

# The complexity of inversion of explicit Goldreich's function by DPLL algorithms

Dmitry Itsykson \* and Dmitry Sokolov \*\*

Steklov Institute of Mathematics at St. Petersburg  
27 Fontanka, 191023, St. Petersburg, Russia  
dmitrits@pdmi.ras.ru  
St. Petersburg Academic University  
8(3) Khlopina, 194021, St.-Petersburg, Russia  
sokolov.dmt@gmail.com

**Abstract.** The Goldreich's function has  $n$  binary inputs and  $n$  binary outputs. Every output depends on  $d$  inputs and is computed from them by the fixed predicate of arity  $d$ . Every Goldreich's function is defined by its dependency graph  $G$  and predicate  $P$ . In 2000 O. Goldreich formulated a conjecture that if  $G$  is an expander and  $P$  is a random predicate of arity  $d$  then the corresponding function is one way. In 2005 M. Alekhnovich, E. Hirsch and D. Itsykson proved the exponential lower bound on the complexity of inversion of Goldreich's function based on linear predicate and random graph by myopic DPLL algorithms. In 2009 J. Cook, O. Etesami, R. Miller, and L. Trevisan extended this result to nonlinear predicates (but for a slightly weaker definition of myopic algorithms). Recently D. Itsykson and independently R. Miller proved the lower bound for drunken DPLL algorithms that invert Goldreich's function with nonlinear  $P$  and random  $G$ . All above lower bounds are randomized.

The main contribution of this paper is the simpler proof of the exponential lower bound of the Goldreich's function inversion by myopic DPLL algorithms. A dependency graph in our construction may be based on an arbitrary expander, particularly it is possible to use an explicit expander; the predicate may be linear or slightly nonlinear. Our definition of myopic algorithms is more general than one used by J. Cook et al. Our construction may be used in the proof of lower bound for drunken algorithms as well.

**Keywords:** DPLL algorithm, expander, one-way function, lower bounds

---

\* Partially supported by Federal Target Programme "Scientific and scientific-pedagogical personnel of the innovative Russia" 2009-2013, RAS Program for Fundamental Research, the president grants NSh-5282.2010.1 and MK-4089.2010.1 and by RFBR.

\*\* Partially supported by Federal Target Programme "Scientific and scientific-pedagogical personnel of the innovative Russia" 2009-2013 and by Yandex Fellowship.

## 1 Introduction

This work continues [1], [2], [3], [4] and is devoted to lower bounds of DPLL (for Davis, Putnam, Logemann, and Loveland) algorithms on satisfiable formulas. DPLL algorithm is a recursive algorithm. On each recursive call it simplifies an input formula  $F$  (without affecting its satisfiability), chooses a variable  $v$  and makes two recursive calls on the formulas  $F[v := 1]$  and  $F[v := 0]$  in some order. It returns the result “Satisfiable” if at least one of recursive calls returns “Satisfiable” (note that it is not necessary to make the second call if the first one was successful). Recursion stops if the input formula becomes trivial. That is, the algorithm is only allowed to backtrack when unsatisfiability in the current branch is proved. A DPLL algorithm is defined by simplification rules and two heuristics: the heuristic **A** chooses a variable for splitting and the heuristic **B** chooses a value that will be investigated first.

The behavior of DPLL algorithms on unsatisfiable formulas is equivalent to tree-like resolution proofs. Therefore lower bounds on DPLL algorithms on unsatisfiable formulas follow from lower bound for resolutions [5]. However the most interesting inputs are satisfiable formulas. Consider for example formulas that code the problem of inversion of one-way function. The most important case for practice is the case there one-way function indeed has preimage. There is no hope of proving a superpolynomial lower bound for all DPLL algorithms on satisfiable formulas since if  $P = NP$ , then the heuristic that chooses the value of a variable that would be investigated first may always choose the correct value.

Exponential lower bounds on running time of myopic and drunken DPLL algorithms on satisfiable formulas were proved in the paper [1]; these two classes of DPLL algorithms cover a lot of known DPLL algorithms. In myopic algorithms heuristics that choose a variable for splitting and that choose a value that will be investigated first have the following restrictions: they can see the formula with erased signs of negations and they also know the number of positive and negative occurrences of every variable and they also can request  $K = n^{1-\varepsilon}$  clauses of the formula to read them precisely. In drunken algorithms the heuristic that chooses variable for splitting may be arbitrary, while the first substituted value is chosen at random with equal probabilities. Lower bounds for myopic algorithms were proved on the formulas that code the system of linear equations over  $\mathbb{F}_2$  based on expander matrices; lower bounds for drunken algorithms were proved on artificial formulas that are based on hard examples for resolution.

The paper [2] gives a cryptographic view on [1]. Namely it was noted in [2] that the lower bound for myopic algorithms [1] was proved on the formulas that code the problem of inversion of Goldreich’s function based on linear predicate. Goldreich’s function [6] has  $n$  binary inputs and  $n$  binary outputs. Every output depends on  $d$  inputs and is computed from them by a fixed predicate of arity  $d$ . Goldreich conjectured that if the dependency graph is an expander and the predicate is random, then the resulting function is one-way. However, linear functions are not interesting from the cryptographic point of view since they can be easily inverted by Gaussian elimination. The main goal of [2] was the proof of lower bound for a function that is potentially hard to invert. J. Cook et al.

consider Goldreich's function based on the predicate  $x_1 + x_2 + \dots + x_{d-2} + x_{d-1}x_d$  and a random graph (a random graph is an expander with high probability). They have proved the exponential lower bound for the weakened<sup>1</sup> variant of myopic algorithms. Recently Itsykson [3] and Miller [4] independently proved the lower bound on the complexity of inversion of Goldreich's function based on random graph and predicate of type  $x_1 + x_2 + \dots + x_{d-k} + Q(x_{d-k+1}, \dots, x_d)$ , where  $Q$  is an arbitrary predicate of arity  $k$  and  $k < d/4$  by drunken algorithms. We should note that the proof from [2] works for this type of predicates as well.

The construction of Goldreich's function in all papers listed above was randomized. In this paper we suggest an explicit construction of Goldreich's function based on expanders (for example the explicit expander from [7] fits our purposes). It is possible to use those formulas in the proof of exponential lower bound for drunken algorithms from [3]. In this paper we demonstrate the lower bound for myopic algorithms. Our proof is technically much simpler than proofs from [1] and [2]. We prove lower bound for the general notion of myopic algorithms (according to the definition from [1]) instead of the weakened variant that was used in [2].

Our Goldreich's function has the following structure: it is the sum of two Goldreich's functions: linear and nonlinear. The linear part is necessary for proving the lower bound for DPLL algorithms while the nonlinear part makes our function hard to invert in practice. The linear part is based on an expander, while nonlinear part may be almost arbitrary but it should depend only on  $n^{\varepsilon/2}$  variables. Of course an adversary may guess the value of variables from nonlinear part and solve the resulting linear system by Gaussian elimination but the running time of such algorithm is  $2^{n^{\varepsilon/2}}$  (still exponential), therefore we believe that there are hard invertible functions among our functions. We actually do not use in the proof the fact that the nonlinear part of predicate is the same for every bit of the output.

The plan of the proof is the following: first of all we slightly modify the expander from the linear part so that its adjacency matrix would have high rank. Since the nonlinear part of our function depends on very few variables we conclude that our Goldreich's function is almost a bijection. In order to prove the lower bound we first of all prove the lower bound for unsatisfiable formulas using lower bound techniques for resolutions from [8]. Using almost linearity and almost bijectivity we prove that with high probability the myopic algorithm makes the formula unsatisfiable during first several steps and we apply the lower bound for unsatisfiable formulas.

Our proof has one disadvantage compared to the proof from [1]; namely our proof works for expanders with degrees, that are large enough, while the technique from [1] works for degrees, that are at least 3. However the proof from [1] of the fact that a myopic algorithm with high probability makes the formula unsatisfiable during first several steps is complicated, while our proof is intuitive

---

<sup>1</sup> In contrast to [1], [2] did not allow the DPLL algorithm to use pure literal simplification rules and also myopic algorithms from [2] may read only constant (opposite to  $n^{1-\varepsilon}$ ) number of clauses per step.

and based on the simple fact from elementary linear algebra: a dimension of a solution space of a satisfiable linear system does not depend on the right hand side.

## 2 Preliminaries

Let  $X = \{x_1, x_2, \dots, x_n\}$  be the set of propositional variables.

A partial substitution is a function  $\rho : X \rightarrow \{0, 1, *\}$ , that maps a variable to its value or leaves it free. The set  $Vars(\rho) = \rho^{-1}(\{0, 1\})$  is the support of the substitution; we denote  $|\rho| = |Vars(\rho)|$ .

If  $\rho_1$  and  $\rho_2$  are two partial substitutions with disjoint support then the substitution  $\rho_1 \cup \rho_2$  can be defined by the natural way.

We say that a string  $y \in \{0, 1\}^n$  is consistent with the partial substitution  $\rho$  (we denote it  $y \sim \rho$ ) if for all  $x_j$  from the support of  $\rho$  the following is satisfied  $y_j = \rho(x_j)$ .

### 2.1 DPLL algorithms

We consider a wide class of SAT algorithms: DPLL (or backtracking) algorithms. A DPLL algorithm is defined by two *heuristics* (procedures): 1) Procedure **A** maps a CNF formula to one of its variables. (This is the variable for splitting). 2) Procedure **B** maps a CNF formula and its variable to  $\{0, 1\}$ . (This value will be investigated at first).

An algorithm may also use some syntactic *simplification rules*. Simplification rules may modify the formula without affecting its satisfiability and may also make substitutions to its variables if their values can be inferred from the satisfiability of the initial formula.

A DPLL algorithm is a recursive algorithm. Its input is a formula  $\varphi$  and a partial substitution  $\rho$ .

**Algorithm 1** *Input: formula  $\varphi$  and substitution  $\rho$*

- Simplify  $\varphi$  by means of simplification rules (assume that simplification rules change  $\varphi$  and  $\rho$ ; all variables that are substituted by  $\rho$  should be deleted from  $\varphi$ ).
- If current formula is empty (that is, all its clauses are satisfied by  $\rho$ ), then return  $\rho$ . If formula contains an empty clause (unsatisfiable), then return “formula is unsatisfiable”.
- $x_j := \mathbf{A}(\varphi)$ ;  $c := \mathbf{B}(\varphi, x_j)$
- Make a recursive call with the input  $(\varphi[x_j := c], \rho \cup \{x_j := c\})$ , if the result is “formula is unsatisfiable”, then make a recursive call with the input  $(\varphi[x_j := 1 - c], \rho \cup \{x_j := 1 - c\})$  and return its result, otherwise return the result of the first recursive call.

**Definition 1.** *Myopic algorithms [1] are DPLL algorithms, where heuristics **A** and **B** have the following restrictions:*

- They can see the whole formula with erased signs of negations.
- For every variable they know the number of its positive and the number of its negative occurrences.
- They may request to read  $K = o(n)$  clauses to read precisely (with negation signs).

*Simplification rules:* 1) Unit clause elimination: if formula contains a clause with only one literal, then make a substitution that satisfies that clause. 2) Pure literals rule: if formula contains a variable that has only positive or only negative occurrences, then substitute it with the corresponding value.

The running time of a DPLL algorithm for a given sequence of random bits is the number of recursive calls.

## 2.2 Expanders

We consider bipartite graphs with each part containing  $n$  vertices. The first part we denote by  $X = \{x_1, x_2, \dots, x_n\}$  and the second we denote by  $Y = \{y_1, y_2, \dots, y_n\}$ . Every vertex from the set  $Y$  has an ordered list of its neighbours from the set  $X$  (repetitions are allowed). All considered graphs are  $d$ -regular: the degree of every vertex from  $Y$  is equal to  $d$ , where  $d$  is a constant.

Every graph has its adjacency matrix over  $\mathbb{F}_2$ . Rows of this matrix correspond to the set  $Y$  and columns correspond to the set  $X$ , the element with coordinates  $(y, x)$  contains the parity of the number of edges between  $y$  and  $x$ .

For set  $A \subseteq Y$  we denote  $\Gamma(A)$  (the set of neighbours of  $A$ ) the set of vertices from  $X$  that are connected with at least one vertex from  $A$ ; we denote  $\delta(A)$  (the boundary of  $A$ ) the set of vertices from  $X$  that have exactly one incoming edge from the set  $A$ .

**Definition 2.** *The graph  $G$  is a  $(r, d, c)$ -expander, if 1) the degree of any vertex in  $Y$  is equal to  $d$ ; 2) for any set  $A \subseteq Y, |A| \leq r$  we have  $|\Gamma(A)| \geq c|A|$ . The graph  $G$  is called a  $(r, d, c)$ -boundary expander if the second condition is replaced by: 2) for any set  $A \subseteq Y, |A| \leq r$  we have  $|\delta(A)| \geq c|A|$ .*

**Lemma 1 (cf. [1], Lemma 1).** *Every  $(r, d, c)$ -expander is also a  $(r, d, 2c - d)$ -boundary expander.*

*Proof.* Let  $A \subseteq Y, |A| \leq r$ , then  $|\Gamma(A)| \geq c|A|$ . The number of edges between  $A$  and  $\Gamma(A)$  may be estimated:  $d|A| \geq |\delta(A)| + 2|\Gamma(A) \setminus \delta(A)| = 2|\Gamma(A)| - |\delta(A)| \geq 2c|A| - |\delta(A)|$ . Finally we get  $|\delta(A)| \geq (2c - d)|A|$ .  $\square$

We need boundary expanders; for this it is enough to have an expander with constant  $c > d/2$ . For example, a random graph is an appropriate expander.

**Lemma 2 ([9], Lemma 1.9).** *For  $d \geq 32$ , for all big enough  $n$  a random bipartite  $d$ -regular graph, where parts  $X$  and  $Y$  contain  $n$  vertices is a  $(\frac{n}{10d}, d, \frac{5}{8}d)$ -expander with probability 0.9, if for every vertex in  $Y$   $d$  edges are chosen independently at random (with repetitions).*

**Corollary 1.** *In terms of Lemma 2 this graph is a  $(\frac{n}{10d}, d, \frac{1}{4}d)$ -boundary expander.*

*Proof.* Follows from Lemma 1. □

There are also explicit constructions of such expanders:

**Lemma 3 ([7]).** *For every constant  $\epsilon > 0$  there is a constant  $d$  such that it is possible to construct a  $(r, d, c)$ -expander in polynomial of  $n$  time, where  $c = (1 - \epsilon)d$ ,  $r = \Omega(n/d)$ .*

**Corollary 2.** *This graph is a  $(\Omega(n/d), d, (1 - 2\epsilon)d)$ -boundary expander.*

### 2.3 Goldreich's function

O. Goldreich in the paper [6] introduces a function  $f : \{0, 1\}^n \rightarrow \{0, 1\}^n$  defined by a graph  $G$  and a predicate  $P : \{0, 1\}^d \rightarrow \{0, 1\}$ . Every string from  $\{0, 1\}^n$  assigns some value to the variables from the set  $X = \{x_1, x_2, \dots, x_n\}$ . The value of  $(f(x))_j$  ( $j$ -th symbol of the string  $f(x)$ ) is computed in the following way: if  $y_j$  has neighbours  $x_{j_1}, x_{j_2}, \dots, x_{j_d}$ , then  $(f(x))_j = P(x_{j_1}, x_{j_2}, \dots, x_{j_d})$ .

### 2.4 Formulas from Goldreich's function

Now we describe the way we code the problem of inversion of Goldreich's function as instance of CNF satisfiability problem.

Let  $g : \{0, 1\}^\ell \rightarrow \{0, 1\}$ , the canonical CNF representation of  $g$  is the following: for every  $c \in \{0, 1\}^\ell$  that satisfies  $g(c) = 0$  we write the clause  $x_1^{c_1} \vee x_2^{c_2} \vee \dots \vee x_\ell^{c_\ell}$ , where  $x_i^0 = x_i$  and  $x_i^1 = \neg x_i$ . The whole formula is the conjunction of all written clauses.

Let  $f$  be the Goldreich's function based on the graph  $G$  and the predicate  $P$ . We represent the equation  $f(x) = b$  in the following way: for every vertex  $y_j \in Y$  that has neighbours  $x_{j_1}, x_{j_2}, \dots, x_{j_d}$  we put down the canonical CNF representation of the equality  $b_j = P(x_{j_1}, x_{j_2}, \dots, x_{j_d})$  using variables  $x_{j_1}, x_{j_2}, \dots, x_{j_d}$ . The conjunctions of all those formulas we denote  $\Phi_{f(x)=b}$ . The part of this formula that corresponds to the vertices from the set  $A \subseteq Y$  we denote  $\Phi_{f(x)=b}^A$ .

**Lemma 4.** *If a function  $g : \{0, 1\}^\ell \rightarrow \{0, 1\}$  is linear on at least two variables, then the canonical CNF representation of  $g$  has exactly  $2^{\ell-1}$  clauses and every variable has an equal number of positive and negative occurrences.*

*Proof.* Let  $g$  have the following  $\mathbb{F}_2$  representation  $g(x_1, x_2, \dots, x_n) = x_1 + x_2 + h(x_3, \dots, x_\ell)$ . Let us denote  $T_0 = h^{-1}(0)$ ,  $T_1 = h^{-1}(1)$ . Then  $g^{-1}(0) = \{00y \mid y \in T_0\} \cup \{11y \mid y \in T_0\} \cup \{01x \mid y \in T_1\} \cup \{10y \mid y \in T_1\}$ . The latter shows that  $|g^{-1}(0)| = 2^{\ell-1}$  and every variable has an equal number of positive and negative occurrences. □

Lemma 4 implies that if a myopic algorithm does not see negation signs, then it can't differ  $g(x) = 0$  from  $g(x) = 1$  when  $g$  is linear on at least 2 variables. Also we note that a canonical CNF formula is still canonical after substitution of the value of a variable.

### 3 Almost bijective Goldreich's function

#### 3.1 Linear function

Let  $G_1$  be a  $d_1$ -regular graph and  $G_2$  be a  $d_2$ -regular graph.  $G_1 + G_2$  is  $(d_1 + d_2)$ -regular graph such that for every vertex from  $Y$  the list of neighbours is a concatenation of lists of neighbours in graph  $G_1$  and graph  $G_2$ . The adjacency matrix of  $G_1 + G_2$  is the sum of adjacency matrices of  $G_1$  and  $G_2$ .

**Proposition 1.** *If graph  $G$  is a  $(r, d, c)$ -expander and  $G'$  is a  $d'$ -regular graph, then  $G + G'$  is a  $(r, d + d', c)$ -expander.*

**Theorem 1.** *Given a graph  $G$  it is possible to construct in polynomial of  $n$  time a 1-regular graph  $T$  such that the rank of adjacency matrix of  $G + T$  is at least  $n - 1$ .*

*Proof.* First of all we prove the auxiliary lemma:

**Lemma 5.** *Let  $a = (\alpha_1, \dots, \alpha_n) \in \mathbb{F}_2^n$ . Then there are at least  $n - 1$  linear independent vectors among  $b_i = (\alpha_1, \dots, \alpha_{i-1}, \alpha_i + 1, \alpha_{i+1}, \dots, \alpha_n)$ , where  $1 \leq i \leq n$ .*

*Proof.* Let us consider the matrix  $A$  of size  $n \times n$ ; all columns of  $A$  are equal to vector  $a$ . Vectors  $b_i$  are columns of the matrix  $A + E$ , where  $E$  is the identity matrix. Since the rank of the sum of matrices is less then or equal to the the sum of ranks we may conclude that  $n = \text{rk } E \leq \text{rk}(A + E) + \text{rk } A$ . All columns of  $A$  are the same, hence  $\text{rk } A \leq 1$  and  $\text{rk}(A + E) \geq n - 1$ .  $\square$

Now we describe the construction of graph  $T$ . We start from an empty set of edges and we will add one edge per step. On the  $i$ -th step for  $1 \leq i \leq n - 1$  we add a neighbour to the vertex  $y_i \in Y$  in such a way that the first  $i$  rows of  $G + T$  are linearly independent. It can be done by the Lemma 5 (we apply the Lemma to the  $i$ -th row of the matrix of graph  $G$ ). We add an arbitrary neighbour to vertex  $y_n$ .  $\square$

**Corollary 3.** *If  $G$  is a  $(r, d, c)$ -expander, then the graph  $G + T$  from the theorem is a  $(r, d + 1, c)$ -expander and the Goldreich's function  $f$  based on  $G + T$  and a linear predicate of arity  $d + 1$  has the following property: for every  $b \in \{0, 1\}^n$  the size of the set  $f^{-1}(b)$  is at most 2.*

#### 3.2 Slightly nonlinear Goldreich's function

Let  $R \subseteq X$  be some subset of  $X$ .  $R$ -graph is a regular graph such that all vertices from  $X \setminus R$  have degree 0.

**Lemma 6.** *Let  $G$  be a  $(d - k)$ -regular graph with an adjacency matrix of rank at least  $n - 1$  and  $H$  be a  $k$ -regular  $R$ -graph. Let  $f$  be Goldreich's function based on  $G + H$  and predicate  $x_1 + x_2 + \dots + x_{d-k} + Q(x_{d-k+1}, \dots, x_d)$ , where  $Q$  is an arbitrary predicate of arity  $k$ . Then for every  $b \in \{0, 1\}^n$  the size of the set  $f^{-1}(b)$  is at most  $2^{|R|+1}$ .*

*Proof.* We consider the system of equalities  $f(x) = b$  and fix values of all variables from the set  $R$ . We get the linear system whose matrix equals to the matrix of graph  $G$  after removing the columns from the set  $R$ . The matrix of  $G$  has rank  $n - 1$ , therefore the resulting system has at most two solutions. Hence the initial system  $f(x) = b$  has at most  $2^{|R|+1}$  solutions.  $\square$

## 4 Lower bound on unsatisfiable formulas

We say that a variable is *sensitive* if by changing its value we change the value of the formula (for every assignment of values of other variables). (The boolean function that corresponds to the formula is linear on all its sensitive variables).

**Theorem 2 ([3]).** *Let  $f$  be a Goldreich's function based on  $G$  and  $P$ , where graph  $G$  is a  $(r, d, c)$ -boundary expander and predicate  $P$  contains at most  $k$  insensitive variables;  $\rho$  is a partial assignment to variables of  $X$  such that the formula  $\Phi_{f(x)=b}|_{\rho}$  is unsatisfiable and for any set of vertices  $A \subseteq Y$ ,  $|A| < \frac{r}{2}$ , the formula  $\Phi_{f(x)=b}^A|_{\rho}$  is satisfiable. Then the running time of any DPLL algorithm (that does not use simplification rules) on the formula  $\Phi_{f(x)=b}|_{\rho}$  is at least  $2^{\frac{(c-k)r}{4} - |\rho| - d}$ .*

## 5 Lower bound on satisfiable formulas

### 5.1 Closure

Let graph  $G$  be a  $(r, d, c)$ -boundary expander. Let  $0 < k < c - 2$  and  $P(x_1, \dots, x_d) = x_1 + \dots + x_{d-k} + Q(x_{d-k+1}, \dots, x_d)$ ; Goldreich's function  $f$  is based on  $G$  and  $P$ .

The next technical definition formalize the following simple idea: suppose that the set  $J$  is removed from the part  $X$  of  $G$  and we want to remove a set  $I$  from  $Y$  (and also  $\Gamma(I)$  from  $X$ ) such that the resulting graph becomes a  $(r/2, d, k + 1)$ -boundary expander. We construct such  $I$  step by step removing sets with small boundary from  $Y$ .

**Definition 3.** *Let  $J \subseteq X$ . The set of vertices  $I \subseteq Y$  is called  $k$ -closure of the set  $J$  if there is a finite sequence of sets  $I_1, I_2, \dots, I_m$  (we denote  $C_\ell = \bigcup_{1 \leq i \leq \ell} I_i$ ,  $C_0 = \emptyset$ ), such that the following properties are satisfied:*

- $I_\ell \subseteq Y$  and  $0 < |I_\ell| \leq \frac{r}{2}$  for all  $1 \leq \ell \leq m$ ;
- $I_i \cap I_j = \emptyset$  for all  $1 \leq i, j \leq m$ ;
- $|\delta(I_\ell) \setminus (\Gamma(C_{\ell-1}) \cup J)| \leq (1 + k)|I_\ell|$ ; for all  $1 \leq \ell \leq m$ ;
- for all  $I' \subseteq Y \setminus C_m$  if  $0 < |I'| \leq \frac{r}{2}$ , then  $|\delta(I') \setminus (\Gamma(C_m) \cup J)| > (1 + k)|I'|$ ;
- $I = C_m$ .

The set of all  $k$ -closures of the set  $J$  we denote as  $Cl^k(J)$ .

**Lemma 7.** *1. For every set  $J \subseteq X$  there exists a  $k$ -closure. 2. Let  $J_1 \subseteq J_2$ , then for every  $I_1 \in Cl^k(J_1)$  there exists  $I_2 \in Cl^k(J_2)$  such that  $I_1 \subseteq I_2$*

**Lemma 8 ([1]).** Let  $|J| < \frac{(c-k-1)r}{2}$ , then for every set  $I \in Cl^k(J)$  the inequality  $|I| \leq (c-k-1)^{-1}|J|$  is satisfied

**Definition 4.** Let  $f : \{0,1\}^n \rightarrow \{0,1\}^n$  be the Goldreich's function based on graph  $G$  and predicate  $P$ ,  $b \in \{0,1\}$ . Partial substitution  $\rho$  is called locally consistent for the equation  $f(x) = b$  if there exists a string  $z \in \{0,1\}^n$  that is consistent to  $\rho$  and a set  $I \in Cl^k(Vars(\rho))$  such that the equality  $f(z)|_I = b|_I$  holds.

**Lemma 9 (cf. [1]).** If the partial substitution  $\rho$  is locally consistent for  $f(x) = b$ , then for all  $Z \subseteq X$ ,  $|Z| \leq \frac{r}{2}$  there exists a string  $z \in \{0,1\}^n$  such that  $z$  is consistent with  $\rho$  and the equality  $f(z)|_Z = b|_Z$  holds.

*Proof.* Proof by contradiction. Consider the minimal  $Z \subseteq Y$  such that  $|Z| \leq \frac{r}{2}$  and for all  $z$  that are consistent to  $\rho$  the nonequality  $f(z)|_Z \neq b|_Z$  holds. Let  $I \in Cl^k(Vars(\rho))$  be from Definition 4. Partial substitution  $\rho$  is locally consistent therefore  $Z \setminus I \neq \emptyset$ .

By the definition of closure  $|\delta(Z \setminus I) \setminus (\Gamma(I) \cup Vars(\rho))| > (k+1)|Z \setminus I|$ , therefore there exists  $y \in Z \setminus I$  such that at least  $k+1$  boundary vertices of set  $Z$  (not from the support of  $\rho$  and not connected with  $I$ ) are connected with  $y$ . The minimality of  $Z$  implies that there exists  $z \in \{0,1\}^n$ , such that  $z \sim \rho$  and  $f(z)|_{Z \setminus \{y\}} = b|_{Z \setminus \{y\}}$ . It is possible also to satisfy the equation corresponding to vertex  $y$  by flipping the  $z$ -value of one of the boundary neighbours of vertex  $y$ . Therefore there exists  $z' \in \{0,1\}^n$  that is consistent with  $\rho$  and  $f(z')|_Z = b|_Z$ . Contradiction.  $\square$

## 5.2 Clever myopic algorithm

We assume that the myopic algorithm runs on the formula  $\Phi_{f(x)=b}$ , where  $f^{-1}(b) \neq \emptyset$ . We describe the clever myopic algorithm. A clever myopic algorithm is allowed to read more clauses precisely (equivalently it may open more bits of  $b$ ). Besides, the clever algorithm doesn't make substitutions that obviously lead to unsatisfiable formulas. It is not hard to see that it is enough to prove the lower bound for clever myopic algorithms; the lower bound for all myopic algorithms will follow.

Now we describe the behavior of clever myopic algorithms more formally. A clever algorithm has a current partial substitution  $\rho$  and a set  $I \in Cl^k(Vars(\rho))$ . At the beginning  $\rho = \emptyset$ ,  $I = \emptyset$ . On each step the clever algorithm simplifies the formula (probably increases  $\rho$  and extends the set  $I$  to the element of  $Cl^k(Vars(\rho))$ ).

If the clever algorithm requests a clause that corresponds to the vertex  $y_j \in Y$  we say that the algorithm opens  $j$ -th bit of output. We assume that all clauses corresponding to  $y_j \in Y$  may be read by a clever algorithm for free.

Consider the heuristic **A** that choose variable  $x$  for splitting. Let  $Z$  be the set of all open bits of output (in particular  $Z$  includes  $K$  bits that were open before  $x$  was chosen). The clever algorithm extends the set  $I$  to the element of

$Cl^k(Vars(\rho) \cup \{x\})$ . The set of open bits is increased:  $Z := Z \cup I$ . The clever algorithm chooses the value of variable  $x$  in order to make the part of formula that corresponds to  $Z$  satisfiable.

**Lemma 10.** *For every clever myopic algorithm  $A$  there exists another clever myopic algorithm  $B$  such that  $B$  does not use pure literal and unit clause elimination rules and the running time of algorithm  $B$  on the formula  $\Phi_{f(x)=b}$  is bounded by polynomial on the running time of algorithm  $A$ .*

*Proof.* If the current predicate in the vertex  $y \in Y$  (taking into account  $\rho$ ) is linear on at least two variables then Lemma 4 implies that there are no pure literals in the formula that corresponds to  $y$ . So predicates in vertices that contain pure literals have at most one linear variable. All such vertices are contained in  $I \in Cl^k(Vars(\rho))$ , hence all corresponding bits of output are open and the algorithm may make a substitution to this pure literal by itself. Similarly, if formula contains a unit clause, then the corresponding vertex is in  $I$  and a clever algorithm may choose the correct substitution by itself.  $\square$

In the following we assume that clever myopic algorithms do not use simplification rules.

Let us denote  $N = \lfloor \frac{(c-k-1)r}{4dK} \rfloor$ .

**Lemma 11 (cf. [1]).** *After  $N$  steps of any clever myopic algorithm the number of open bits is at most  $\frac{r}{2}$ .*

*Proof.* The number of open bits is at most  $K \frac{(c-k-1)r}{4dK} + |Cl^k(Vars(\rho))|$ , where  $\rho$  is the current substitution. By Lemma 8  $|Cl^k(Vars(\rho))| \leq \frac{|Vars(\rho)|}{c-k-1}$ . Since  $|Vars(\rho)| \leq \frac{(c-k-1)r}{4}$  we may conclude  $K \frac{(c-k-1)r}{4dK} + |Cl^k(Z)| \leq \frac{(c-k-1)r}{4d} + \frac{r}{4} \leq \frac{r}{2}$ .  $\square$

**Corollary 4.** *During the first  $N$  steps a clever myopic algorithm does not backtrack (backtracking corresponds to a leaf of the splitting tree) and  $\rho$  is locally consistent.*

*Proof.* During  $N$  steps the number of open bits is at most  $\frac{r}{2}$ . We prove by induction that the current substitution is locally consistent. It is trivial for the beginning. Induction step follows from the fact that the value of the variable is chosen in such a way that  $\Phi_{f(x)=b}|_I$  is satisfiable. This is possible by Lemma 9 and by induction hypothesis.  $\square$

Our goal is to show that after  $N$  steps of a clever myopic algorithm the current formula will be unsatisfiable with high probability.

From this point we assume that graph  $G$  has the type  $G_L + H$ , where  $G_L$  is a  $(d-k)$ -regular and  $H$  is a  $k$ -regular  $R$ -graph; the rank of adjacency matrix of  $G_L$  is at least  $n-1$ . The Goldreich's function  $f$  based on  $G$  and  $P$  is linear on variables  $X \setminus R$ .

**Lemma 12.** Let  $b \in \{0, 1\}^n$  and  $J \subseteq X$ . Let  $y \in \{0, 1\}^n$  and  $Z \subseteq Y$ , we define set  $X_y = \{x \in \{0, 1\}^n \mid \forall j \in (X \setminus J) x_j = y_j\}$  and set  $S_y = \{x \in X_y \mid f(x)|_Z = b|_Z\}$ . Then either  $|S_y| \geq 2^{|J|-|Z|-|J \cap R|}$  or  $|S_y| = 0$ .

*Proof.* We have to estimate the number of  $x \in X_y$  that satisfies the system of equalities  $f(x)|_Z = b|_Z$ . If we fix the values for variables  $x_j$  for  $j \in J \cap R$  then the system becomes linear over variables  $x_j$  for  $j \in (J \setminus R)$ . The rank of the system does not exceed  $|Z|$  and the number of variables is at least  $|J| - |J \cap R|$  (it is not necessary for all those variables to have explicit occurrences in the system). Thus if a solution exists then the dimension of the solution space is at least  $|J| - |J \cap R| - |Z|$ . Since our system is over field  $\mathbb{F}_2$  the number of solutions is at least  $2^{|J|-|Z|-|J \cap R|}$  even for fixed values of  $x_j$ ,  $j \in J \cap R$ .  $\square$

Let  $Z$  be the set of open bits  $b$  in the equation  $f(x) = b$ ,  $\rho$  be some partial substitution; we denote  $C_{\rho, Z, b}$  the set of  $x \in \{0, 1\}^n$  that are consistent with  $\rho$  and satisfy  $f(x)|_Z = b|_Z$ . Formally  $C_{\rho, Z, b} = \{x \mid f(x)|_Z = b|_Z, x \sim \rho\}$ .

**Lemma 13.** Let  $Z \subseteq Y$ ,  $|Z| < \frac{r}{2}$ ,  $J \subseteq X$ . Then for every two locally consistent substitutions  $\rho_1, \rho_2$  with  $\text{Vars}(\rho_1) = \text{Vars}(\rho_2) = J$  and for every  $b \in \{0, 1\}^n$  the following is satisfied:  $\frac{|C_{\rho_1, Z, b}|}{|C_{\rho_2, Z, b}|} \leq 2^{|R|}$ .

*Proof.*

$$\frac{|C_{\rho_1, Z, b}|}{|C_{\rho_2, Z, b}|} = \frac{\sum_{\sigma} |C_{\rho_1 \cup \sigma, Z, b}|}{\sum_{\sigma} |C_{\rho_2 \cup \sigma, Z, b}|},$$

where the sum in both cases is over partial substitutions  $\sigma$  with support  $\text{Vars}(\sigma) = R \setminus J$ .

We show that the size of the set  $C_{\rho_i \cup \sigma, Z, b}$  is either 0 or some fixed value and not dependant on  $\sigma$  and  $i \in \{1, 2\}$ .

The size of the set  $C_{\rho_i \cup \sigma, Z, b}$  equals the number of solutions of the system of equations  $f(x)|_Z = b|_Z$  if some bits of  $x$  are fixed by substitution  $\rho_i \cup \sigma$ . This fixation makes the system linear. Note that the rank of this system does not depend on substitutions  $\rho_i$  and  $\sigma$  (since  $\rho_i$  and  $\sigma$  influence only the column of constants in the system). Therefore, if such system has a solution then the number of solutions does not depend on  $i$  and  $\sigma$ .

Since the substitution  $\rho_i$  is locally consistent and  $|Z| < \frac{r}{2}$ , Lemma 9 implies that there exists such substitution  $\sigma_i$  with support  $\text{Vars}(\sigma_i) = R \setminus J$  that  $C_{\rho_i \cup \sigma_i, Z, b} \neq \emptyset$ .

$$\frac{|C_{\rho_1, Z, b}|}{|C_{\rho_2, Z, b}|} \leq \frac{2^{|R|} |C_{\rho_1 \cup \sigma_1, Z, b}|}{|C_{\rho_2 \cup \sigma_2, Z, b}|} = 2^{|R|}.$$

$\square$

**Theorem 3.** Assume  $|R| = o(\frac{n}{K})$  and  $\rho$  is the current substitution after  $N$  steps of a clever myopic algorithm running on the formula  $\Phi_{f(x)=b}$  for some  $b \in f(\{0, 1\}^n)$  and  $Z$  is the set of open bits. Then  $\Pr_{y \leftarrow U(\{0, 1\}^n)}[\exists x : x \sim \rho, f(x) = f(y) \mid f(y)|_Z = b|_Z] \leq 2^{-\Omega(\frac{n}{K})}$ .

Before giving a formal prove we informally describe the main idea. For simplicity we assume that  $R = \emptyset$  therefore the predicate  $P$  is linear . We consider a clever myopic algorithm after  $N$  steps (i.e.  $|\rho| = N$ ). In this moment the size of the set  $I \in Cl^k(Vars(\rho))$  does not exceed  $(1 - \varepsilon)N$  for some positive  $\varepsilon$  by Lemma 8. We apply Lemma 12 for  $J = Vars(\rho)$  and  $Z = I$ , Lemma 12 states that the number of locally consistent substitutions is at least  $2^{|Vars(\rho)| - |I|} = 2^{\Omega(N)}$ .

Lemma 9 and Lemma 11 imply that every local consistent partial substitution may be extended to the full substitution that is consistent with open bits of the right hand side. Lemma 13 states that the number of such extensions is the same for every locally consistent substitution if  $R = \emptyset$ . A myopic algorithm has no chance to find one substitution among all locally consistent substitutions since they all have equal chances to be correct. Since our linear system has at most two solutions (if  $R = \emptyset$ ), there are at most two locally consistent substitutions that can be extended to the solution of the system. Therefore the probability of correct substitution is at most  $2^{-\Omega(N)} = 2^{-\Omega(\frac{N}{K})}$ .

*Proof (Theorem 3).* Corollary 4 implies that during  $N$  steps the algorithm does not backtrack and  $|\rho| = N$ .

We apply Lemma 12 for  $J = Vars(\rho)$  and  $Z = I$ , where  $I \in Cl^k(Vars(\rho))$  is from definition of a clever myopic algorithm after step  $N$ . Since  $b \in f(\{0, 1\}^n)$  there exists  $y \in \{0, 1\}^n$  such that  $S_y \neq \emptyset$  ( $S_y$  is defined in the Lemma 12) and the inequality  $|S_y| \geq 2^{|Vars(\rho)| - |I| - |R|}$  holds. Therefore at least  $2^{|Vars(\rho)| - |I| - |R|}$  substitutions with support  $Vars(\rho)$  are locally consistent.

$$\begin{aligned} & \Pr_{y \leftarrow U(\{0,1\}^n)} [\exists x : x \sim \rho, f(x) = f(y) \mid f(y)|_Z = b|_Z] \\ &= \Pr_{y \leftarrow U(\{0,1\}^n)} [f^{-1}(f(y)) \cap C_{\rho, Z, b} \neq \emptyset \mid f(y)|_Z = b|_Z] \\ &\leq \max_y |f^{-1}(f(y))| \cdot \Pr_{y \leftarrow U(\{0,1\}^n)} [y \in C_{\rho, Z, b} \mid f(y)|_Z = b|_Z] \end{aligned}$$

By Lemma 6 the first term may be estimated as  $\max_y |f^{-1}(f(y))| \leq 2^{|R|+1}$ . Let us estimate the second term:  $\Pr_{y \leftarrow U(\{0,1\}^n)} [y \in C_{\rho, Z, b} \mid f(y)|_Z = b|_Z] \leq \frac{\max_{\sigma} |C_{\sigma, Z, b}|}{\sum_{\sigma} |C_{\sigma, Z, b}|}$ , where  $\sigma$  goes through all locally correct substitutions with the support  $Vars(\rho)$ . By Lemma 13  $\frac{\max_{\sigma} |C_{\sigma, Z, b}|}{\sum_{\sigma} |C_{\sigma, Z, b}|} \leq 2^{|R|} \frac{\min_{\sigma} |C_{\sigma, Z, b}|}{2^{|Vars(\rho)| - |I| - |R|} \min_{\sigma} |C_{\sigma, Z, b}|} = 2^{2|R| + |I| - |Vars(\rho)|}$ .

Altogether:

$$\Pr_{y \leftarrow U(\{0,1\}^n)} [\exists x : x \sim \rho, f(x) = f(y) \mid f(y)|_Z = b|_Z] \leq 2^{3|R| + |I| - |Vars(\rho)| + 1}.$$

Since  $I \in Cl^k(Vars(\rho))$  the Lemma 8 implies  $|I| \leq (c - k - 1)^{-1} |Vars(\rho)|$ . The statement of the theorem follows from  $Vars(\rho) = \Omega(\frac{N}{K})$  and  $c > k + 2$ .  $\square$

**Theorem 4.** Let  $|R| = o(\frac{n}{K})$ , then for every myopic algorithm  $A$  the following inequality holds:  $\Pr_{y,s}[t_A(\Phi_{f(x)=f(y)}) \geq 2^{\Omega(n)}] \geq 1 - 2^{-\Omega(\frac{n}{K})}$ , where  $t_A(x)$  denotes the running time of  $A$  on input  $x$  and  $s$  is a string of random bits used by  $A$ .

*Proof.* Lemma 10 implies that it is enough to prove the Theorem for clever myopic algorithms that do not use simplification rules.

We fix the string of random bits  $s$  and prove that for algorithms that use  $s$  instead of random bits the following holds:  $\Pr_y[t_A(\Phi_{f(x)=f(y)}) \geq 2^{\Omega(n)}] \geq 1 - 2^{-\Omega(\frac{n}{K})}$ , and the theorem follows.

We consider a clever myopic algorithm after  $N$  steps on the formula  $\Phi_{f(x)=f(y)}$ . Let  $Z_y$  be the set of open bits of output by this moment. Note that for a fixed string  $s$  the behavior of algorithm during the first  $N$  steps is the same for all  $y' \in \{0,1\}^n$  such that  $f(y')|_{Z_y} = f(y)|_{Z_y}$  (in this case  $Z_{y'} = Z_y$ ). Thus the set of all  $y \in \{0,1\}^n$  may be split on the finite number of classes of equivalence  $S_1, S_2, \dots, S_m$  such that for all  $y$  and for all  $y' \in S_y$  the values of  $Z'_y$  are the same and the values of  $f(y)|_{Z_y}$  are the same, and this is not true for different classes.

$$\Pr_y[t_A(\Phi_{f(x)=f(y)}) \geq 2^{\Omega(n)}] = \sum_{i=1}^m \Pr_y[t_A(\Phi_{f(x)=f(y)}) \geq 2^{\Omega(n)} \mid y \in S_i] \Pr_y[y \in S_i].$$

By Theorem 3 after  $N$  steps of a clever myopic algorithm

$$\Pr_y[\Phi_{f(x)=f(y)}|_{\rho} \text{ is unsatisfiable} \mid y \in S_i] \geq 1 - 2^{-\Omega(\frac{n}{K})},$$

where  $\rho$  is the current substitution that is locally consistent. Finally Theorem 2 implies that  $\Pr_y[t_A(\Phi_{f(x)=f(y)}) \geq 2^{\Omega(n)} \mid y \in S_i] \geq 1 - 2^{-\Omega(\frac{n}{K})}$ . The theorem follows from the last inequality.  $\square$

In conclusion we describe the construction of Goldreich's function that suits the previous theorem.

We choose  $\epsilon = \frac{1}{4k+1}$  and for given  $\epsilon$  we construct an  $(r, d, (1-\epsilon)c)$ -expander  $H$  by Lemma 3. The constant  $d$  satisfies the inequality  $d \geq 4k+1$ . By Theorem 1 we add to the constructed graph such 1-regular graph  $T$  that the resulting graph  $H+T$  has the adjacency matrix with rank at least  $n-1$ . The resulting graph is a  $(r, d+1, (1-\frac{1}{4k+1})d)$ -expander. We choose the subset  $R \subseteq X$  of size  $o(n/K)$  and  $k$ -regular  $R$ -graph  $F$ . We define  $G = H+T+F$ ; graph  $G$  is a  $(r, d+1+k, (1-\frac{1}{4k+1})d)$ -expander and hence a  $(r, d+1+k, d(1-\frac{1}{4k+1})-k-1)$ -boundary expander. For  $k > 1$  the inequality  $d(1-\frac{1}{4k+1})-k-1 > k+2$  holds. For such graph any predicate of the type  $x_1 + \dots + x_{d-k} + Q(x_{d-k+1}, \dots, x_d)$  is suitable, where  $Q$  is arbitrary predicate of arity  $k$ . It may be easily verified that we do not use the fact that the predicate  $Q$  is the same for all vertices of the set  $Y$ .

## References

1. Michael Alekhnovich, Edward A. Hirsch, and Dmitry Itsykson. Exponential lower bounds for the running time of DPLL algorithms on satisfiable formulas. *J. Autom. Reason.*, 35(1-3):51–72, 2005.
2. James Cook, Omid Etesami, Rachel Miller, and Luca Trevisan. Goldreich’s one-way function candidate and myopic backtracking algorithms. In *proceedings of TCC*, pages 521–538. Springer-Verlag, 2009.
3. D. Itsykson. Lower bound on average-case complexity of inversion of Goldreich function by ”drunken” backtracking algorithms. In *Proceedings of International Computer Science Symposium in Russia*, volume 6072 of *Lecture Notes in Computer Science*, pages 204–215. Springer, 2010.
4. Rachel Miller. Goldreich’s one-way function candidate and drunken backtracking algorithms. Master’s thesis, University of Virginia, 2009. Distinguished Majors Thesis.
5. G. S. Tseitin. On the complexity of derivation in the propositional calculus. *Zapiski nauchnykh seminarov LOMI*, 8:234–259, 1968. English translation of this volume: Consultants Bureau, N.Y., 1970, pp. 115–125.
6. Oded Goldreich. Candidate one-way functions based on expander graphs. Technical Report 00-090, Electronic Colloquium on Computational Complexity, 2000.
7. M. Capalbo, O. Reingold, S. Vadhan, and A. Wigderson. Randomness conductors and constant-degree expansion beyond the degree/2 barrier. In *Proceedings of the 34th Annual ACM Symposium on Theory of Computing*, pages 659–668, 2002.
8. E. Ben-Sasson and A. Wigderson. Short proofs are narrow — resolution made simple. *Journal of ACM*, 48(2):149–169, 2001.
9. S. Hoory, N. Linial, and A. Wigderson. Expander graphs and their applications. *Bulletin of the American Mathematical Society*, 43:439–561, 2006.