

Improved Agreeing-Gluing Algorithm

Igor Semaev

Department of Informatics, University of Bergen, Norway
igor@ii.uib.no

Abstract. A system of algebraic equations over a finite field is called sparse if each equation depends on a low number of variables. Finding efficiently solutions to the system is an underlying hard problem in the cryptanalysis of modern ciphers. In this paper a deterministic Improved Agreeing-Gluing Algorithm is introduced. The expected running time of the new Algorithm on uniformly random instances of the problem is rigorously estimated. The estimate is at present the best theoretical bound on the complexity of solving average instances of the problem. In particular, this is a significant improvement over those in our earlier papers [10, 11]. In sparse Boolean equations a gap between the worst case and the average time complexity of the problem has significantly increased.

1 Introduction

1.1 The problem and motivation

Let (q, l, n, m) be a quadruple of natural numbers, where q is a prime power. Then F_q denotes a finite field with q elements and $X = \{x_1, x_2, \dots, x_n\}$ is a set of variables from F_q . By X_i , $1 \leq i \leq m$ we denote subsets of X of size $l_i \leq l$. The system of equations

$$f_1(X_1) = 0, \dots, f_m(X_m) = 0 \tag{1}$$

is considered, where f_i are polynomials over F_q and they only depend on variables X_i . Such equations are called l -sparse. A solution to (1) over F_q is an assignment in F_q to all n variables X that satisfies all equations (1). That is a vector of length n over F_q provided the variables X are somehow ordered. The main goal is to find all solutions over F_q . We suggest a deterministic Improved Agreeing-Gluing (IAG) Algorithm. It is presented by two variations. The expected complexity of one variation is rigorously estimated assuming uniform distribution on the problem instances; see Section 1.2 for detail. The results provide a significant improvement over earlier average time complexity estimates in [10, 11].

The approach, which exploits the sparsity of equations and doesn't depend on their algebraic degree, was studied in [14, 8, 10, 11]. These are guess-and-determine algorithms. In sparse equations the number of guesses on a big enough variable set Y and the time to produce them is much lower than $q^{|Y|}$ due to the Search Algorithm; see Section 6. Previously, no preference was made on which variables to guess. We now argue that guessing values of some particular variables leads to better asymptotic complexity bounds.

Gröbner bases algorithm was designed to solve general algebraic equation systems; see for instance [5]. In Boolean case, where $q = 2$ and $m = n$, the conjectural average complexity bound is higher than 2^n bit operations except for quadratic equations [1]. The best heuristic bound is then of order 1.7^n [13]. In contrast, our estimates are rigorous mathematical statements and very

low exponential functions themselves even in non quadratic case. Sparse equations may be encoded by a CNF formula and solved with a Sat-solving software. The asymptotical complexity of modern Sat-solvers, as MiniSat, is unknown, though they may be fast in practice [3] for relatively low parameters.

The article was motivated by applications in cryptanalysis. Modern ciphers are product, the mappings they implement are compositions of not so many functions in a low number of variables. The similar is true for asymmetric ciphers. Intermediate variables are introduced to simplify equations, describing the cipher, and to get a system of sparse equations. For a more general type of sparse equations, Multiple Right Hand Side linear equations describing in particular AES; see [9]. An efficient solution of the equations breaks the cipher.

Let Y be an ordered string of variables and a be an F_q -vector of the same length. We say that a is a vector in variables Y , or Y -vector, if the entries of a may be assigned to the variables Y , for instance, in case of fixation.

1.2 Probabilistic model

We look for the set of all solutions to (1) over F_q , so we only consider for f_i polynomials of degree at most $q - 1$ in each variable. Obviously, the equation $f_i(X_i) = 0$ is determined by the pair (X_i, V_i) , where V_i is the set of X_i -vectors, where f_i is zero. Given q, n, m , and $l_1, \dots, l_m \leq l$, uniform distribution on instances is assumed. As any particular information on equations is beforehand assumed unknown, this looks the most fair probabilistic model to compute expected complexities. The uniformity means

1. the equations in (1) are independently generated. Each equation $f_i(X_i) = 0$ is determined by
2. the subset X_i of size l_i taken uniformly at random from the set of all possible l_i -subsets of X , that is with the probability $\binom{n}{l_i}^{-1}$,
3. and the polynomial f_i taken uniformly at random and independently of X_i from the set of all polynomials of degree $\leq q - 1$ in each of variables X_i . In other words, with the equal probability $q^{-q^{l_i}}$.

Running time of any deterministic solving algorithm is a random variable under that model. We assume that m/n tends to $d \geq 1$ as q and l are fixed and n tends to infinity.

Table 1. Algorithms' running time: $q = 2$ and $m = n$.

	l	3	4	5	6
the worst case, [6]		1.324^n	1.474^n	1.569^n	1.637^n
Gluing1, expectation, [10]		1.262^n	1.355^n	1.425^n	1.479^n
Gluing2, expectation, [10]		1.238^n	1.326^n	1.393^n	1.446^n
Agreeing-Gluing, expectation, [11]		1.113^n	1.205^n	1.276^n	1.334^n
	r	2	3	3	4
Weak Improved Agreeing-Gluing, expectation		1.029^n	1.107^n	1.182^n	1.239^n

2 Previous Ideas and the New Approach

One earlier method [10] is based on subsequent computing solutions U_k to the equation subsystems: $f_1(X_1), \dots, f_k(X_k)$ for $k = 1, \dots, m$. Gluing procedure extends instances U_k to instances U_{k+1} by walking throughout a search tree. In the end, all system solutions are U_m . The running time is determined by the maximal of $|U_k|$. Gluing2 is a time-memory trade-off variation of the basis Gluing1 Algorithm. See Table 1 for their running time expectation in case of n Boolean equations in n variables and a variety of l .

In Agreeing-Gluing Algorithm [11] we only extend those intermediate solutions from U_k that do not contradict with the rest of the equations $f_{k+1}(X_{k+1}) = 0, \dots, f_m(X_m) = 0$. That makes lots of search tree branches cut and implies a better average time complexity.

Let Z_r denote variables that occur in at least r equations (1). The new method has two variations. In the Strong IAG Algorithm the largest r , where Z_r is not empty, is taken. Then Z_r -vectors that do not contradict any of (1) are generated by the Search Algorithm; see Section 6. We denote them W_r . For each $a \in W_r$ the variables Z_r are substituted by the entries of a . New l -sparse equations in a smaller variable set $X \setminus Z_r$ are to solve. One then recursively computes W_{r-1} and so on. It is enough to obtain W_2 as all system solutions are then easy to deduce; see Lemma 1 below.

In the Weak IAG Algorithm r is a parameter. The vectors W_r are generated by the Search Algorithm. The variables Z_r are substituted by the entries of $a \in W_r$. New equations, in case $r \geq 3$, are encoded by a CNF formula and local search algorithm is applied to find all solutions. Two last lines in Table 1 show expected complexity of the Weak IAG Algorithm and the optimal value of r . The expected complexity of the Strong IAG Algorithm is not presented in this article. The Agreeing-Gluing Algorithm [11] is a particular case of the present method for $r = 1$.

3 Trivially unsolvable equations

The probability that a randomly chosen equation in l variables is solvable over F_q , i.e., admits at least one solution over F_q , is $1 - (1 - \frac{1}{q})^{q^l}$. So the probability the equation system (1) is trivially unsolvable(at least one of the equations has no solutions over F_q) is $1 - \left[1 - (1 - \frac{1}{q})^{q^l}\right]^m$. This value tends to 1 as l and q are fixed and $m = dn$ tends to infinity. It is very easy to recognize, with some average complexity R , a trivially unsolvable equation system. However, for small d that only gives a negligible contribution to the average complexity estimate while it is exponential. Let Q denote average complexity of a deterministic algorithm on all instances of (1). Let Q_1 denote average complexity of the algorithm on the instances of (1) which are not trivially unsolvable, i.e., each equation has at least one solution over F_q . In both cases uniform distribution is assumed. By the conditional expectation formula,

$$Q = \left[1 - (1 - \frac{1}{q})^{q^l}\right]^{dn} Q_1 + \left(1 - \left[1 - (1 - \frac{1}{q})^{q^l}\right]^{dn}\right) R.$$

Therefore, $Q_1 < \left[1 - (1 - \frac{1}{q})^{q^l}\right]^{-dn} Q$. For $q = 2$ and $d = 1$ that will affect the bound at $l = 3$: in case of the Weak IAG Algorithm, Q_1 becomes bounded by 1.033^n . For all other l the influence is negligible: estimates for Q and Q_1 are almost identical. For larger $d = 1 + \delta$ the contribution is larger, but Q becomes sub-exponential fast. So Q_1 remains bounded by a very low exponential

function at least for low δ . In fact, we believe that Q_1 becomes sub-exponential too, though it is not proved here.

4 Notation and Example

Let $r \geq 1$ and $Z_r(k)$ denote the set of variables that appear in at least r of X_1, \dots, X_k , so $Z_r(m) = Z_r$. Let a be a $Z_r(k)$ -vector. Assume (X_i, V_i) is one of the equations. If $X_i \not\subseteq Z_r(k)$, then $V_i(a)$ denotes projections to variables $X_i \setminus Z_r(k)$ of $b \in V_i$, where b and a agree on common variables $X_i \cap Z_r(k)$. In other words, $V_i(a)$ are solutions in variables $X_i \setminus Z_r(k)$ to $f_i(X_i) = 0$ after the $Z_r(k)$ get substituted by the entries of a . If $X_i \subseteq Z_r(k)$ and the projection of a to X_i does not appear in V_i , then we write $V_i(a) = \emptyset$. Otherwise, $V_i(a) \neq \emptyset$. It is obvious that $V_i(a) = \emptyset$ iff the fixation of $Z_r(k)$ by constants a contradicts the equation (X_i, V_i) .

Let $W_r(k)$ be the set of $Z_r(k)$ -vectors a such that $V_i(a) \neq \emptyset$ for all $i = 1, \dots, m$. In Section 5 a procedure, called Search Algorithm, that extends instances $W_r(k)$ to $W_r(k+1)$ is described. Its output is $W_r = W_r(m)$. Three cases should be studied separately.

Case $r = 1$. Then $U_m = W_1$, all system solutions in variables $X_1 \cup \dots \cup X_m$. Extending $W_1(k)$ to $W_1(k+1)$ by walking over a search tree is the Agreeing-Gluing Algorithm [11].

Case $r = 2$. Remark that variables in $X_i \setminus Z_2$ are different for different $1 \leq i \leq m$.

Lemma 1. *Let $X_{i_1}, X_{i_2}, \dots, X_{i_s}$ be all variable sets such that $X_{i_j} \not\subseteq Z_2$. After reordering of variables it holds that*

$$U_m = \bigcup_{a \in W_2} \{a\} \times V_{i_1}(a) \times V_{i_2}(a) \dots \times V_{i_s}(a). \quad (2)$$

Example Let the system of three Boolean equations be given:

$$\begin{array}{c|c|c} x_1 & x_2 & x_3 \\ \hline 0 & 0 & 1 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 0 & 1 \end{array}, \quad \begin{array}{c|c|c} x_3 & x_4 & x_5 \\ \hline 0 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{array}, \quad \begin{array}{c|c|c} x_5 & x_6 & x_7 \\ \hline 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{array}.$$

Then $Z_2(2) = \{x_3\}$ and $Z_2(3) = \{x_3, x_5\}$ and the directed products (2) are:

$$\begin{array}{c|c|c|c|c|c} x_3 & x_5 & x_1 & x_2 & x_4 & x_6 & x_7 \\ \hline 0 & 0 & 1 & 0 & \times & 0 & \times & 0 & 0 \\ & & & & & & & 1 & 1 \\ \hline 0 & 1 & 1 & 0 & \times & 0 & \times & 1 & 0 \\ & & & & & & & 0 & 1 \\ \hline 1 & 1 & 0 & 0 & \times & 0 & \times & 1 & 0 \\ & & 1 & 1 & & 1 & & 0 & 1 \\ & & 1 & 0 & & & & & \end{array}.$$

So 16 solutions to the system are represented by three strings of length 2, that is by 0, 0, and 0, 1, and 1, 1 related to variables x_3, x_5 .

Case $r \geq 3$. The Search Algorithm returns some $a \in W_r$. The variables Z_r are substituted by the entries of a . The problem is represented in a conjunctive normal form (CNF) with the clause length of at most $l_1 = \lceil \log_2 q \rceil l$ and in $n_1 = \lceil \log_2 q \rceil |X \setminus Z_r|$ Boolean variables. Local search

algorithm, described in [4], is used to find all solutions. In worst case, that takes $O(N(2 - \frac{2}{l_1+1})^{n_1})$ bit operations to find N solutions. Other estimates, e.g. presented in the first line of Table 1, do not improve the overall results for $q = 2$.

5 Weak IAG Algorithm

We define a rooted search tree. The tree has at most $m + 1$ levels numbered $0, 1, \dots, m$. The root at level 0 is labeled by \emptyset . Vertices at level k are labeled by vectors $W_r(k)$ or \emptyset if $Z_r(k) = \emptyset$. Let a be a level k vertex. It is connected to a level $k + 1$ vertex b whenever a is a sub-vector of b . Remark that $Z_r(k) \subseteq Z_r(k + 1)$. We now describe the Algorithm.

Stage 1 (Search Algorithm) It starts at the root. Let the Algorithm be at a level k vertex a which is extended to b with a projection of $V_{k+1}(a)$ to variables $Z_r(k + 1) \setminus Z_r(k)$. If b does not contradict to any of the equations, then b is a level $k + 1$ vertex. The Algorithm walks to that. Otherwise, another projection is taken to extend a . If all the projections are exhausted, the Algorithm backtracks to level $k - 1$. This stage output is $W_r = W_r(m)$.

Stage 2 Let the Algorithm achieve a vertex $a \in W_r$. If $r = 1$, then a is a system solution. If $r = 2$, then the system solutions are deduced with (2). If $r \geq 3$, a system of l -sparse equations in variables $X \setminus Z_r$ after substituting Z_r by constants a is solved with local search. If no vertex at level m is hit, then the system has no solution.

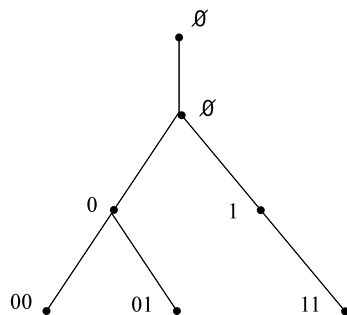


Fig. 1. The search tree.

Theorem 1. Let $r < 3$, then the algorithm running time is $O\left(m \sum_{i=r}^{m-1} |W_r(i)|\right)$ operations with vectors over F_q of length at most n . Let $r \geq 3$, then running time is $O\left(m \sum_{i=r}^{m-1} |W_r(i)| + |W_r| c^{n-|Z_r|}\right)$ operations, where $c = \left(2 - \frac{2}{l \lceil \log_2 q \rceil + 1}\right)^{\lceil \log_2 q \rceil}$.

$|W_r(k)|$ are random variables that depend on sets X_1, \dots, X_m and the polynomials f_1, \dots, f_m . In Section 8 we estimate the maximal of their expectations. The expectation of $|W_r|c^{n-|Z_r|}$ is estimated in Section 9. For a range of r the estimates are computed with an optimization software like MAPLE. Remark that the computation does not depend on n . One then finds r such that the running time expectation is minimal.

The search tree for the example system is presented in Fig. 1, where $r = 2$. Level 2 vertices are labeled by $W_2(2) = \{(0), (1)\}$, vectors in variables $Z_2(2) = \{x_3\}$. The vertices at level 3 are labeled by $W_2(3) = \{(0, 0), (0, 1), (1, 1)\}$, vectors in variables $Z_2(3) = \{x_3, x_5\}$.

6 General Search Algorithm

Given $Y \subseteq X$, find all Y -vectors over F_q that do not contradict any of equations (1). Take a subset sequence $Y_1 \subseteq Y_2 \subseteq \dots \subseteq Y_s = Y$. That defines a search tree. The root is labeled by \emptyset , the vertices at level $1 \leq k \leq s$ are labeled by Y_k -vectors V_k that do not contradict any of (1). Vertices a and b at subsequent levels are connected if a is a sub-vector of b . The algorithm walks throughout the tree by constructing instances V_k with backtracking. The running time is proportional to $|V_1|q^{|Y_2 \setminus Y_1|} + |V_2|q^{|Y_3 \setminus Y_2|} + \dots + |V_{s-1}|q^{|Y_s \setminus Y_{s-1}|}$ operations. One may take a sequence of subsets that minimizes the running time. In IAG Algorithms the sequence is $Z_r(r) \subseteq Z_r(r+1) \subseteq \dots \subseteq Z_r(m) = Z_r$. In practice, one may check whether Y_k -vector a contradicts the whole system (1) but not only each of the equations taken separately. One then runs the Agreeing Algorithm [9, 12] after the variables Y_k get substituted by constants a .

7 Tools

In this Section we collect miscellaneous auxiliary statements. Let $\eta = \eta(x, y)$ be any variable that depends on two independent discrete random variables x and y . Then $\mathbf{E}_y \eta = \mathbf{E}_y \eta(x, y)$ denotes the expectation of η , where y is generated to its initial distribution.

Lemma 2. [11] *For the full expectation of $\eta = \eta(x, y)$ we have*

$$\mathbf{E}_{x,y} \eta = \mathbf{E}_x(\mathbf{E}_y(\eta)).$$

Random Allocations Theory studies random allocations of particles(balls) into boxes, see [7]. Let k complexes of particles be independently and uniformly allocated into n boxes, $l_i \leq n$ particles at the i -th allocation. This means that at the i -th allocation any l_i boxes are occupied with the equal probability $\binom{n}{l_i}^{-1}$. This is how variable sets X_1, \dots, X_m are generated according to Section 1.2. Let ν_1, \dots, ν_n be the string of box frequencies, that is ν_i is the number of particles in the i -th box after k such allocations. Let $A = A(\nu_1, \dots, \nu_n)$ be any event depending on the frequencies ν_i . Let also $\mathbf{Pr}(A | l_1, \dots, l_k)$ denote the probability of the event A under the allocation by complexes.

Lemma 3.

$$\mathbf{Pr}(A | l_1, \dots, l_k) \leq \frac{\mathbf{Pr}(A | 1, \dots, 1)}{\prod_{i=1}^k (1 - 1/n) \dots (1 - (l_i - 1)/n)},$$

where $\mathbf{Pr}(A | 1, \dots, 1)$ is the probability of A under condition that $L = l_1 + \dots + l_k$ particles are allocated one after the other, i.e., L of 1's in the expression $\mathbf{Pr}(A | 1, \dots, 1)$.

Proof. Let L particles be independently allocated into n boxes one after the other. Let B denote the event that the first l_1 particles were allocated into different boxes, the following l_2 were allocated into different boxes etc, until the last l_k particles were allocated into different boxes. In other words, the event B occurs if the particles are allocated by complexes of size l_1, l_2, \dots, l_k . Then $\Pr(B) = \prod_{i=1}^k (1 - 1/n) \dots (1 - (l_i - 1)/n)$ as the particles were allocated independently. By the complete probability formula we get

$$\begin{aligned} \Pr(A|1, \dots, 1) &= \Pr(B) \Pr(A|B) + \Pr(\bar{B}) \Pr(A|\bar{B}) \\ &\geq \Pr(B) \Pr(A|B) = \Pr(B) \Pr(A|l_1, \dots, l_k) \end{aligned}$$

as $\Pr(A|B) = \Pr(A|l_1, \dots, l_k)$.

Let $f^n(z) = \sum_{k=0}^{\infty} a_{n,k} z^k$, where $a_{n,k} \geq 0$, be an analytic function for any natural n .

Lemma 4. For any real $z_0 > 0$

$$a_{n,k} \leq \frac{f^n(z_0)}{z_0^k}.$$

Proof. The expansion of f^n has only nonnegative coefficients, so $a_{n,k} z_0^k \leq f^n(z_0)$.

To minimize the estimate one takes a positive root z_0 to

$$\frac{\partial(n \ln f(z) - k \ln z)}{\partial z} = 0$$

if there exist any. In case there is only one root, the Lemma estimate is proportional to the main term of the asymptotic expansion for $a_{n,k}$ with saddle point method as n and k tend to infinity; see [2]. Lemma 4 estimate is then asymptotically close to the real value of $a_{n,k}$. We use this observation in Lemmas 5, 6 and 10.

Let $\mu_r = \mu_r(t, n)$ be the number of boxes with just r particle after uniform allocation of t particles one after the other into n boxes. Let $\mathbf{E}(x_1^{\mu_{r_1}} \dots x_s^{\mu_{r_s}})$ be the expectation of the random variable $x_1^{\mu_{r_1}} \dots x_s^{\mu_{r_s}}$, where x_1, \dots, x_s are any variables. By definition,

$$\mathbf{E}(x_1^{\mu_{r_1}} \dots x_s^{\mu_{r_s}}) = \sum_{k_1, \dots, k_s} \Pr(\mu_{r_1} = k_1, \dots, \mu_{r_s} = k_s) x_1^{k_1} \dots x_s^{k_s}.$$

Theorem 2 in Chapter 2, Section 1 of [7] states

$$\sum_{t=0}^{\infty} \frac{n^t z^t}{t!} \mathbf{E}(x_1^{\mu_{r_1}} \dots x_s^{\mu_{r_s}}) = \left(e^z + \frac{z^{r_1}}{r_1!} (x_1 - 1) + \dots + \frac{z^{r_s}}{r_s!} (x_s - 1) \right)^n. \quad (3)$$

In particular, we get

$$\sum_{t=0}^{\infty} \frac{n^t z^t}{t!} \mathbf{E}(x_0^{\mu_0} \dots x_{r-1}^{\mu_{r-1}}) = \left(e^z + (x_0 - 1) + \dots + \frac{z^{r-1}}{(r-1)!} (x_{r-1} - 1) \right)^n.$$

We there put $x_0 = \dots = x_{r-1} = 0$ and get

$$\left(e^z - 1 - z \dots - \frac{z^{r-1}}{(r-1)!} \right)^n = \sum_{t=nr}^{\infty} \frac{n^t z^t}{t!} \Pr(\mu_0 = 0, \dots, \mu_{r-1} = 0)$$

as $\Pr(\mu_0 = 0, \dots, \mu_{r-1} = 0) = 0$ for $t < nr$. Let $g(z) = e^z - 1 - z \dots - \frac{z^{r-1}}{(r-1)!}$.

Lemma 5. Let $r \geq 1$. For any natural number $t \geq nr$

$$\Pr(\boldsymbol{\mu}_0(t, n) = 0, \dots, \boldsymbol{\mu}_{r-1}(t, n) = 0) \leq \frac{g^n(x_0) t!}{x_0^t n^t},$$

where x_0 is the only nonnegative root of the equation $n \frac{x(e^x - 1 - x \dots - \frac{x^{r-2}}{(r-2)!})}{e^x - 1 - x \dots - \frac{x^{r-1}}{(r-1)!}} = t$.

Proof. Let $t > nr$. Then the equation has the only positive solution x_0 . The statement is true by Lemma 4. Let $t = nr$, then $x_0 = 0$. One sees that $\frac{g^n(x_0) t!}{x_0^t n^t}$ is defined at $x_0 = 0$ and equal to $\frac{(nr)!}{(r!)^n n^{nr}}$. On the other hand, one directly computes

$$\Pr(\boldsymbol{\mu}_0(nr, n) = 0, \dots, \boldsymbol{\mu}_{r-1}(nr, n) = 0) = \frac{(nr)!}{(r!)^n n^{nr}}.$$

It follows from (3) that

$$\sum_{t=0}^{\infty} \frac{n^t z^t}{t!} \mathbf{E}(x_1^{\boldsymbol{\mu}_1} \dots x_{r-1}^{\boldsymbol{\mu}_{r-1}}) = \left(e^z + z(x_1 - 1) + \dots + \frac{z^{r-1}}{(r-1)!} (x_{r-1} - 1) \right)^n. \quad (4)$$

Substitute $x_i = x^i$ for $i = 1, \dots, r-1$. Then

$$\begin{aligned} & \sum_{t \geq k} \frac{n^t z^t x^k}{t!} \Pr(\boldsymbol{\mu}_1 + 2\boldsymbol{\mu}_2 + \dots + (r-1)\boldsymbol{\mu}_{r-1} = k) \\ &= \sum_{t=0}^{\infty} \frac{n^t z^t}{t!} \mathbf{E}(x^{\boldsymbol{\mu}_1 + 2\boldsymbol{\mu}_2 + \dots + (r-1)\boldsymbol{\mu}_{r-1}}) \\ &= \left[e^z - \left(z + \dots + \frac{z^{r-1}}{(r-1)!} \right) + \left(zx + \dots + \frac{(zx)^{r-1}}{(r-1)!} \right) \right]^n. \end{aligned}$$

because $\Pr(\boldsymbol{\mu}_1 + 2\boldsymbol{\mu}_2 + \dots + (r-1)\boldsymbol{\mu}_{r-1} = k) = 0$ if $t < k$. We again denote zx by x , then

$$\begin{aligned} & \sum_{t \geq k} \frac{n^t z^{t-k} x^k}{t!} \Pr(\boldsymbol{\mu}_1 + 2\boldsymbol{\mu}_2 + \dots + (r-1)\boldsymbol{\mu}_{r-1} = k) \\ &= \left[e^z - \left(z + \dots + \frac{z^{r-1}}{(r-1)!} \right) + \left(x + \dots + \frac{x^{r-1}}{(r-1)!} \right) \right]^n. \end{aligned}$$

We now put $z = 0$. Therefore,

$$\left(1 + x + \dots + \frac{x^{r-1}}{(r-1)!} \right)^n = \sum_{t=0}^{(r-1)n} \frac{n^t x^t}{t!} \Pr(\boldsymbol{\mu}_1 + 2\boldsymbol{\mu}_2 + \dots + (r-1)\boldsymbol{\mu}_{r-1} = t).$$

We remark that the probability is zero if $t > (r-1)n$. Let $h(x) = 1 + x \dots + \frac{x^{r-1}}{(r-1)!}$.

Lemma 6. Let $r \geq 1$. For any natural number t such that $(r-1)n \geq t \geq 0$ we have

$$\Pr(\mu_1 + 2\mu_2 + \dots + (r-1)\mu_{r-1} = t) \leq \frac{h^n(y_0)}{y_0^t} \frac{t!}{n^t},$$

where y_0 is the only nonnegative root (including ∞) of the equation $n \frac{x(1+x+\dots+\frac{x^{r-2}}{(r-2)!})}{1+x+\dots+\frac{x^{r-1}}{(r-1)!}} = t$.

Proof. Let $(r-1)n > t > 0$. Then the equation has the only positive solution y_0 . The estimate is true by Lemma 4. Let $t = 0$, then $y_0 = 0$ and the Lemma is true as both the sides of the inequality are 1. Let $t = (r-1)n$, then $y_0 = \infty$. The right hand side of the inequality is defined at $y_0 = \infty$ and equal to $\frac{t!}{((r-1)!)^n n^t}$. On the other hand,

$$\Pr(\mu_1 + 2\mu_2 + \dots + (r-1)\mu_{r-1} = t) = \Pr(\mu_{r-1}(t, n) = n) = \frac{t!}{((r-1)!)^n n^t}.$$

Lemma 7. For every integer number $k \geq 0$ it holds that

$$k^k e^{-k} \leq k! \leq k^k e^{-k} \sqrt{2\pi(k+1)}.$$

8 Complexity Estimate. Stage 1

We now estimate the expectation of $|W_r(k)|$. Its maximum in k will be estimated with (13).

Theorem 2. Let the variable sets X_1, \dots, X_m be fixed while the polynomials f_1, \dots, f_m generated according to the probabilistic model. Then

$$\mathbf{E}_{f_1, \dots, f_m} |W_r(k)| = q^{|Z_r(k)|} \prod_{i=1}^m \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_i \setminus Z_r(k)|}}\right).$$

Proof. We fix an F_q -vector a in variables $Z_r(k)$ and compute $\Pr(a \in W_r(k))$, the probability that a doesn't contradict any of $f_i(X_i) = 0$. As f_i are independent,

$$\Pr(a \in W_r(k)) = \prod_{i=1}^m \Pr(V_i(a) \neq \emptyset).$$

One sees $\Pr(V_i(a) = \emptyset) = \left(1 - \frac{1}{q}\right)^{q^{|X_i \setminus Z_r(k)|}}$ doesn't depend on a . So $\mathbf{E}_{f_1, \dots, f_m} |W_r(k)| =$

$$\sum_a \Pr(a \in W_r(k)) = q^{|Z_r(k)|} \Pr(a \in W_r(k)) = q^{|Z_r(k)|} \prod_{i=1}^m \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_i \setminus Z_r(k)|}}\right).$$

By Lemma 2, $\mathbf{E}|W_r(k)| = \mathbf{E}_{X_1, \dots, X_m} (\mathbf{E}_{f_1, \dots, f_m} |W_r(k)|)$. So

$$\mathbf{E}|W_r(k)| = \mathbf{E}_{X_1, \dots, X_m} \left(q^{|Z_r(k)|} \prod_{i=1}^m \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_i \setminus Z_r(k)|}}\right) \right). \quad (5)$$

According to the probabilistic model, X_1, \dots, X_m are uniformly allocated into the whole variable set X of size n . So we use the language of particle allocation into n boxes from now. In particular, $Z_r(k)$ is the set of boxes with at least r particles after uniform allocation by complexes of size l_1, \dots, l_k . We split the last product

$$\mathbf{E}|W_r(k)| = \mathbf{E}_{X_1, \dots, X_m} \left(q^{|Z_r(k)|} \prod_{i=1}^k \left(1 - \left(1 - \frac{1}{q} \right)^{q^{|X_i \setminus Z_r(k)|}} \right) \prod_{j=k+1}^m \left(1 - \left(1 - \frac{1}{q} \right)^{q^{|X_j \setminus Z_r(k)|}} \right) \right).$$

We use the conditional expectation formula under the condition $A = A(U, t_1, \dots, t_k)$. The event A occurs if $Z_r(k) = U$ and $|X_i \setminus U| = t_i$, where $i = 1, \dots, k$. We get

$$\mathbf{E}|W_r(k)| = \sum_U \sum_{t_1, \dots, t_k} q^{|U|} \prod_{i=1}^k \left(1 - \left(1 - \frac{1}{q} \right)^{q^{t_i}} \right) \mathbf{E}(A) \Pr(A), \quad (6)$$

where U runs over all subsets of X and $0 \leq t_i \leq l_i$. We denoted

$$\mathbf{E}(A) = \mathbf{E}_{X_1, \dots, X_m} \left(\prod_{j=k+1}^m \left(1 - \left(1 - \frac{1}{q} \right)^{q^{|X_j \setminus Z_r(k)|}} \right) \middle| A \right).$$

We remark that for any fixed $A = A(U, t_1, \dots, t_k)$ the probability $\Pr(A)$ and the expectation $\mathbf{E}(A)$ are increasing if all l_i become l . So it is enough to upper bound $\mathbf{E}|W_r(k)|$ in case $l_i = l$ only. We now estimate the probability of the event A , where $|U| = u$.

Lemma 8. *Let $L = lk$ and $T = t_1 + \dots + t_k$. Then*

$$\Pr(A) \leq \frac{\left(\frac{u}{n} \right)^{L-T} \left(\frac{n-u}{n} \right)^T P_1(L-T, u) P_2(T, n-u) \prod_{i=1}^k \binom{l}{t_i}}{\prod_{i=1}^{l-1} \left(1 - \frac{i}{n} \right)^k},$$

where

$$\begin{aligned} P_1(L-T, u) &= \Pr(\mu_0(L-T, u) = 0, \dots, \mu_{r-1}(L-T, u) = 0), \\ P_2(T, n-u) &= \Pr(\mu_1(T, n-u) + 2\mu_2(T, n-u) + \dots + (r-1)\mu_{r-1}(T, n-u) = T). \end{aligned}$$

Proof. Assume $0 < u < n$ and $r \geq 2$, otherwise the statement is easy with Lemma 3. Let the event B occur if $|X_i \setminus U| = t_i$ for $i = 1, \dots, k$. Then $\Pr(A) = \Pr(B)\Pr(A|B)$.

$$\begin{aligned} \Pr(B) &= \prod_{i=1}^k \Pr(|X_i \setminus U| = t_i) = \prod_{i=1}^k \frac{\binom{u}{l-t_i} \binom{n-u}{t_i}}{\binom{n}{l}} = \\ &= \prod_{i=1}^k \binom{l}{t_i} \left(\frac{u}{n} \right)^{l-t_i} \left(\frac{n-u}{n} \right)^{t_i} \frac{\left(1 - \frac{1}{u} \right) \dots \left(1 - \frac{l-t_i-1}{u} \right) \left(1 - \frac{1}{n-u} \right) \dots \left(1 - \frac{t_i-1}{n-u} \right)}{\left(1 - \frac{1}{n} \right) \dots \left(1 - \frac{l-1}{n} \right)} \\ &= \frac{\left(\frac{u}{n} \right)^{L-T} \left(\frac{n-u}{n} \right)^T \prod_{i=1}^k \binom{l}{t_i}}{\prod_{i=1}^{l-1} \left(1 - \frac{i}{n} \right)^k} \prod_{i=1}^k \left(1 - \frac{1}{u} \right) \dots \left(1 - \frac{l-t_i-1}{u} \right) \left(1 - \frac{1}{n-u} \right) \dots \left(1 - \frac{t_i-1}{n-u} \right). \end{aligned}$$

The event $A|B$ occurs if and only if the following two events A_1 and A_2 occur simultaneously. First, the complexes of $l-t_1, \dots, l-t_k$ particles are allocated into $|U| = u$ boxes, where each box

is occupied by at least r particles. Second, the complexes of t_1, \dots, t_k particles are allocated into $|X \setminus U| = n - u$ boxes, where each box is occupied by at most $r - 1$ particles. These are independent events. Therefore $\Pr(A|B) = \Pr(A_1)\Pr(A_2)$.

Let $\mu'_s(t_1, \dots, t_k, n)$ be the number of boxes with exactly s particles after k uniform allocations into n boxes by complexes of t_1, \dots, t_k particles. The event A_1 occurs if and only if $\mu'_i(l - t_1, \dots, l - t_k, u) = 0$ for $i = 0, \dots, r - 1$. The event A_2 occurs if and only if $\mu'_i(t_1, \dots, t_k, n - u) = 0$ for $i \geq r$. The latter is equivalent to

$$\mu'_1(t_1, \dots, t_k, n - u) + 2\mu'_2(t_1, \dots, t_k, n - u) + \dots + (r - 1)\mu'_{r-1}(t_1, \dots, t_k, n - u) = T.$$

By Lemma 3,

$$\Pr(A_1) \leq \frac{P_1(L - T, u)}{\prod_{i=1}^k \left(1 - \frac{1}{u}\right) \dots \left(1 - \frac{l - t_i - 1}{u}\right)}$$

and

$$\Pr(A_2) \leq \frac{P_2(T, n - u)}{\prod_{i=1}^k \left(1 - \frac{1}{n - u}\right) \dots \left(1 - \frac{t_i - 1}{n - u}\right)}$$

So $\Pr(A) = \Pr(B)\Pr(A|B) =$

$$= \Pr(B)\Pr(A_1)\Pr(A_2) \leq \frac{\left(\frac{u}{n}\right)^{L-T} \left(\frac{n-u}{n}\right)^T P_1(L - T, u) P_2(T, n - u) \prod_{i=1}^k \binom{l}{t_i}}{\prod_{i=1}^{l-1} \left(1 - \frac{i}{n}\right)^k}.$$

Lemma 9. $\mathbf{E}(A) = \mathbf{E}(u)$, where

$$\mathbf{E}(u) = \prod_{j=k+1}^m \mathbf{E}_{X_j} \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_j \setminus U|}}\right)$$

does not depend on the set U but rather on its size u .

Proof. For any X_1, \dots, X_k , where A occurs

$$\mathbf{E}_{X_{k+1}, \dots, X_m} \left(\prod_{j=k+1}^m \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_j \setminus Z_r(k)|}}\right) \middle| A \right) = \prod_{j=k+1}^m \mathbf{E}_{X_j} \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_j \setminus U|}}\right).$$

That doesn't depend on X_1, \dots, X_k . By Lemma 2, $\mathbf{E}(A) = \mathbf{E}_{X_1, \dots, X_k}(\mathbf{E}_{X_{k+1}, \dots, X_m}(\dots)) = \mathbf{E}(u)$. The value depends on the size of the set U and not on the set itself.

From (6) we get

$$\mathbf{E}|W_r(k)| = \sum_{u=0}^n \binom{n}{u} q^u \mathbf{E}(u) \sum_{t_1, \dots, t_k} \prod_{i=1}^k \left(1 - \left(1 - \frac{1}{q}\right)^{q^{t_i}}\right) \Pr(A) \quad (7)$$

as $\Pr(A)$ only depends on u, t_1, \dots, t_k . From (7) by Lemmas 8 and 9,

$$\begin{aligned} \mathbf{E}|W_r(k)| &\leq \frac{1}{\prod_{i=1}^{l-1} \left(1 - \frac{i}{n}\right)^k} \sum_{u=0}^n \binom{n}{u} q^u \mathbf{E}(u) \\ &\times \sum_{T=0}^L C_T \left(\frac{u}{n}\right)^{L-T} \left(\frac{n-u}{n}\right)^T P_1(L - T, u) P_2(T, n - u), \end{aligned} \quad (8)$$

where $C_T = \sum_{t_1+\dots+t_k=T} \prod_{i=1}^k \binom{l}{t_i} \left(1 - (1 - \frac{1}{q})^{q^{t_i}}\right)$. Let $f(z) = \sum_{t=0}^l \binom{l}{t} \left(1 - (1 - \frac{1}{q})^{q^t}\right) z^t$. It is obvious $f^k(z) = \sum_{T=0}^{lk} C_T z^T$.

Lemma 10. For every $0 \leq T \leq lk$ we have $C_T \leq \frac{f^k(z_0)}{z_0^T}$, where z_0 is the only nonnegative root (including ∞ for $T = lk$) to the equation $k \frac{\sum_{t=1}^l t \binom{l}{t} \left(1 - (1 - \frac{1}{q})^{q^t}\right) z^t}{\sum_{t=0}^l \binom{l}{t} \left(1 - (1 - \frac{1}{q})^{q^t}\right) z^t} = T$.

Let $u = \beta n$, then

$$\mathbf{E}_{X_i} \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_i \setminus U|}}\right) = 1 - \sum_{t=0}^l \frac{\binom{u}{l-t} \binom{n-u}{t}}{\binom{n}{l}} \left(1 - \frac{1}{q}\right)^{q^t} = 1 - \sum_{t=0}^l \frac{\binom{\beta n}{l-t} \binom{n-\beta n}{t}}{\binom{n}{l}} \left(1 - \frac{1}{q}\right)^{q^t}.$$

By taking $\lim_{n \rightarrow \infty}$, we get

Lemma 11. Let $|U| = \beta n$, where $0 \leq \beta \leq 1$ as n tends to ∞ , then

$$\mathbf{E}_{X_i} \left(1 - \left(1 - \frac{1}{q}\right)^{q^{|X_i \setminus U|}}\right) = F(\beta) + O\left(\frac{1}{n}\right),$$

where $F(\beta) = 1 - \sum_{t=0}^l \binom{l}{t} \beta^{l-t} (1 - \beta)^t \left(1 - \frac{1}{q}\right)^{q^t}$ and $O\left(\frac{1}{n}\right)$ is uniformly bounded in β .

Lemmas 9 and 11 imply $\mathbf{E}(u) \leq (F(\beta) + \epsilon)^{m-k}$, where ϵ is any positive number and n is big enough. Let $L = \alpha n$ and $T = \gamma n$. So $\frac{m-k}{n} = \frac{m}{n} - \frac{\alpha}{l}$ and

$$\mathbf{E}(u) \leq (F(\beta) + \epsilon)^{(d - \frac{\alpha}{l})n} \quad (9)$$

for any positive ϵ as n tends to ∞ . By Lemma 5, $P_1(L - T, u) \leq \frac{g^u(x_0) (L - T)!}{x_0^{L-T} u^{L-T}}$. Therefore,

$$P_1(L - T, u) \leq \left[\frac{g^\beta(x_0)}{x_0^{\alpha - \gamma}} \left(\frac{\alpha - \gamma}{\beta e}\right)^{\alpha - \gamma} + \epsilon \right]^n, \quad (10)$$

for any positive ϵ and big enough n , where x_0 is the only nonnegative root of the equation $\beta \frac{x \left(e^x - 1 - x \dots - \frac{x^{r-2}}{(r-2)!}\right)}{e^x - 1 - x \dots - \frac{x^{r-1}}{(r-1)!}} = \alpha - \gamma$. By Lemma 6, $P_2(T, n - u) \leq \frac{h^{n-u}(y_0) T!}{y_0^T (n-u)^T}$. Therefore,

$$P_2(T, n - u) \leq \left[\frac{h^{1-\beta}(y_0)}{y_0^\gamma} \left(\frac{\gamma}{(1-\beta)e}\right)^\gamma + \epsilon \right]^n, \quad (11)$$

for any positive ϵ and big enough n , where y_0 is a nonnegative root of the equation $(1-\beta) \frac{v \left(1 + y + \dots + \frac{y^{r-2}}{(r-2)!}\right)}{1 + y + \dots + \frac{y^{r-1}}{(r-1)!}} = \gamma$. By Lemma 10,

$$C_T \leq \left(\frac{f^{\frac{\alpha}{l}}(z_0)}{z_0^\gamma} \right)^n, \quad (12)$$

where z_0 is the only nonnegative root to $\frac{\alpha}{l} \frac{\sum_{t=1}^l t \binom{l}{t} \left(1 - (1 - \frac{1}{q})^{q^t}\right) z^t}{\sum_{t=0}^l \binom{l}{t} \left(1 - (1 - \frac{1}{q})^{q^t}\right) z^t} = \gamma$. From (8) with (10), (11), (12) and (9) we now get

$$\mathbf{E}|W_r(k)| \leq n(lm + 1) \max \left[\frac{q^\beta f^{\frac{\alpha}{l}}(z_0) g^\beta(x_0) h^{1-\beta}(v_0) (\alpha - \gamma)^{\alpha - \gamma} \gamma^\gamma}{\beta^\beta (1 - \beta)^{1-\beta} (z_0 v_0)^\gamma x_0^{\alpha - \gamma} e^\alpha} F(\beta)^{d - \frac{\alpha}{l}} + \epsilon \right]^n,$$

for any positive ϵ and big enough n , where Lemma 7 was used to bound the binomial coefficient $\binom{n}{u}$. Therefore,

$$\mathbf{E}|W_r(k)| \leq \left[\max \left(\frac{q^\beta f^{\frac{\alpha}{r}}(z_0) g^\beta(x_0) h^{1-\beta}(y_0) (\alpha - \gamma)^{\alpha-\gamma} \gamma^\gamma}{\beta^\beta (1 - \beta)^{1-\beta} (z_0 y_0)^\gamma x_0^{\alpha-\gamma} e^\alpha} F(\beta)^{d-\frac{\alpha}{r}} \right) + \epsilon \right]^n \quad (13)$$

for any positive ϵ and big enough n . The maximum is over $0 \leq \beta \leq 1$ and $0 \leq \gamma \leq \alpha$. We remark that the parameters α, β, γ should satisfy $r\beta \leq \alpha - \gamma$ and $(r-1)(1-\beta) \leq \gamma$, otherwise $P_1(L-T, u) = 0$ or $P_2(T, n-u) = 0$. The complexity of the first stage is upper bounded by the maximum of (13) over $0 \leq \alpha \leq dl$. For any q, l, d, r that maximum may be computed with an advanced optimization package like MAPLE.

9 Complexity Estimate. Stage 2

Let $r \geq 3$. Let $W_r = W_r(m)$ and $Z_r = Z_r(m)$. Let X_1, \dots, X_m be fixed and f_1, \dots, f_m randomly generated. Then one proves that $\mathbf{E}_{f_1, \dots, f_m} (|W_r| c^{n-|Z_r|})$ is the expected complexity to compute all solutions, where c is defined in Theorem 1. Similarly to Theorem 2,

$$\mathbf{E} \left(|W_r| c^{n-|Z_r|} \right) = \mathbf{E}_{X_1, \dots, X_m} \left(q^{|Z_r|} c^{n-|Z_r|} \prod_{i=1}^m \left(1 - \left(1 - \frac{1}{q} \right)^{q^{|X_i \setminus Z_r|}} \right) \right).$$

Let $L = lm$. Similarly to (8),

$$\begin{aligned} \mathbf{E} \left(|W_r| c^{n-|Z_r|} \right) &\leq \frac{1}{\prod_{i=1}^{l-1} \left(1 - \frac{i}{n} \right)^m} \sum_{u=0}^n \binom{n}{u} q^u c^{n-u} \\ &\times \sum_{T=0}^L C_T \left(\frac{u}{n} \right)^{L-T} \left(\frac{n-u}{n} \right)^T P_1(L-T, u) P_2(T, n-u), \end{aligned}$$

where $C_T = \sum_{t_1 + \dots + t_m = T} \prod_{i=1}^m \binom{l}{t_i} \left(1 - \left(1 - \frac{1}{q} \right)^{q^{t_i}} \right)$. Therefore,

$$\mathbf{E} \left(|W_r| c^{n-|Z_r|} \right) \leq \left[\max \left(\frac{q^\beta c^{1-\beta} f^d(z_0) g^\beta(x_0) h^{1-\beta}(y_0) (dl - \gamma)^{dl-\gamma} \gamma^\gamma}{\beta^\beta (1 - \beta)^{1-\beta} (z_0 y_0)^\gamma x_0^{dl-\gamma} e^{dl}} \right) + \epsilon \right]^n \quad (14)$$

for any positive ϵ and big enough n . The maximum is over $0 \leq \beta \leq 1$ and $0 \leq \gamma \leq dl$. We remark that the parameters α, β, γ should satisfy $r\beta \leq dl - \gamma$ and $(r-1)(1-\beta) \leq \gamma$, otherwise $P_1(L-T, u) = 0$ or $P_2(T, n-u) = 0$.

References

1. M. Bardet, J.-C. Faugère, and B. Salvy, *Complexity of Gröbner basis computation for semi-regular overdetermined sequences over F_2 with solutions in F_2* , Research report RR-5049, INRIA, 2003.
2. E.T. Copson, *Asymptotic expansions*, Cambridge University Press, 1965.
3. N. T. Courtois and G. V. Bard, *Algebraic Cryptanalysis of the Data Encryption Standard*, Crypt. ePrint Arch., report 2006/402.

4. E. Dantsin, A. Goerdt, E. A. Hirsch, R. Kannan, J. M. Kleinberg, C. H. Papadimitriou, P. Raghavan, U. Schning, *A deterministic $(2 - 2/(k + 1))^n$ algorithm for k -SAT based on local search*. Theor. Comput. Sci. 289(2002), pp.69–83.
5. J.-C. Faugère, *A new efficient algorithm for computing Gröbner bases without reduction to zero (F5)*, in ISSAC 2002, pp. 75 – 83, ACM Press, 2002.
6. K. Iwama, *Worst-Case Upper Bounds for k SAT*, The Bulletin of the EATCS, vol. 82(2004), pp. 61–71.
7. V. Kolchin, A. Sevast'yanov, and V. Chistyakov, *Random allocations*, John Wiley & Sons, 1978.
8. H. Raddum, *Solving non-linear sparse equation systems over $GF(2)$ using graphs*, University of Bergen, preprint, 2004.
9. H. Raddum, I. Semaev, *Solving Multiple Right Hand Sides linear equations*, Des. Codes Cryptogr., vol.49 (2008), pp.147–160.
10. I. Semaev, *On solving sparse algebraic equations over finite fields*, Des. Codes Cryptogr., vol. 49 (2008), pp.47–60.
11. I. Semaev, *Sparse algebraic equations over finite fields*, SIAM J. on Comp., vol. 39(2009), pp. 388-409.
12. I. Semaev, *Sparse Boolean equations and circuit lattices*, Des. Codes Cryptogr.,(2010), to appear.
13. B.-Y. Yang, J-M. Chen, and N.Courtois, *On asymptotic security estimates in XL and Gröbner bases-related algebraic cryptanalysis*, LNCS 3269, pp. 401–413, Springer-Verlag, 2004.
14. A. Zakrevskij, I. Vasilkova, *Reducing large systems of Boolean equations*, 4th Int. Workshop on Boolean Problems, Freiberg University, September, 21–22, 2000.