

## On Secure Embedded Token Design

Simon Hoerder, Kimmo Järvinen, Daniel Page

► **To cite this version:**

Simon Hoerder, Kimmo Järvinen, Daniel Page. On Secure Embedded Token Design. Lorenzo Cavallaro; Dieter Gollmann. 7th International Workshop on Information Security Theory and Practice (WISTP), May 2013, Heraklion, Greece. Springer, Lecture Notes in Computer Science, LNCS-7886, pp.112-128, 2013, Information Security Theory and Practice. Security of Mobile and Cyber-Physical Systems. <10.1007/978-3-642-38530-8\_8>. <hal-01485937>

**HAL Id: hal-01485937**

**<https://hal.inria.fr/hal-01485937>**

Submitted on 9 Mar 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# On secure embedded token design

## Quasi-looped Yao circuits and bounded leakage

S. Hoerder<sup>1</sup> and K. Järvinen<sup>2</sup> and D. Page<sup>1</sup>

<sup>1</sup> University of Bristol, Department of Computer Science,  
Merchant Venturers Building, Woodland Road, Bristol, BS8 1UB, UK.  
`{hoerder,page}@cs.bris.ac.uk`

<sup>2</sup> Aalto University, Department of Information and Computer Science,  
P.O. Box 15400, FI-00076 Aalto, Finland.  
`kimmo.jarvinen@aalto.fi`

**Abstract.** Within a broader context of mobile and embedded computing, the design of practical, secure tokens that can store and/or process security-critical information remains an ongoing challenge. One aspect of this challenge is the threat of information leakage through side-channel attacks, which is exacerbated by any resource constraints. Along these lines, this paper extends previous work on use of Yao circuits via two contributions. First, we show how careful analysis can fix the maximum number of leakage occurrences observed during a DPA attack, effectively bounding leakage from a Yao-based token. To achieve this we use modularised Yao circuits, which also support our second contribution: the first Yao-based implementation of a secure authentication payload, namely HMAC based on SHA-256.

## 1 Introduction

A vast range of challenges and opportunities has emerged as a result of the (ongoing) proliferation of mobile and embedded computing. We now routinely and fundamentally depend on a broad range of complex, highly integrated mobile devices and applications and supporting technologies and techniques need to keep pace with such developments. We consider a motivating example within this context, namely the establishment of a secure communication channel between some token and another party. Although we limit our remit to tokens that are more capable than a (basic) smart-card (e.g., a mobile telephone or USB token, with concrete exemplars including the IBM ZTIC<sup>3</sup>), the secure, efficient implementation of such a device is clearly of central importance.

Although well studied cryptographic techniques can satisfy varied requirements at a high level, real-world security guarantees are still difficult to achieve: the field of physical security in particular, including passive side-channel attacks, represents a central concern. In general, a typical side-channel adversary acquires execution profiles by observing a device and, in an ensuing analysis phase, has

---

<sup>3</sup> <http://www.zurich.ibm.com/ztic/>

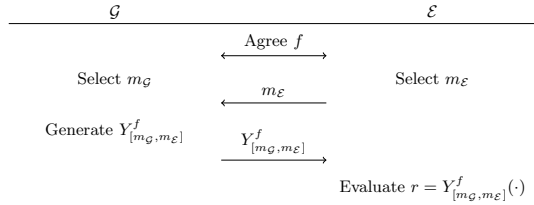
to recover a security-critical target value (e.g., key material) from the profiles. Although practical bounds on the number of profiles collected or processed may exist, attacker capability in this respect scales over time (e.g., with Moore’s law); ideally this should be catered for by any countermeasure.

Our focus is the threat of Simple (SPA) and Differential Power Analysis (DPA) [10], including variations thereof such as electro-magnetic emission analysis. We cater for timing side-channels, but explicitly deem (semi-)invasive or active attacks as outside the scope of this paper. In our scenario, profiles acquired take the form of traces that describe power consumption by the token. One broad class of countermeasures aims to produce an implementation of some primitive, and secure it by applying attack-specific countermeasures at an algorithmic, software and hardware level. Selected examples include schemes for hiding [16, Chapter 7] and masking [16, Chapter 10] input and/or intermediate values. An alternative class aims to formulate then implement a primitive which is secure-by-design. Following a philosophy that says security should be treated as a first-class goal [9, 23], this is an attractive direction in that robust guarantees can be provided. However, as the emerging field of leakage-resilient primitives (see [25] for example) illustrates, difficulties remain. Most importantly, leakage-resilient cryptography has focused on assuring security provided leakage entropy remains below a certain bound; unfortunately, no practical means of (provably) verifying such a bound for a given device is currently available.

We extend work on Yao circuits [28, 29, 15, 1, 5, 6, 14], especially their implementation on tokens [8, 7] as a means of performing leakage-resilient computation within the motivating scenario above. Our focus is practicality: we aim to keep the entire architecture as simple as possible in order to reduce manufacturing and verification cost. Careful analysis of said architecture places a bound  $\tau$  on the number of useful leakage occurrences an attacker can observe. For a given signal-to-noise ratio, this forces an attacker to develop more effective attacks rather than simply using more traces; when combined with conventional countermeasures, it can effectively prevent such attacks. One can view this process as playing a similar role to key refresh [18, 17], but without the need for synchronization. In addition, we add the first secure authentication payload, HMAC [20], to the list of applications implemented as Yao circuits. The overall result are leakage-*bounded* implementations of both AES<sub>128</sub> and HMAC.

## 2 Background

An important aspect of formalising the ability of a side-channel attack(er) is the selection of a model that describes the form of leakage from a token (and thus exploitable by the attacker). The model proposed by Standaert et al. [24] can be directly applied to our scenario with just one minor modification. In said model, adversarial success depends, among other factors, on the number of oracle queries allowed per primitive independent of updates to secret values (e.g., use of key refresh schemes). We therefore replace the number of oracle queries with the number of observable leakage occurrences per secret value.



**Fig. 1:** A high-level, generic description of Yao circuit generation and evaluation. Note that all inputs and outputs are implicitly defined wrt.  $f$ , and that depending on the scenario a) one or more of  $f$ ,  $m_G$  and  $m_E$  could alternatively be provided as input to the protocol, and b) subsequent use of  $r$  could be included.

## 2.1 Yao circuits at a high level

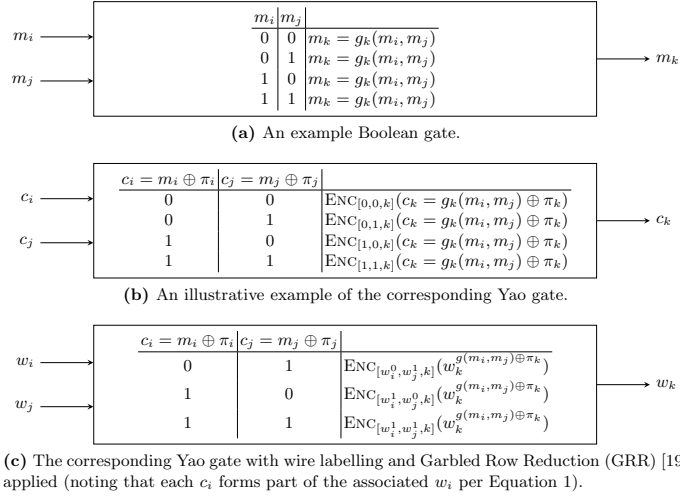
An implementation of a functionality  $f$  as a standard Boolean circuit, say  $B^f$ , can be evaluated using an input  $x$  to give an output  $r = B^f(x)$ . We use  $B_x^f$  to denote an implementation with an embedded fixed input and retain a similar form for evaluation  $r = B_x^f(\cdot)$ . At a high level, Yao circuits simply represent a non-standard implementation of  $f$ , say  $Y^f$ , which allows the associated evaluation to be secure.

Use of a given Yao circuit can be described formally as a secure two-party computation protocol. The parties involved are a circuit generator  $\mathcal{G}$  who (given  $f$  and  $x$ ) produces  $Y_x^f$ , and a circuit evaluator  $\mathcal{E}$  who (given  $Y_x^f$ ) computes  $r = Y_x^f(\cdot)$ ; both are illustrated in Figure 1. Thus the side-channel attack surface is shifted away from individual primitives  $f$  onto the task of generating a Yao circuit  $Y_{[m_G, m_E]}^f$  where we can give strong bounds on leakage. In contrast to the original usage as a two-party computation protocol, the token is trusted and we do not need oblivious transfer to communicate the circuit inputs  $m_E, m_G$ . Note that in theory  $\mathcal{G}$  could evaluate  $Y_{[m_G, m_E]}^f$  as well as generating it, but we deem it more economic in most cases to let  $\mathcal{E}$  do the evaluation.

Imagine  $g_k$  refers to some  $k$ -th truth table wlog. describing a 2-input, 1-output Boolean function or gate instance within  $B^f$ . Both the inputs to and outputs from said gate are provided on wires indexed using a unique wire identifier (or wire ID): we write  $m_i$  for the value carried by the wire with wire ID  $i$ . Note that the output wire ID can act to identify a given gate instance (i.e., acts as a gate ID). Figure 2a is a trivial starting point outlining how such a gate can be evaluated.

Yao circuits are constructed by forming a “garbled” Yao gate instance in  $Y^f$  for each Boolean gate instance in  $B^f$ . Both inputs to and outputs from the Yao gate are altered to mask their underlying value: this means each  $m_i$  is replaced by  $c_i$ , an encryption of the former. Given  $\text{ENC}_x(y)$  denotes the encryption of  $y$  under the key  $x$  using some symmetric cipher (with a  $\kappa$ -bit key and  $\beta$ -bit block size), Figure 2b illustrates a Yao gate corresponding to the above. Note that

- the public  $c_i$  and  $c_j$  inputs (whose secret underlying values are  $m_i$  and  $m_j$ ) are provided on wires with public indices  $i$  and  $j$ ,
- the public  $c_k$  output (whose secret underlying value is  $m_k$ ) is provided on a wire with public index  $k$ ,



**Fig. 2:** A step-by-step comparison of a 2-input, 1-output Boolean gate (whose function is described by  $g_k$ ), and the associated Yao gate construction before and after optimisation.

- the standard gate functionality is  $g_k$ , and
  - $\pi_i$ ,  $\pi_j$  and  $\pi_k$  act as secret permutation bits on the rows of the truth table.
- During evaluation of the gate,  $\mathcal{E}$  gets  $c_i$  and  $c_j$  meaning it cannot recover the underlying values of  $m_i$  and  $m_j$  since it does not know  $\pi_i$  and  $\pi_j$ . However,  $c_i$  and  $c_j$  index *one* entry of the truth table and allow *only* this entry to be decrypted (since they determine the associated cipher key) and yield  $c_k$ . The central idea is that a Yao gate reveals nothing about a) the gate functionality nor b) underlying values, iff. it is evaluated at most once.

## 2.2 Abstract realisation of Yao circuits

**Wire labels** The illustrative example above has a major shortcoming: the effective cipher key size is just two bits (since all wire IDs are public), meaning the key is inherently susceptible to exhaustive search. To combat this, given a security parameter  $\lambda$  one replaces the Boolean value communicated on each wire with a randomised  $\lambda$ -bit wire label. Let

$$w_i^{c_i} = \rho_i \parallel (m_i \oplus \pi_i) = \rho_i \parallel c_i \quad (1)$$

denote the  $i$ -th such wire label in the general case where  $m_i \oplus \pi_i = c_i$  for  $c_i \in \{0, 1\}$ ,  $\rho_i \xleftarrow{\$} \mathbb{Z}_2^{\lambda-1}$  and  $\pi_i \xleftarrow{\$} \mathbb{Z}_2$ . Note that the combination of  $\rho_i$  and  $\pi_i$  could be viewed as a  $\lambda$ -bit ephemeral key, implying  $\kappa = 2\lambda + \epsilon$  where  $\epsilon$  is the number of bits used to encode wire IDs, and that Equation 1 caters for two optimisations outlined below.

*The Garbled Row Reduction (GRR) optimisation* was introduced by Naor et al. [19]. In short, by carefully selecting

$$w_k^{c_k} = w_k^{g(m_i, m_j) \oplus \pi_k} = \text{ENC}_{[w_i^0, w_j^0, k]}(\gamma) \quad (2)$$

for  $c_i = c_j = 0$  and a suitable public constant  $\gamma$ , the first truth table entry can be eliminated (as illustrated by Figure 2c). This is attractive since it permits up to a 25% reduction in communication of the Yao circuit from  $\mathcal{G}$  to  $\mathcal{E}$ , plus reduces the amount of storage required.

The “free XOR” optimisation introduced by Kolesnikov and Schneider [11] realises XOR gates (almost) for free. Given a global (i.e., one for each instance of  $Y^f$ ), secret constant  $\Delta \in \mathbb{Z}_2^{\lambda-1}$  they select  $w_i^1$  as in Equation 1, then define

$$w_i^0 = w_i^1 \oplus (\Delta \parallel 1) = (\rho_i \oplus \Delta) \parallel (\pi_i \oplus (m_i \oplus 1)) \quad (3)$$

to allow computation of XOR gates via the relationships

$$\begin{aligned} w_i^0 \oplus w_i^1 &= w_i^1 \oplus (\Delta \parallel 1) \oplus w_j^1 &= w_i^1 \oplus w_j^0 &= w_k^1 \\ w_i^0 \oplus w_i^0 &= w_i^1 \oplus (\Delta \parallel 1) \oplus w_j^1 \oplus (\Delta \parallel 1) &= w_i^1 \oplus w_j^1 &= w_k^0 = w_k^1 \oplus (\Delta \parallel 1) \end{aligned}$$

### 2.3 Previous Token implementations of Yao circuits

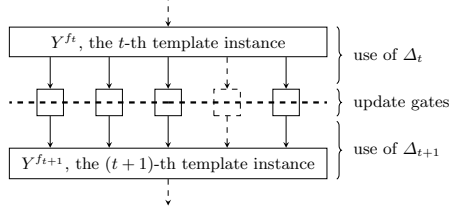
As far as cryptographic primitives are concerned, previous work (with the exception of [12]) focus on use of  $\text{AES}_{128}$  as the functionality  $f$ . Other functionalities considered relate instead to higher level applications, e.g., database search [5, 14] or bioinformatics [6]. Pinkas et al. [21] provide the first feasibility results (using software) of Yao-based constructions; since they relate more directly to the chosen scenario, we detail work by Järvinen et al. [8, 7] below which both implement Yao circuits on tokens but do not use modularisation.

*Secure computation via One-Time Programs (OTPs):* In [8], Järvinen et al. consider a scenario wherein  $\mathcal{E}$  is a hardware token, and  $\mathcal{G}$  is a trusted party during a setup stage. The idea is that  $\mathcal{G}$  as token issuer stores One-Time Memories (OTMs) for a fixed number  $x$  of OTPs represented by  $Y^f$  on the token. At run-time, the token owner uses one set of OTMs to form the input labels corresponding to his data, and the token evaluates the Yao circuit before finally releasing the result.

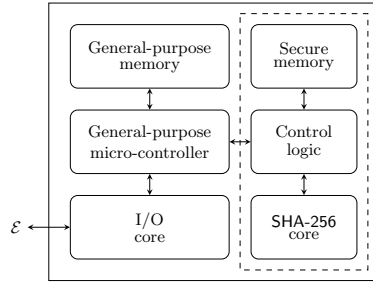
The advantage of this scenario is that very little protection against side-channel attacks is required. However, the major disadvantage is the limited number of Yao circuits: in most real-life scenarios this is unacceptable. In addition, a generic framework without this disadvantage is given at the cost of losing the leakage-resilient circuit generation. Our work can be seen as a leakage-resilient, more flexible version of the framework.

*Secure computation via out-sourcing:* In [7], Järvinen et al. consider a scenario wherein  $\mathcal{E}$  is a cloud computing provider; the role of  $\mathcal{G}$  is split between a secure hardware token  $\mathcal{G}_S$  and some other party  $\mathcal{G}_U$  (e.g., a desktop workstation). The idea is for  $\mathcal{G}_U$  to generate  $B^f$ , which is passed to  $\mathcal{G}_S$  and translated (securely) into  $Y^f$ . The  $Y^f$  can then be evaluated by  $\mathcal{E}$ , with the overall effect of securely out-sourcing computation from  $\mathcal{G}_U$  to  $\mathcal{E}$  (and the token  $\mathcal{G}_S$ ).

This scenario is advantageous in the sense it allows a flexible choice of  $f$  (wrt. the token) and is very speed- and memory-efficient. However, it has a relatively high hardware requirement: in addition to the  $\text{SHA-256}$  core it requires at least one  $\text{AES}_{128}$  core, [7] uses two, which all have to be free from leakage.



**Fig. 3:** An example of updating  $\Delta$  between two template instances, with the heavy dashed line denoting the boundary between use of  $\Delta_t$  and  $\Delta_{t+1}$ .



**Fig. 4:** A block diagram of our proposed token architecture. Only the Yao core (within the dashed box) has requirements for conventional side channel countermeasures.

### 3 Supporting alterations to traditional Yao circuits

To support our design in Section 4, we first outline two supporting concepts: detail relating to their realisation and utility is deferred until later. While neither represents a significant change to underlying theory, we posit that both significantly ease the practical task of constructing and using Yao circuits.

#### 3.1 Circuit modularisation

**Existing limitations** Consider a typical iterative block cipher design, with  $s$  rounds in total (e.g.  $\text{AES}_{128}$  with  $s = 10$ ). The functionality for round  $i$  is described by  $f_i$ , which implemented as a Boolean circuit is  $B^{f_i}$ . The  $s$  different round functionalities can be the same or different as required, with the overall cipher thus described by the functionality

$$f = f_{s-1} \circ \dots \circ f_1 \circ f_0.$$

One can implement this either combinatorially, whereby instances of each  $B^{f_i}$  are “unrolled” to form the resulting monolithic circuit, or iteratively, whereby instances of each unique  $B^{f_i}$  are “looped over” with a need for only one circuit instance and some control logic. Traditional Yao circuits *must* adopt the former approach: no sequential elements (e.g., latches, clock signals) are allowed because each gate can be evaluated at most once (to ensure security). One cannot reuse the resulting Yao circuit unless rerandomisable constructions [4] are considered; typically these incur a prohibitive overhead.

One impact of this restriction is that previous work almost exclusively focuses on  $\text{AES}_{128}$  (as far as cryptographic primitives are concerned) which a) has a fixed  $s$  and can hence be unrolled, and b) has a fairly compact hardware implementation. We know of only one implementation of another cryptographic primitive [12], wherein an (insecure) implementation of 256-bit RSA is described. Even with such small operands, the resulting Yao circuit is  $\approx 8500$  times larger than their  $\text{AES}_{128}$  circuit, in part as a result of the requirement to unroll the loop representing a binary, modular exponentiation.

**Modularised Yao circuits** To address this issue, we make use of modular Yao circuits. Similar concepts have been used recently<sup>4</sup> in [6, 14] to achieve efficiency gains for large circuits, but only [14] mentions the possibility of using run-time parameters to control circuit assembly.

Traditionally (right) each monolithic Boolean circuit  $B^f = B^{f_{s-1} \circ \dots \circ f_1 \circ f_0}$  is stored and must be translated into the corresponding Yao circuit by  $\mathcal{G}$  each time the latter needs to be evaluated. In our alternative (left),  $\mathcal{G}$  holds only the static description of each template  $B^{f_i}$  (and associated meta-data), instantiating them to form  $Y^f$  without holding the entirety of the (potentially large)  $B^f$ . This technique permits a quasi-loop: given  $s$ , e.g. the number of blocks to be hashed, at run-time (rather than being fixed), the token unrolls one circuit template  $B^{f_i}$ , to form the resulting Yao circuit. The resulting reduction in resource requirement and increased flexibility allows us to implement HMAC. Generation of the corresponding  $Y^f$  can be streamed in that  $\mathcal{G}$  communicates one part at a time to  $\mathcal{E}$ .

### 3.2 Updating $\Delta$ between template instances

In previous works [21, 8, 7] using the free XOR trick [11], the authors argue that for correctness  $\Delta$  must remain constant within a given Yao circuit. Indeed, if an XOR gate is evaluated using constants  $\Delta_1 \neq \Delta_2$ , the result is incorrect as

$$\begin{aligned} w_i^0 \oplus w_j^1 &= (w_i^1 \oplus (\Delta_1 \parallel 1)) \oplus w_j^1 \neq w_i^1 \oplus (w_j^1 \oplus (\Delta_2 \parallel 1)) = w_i^1 \oplus w_j^0 \\ w_i^0 \oplus w_j^0 &= (w_i^1 \oplus (\Delta_1 \parallel 1)) \oplus (w_j^1 \oplus (\Delta_2 \parallel 1)) \neq w_i^1 \oplus w_j^1 \end{aligned}$$

shows. Despite this, a crucial observation is that  $\Delta$  *can* be changed *if* said change is applied consistently. Although doing so has no functional benefit, we use it to directly formulate a leakage bound within Section 5.

In some modularised Yao circuit, imagine the  $t$ -th template instance uses  $\Delta_t$ . The next,  $(t+1)$ -th instance can then use  $\Delta_{t+1}$  (as illustrated by Figure 3) iff each connecting wire is updated appropriately. The simplest approach is to consider a dedicated 1-input, 1-output update gate with the identity functionality, i.e.,  $g_k(x) = x$ . Given

$$w_k^{c_k} = w_k^{m_i \oplus \pi_k} = \text{ENC}_{[w_i^0, w_i^0, k]}(\gamma)$$

by definition, and that  $m_i = m_k$  since the gate updates  $\Delta$  rather than change the underlying input value, Equation 2 means the associated wire labels are

$$\begin{aligned} w_i^0 &= w_i^1 \oplus (\Delta_t \parallel 1) \\ w_k^0 &= w_k^1 \oplus (\Delta_{t+1} \parallel 1) \end{aligned}$$

However, this suggests that *any* non-XOR gate can be used to perform the update without cost: the existing GRR-optimised Yao gate truth table only needs to have  $\Delta_{t+1}$  folded into the label  $w_k^{g_k(m_i, m_j) \oplus \pi_k}$  instead of  $\Delta_t$  where appropriate.

<sup>4</sup> A year earlier, [5] proposed a different kind of modularity, namely mixing Yao circuits and homomorphic operations. The cost associated with additional hardware required to support homomorphic operations means we do not adopt this approach.



## 4 Token design

We consider a scenario wherein  $\mathcal{G}$  is a hardware token that needs to be secured against SPA and DPA attacks, and  $\mathcal{E}$  is some other party to which computation is outsourced. The idea is to put a bound on the number of times secret values are used for any computation and leakage can be observed before the value – akin to key refresh – is updated. The known, residual leakage can then be accommodated by careful use of conventional countermeasures.

### 4.1 Cipher construction

To instantiate the symmetric cipher required to encrypt wire labels, we follow existing work [21, 7, 8] by using a one-time pad like construction

$$\text{ENC}_{[w_i^{c_i}, w_j^{c_j}, k]}(w_k^{c_k}) = \text{SHA-256}(w_i^{c_i} \parallel w_j^{c_j} \parallel k) \oplus w_k^{c_k}$$

where the SHA-256 output is implicitly truncated to  $\lambda$  bits to match the wire label size. We reuse SHA-256 as a secure Pseudo-Random Number Generator (PRNG), splitting the SHA-256 into two 128-bit values

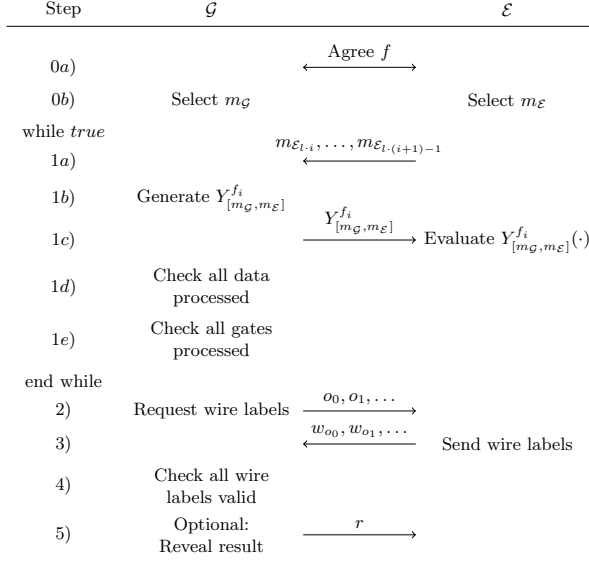
$$\begin{aligned} [x, y] &\leftarrow \text{SHA-256}(\text{seed}_{t-1}) \\ \text{seed}_t &\leftarrow \text{seed}_{t-1} \oplus y \\ \text{rand} &\leftarrow x \end{aligned}$$

so  $x = \text{rand}$  is used as a wire label for example, while  $y$  is used exclusively to update the seed. Note SHA-256 is therefore the *only* significant cryptographic core required by the token. This construction is certainly secure if the PRNG is modelled as a random oracle which is a weaker model than the one we have for our Yao circuits. Intuitively however, some form of correlation robustness should be sufficient. Research on the correlation robustness of hash functions is still developing, see [11, 3] for example, and therefore we defer the exact security requirements to future work.

Previous work has used domain-specific languages (e.g. SFDL [15, 1]) to implement the payloads. While this *may* ease implementation, it reduces control over the actual circuit layout. Standard Hardware Description Languages (HDLs) are, in fact, well suited to payload description: they simply lack the protocol aspect of Yao circuits. Therefore we aim to reap benefits of familiarity, design and code portability by using a VHDL dialect and associated compiler, both of our own design, to describe structural and behavioural circuit templates.

### 4.2 Operational protocol and token architecture

The communication between the token  $\mathcal{G}$  and the evaluator  $\mathcal{E}$  is shown in Figure 5.  $\mathcal{E}$  can, for example, be a local untrusted work station or an untrusted but more powerful chip within a mobile phone; in most cases it will not be



**Fig. 5:** The two-party computation protocol reflecting circuit modularisation. Note that  $\mathcal{E}$  does not communicate  $m_{\mathcal{E}}$  in one block, but rather in multiple  $l$ -bit blocks. Output wires of the Yao circuit (i.e. the wires carrying results from the functionality) have wire IDs  $o_0, o_1, \dots$

the authentication partner. We assume a physical connection between the parties, and hence focus on optimising their workload rather than the number of communication rounds.

Initially, in step 0,  $\mathcal{E}$  requests a functionality  $f$  (e.g., HMAC or AES<sub>128</sub>) and both parties need to have (or generate) corresponding inputs  $m_{\mathcal{G}}$  and  $m_{\mathcal{E}}$ . Step 1, from a theoretic perspective, is the same as the monolithic communication in Figure 1 despite now supporting the modularised approach. Specifically,  $\mathcal{G}$  generates the Yao circuit  $Y^f$  based on the circuit templates  $B^{f_i}$  and sends it step-by-step to be evaluated. Note that modularisation forces three important checks:

1. step 1d checks if all input values (from both  $m_{\mathcal{G}}$  and  $m_{\mathcal{E}}$ ) have been used,
2. step 1e checks if all gates have evaluated, and
3. step 4 checks if all output values are valid, i.e., if  $w_{o_j} \in \{w_{o_j}^0, w_{o_j}^1\} \forall o_j$ .

When a check condition can not be satisfied the token aborts immediately, meaning in particular that it does not reach step 5 where the result  $r = f(m_{\mathcal{G}}, m_{\mathcal{E}}) = Y_{[m_{\mathcal{G}}, m_{\mathcal{E}}]}^f$  is revealed, and that *seed* values of the PRNG are not accidentally reused.

To support the protocol outlined above, Figure 4 outlines a proposed token architecture. The main components and their roles are as follows:

- A general-purpose micro-controller manages the communication protocol in Figure 5, controlling assembly of  $Y^f$  from the  $B^{f_i}$  circuit templates and performing other tasks (such as message padding for HMAC).
- A general-purpose memory holds the  $B^{f_i}$ , micro-controller program and other public run-time variables values which only require correctness.

- Storage of and computation using secret values is limited to the Yao core, which consists of:
  - A SHA-256 core, used to encrypt wire labels and also as a PRNG.
  - Control logic used for auxiliary operations such as wire label generation.
  - A secure memory, split into two parts: a non-volatile part holds  $m_G$  and *seed*, while a volatile part holds  $\Delta_t, \Delta_{t+1}$  and, for each wire  $i$ , the tuple  $\{w_i^0, w_i^1\}$ . Note that  $\Delta_{t+1} = 0$  unless the token is processing update gates, and that a crucial role of the secure memory as a whole is to prevent read-out or other leaks of values such as  $m_G$  and *seed*.

## 5 Analysis and results

### 5.1 Security analysis

This section attempts to explore security aspects of the proposal outlined in Section 4. After a statement of general assumptions, we deal specifically with potential attack vectors exploited during an SPA or DPA attack, or by a malicious adversary within the operational protocol.

**Security assumptions** To be successful, the adversary has to recover the input  $x_G$  held by the token for which he can either attempt to recover  $x_G$  directly (e.g., via a DPA attack), or try to “ungarble” the Yao circuit  $Y_{x_G, x_E^f}$  (or part thereof). In showing neither strategy is viable, we make some important assumptions:

- A-1** An authentication protocol that prevents man-in-the-middle attacks against  $m_E$  in step 1d of the Yao protocol must be selected (if this threat is relevant).
- A-2** The control-flow of the token, managed by the micro-controller, is tamper-proof which implies integrity of the general-purpose memory. Although this is a strong requirement, it is common for embedded systems.
- A-3** The hash function used for encryption of wire labels (in our case SHA-256) must be circular-2-correlation robust (see Choi et al. [3]).
- A-4** The token cannot be reset, and no randomness reused: this prevents an adversary forcing the token to regenerate the same Yao circuit with the same randomness, then reevaluating it with different inputs.

**Power analysis adversaries** SPA attacks attempt to recover the target value using one (or at least very few) traces and attack capabilities do not scale over time; examination of data-dependent control-flow is one example. There are two possible attack vectors:

- SPA-1** For each  $i$ -th input wire, the token must send either  $w_i^0$  or  $w_i^1$  to the evaluator depending on the underlying value of  $m_i$ . To succeed, the adversary must be able to determine whether  $m_i$  is 0 or 1 (for all  $m_i \in x_G$ ).
- SPA-2** Gates such as AND and OR are biased towards 0 or 1 in their output: if the adversary determines during computation of

$$\text{ENC}_{[w_i^{c_i}, w_j^{c_j}, k]}(w_k^{g(m_i, m_j) \oplus \pi_k})$$

which truth table row contains the minority output, they can reverse the

permutations (i.e.,  $\pi_i$  and  $\pi_j$ ) and recover the underlying values of almost all output wire labels from non-XOR gates.

For hardware implementations, both SPA attack vectors will be implemented using multiplexers whose data dependency of the power consumption is usually already hidden well enough without countermeasures (e.g., [16, Appendix A.3]). Even if this is not the case, traditional countermeasures (e.g., random masking of the select signal with corresponding permutation of the inputs) are efficient.

In contrast, DPA attacks attempt to recover the security-critical target value by applying statistical distinguishers to a large set of traces; issues of signal-to-noise ratio, as well as explicit countermeasures, determine the exact number. More formally, let  $k$  be the target value,  $v$  be a variable value and  $r$  the result of some generic operation  $\odot$ . A DPA adversary collects traces relating to execution of  $r_i \leftarrow k \odot v_i$  for many different  $v_1, v_2, \dots, v_\sigma$ . The potency of a DPA attack is then judged by  $\sigma$ , the number of traces required to be reliably recover  $k$ . Our approach is to have a design-time constant bound  $\tau \ll \sigma$  instead of allowing the adversary to control it. Put another way, we bound the leakage such an adversary can utilise in a DPA attack: if the application of conventional countermeasures can prevent attacks with said leakage level, the token is secure.

**DPA-1** The token must compute  $w_i^0 \leftarrow w_i^1 \oplus (\Delta_t \parallel 1)$  for every wire. As such we have

$$\tau_{\text{DPA-1}} = \max(\delta_1, \delta_2, \dots)$$

where  $\delta_t$  denotes the number of wires using  $\Delta_t$ . If the technique in Section 3.2 is used correctly,  $\tau_{\text{DPA-1}}$  is a constant determined by the token designer.

**DPA-2** For each gate, the four values of  $w_{\{i,j\}}^{\{0,1\}}$  are each used twice as input to SHA-256( $w_i^{\{0,1\}} \parallel w_j^{\{0,1\}} \parallel k$ ). Focusing on one label, wlog.  $w_i$  say, and one external value, wlog. 0 say, the attacker gets two traces for each gate where  $w_i^0$  is used as input. Therefore, we have

$$\tau_{\text{DPA-2}} = 2 \cdot \max_{\forall k} (G_k)$$

where  $G_k$  represents the fan-out of the  $k$ -th gate (and input wires are also considered as being driven by imaginary gates). Note that a similar attack vector exists when the token processes an XOR gate. Such a gate must compute  $w_k^0$  and  $w_k^1$ , and one possible approach is to compute

$$\begin{aligned} w_k^0 &\leftarrow w_i^0 \oplus w_j^0 \\ w_k^1 &\leftarrow w_i^0 \oplus w_j^1 \end{aligned} \tag{4}$$

in which case  $\tau_{\text{DPA-2}}$  conveniently covers this attack vector as well.

Concrete, non-optimised examples for these bounds are given in Section 5.2. If our design is used to protect against DPA attacks, functionalities that were inherently secure against SPA clearly inherit any SPA vulnerabilities of the underlying Yao circuit approach. We suggest that preventing SPA attacks on our design using conventional countermeasures is, broadly speaking, easier than preventing DPA attacks on the functionality in question: the cost of preventing the former is easily justified by the improvement offered wrt. the latter.

**Catering for timing analysis adversaries** The execution time associated with generating of a Yao circuit is inherently independent of the inputs to that Yao circuit: it depends *only* on the circuit size. As far as the architecture is concerned, we do not use a cache for the micro-processor in order to avoid cache-based timing attacks. Working without a cache is a common decision for cryptographic tokens and therefore not an exceptional burden of our design.

**Catering for malicious adversaries** One advantage of Yao circuits is the availability of related security proofs. For semi-honest adversaries, Figure 5 preserves proofs already given by [11, 21, 7, 8, 3]. This is a direct result of steps 1a to 1e being equivalent to the single generate-evaluate step from previous protocols.

However, we also need to consider malicious adversaries. Lindell et al. [13] show how a two-party computation protocol using Yao circuits can cover the case of malicious  $\mathcal{G}$  using a cut-and-choose approach. Our scenario is far less complex, since  $Y^f$  is generated by the token which is implicitly trusted: we disallow a malicious  $\mathcal{G}$ . Therefore we only have to consider a malicious  $\mathcal{E}$ , and show it cannot learn anything about  $x_{\mathcal{G}}$  not implied by the result  $r = f(x_{\mathcal{G}}, x_{\mathcal{E}})$  by deviating from the protocol.  $\mathcal{E}$  has two options (which we expand on below): it can either a) provide faulty input or b) perform variants on early termination.

*Faulty data.* The only steps where  $\mathcal{E}$  can provide faulty data are 1a and 3. As  $m_{\mathcal{E}}$  in step 1a is the input of the functionality, sending a faulty  $m_{\mathcal{E}}$  has no impact on the security of the Yao protocol: it can only influence the output of the functionality  $f$ . In contrast, sending faulty  $w_{o_j}$  in step 3 is potentially a problem if the adversary sends labels from intermediate wires instead of the output labels. However, this is prevented by the check in step 4 which ensures that for each output wire  $o_j$ , exactly one wire label  $w_{o_j}^{c_{o_j}} \in \{w_{o_j}^0, w_{o_j}^1\}$  has been sent before the result is revealed.

*Early termination.* Since  $\mathcal{E}$  can not learn anything from one of the partial circuits (i.e., a given  $Y_{[m_{\mathcal{G}}, m_{\mathcal{E}}]}^{f_i}$ ) until the protocol is finished (i.e., until  $\mathcal{G}$  reveals the result),  $\mathcal{E}$  cannot profit from straight early termination. However, if the functionality  $f$  requires  $s$  iterations of a loop to form  $f = f_{s-1} \circ \dots \circ f_1 \circ f_0$ , per the description of  $\text{AES}_{128}$  in Section 3.1 for instance, the adversary could potentially gain information from terminating the loop early, i.e., to get  $f' = f_{s'-1} \circ \dots \circ f_1 \circ f_0$  for  $0 < s' < s$ : this would be analogous to a reduced-round attack. To prevent this, we require the token to check (in steps 1d and 1e) whether the Yao circuit for  $f$  has been completely generated or whether some  $Y_{[m_{\mathcal{G}}, m_{\mathcal{E}}]}^{f_i}$  is missing.

## 5.2 Experimental results and analysis

Our goal is to study gross, indicative metrics and trade-offs rather than focus on absolute figures that could be improved via incremental optimisation. For the evaluation of our proposed design, we implemented a VHDL compiler (per Section 4.1), a token simulator  $\mathcal{G}$ , an evaluator  $\mathcal{E}$  as well as two payload functionalities, namely  $\text{AES}_{128}$  and  $\text{HMAC}_{\text{SHA-256}}$ . For each payload, we considered variants that differ in their frequency of  $\Delta$  update: for  $\text{AES}_{128}$  three variants are used, for  $\text{HMAC}_{\text{SHA-256}}$  two variants. The variants are as follows:

	#blocks	# $\Delta$	#XOR	#non-XOR	#SHA-256	RAM	$ B^i $	$ m_{G,sec} $	$ m_{G,pub} $	$\tau_{\text{DPA-1}}$	$\tau_{\text{DPA-2}}$
AES	1	1	19088	5760	24578	245.9kB	12318B	176B	-	7296	11
AES_U1	1	2	19088	5888	25091	263.4kB	13628B	176B	-	3776	11
AES_U9	1	10	19088	6912	29195	262.3kB	12845B	176B	-	960	11
HMAC	1	1	148080	129680	556866	1883.6kB	121942B	64B	32B	167824	19
	2	1	222120	194520	835170	2489.8kB	121942B	64B	32B	251608	19
	3	1	296160	260384	1113474	3069.0kB	121942B	64B	32B	335392	19
	4	1	370200	324200	1391778	3671.4kB	121942B	64B	32B	419176	19
HMAC_U	1	3	148080	130192	558916	1911.0kB	122981B	64B	32B	84040	19
	2	4	222120	195288	835170	2500.8kB	122981B	64B	32B	84040	19
	3	5	296160	260384	1117574	3113.4kB	122981B	64B	32B	84040	19
	4	6	370200	325480	1396903	3724.6kB	122981B	64B	32B	84040	19

**Table 1:** Efficiency metrics and leakage bounds for our token design and a range of payload implementations. The block size for AES<sub>128</sub> is 128 bits, for HMAC it is 512 bits (including the padding in the last block).

**AES** Baseline AES<sub>128</sub> implementation without updating of  $\Delta$ .

**AES\_U1** AES<sub>128</sub> with a  $\Delta$  update after the fifth iteration of the round function.

**AES\_U9** AES<sub>128</sub> with a  $\Delta$  update after every iteration of the round function.

**HMAC** Baseline HMAC implementation without updating of  $\Delta$ .

**HMAC\_U** HMAC with  $\Delta$  updates after every iteration of the compression function.

Table 1 details efficiency metrics for implementation of these variants on the platform described and shows the two associated bounds  $\tau_{\text{DPA-1}}$  and  $\tau_{\text{DPA-2}}$ . The first three columns specify the payload, the number of input blocks from  $\mathcal{E}$  and the number of  $\Delta$  values being used at run-time.

**Efficiency** The columns #XOR and #non-XOR in Table 1 give the number of gates in the resultant Yao circuit. Compared to [21, 8] we have considerably smaller AES<sub>128</sub> circuits, which is mainly due to omission of key scheduling and, to a less extent, use of more optimised S-box formulas of Boyar et al. [2]. The omission of key scheduling implies a small penalty of having to store all round keys  $m_{G,sec}$  in secure ROM.

The column #SHA-256 shows the number of distinct uses of the SHA-256 core, each a one-block hash. Ignoring the absolute simulation time, we feel this metric best represents the execution time of a concrete token since the SHA-256 core will most likely be the throughput bottleneck. [7] use a SHA-256 core which requires 67 cycles per 512 bit block at 66 MHz. Based on these numbers, a crude time estimation (based only on calls to the SHA-256 core) is 24ms for AES and 1418ms for HMAC\_U with 4 message blocks.

A significant issue is the amount of RAM required at run-time. To assess this, we measured the simulator heap and stack usage using the Valgrind `massif` tool [27]. We note that the tool itself is not perfect, and that the result includes overhead of up to 20% relating to performance and security counters. Even so, the indicative RAM requirement is large: it remains within the capability of devices in our remit, but clearly beyond smart-cards or RFID tokens for example. The requirement stems in the most part from storing *all* wire labels  $\{w_i^0, w_i^1\}$  in RAM. One possible trade-off would be to store only  $w_i^0$  and recompute  $w_i^1$  when needed, reducing the RAM usage by a factor of two, but increasing the number of traces available by a factor similar to the maximum fan-out. [7] chose a keyed PRNG which allows recomputation of  $w_i^0$  when needed, thus reducing the RAM requirements drastically. However, any keyed PRNG is vulnerable to DPA attacks with unlimited  $\tau$  which negates our aim of bounding the leakage.

An interesting observation can be made about the RAM usage of AES\_U1 and AES\_U9. Intuitively, one would expect the RAM usage to always grow in line with the number of  $\Delta_i$  used. In this case, the opposite happens because AES\_U1 applies the  $\Delta$  updating within the top-level entity (which also accounts for the larger  $|B^{f_i}|$ ), requiring more wires for which RAM is allocated during the entire run-time. AES\_U9 performs the updating at the end of the round function entity instead, and the RAM for additional wires can be deallocated as soon as each round function instance of has been completed.

The size of the templates,  $|B^{f_i}|$  (stored in unsecured ROM), profits directly from modularisation. As predicted, the size of  $|B^{f_i}|$  for HMAC does not depend on the message size as it would have for the traditional approach.

**Security** Having explained  $\tau_{\text{DPA-1}}$  and  $\tau_{\text{DPA-2}}$  in Section 5.1, we note that our  $\Delta$  updating technique limits  $\tau_{\text{DPA-1}}$  as predicted; note esp. the HMAC\_U payload, where updating  $\Delta$  fixes previously unlimited leakage to a constant chosen by the token designer.

The result for  $\tau_{\text{DPA-2}}$  is an absolute upper bound, i.e., for all output wires we counted how often it gets used while processing the follow-up gates. As explained in Section 5.1, if a wire is used as input to a non-XOR gate each label gets used twice; for XOR gates Equation 4 gives the numbers relevant to our implementation. For an attacker it will be very difficult to combine traces from two different operations like this but we prefer to err on the side of security by overestimating the attacker. With numbers this low, **DPA-2** is almost irrelevant as an attack vector. But having a low  $\tau_{\text{DPA-2}}$  was an explicit aim of our work:  $\tau_{\text{DPA-2}}$  is the only possible attack vector on the SHA-256 core, and therefore  $\tau_{\text{DPA-2}}$  is crucial to determine the level of conventional countermeasures needed to protect the SHA-256 core. Compared to the SHA-256 core, protecting the XOR from **DPA-1** to match a much higher  $\tau_{\text{DPA-1}}$  is inexpensive.

As a reference one may look at the Power-Trust micro-processor of Tillich et al. [26], which has parts of the ALU implemented within a secure zone. For evaluation purposes they implemented the secure zone in three different logic styles (namely CMOS, iMDPL [22] and DWDDL [30]) and performed a DPA attack against an AES<sub>128</sub> software implementation using the secure zone. While it is difficult to directly extrapolate from a design as different from ours, this at least gives an estimate: there is no reason why secure logic styles such as iMDPL and DWDDL should fare worse for our token. For the DPA attack on the secure zone to be successful, Tillich et al. required 130,000 traces against the (unprotected) CMOS implementation, 260,000 traces against the iMDPL implementation and 675,000 traces against the DWDDL implementation. With  $\tau_{\text{DPA-1}} = 7296$  in the worst case for AES<sub>128</sub> and  $\tau_{\text{DPA-1}} = 84040$  for HMAC\_U we surmise that both iMDPL and DWDDL would have successfully thwarted the DPA attack from Tillich et al. against an implementation of our token. It is important to note, that for both bounds we did not yet try to find the absolute minimum. For example it is possible to add additional gates to achieve fan-out = 2 and thus  $\tau_{\text{DPA-2}} \leq 4$  while  $\tau_{\text{DPA-1}}$  can be easily reduced by updating  $\Delta$  more often within the round resp. compression functions, not just at their end.

## 6 Conclusions

In essence, this paper has demonstrated that an embedded token can be designed which gives strong bounds on the number of useful traces a power analysis adversary can collect. Our design methodology

1. is generic in that it works for all payloads and use-cases (cf. PIN block),
2. does not impose limits on the token lifetime,
3. does not require synchronization (cf. key update schemes),
4. is easily verifiable, and
5. successfully limits attack capabilities for side-channel adversaries.

In relation to the former point, we have already extended previous work through support for a Yao circuit for HMAC. Exploration of further primitives based on modularisation (including methods and trade-offs to further reduce the leakage bound), plus incremental optimisation of both the token design and operational protocol (especially the RAM requirement) are interesting avenues for further work. In relation to the latter point, the clear next step is to produce experimental results from a concrete implementation of the token. This would, for example, allow investigation of the concrete leakage and whether implementation specific high order moments occur which increase  $\tau$  for higher order attacks.

*Acknowledgements* The work described in this paper has been supported in part by EPSRC grant EP/H001689/1 and by Academy of Finland, project #138358. We would like to thank Elisabeth Oswald for her valuable comments.

## References

1. A. Ben-David, N. Nisan, and B. Pinkas. FairplayMP - A Secure Multi-Party Computation System. In *CCS*, pages 257–266, 2008.
2. J. Boyar and R. Peralta. A new combinational logic minimization technique with applications to cryptology. In *SEA*, volume LNCS 6049, pages 178–189, 2010.
3. S.G. Choi, J. Katz, and R. Kumaresan. On the Security of the “Free-XOR” Technique. In *TCC*, volume LNCS 7194, pages 39–53, 2012.
4. C. Gentry, S. Halevi, and V. Vaikuntanathan. i-Hop Homomorphic Encryption and Rerandomizable Yao Circuits. In *CRYPTO*, volume LNCS 6223, pages 155–172, 2010.
5. W. Henecka, S. Kögl, A.-R. Sadeghi, T. Schneider, and I. Wehrenberg. TASTY: Tool for Automating Secure Two-party computations. In *CCS*, pages 451–462, 2010.
6. Y. Huang, D. Evans, J. Katz, and L. Malka. Faster Secure Two-Party Computation Using Garbled Circuits. In *USENIX Security Symposium*, 2011.
7. K. Järvinen, V. Kolesnikov, A.-R. Sadeghi, and T. Schneider. Embedded SFE: Offloading Server and Network using Hardware Tokens (short version). In *Financial Cryptography*, volume LNCS 6052, pages 207–221, 2010.
8. K. Järvinen, V. Kolesnikov, A.-R. Sadeghi, and T. Schneider. Garbled Circuits for Leakage-Resilience: Hardware Implementation and Evaluation of One-Time Programs. In *CHES*, volume LNCS 6225, pages 383–397, 2010.
9. P. Kocher, R. Lee, G. McGraw, A. Raghunathan, and S. Ravi. Security as a New Dimension in Embedded System Design. In *DAC*, pages 753–760, 2004.



10. P.C. Kocher, J. Jaffe, and B. Jun. Differential Power Analysis. In *CRYPTO*, volume LNCS 1666, pages 388–397, 1999.
11. V. Kolesnikov and T. Schneider. Improved Garbled Circuit: Free XOR Gates and Applications. In *ICALP*, volume LNCS 5126, pages 486–498, 2008.
12. B. Kreuter, A. Shelat, and C. Shen. Billion-gate secure computation with malicious adversaries. In *USENIX Security Symposium*, 2012.
13. Y. Lindell and B. Pinkas. An Efficient Protocol for Secure Two-Party Computation in the Presence of Malicious Adversaries. In *EUROCRYPT*, volume LNCS 4515, pages 52–78, 2007.
14. L. Malka and J. Katz. VMCrypt – Modular Software Architecture for Scalable Secure Computation. In *CCS*, pages 715–724, 2011.
15. D. Malkhi, N. Nisan, B. Pinkas, and Y. Sella. Fairplay - A Secure Two-Party Computation System. In *USENIX Security Symposium*, pages 287–302, 2004.
16. S. Mangard, E. Oswald, and T. Popp. *Power Analysis Attacks: Revealing the Secrets of Smart Cards*. Springer-Verlag, 2007.
17. M. Medwed, C. Petit, F. Regazzoni, M. Renauld, and F.-X. Standaert. Fresh re-keying II: Securing multiple parties against side-channel and fault attacks. In *CARDIS*, pages 115–132. LNCS 7079, 2011.
18. M. Medwed, F.-X. Standaert, J. Großschädl, and F. Regazzoni. Fresh re-keying: Security against side-channel and fault attacks for low-cost devices. In *AFRICACRYPT*, pages 279–296. LNCS 6055, 2010.
19. M. Naor, B. Pinkas, and R. Sumner. Privacy preserving auctions and mechanism design. In *Electronic Commerce*, pages 129–139, 1999.
20. National Institute of Standards and Technology (NIST). The Keyed-Hash Message Authentication Code (HMAC). Federal Information Processing Standards Publication 198-1, Jul. 2008.
21. B. Pinkas, T. Schneider, N.P. Smart, and S.C. Williams. Secure Two-Party Computation is Practical. In *ASIACRYPT*, volume LNCS 5912, pages 250–267, 2009.
22. T. Popp, M. Kirschbaum, T. Zefferer, and S. Mangard. Evaluation of the Masked Logic Style MDPL on a Prototype Chip. In *CHES*, volume LNCS 4727, pages 81–94, 2007.
23. S. Ravi, A. Raghunathan, P.C. Kocher, and S. Hattangady. Security in Embedded Systems: Design Challenges. *TECS*, 3(3):461–491, 2004.
24. F.-X. Standaert, T.G. Malkin, and M. Yung. A Unified Framework for the Analysis of Side-Channel Key Recovery Attacks. In *EuroCrypt*, volume LNCS 5479, pages 443–461, 2009.
25. F.-X. Standaert, O. Pereira, Y. Yu, J.-J. Quisquater, M. Yung, and E. Oswald. Leakage Resilient Cryptography in Practice. In *Towards Hardware-Intrinsic Security*, pages 99–134. 2010.
26. S. Tillich, M. Kirschbaum, and A. Szekely. Implementation and Evaluation of an SCA-Resistant Embedded Processor. In *CARDIS*, volume LNCS 7079, pages 151–165, 2011.
27. Valgrind Project. Massif User Manual. <http://valgrind.org/docs/manual/ms-manual.html>.
28. A.C. Yao. Protocols for secure computations. In *Foundations of Computer Science*, pages 160–164, 1982.
29. A.C. Yao. How to generate and exchange secrets. In *Foundations of Computer Science*, pages 162–167, 1986.
30. P. Yu and P. Schaumont. Secure FPGA circuits using controlled placement and routing. In *CODES+ISSS*, pages 45–50, 2007.