 2013/436. <http://eprint.iacr.org/2013/436> (2013)
- [BFJ+09] Bouillaguet, C., Fouque, P.-A., Joux, A., Treger, J.: A family of weak keys in HFE (and the corresponding practical key-recovery). IACR Cryptology ePrint Archive, Report 2009/619. <http://eprint.iacr.org/2009/619> (2009)
- [BFS+13] Bardet, M., Faugère, J.-C., Salvy, B., Spaenlehauer, P.-J.: On the complexity of solving quadratic boolean systems. *J. Complex.* **29**(1), 53–75 (2013)
- [BGP06] Berbain, C., Gilbert, H., Patarin, J.: QUAD: a practical stream cipher with provable security. In: Vaudenay, S. (ed.) EUROCRYPT 2006. LNCS, vol. 4004, pp. 109–128. Springer, Heidelberg (2006)
- [CBW08] Courtois, N.T., Bard, G.V., Wagner, D.: Algebraic and slide attacks on KeeLoq. In: Nyberg, K. (ed.) FSE 2008. LNCS, vol. 5086, pp. 97–115. Springer, Heidelberg (2008)
- [CGP02] Courtois, N., Goubin, L., Patarin, J.: SFLASH, A Fast Asymmetric Signature Scheme for Low-Cost Smartcards: Primitive Specification, Second Revised Version. <https://www.cosic.esat.kuleuven.be/nessie/tweaks.html> (2002)
- [CKP+00] Courtois, N.T., Klimov, A.B., Patarin, J., Shamir, A.: Efficient algorithms for solving overdefined systems of multivariate polynomial equations. In: Preneel, B. (ed.) EUROCRYPT 2000. LNCS, vol. 1807, pp. 392–407. Springer, Heidelberg (2000)
- [Fau02] Faugère, J.-C.: A new efficient algorithm for computing Gröbner bases without reduction to zero (F_5). In: International Symposium on Symbolic and Algebraic Computation – ISSAC 2002, July 2002, pp. 75–83. ACM Press (2002)
- [Pat96] Patarin, J.: Hidden fields equations (HFE) and isomorphisms of polynomials (IP): two new families of asymmetric algorithms. In: Maurer, U.M. (ed.) EUROCRYPT 1996. LNCS, vol. 1070, pp. 33–48. Springer, Heidelberg (1996)
- [PCG01] Patarin, J., Courtois, N.T., Goubin, L.: QUARTZ, 128-bit long digital signatures. In: Naccache, D. (ed.) CT-RSA 2001. LNCS, vol. 2020, pp. 282–297. Springer, Heidelberg (2001)
- [UG384] Spartan-6 FPGA Configurable Logic Block — User Guide. v1.1 UG384. Xilinx Inc., Feb. 2010
- [YCC04] Yang, B.-Y., Chen, J.-M., Courtois, N.T.: On asymptotic security estimates in XL and Gröbner bases-related algebraic cryptanalysis. In: López, J., Qing, S., Okamoto, E. (eds.) ICICS 2004. LNCS, vol. 3269, pp. 401–413. Springer, Heidelberg (2004)