

# Evolving S-boxes Based on Cellular Automata with Genetic Programming

Stjepan Picek<sup>1</sup>, Luca Mariot<sup>2</sup>, Domagoj Jakobovic<sup>3</sup>, and Alberto  
Leporati<sup>1</sup>

<sup>1</sup>KU Leuven, imec-COSIC, Kasteelpark Arenberg 10 3001 Leuven, Belgium ,  
[stjepan@computer.org](mailto:stjepan@computer.org)

<sup>2</sup>Dipartimento di Informatica, Sistemistica e Comunicazione, Università  
degli Studi di Milano-Bicocca, Viale Sarca 336, 20126 Milano, Italy ,  
[{luca.mariot, leporati}@disco.unimib.it](mailto:{luca.mariot, leporati}@disco.unimib.it)

<sup>3</sup>University of Zagreb, Unska 3, 10000 Zagreb, Croatia ,  
[domagoj.jakobovic@fer.hr](mailto:domagoj.jakobovic@fer.hr)

## Abstract

The design of cryptographically strong Substitution Boxes (S-boxes) is an interesting problem from both a cryptographic perspective as well as the combinatorial optimization one. Here we introduce the concept of evolving cellular automata rules that can be then translated into S-boxes. With it, we are able to find optimal S-boxes for sizes from  $4 \times 4$  up to  $7 \times 7$ . As far as we know, this is the first time a heuristic approach is able to find optimal S-boxes for sizes larger than 4.

**Keywords** Substitution boxes, Genetic Programming, Cellular automata, Cryptography

## 1 Introduction

In the process of designing block ciphers, one well explored direction is to build the so-called Substitution-Permutation Network (SPN) cipher. Such ciphers usually consist of an XOR operation with the key/subkeys, a linear layer, and a substitution layer [3]. A standard way to build the substitution layer is to use one or more Substitution Boxes (S-boxes) where a number of properties need to be satisfied for an S-box to be useful in practice. An S-box, or  $(n, m)$  function, is a mapping from  $n$  inputs into  $m$  outputs.

From the cryptographic properties perspective, the bare minimum one would need to consider when designing S-boxes with the same number of inputs and outputs (as we consider in this paper) is for them to be

bijjective, with high nonlinearity, and low differential uniformity. When utilizing heuristics in the design of S-boxes, a common approach is to use the permutation encoding since it ensures the bijectivity property. However, already for the size  $5 \times 5$  heuristics are not able to reach the optimal values of nonlinearity and differential uniformity and therefore algebraic constructions are the common method of choice [5]. However, there are several ciphers that use S-boxes not directly obtained by algebraic constructions or heuristics, but where the S-boxes are actually cellular automata (CA). The best known example is the Keccak sponge construction that is now part of the SHA-3 standard [1].

In this paper, we focus on the investigation of CA rules that are able to produce S-boxes with optimal cryptographic properties. In particular, we use Genetic Programming (GP) to evolve CA local rules in the form of Boolean functions, by viewing the corresponding CA as S-boxes. With our approach we are able to produce large quantities of S-boxes with optimal cryptographic properties, dependent on the chosen optimization criteria.

## 2 S-boxes and Cellular Automata

Let  $n, m$  be positive integers, i.e.,  $n, m \in \mathbb{N}^+$ . The set of all  $n$ -tuples of elements in the field  $\mathbb{F}_2$  is denoted as  $\mathbb{F}_2^n$ , where  $\mathbb{F}_2$  is the Galois field with two elements. The inner product of two vectors  $a$  and  $b$  from  $\mathbb{F}_2^n$  equals  $a \cdot b = \bigoplus_{i=1}^n a_i b_i$  where “ $\bigoplus$ ” represents addition modulo two. An  $(n, m)$ -function is any mapping  $F$  from  $\mathbb{F}_2^n$  to  $\mathbb{F}_2^m$ . For any set  $S$ , we denote  $S \setminus \{0\}$  by  $S^*$ .

An  $(n, m)$ -function  $F$  is *balanced* if it takes every value of  $\mathbb{F}_2^m$  the same number  $2^{n-m}$  of times.

The *nonlinearity*  $N_F$  of an  $(n, m)$ -function  $F$  equals the minimum nonlinearity of all its component functions  $v \cdot F$ , where  $v \in \mathbb{F}_2^{m*}$  [2]:

$$N_F = 2^{n-1} - \frac{1}{2} \max_{\substack{a \in \mathbb{F}_2^n \\ v \in \mathbb{F}_2^{m*}}} |W_F(a, v)|. \quad (1)$$

Here,  $W_F(a, v)$  is the Walsh-Hadamard transform of an  $(n, m)$ -function  $F$ , defined as:

$$W_F(a, v) = \sum_{x \in \mathbb{F}_2^n} (-1)^{v \cdot F(x) \oplus a \cdot x}, \quad a, v \in \mathbb{F}_2^m. \quad (2)$$

Let  $F: \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$ ,  $a \in \mathbb{F}_2^n$  and  $b \in \mathbb{F}_2^m$ . We denote:

$$D_F(a, b) = \{x \in \mathbb{F}_2^n : F(x) + F(x + a) = b\}. \quad (3)$$

The entry at the position  $(a, b)$  corresponds to the cardinality of the delta difference table  $D_F(a, b)$ . The *differential uniformity*  $\delta_F$  is then defined as [4]:

$$\delta_F = \max_{a \neq 0, b} D_F(a, b). \quad (4)$$

Cellular automata (CA) are a parallel computational model which can be represented as a regular grid of cells, such that each cell synchronously updates its state according to a local rule, which depends on a specific number of neighboring cells. In this paper we consider a CA model where the cells are arranged on a finite periodic one-dimensional array and are described by a binary state, 0 or 1.

Formally, let  $m, n \in \mathbb{N}$  with  $m \geq n$ , and let  $f : \mathbb{F}_2^n \rightarrow \mathbb{F}_2$  be a Boolean function. A *cellular automaton* (CA) of length  $m$  with *local rule*  $f$  and *periodic boundary condition* is a vectorial Boolean function  $F : \mathbb{F}_2^m \rightarrow \mathbb{F}_2^m$  defined for all  $x = (x_1, \dots, x_m) \in \mathbb{F}_2^m$  as:

$$F(x_1, \dots, x_m) = (f(x_1, x_2, \dots, x_n), \dots, f(x_m, x_1, \dots, x_{n-1})) . \quad (5)$$

In the rest of the paper we assume  $m = n$ , and we identify  $F$  (the CA) with an S-box.

### 3 GP Approach and Results

We use GP to evolve Boolean functions of  $n$  variables, in the form of trees, which are used as CA local rules. In this process, we assume that the state of a CA is represented as a periodic one-dimensional binary array of size  $n$ . The elements of the binary array are used as GP terminals, where the terminal  $c_0$  denotes the value that is being updated. The terminals  $c_1, \dots, c_{n-1}$  denote the cells to the right of the current cell. In our experiments, the neighborhood of a cell is formed by the cell itself and the  $n - 1$  cells to its right, so each value in the current state can be used in a local update rule. A candidate Boolean function obtained with GP is evaluated in the following manner: all the possible  $2^n$  input states are considered, and for each state the same rule is applied in parallel to each of the bits to determine the next state. The obtained global rule represents a candidate S-box that is then evaluated according to the desired cryptographic criteria.

We use the following function set: NOT, which inverts its argument, XOR, NAND, NOR, each of which takes two input arguments. Additionally, we use the function IF, which takes three arguments and returns the second one if the first one evaluates to *true*, and the third one otherwise. In the evolution process, GP uses a 3-tournament selection and mutation with a probability of 0.5. The variation operators are simple tree crossover, uniform crossover, size fair, one-point, and context preserving crossover [6] (selected at random) and subtree mutation. The initial population is created at random with ramped half-and-half initialization; every experiment is repeated 50 times and the population size equals 2000.

In the fitness function first the balancedness is verified, and if an S-box is balanced, we give it a value of zero, otherwise the value equals -1; this is denoted with the label *BAL*. If the S-box is balanced, we calculate the

Table 1: Statistical results and comparison.

| S-box size   | $T_{max}$ | GP         |        |       | Ours      |            | Related work |            |
|--------------|-----------|------------|--------|-------|-----------|------------|--------------|------------|
|              |           | Max        | Avg    | SD    | $N_F$     | $\delta_F$ | $N_F$        | $\delta_F$ |
| $4 \times 4$ | 16        | <b>16</b>  | 16     | 0     | 4         | 4          | 4            | 4          |
| $5 \times 5$ | 42        | <b>42</b>  | 41.73  | 1.01  | <b>12</b> | <b>2</b>   | 10           | 4          |
| $6 \times 6$ | 86        | <b>84</b>  | 80.47  | 4.72  | <b>24</b> | <b>4</b>   | 22           | 6          |
| $7 \times 7$ | 182       | <b>182</b> | 155.07 | 8.86  | <b>56</b> | <b>2</b>   | 48           | 6          |
| $8 \times 8$ | 364       | 318        | 281.87 | 13.86 | 82        | 20         | <b>104</b>   | <b>8</b>   |

nonlinearity and differential uniformity (which is subtracted from the value  $2^n$ , since we aim to minimize the value of that property) and **maximize** the resulting value:

$$fitness = BAL + \Delta_{BAL,0}(N_F + (2^n - \delta_F)). \quad (6)$$

$\Delta_{BAL,0}$  represents the Kronecker delta function that equals one when the function is balanced (i.e.,  $BAL = 0$ ) and zero otherwise.

In Table 1 we give the statistical results and compare the best values obtained here with the state-of-the-art results obtained with EC [5], where better values are in bold style. Statistical results are averaged over the best obtained values for each run, and column  $T_{max}$  denotes the theoretical maximal value of the fitness (or one that would correspond to the assumed theoretical maximal value). For the  $4 \times 4$  size both our technique and related work can easily reach optimal values. However, for sizes  $5 \times 5$  till  $7 \times 7$  our approach yielded significantly improved results. On the other hand, for the  $8 \times 8$  size, related work managed to obtain better nonlinearity and differential uniformity. Still, we emphasize that even those results are far from the best one where nonlinearity equals 112 and differential uniformity equals 4 (note that it is only assumed but not proven that the maximal nonlinearity equals 112).

## 4 Conclusions

Our approach shows great potential and marks the first time that heuristic techniques are able to find optimal S-boxes for sizes larger than  $4 \times 4$ . Moreover, our approach transforms a problem that has been up to now of extreme difficulty into a simpler problem for certain S-box sizes. Naturally, since not all S-boxes can be represented as a CA, our technique cannot be used to design all optimal S-boxes of the corresponding size. However, we believe that the corpus of obtainable functions is still large enough to give a sufficient diversity for future block cipher designs.

## Acknowledgments

This work has been supported in part by Croatian Science Foundation under the project IP-2014-09-4882.

## References

- [1] G. Bertoni, J. Daemen, M. Peeters, and G. V. Assche. The KECCAK reference, January 2011. <http://keccak.noekeon.org/>.
- [2] C. Carlet. Vectorial Boolean Functions for Cryptography. In Y. Crama and P. L. Hammer, editors, *Boolean Models and Methods in Mathematics, Computer Science, and Engineering*, pages 398–469. Cambridge University Press, New York, USA, 1st edition, 2010.
- [3] L. R. Knudsen and M. Robshaw. *The Block Cipher Companion*. Information Security and Cryptography. Springer, 2011.
- [4] K. Nyberg. Perfect Nonlinear S-Boxes. In *Advances in Cryptology - EUROCRYPT '91, Workshop on the Theory and Application of Cryptographic Techniques, Brighton, UK, April 8-11, 1991, Proceedings*, volume 547 of *Lecture Notes in Computer Science*, pages 378–386. Springer, 1991.
- [5] S. Picek, M. Cupic, and L. Rotim. A New Cost Function for Evolution of S-boxes. *Evolutionary Computation*, 2016.
- [6] R. Poli, W. B. Langdon, and N. F. McPhee. *A field guide to genetic programming*. Published via <http://lulu.com> and freely available at <http://www.gp-field-guide.org.uk>, 2008. (With contributions by J. R. Koza).