

Exploiting Bit-level Parallelism in GPGPUs: a Case Study on KEELOQ Exhaustive Key Search Attack

Giovanni Agosta, Alessandro Barenghi and Gerardo Pelosi

Dipartimento di Elettronica e Informazione (DEI)

Politecnico di Milano, 20133 Milano (MI), Italy

Email: {agosta, barenghi, pelosi}@elet.polimi.it

Abstract—Graphic Processing Units (GPU) are increasingly popular in the field of high-performance computing for their ability to provide computational power for massively parallel problems at a reduced cost. However, the programming model exposed by the GPGPU software development tools is often insufficient to achieve full performance, and a major rethinking of algorithmic choices is needed. In this paper, we showcase such an effect on a case study drawn from the cryptography application domain. The pervasive use of cryptographic primitives in modern embedded systems is a growing trend. Small, efficient cryptosystems have been effectively employed to design and implement keyless password-based access control systems in various wireless authentication applications. The security margin provided by these lightweight ciphers should be accurately examined in light of the speed and area constraints imposed by the target environment. We present a re-design of the ASIC-oriented KEELOQ implementation to perform efficient exhaustive key search attacks while fitting tightly the parallel programming model exposed by modern GPUs. Indeed, the *bitslicing* technique allows the intrinsic parallelism offered by word-oriented SIMD computations to be effectively exploited. Through proper adaptation of the algorithm implementation to a platform radically different from the one it was designed for, we achieved a $\times 40$ speedup in the computation time with respect to a single-core CPU brute-force attack, employing only consumer grade hardware. The outstanding speedup obtainable points to a significant weakening of the cipher security margin, since it proves that anyone with off-the-shelf hardware is able to circumvent the security measures in place.

I. INTRODUCTION

In the last years, Graphics Processing Units (GPUs) have raised wide interest as sources of computational power for non-graphical applications, due to the availability of programming models such as CUDA and OpenCL that are vastly more accessible to experts of other domains than graphics rendering APIs (OpenGL and DirectX) [8]. A major strength of GPGPU-based platform are their appealing cost-performance figures of merit. In recent times even in the field of High Performance Computing there have been major investments to build GPGPU-based supercomputers. However, there are also factors that hinder the expansion of GPGPU computing, especially the difficulty of programming efficient applications using the available programming models. Special attention must be placed to tailor the application and its algorithmic components to the specific needs of the parallel hardware, e.g. by minimizing control flow divergence and exposing as much parallelism as possible while minimizing synchronization overheads [8]. In this paper, we show how the use

of specialized techniques can lead to large speedups, thus allowing the GPU to contend on an equal or favorable base (in terms of computation throughput per euro) with solutions based on CPUs or reconfigurable hardware. The field of cryptography has been explored since the first GPGPU attempts using graphics rendering APIs [15]. Especially, code breaking is attractive [2], [6], [7], because it requires vast amounts of computational power. We use as a case study the KEELOQ algorithm [12], which is used in remote keyless entry systems (e.g., vehicle doors or building entrances) or as authentication mechanism in wireless protocols. Remote keyless entry systems are based on a password based access control mechanism realized through the unidirectional transmission between a secure token (*encoder*) and a receiver (*decoder*). Unauthorized accesses are possible when the encoded password (*access code*) is fixed or it is derived from a relatively low number of possible combinations. In order to prevent this kind of threat, KEELOQ is employed in the so-called *rolling code* (also known as *hopping code*) mode of operation. The basic idea is to have the access code change each time it is used through picking it from a sequence of codewords that cannot be predicted even knowing a very large number of previously used ones. The generation of such a sequence is based on the definition of both a uni-directional command transfer protocol and an encryption engine to provide the codewords to be transmitted. From an operational point of view, the information transmitted by the encoder is composed by two parts: the code-hopping part (which changes each time the encoder is activated) and a second un-encrypted part, principally containing the encoder serial number, used for identifying the transmitter at a receiving decoder. To this end, the receiver decrypts the codeword, and compares the recovered counter value with its internal one, and the recovered serial number with the one received along with the codeword. If both values match, the token is granted access. Algorithms such as KEELOQ are designed for dedicated hardware implementation, since the target devices (remote controllers) are manufactured as very low cost ASICs. So, their direct implementation in software has much lower performances – which, in principle, makes it easier to carry out an attack using configurable hardware such as FPGAs. However, we show how the introduction of a level of parallelism not commonly seen in GPGPU algorithm design, *bit-level* parallelism, can lead to a $\times 40$ speedup over a CPU core. The rest of this paper is organized as follows.

Section II introduces the KEELOQ cipher, while Section III reviews the characteristics of the NVIDIA GPU families target in this study, as well as the programming model implemented by the CUDA development tools. Section IV describes the design of our solution and Section V provides the experimental evaluation on the case study. Finally, Section VI outlines the most closely related works, while Section VII draws some conclusions.

II. THE KEELOQ CIPHER

KEELOQ is the most scrutinized encryption engine used in remote keyless entry systems. It is a proprietary hardware-dedicated block cipher designed as a pair of Feedback Shift Registers (FSR) coupled with a Non Linear function (NL). Figure 1 shows the internal structure of the KEELOQ cipher: the secret key is stored in the red register on the left and is at most 64-bit wide. The key register is a FSR, and the key is mixed with the output of the state one bit per clock cycle. The 32-bit long Non Linear Feedback Shift Register (NLFSR) on the right hand side constitutes the nonlinear component of the cipher providing its effective security margin. Five bits of the NLFSR are combined together by means of a non linear function described by an equation over \mathbb{Z}_2 among five bits of the status register. The non linear function outputs a single bit per clock cycle, which is added to the aforementioned key bit and to b_{16} and b_0 , and employed as the feedback bit of the NLFSR. To encrypt a 32-bit plaintext block, the NLFSR is initialized with the value of the plaintext, and subsequently the entire system is clocked 528 times. After the 528 updates of both registers, the content of the NLFSR is the final ciphertext. The most common mode of operation for KEELOQ is the so-called *hopping code*, in a scenario where a remote *encoder* transmits a codeword to the authorizing *decoder* (receiver). This mode of operation involves encrypting a plaintext built out of a counter and a unique identifier (ID) of the encoding device. Every time a new 32-bit codeword (i.e. a ciphertext block) must be generated, the counter is incremented and the new plaintext is encrypted. Then, the codeword is transmitted along with the encoding device ID. The secret 64-bit key of any encoder is generated through the decoder engine as a pair of 32-bit codewords. Such a procedure implies that the decoder is able to generate the secret keys for a number of encoders starting from: (i) an embedded 64-bit master key (which is fixed by the manufacturer of the keyless entry system), (ii) the ID of the encoding device, (iii) and a random seed composed by 32, 48 or 60 bits.

A potential attacker may retrieve the master key from the decoding device (receiver) and eavesdrop the ID of an encoder when it is transmitted along with a codeword. Therefore, the use of a secret random seed in the secret key generation phase avoid the leakage of the secret key of the targeted encoder. A brute-forcing attack aimed at recovering the secret key of the transmitting encoder (EK) employs two consecutively transmitted codewords, each of which is bound to the encoder ID. The attacker computes a candidate 64-bit value for EK through guessing on the bits of the random seed, while

the value of the remaining part of the secret key is easily derived from the specification of the key generation protocol. Subsequently, she checks the ID value resulting from the decryption of the first codeword, and whether a match is found, the output derived from the decryption of the second codeword (employing the same EK) is used as a confirmatory step.

III. GENERAL PURPOSE COMPUTING WITH GPUS

The GPGPU devices targeted in this work are based on the NVIDIA GT200 and Fermi architectures. Figure 2 shows a sketch of the NVIDIA GTX470 (Fermi) streaming processor array. A streaming multiprocessor (SM) contains 32 streaming processors, four special functional units and a multithreaded instruction issue unit (respectively indicated as SP, SFU and MT-Issue in Figure 2. This is a fourfold increase over the GT200 SMs. A streaming multiprocessor concurrently executes two groups of 32 threads called *warps*, for a total of 64 concurrent threads. Since each thread in a warp has its own control flow, their execution paths may diverge due to the independent evaluation of conditional statements; when this happens, the warp serially executes each path. Each multiprocessor executes warps much like the *Single Instruction Multiple Data* (SIMD) paradigm, as every thread is assigned to a different SP and every active thread executes the same instruction on different data. Finally, the Fermi architecture includes both L1 and L2 cache memories, with the L1 configurable between cache and shared memory behavior and shared by the SPs in a single SM, and the L2 shared among all SMs in the device. The earlier GT200 only has a fast shared memory shared within each SM. GPGPU computing requires the programmer to manage a heterogeneous system (CPU host plus GPU device) as well as to handle the massive parallelism exposed by the GPU hardware. The Compute Unified Device Architecture (CUDA) [9], [14], proposed by NVIDIA for its graphics processors starting with the G80 series [5], exposes a programming model that integrates host and GPU code in the same C++ source files. On the GPU device side, a *Single Instruction, Multiple Threads* (SIMT) programming model is exposed, where a single *kernel* is executed by a user-specified number of threads. Every CUDA kernel is explicitly invoked by host code and executed by the device, while the host-side code continues the execution asynchronously after instantiating the kernel. On the host side, a specific synchronizing function call is provided to wait for the completion of the active asynchronous kernel computation. The CUDA programming model abstracts the actual parallelism implemented by the hardware architecture, providing the concepts of *block* and *thread* to express concurrency in algorithms. A block captures the notion of a group of concurrent threads. Blocks are required to execute independently, so that it has to be possible to execute them in any order (in parallel or in sequence). Therefore, the synchronization primitives semantically act only among threads belonging to the same block. Intra-block communications among threads use the *logical shared memory* associated with that block. Since the architecture does not provide support for message-

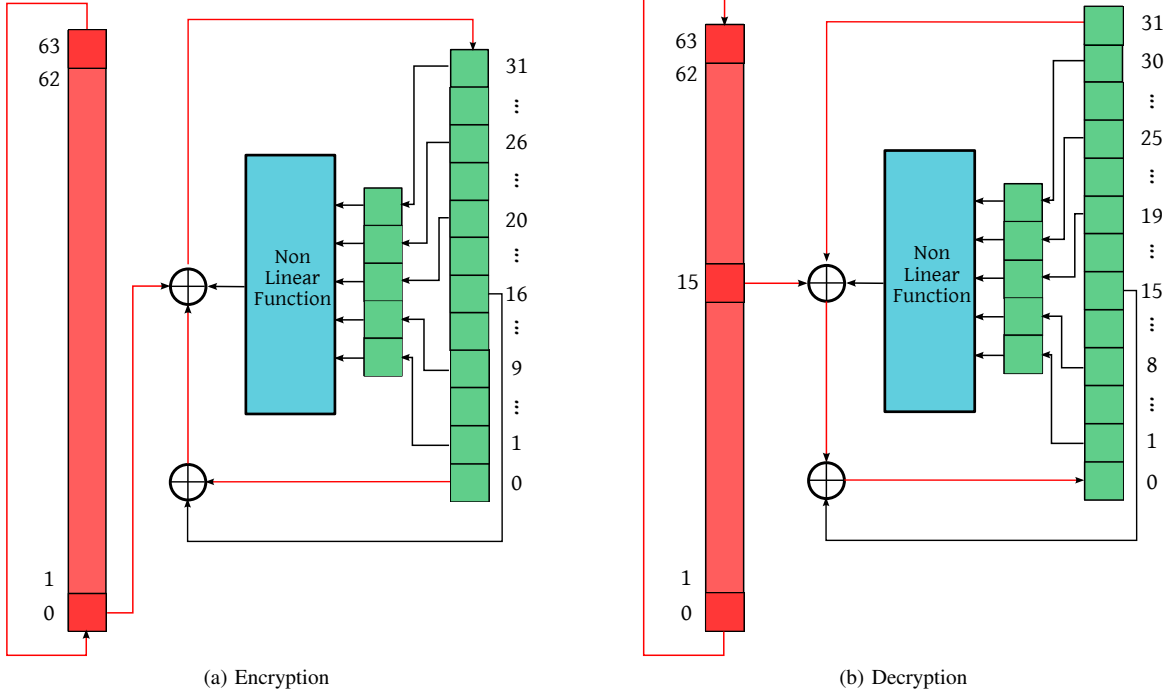


Figure 1. KEELOQ Cipher

passing, threads belonging to different blocks must communicate through *global memory*. Note that while the OpenCL language and API [10] are gaining momentum as the industry standard in programming heterogeneous platforms composed of host CPUs and programmable accelerators, including GPGPUs, the implementations provided are still not mature enough to compete, on NVIDIA devices, with the vendor-specific software development tools. However, the programming model provided in OpenCL is, as far as GPGPU programming goes, essentially based on the same principles as the SIMT model exposed in CUDA, so the techniques and results shown in this work can be easily extended to OpenCL-driven devices.

IV. ADAPTATION TO PARALLEL ARCHITECTURES

Many-core architectures offer large amount of parallel computing power by supplying the developer with hundreds of processing cores, each endowed with limited resources. In GPGPU, key resource limitations include:

Control flow divergence as multiple divergent control flows can be handled safely from the point of view of functionality, but with major performance losses as parallelism is inhibited along the different control flows – essentially, divergent flows of control are serialized, regardless of the data dependences among the divergent threads (which may well be non-existent). This limitation is due to the hardware design of GPGPU, where the processors in a multiprocessor unit are bound to the same program counter.

Local memory availability as a limited amount of very fast local memory must be shared among numerous processing elements. While the sharing allows fast communication among

the processing elements, the local memory is much more useful when used in a read-only way, or partitioned for local use by each processing element, since true shared accesses still require costly synchronization operations, and are often difficult to code. To exploit such parallel computing power, the critical issue is to be able to express a given application or algorithm in a form amenable to parallel execution on the target device. The literature reports three main sources of parallelism, which can be exploited with different degrees of success on various types of parallel architectures:

Thread-level parallelism is obtained when two or more tasks (regions of code with independent control flow) can be executed in parallel with few or no data dependencies (in the former case, synchronizations will be needed within each task, in the latter the synchronization point will be the end of the tasks). Thread-level parallelism is exposed by complex applications, where multiple independent tasks are performed, and is best exploited on symmetric multiprocessors, where each processor is endowed with sufficient resources to executed its assigned task. It is not suited for GPGPUs, since control flow divergence is a major factor for performance reduction in these architectures.

Loop-level parallelism is found in parallel loop constructs, where each iteration of the loop is data-independent from the others (or has limited synchronization requirements). Loop level parallelism is an excellent fit for vector processors, SIMD processors and GPGPUs, since control is fixed and identical for all iterations (barring nested conditionals, which can often be transformed to predicated code).

Instruction-level parallelism is achieved at the finest of

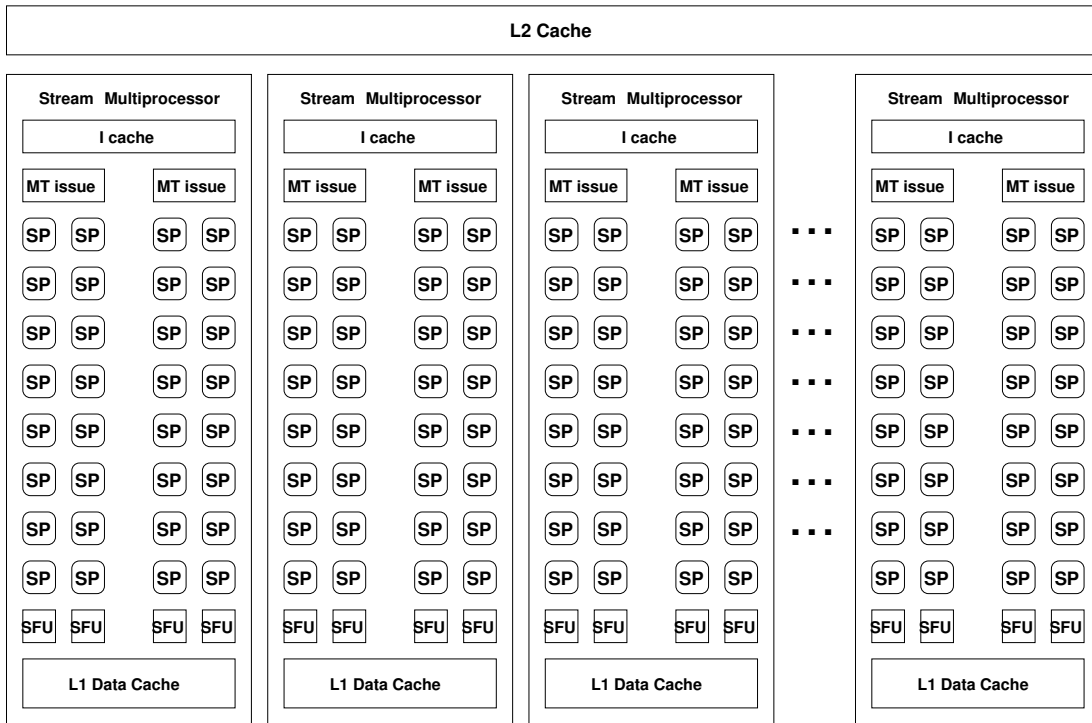


Figure 2. Overview of the NVIDIA GTX470 (Fermi) streaming processors architecture: each stream multiprocessor (SM) contains 32 streaming processors (SP), plus four special function units (SFU). A configurable L1 cache/shared memory is local to each stream multiprocessor, while L2 cache is shared among the entire set of SM. Up to 16 SM can be present in a single unit.

the three common granularities, where independent instructions can be parallelized. It is commonly exploited by super-scalar and Very Long Instruction Word architectures, but, like Thread-level parallelism, it is unsuitable for GPGPU due to the need to executed different instructions in parallel, rather than the same instruction of different data. It would therefore seem that Loop-level parallelism is the only viable choice for GPGPUs, but this model is not exposed by many types of codes. A typical example are encryption primitives designed for hardware implementation. In this case, parallelism is rarely available, but this is not an issue, since the implementation is performed through dedicated ASIC, and may be even considered a benefit, since software implementations are often aimed at *breaking* the encryption through brute force attacks. The usage of GPGPUs to perform brute force attacks is well-documented, but is often limited to mere juxtaposition of several encryption operations with different keys. However, it is possible to push the parallelization further, by introducing an entirely different level of parallelism, *Bit-level parallelism*. Here, the goal is to parallelize operations at the single bit level, thereby obtaining remarkably uniform parallel operations. This technique is know as *bitslicing* [4].

Bitslicing refers to a software technique of using a general purpose CPU to implement Single Instruction Multiple Data (SIMD) operations. The strategy consists of packing the bit values belonging to different operands within a single register and of using general-purpose arithmetic/logic instructions as specialized virtual processing elements designed for SIMD

operations at bit level. Most of the symmetric cryptographic primitives are designed to process input data at bit level. Therefore, the software implementations of such algorithms on not-specialized architectures may greatly benefit from the application of the bit-slicing strategy as long as the underlying hardware resources in terms of number of registers are easily available. Figure 3 reports a public domain, plain C implementation of KEELOQ from [1], while Figure 4 reports our bitsliced CUDA implementation. In the case of KEELOQ breaking, the bitslicing technique is employed through the decryption of the same 32-bit value using all possible keys. In the original code, operations on individual bits of the input are performed by means of *shift* and *mask* operations. In the optimized code, all operations work on full 32-bit words. To this end, the 32-bit input data is expanded (on the CPU host side) to a 32-word array, `bitsl_data_in`, where the i -th word in the array is `0xFFFFFFFF` if the i -th bit of the original text is set, or `0x00000000` otherwise. Input data is identical for each thread, but the same is not true for the key and decrypted output, which are stored separately for each of the $n_{threads} \times n_{blocks}$ threads. The bitsliced keys (each of 64 bits) are generated in blocks of 32 keys, starting from zero and progressively increasing its value. Each 64-word array `bitsl_key[bI][tI]` generated in this way has the five last words (corresponding to the lower bits of the original keys) always equal to the encoding of the same 32 values, which are added to a “base” key value. The base key value, in turn, increases in steps of 32. Thus, the number

```

1 #define NLF 0x3A5C742E
2 #define bit(x,n) ((x)>>(n))&1
3 #define g5(x,a,b,c,d,e) \
4 (bit(x,a)+bit(x,b)*2+bit(x,c)*4+bit(x,d)*8+bit(x,e)*16)
5
6 uint32_t KeeLoq_Decrypt (uint32_t data, uint64_t
7 key) {
8     uint32_t x=data, r;
9     for (r=0; r<528; r++)
10         x= (x<<1)^bit(x,31)^bit(x,15)^(u32)bit(key,(15-r)
11             &63)
12             ^bit(NLF,g5(x,0,8,19,25,30));
13     return x;
14 }

```

Figure 3. Plain C KEELOQ implementation [1].

```

1 __device__ uint32_t bitsl_data_in[32];
2 __device__ uint32_t bitsl_data_out[NBLOCKS][NTHREADS
3 ][32];
4 __device__ uint32_t bitsl_key[NBLOCKS][NTHREADS
5 ][64];
6
7 __global__ void KeeLoq_Decrypt_Bitslice() {
8     int bI = blockIdx.x, tI = threadIdx.x;
9     uint32_t data[32];
10    #pragma unroll 32
11    for (int j=0; j<32; j++) data[j]=bitsl_data_in[j];
12
13    uint32_t key_r, nlf, rs_data;
14    for (int r=0; r<528; r++) {
15        key_r = bitsl_key[bI][tI][(15 - r) & 0x3F];
16        nlf = NonLinearFunction(data, 0, 8, 19, 25, 30);
17        rs_data = data[31] ^ data[15] ^ key_r ^ nlf;
18        for (int i=31; i>0; i--) data[i] = data[i - 1];
19        data[0] = rs_data;
20    }
21
22    #pragma unroll 32
23    for (int j=0; j<32; j++) bitsl_data_out[bI][tI][j]=
24        data[j];
25 }

```

Figure 4. Bitsliced KEELOQ CUDA kernel. The plaintext, key and ciphertext are stored in the GPU main memory at the beginning of the kernel execution, then copied to registers during the kernel itself.

of parallel encryption runs is 32 per thread, as shown by the kernel in Figure 4, with configurable number $n_{threads}$ of threads per each CUDA block. Overall, a grand total of $32 \times n_{threads} \times n_{blocks}$ encryption runs are performed at every time by the GPU. It is worth noting that the shared memory is not used, since the Fermi architecture provides a large number of registers. A full analysis of the tradeoffs will be shown in the next section.

V. EXPERIMENTAL RESULTS

We implemented a fully bitsliced version of the KEELOQ cipher both employing the CUDA programming model and pure C. The pure C version has been run on the host CPU to provide a reference implementation as far as throughput goes. The running environment where the bruteforcing speed tests were performed is an Intel Core i7 920 based system with 12Gb DDR3 DRAM, running Gentoo Linux AMD64. All the

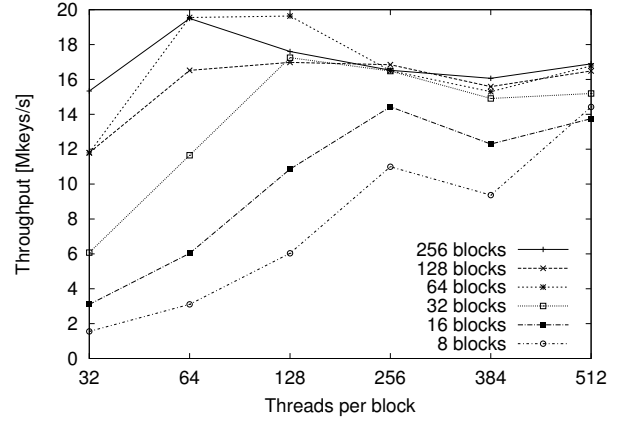


Figure 5. Throughput of the bitsliced implementation of the KEELOQ breaker on the GeForce GTX470 card, related to the number of threads per block and the number of blocks per CUDA kernel invocation

GPU binaries were compiled employing `nvcc` 4.0 from nVidia CUDA toolkit 4.0, while the CPU baseline versions were compiled with `gcc` 4.4.6. The bitsliced implementation of the cipher has been tested on two different GPUs, which have been mounted as the only device on the 16 lane PCI-Express 2.0 port available on the motherboard in order to test the difference in performances. The first GPU card is a GeForce GTX 260 equipped with 894 Mb of GDDR5 video RAM and 192 CUDA cores, while the second card employed for testing is a GeForce GTX 470 with 448 CUDA cores and 1280 MB of GDDR5 video RAM. An important step in the evaluation of the performances of our bitsliced implementation of KEELOQ on CUDA is the exploration of two parameters: the number of threads composing a CUDA block and the number of blocks constituting a CUDA kernel call. The first parameter regulates the level of register pressure on the shared register file of the streaming multiprocessor and the number of warps into which a CUDA block is split. Since the basic execution unit of a streaming multiprocessor is a single warp, the choice of the number of threads should consider only multiples of 32 to achieve the best fit. The level of register pressure on the Fermi architecture is dictated by the fact that the 32768 registers are shared among the contexts of up to 3 different blocks which can be scheduled on the same streaming multiprocessor. In addition to this, the SMP issue unit of the Fermi architecture is able to dual issue warps, thus it is necessary to keep twice the contexts in the registers. Combining these data with the fact that a single bitsliced KEELOQ breaking thread employs at most 45 registers, we obtain a SMP register pressure which can be computed as $270 \times n_{threads}$. The second parameter to be chosen regulates the level of global computational load imposed on the GPU. The main point in choosing this parameter is provide at least enough computations to the GPU so that no SMPs remain idle. Moreover, since the SMP issue unit is able to interleave different blocks in order to hide global memory access latencies, it is wise to provide extra workload to the GPU to exploit this feature. These two considerations pointed to the creation of a CUDA kernel as large

Table I
EXPECTED TIMINGS AND MEASURED THROUGHPUT FOR THE EXHAUSTIVE SEARCH OF THE KEELOQ KEY GENERATION SEED.

Seed Length	Single Core	Four Cores	Single GPU	
	Core i7 920 [h]	Core i7 920 [h]	GTX260 [h]	GTX470 [h]
32	2.6	0.73	0.14	0.04
48	$1.73 \cdot 10^5$	$4.84 \cdot 10^4$	$9.45 \cdot 10^3$	$3.98 \cdot 10^3$
60	$7.08 \cdot 10^8$	$1.98 \cdot 10^8$	$3.81 \cdot 10^7$	$1.63 \cdot 10^7$
Throughput [key/s]	$4.51 \cdot 10^5$	$1.61 \cdot 10^6$	$8.27 \cdot 10^6$	$1.96 \cdot 10^7$

as possible with architectures up to the GT200, since the static scheduling of the blocks on the SMPs did not account for extra time overhead. With the introduction of a new scheduler for multiple kernels on the Fermi architecture, this consideration may not be still valid. Figure 5 reports the results of the exploration of the implementation parameter space: coherently with the previous considerations, the best solution is reached with 128 threads per block (34560 employed registers), when the number of blocks per SMP is enough to fill all the issue queues completely. Raising further the number of blocks per kernel leads to a decrease in performances which can be ascribed to the extra context switching effort imposed on the new scheduler. As expected also raising further the number of threads per block leads to a significant decrease in throughput due to the hindering of context switches caused by the frequent register spills and fills. An analogous exploration campaign has been lead also on the GTX260 card, yielding 64 threads per block as the best performing choice of the parameter. This choice is coherent with the fact that the shared register file of the GTX260 is 16384 since 64 threads per block allow the issue unit of the streaming multiprocessor to perform the context switching between the three blocks in queue without the need to spill part of the register file to the global memory. In this case, however, increasing arbitrarily the number of blocks per kernel did not induce any performance penalty as expected from the GT200 architecture. After choosing the optimal number of threads per block and blocks per kernel invocation, we evaluate the effective time needed in order to break the KEELOQ key generation mechanism, with respect to the length of the employed seed. Table I reports the expected running times of an attack, depending on the chosen platform to perform the exhaustive search. Taking as a reference value the throughput obtained by the bitsliced implementation of KEELOQ running on the host CPU (419430 keys/s), we notice that employing a 32 bit seed for the key generation does not yield a sufficient security margin, as the remote key can be recovered in 3 hours of computation. The bitsliced implementations running on the GTX260 and GTX470 GPUs achieve a $\times 20.5$ and a $\times 43.5$ speedup respectively, allowing a possible attacker to breach even the security of the 48 bit seed key generation mechanism in a few months. Since the exhaustive search can be split over multiple GPUs, it is possible to lower the attack time to a single week, while keeping the cost envelope of the equipment below \$10000, as this budget allows an attacker to build a 20 GTX 470 cluster

with the current market prices.

VI. RELATED WORK

The first cryptanalysis of KEELOQ is presented in [3]. The attack is based on the *slide technique* and a linear approximation of the non-linear Boolean function used in the cryptographic engine. The attack requires 2^{52} encryptions, 16GB of storage and the entire codebook, i.e., 2^{32} known plaintexts. In [16] the authors introduce a specific key recovery attack against KEELOQ which combines the technique of slide attacks with a novel meet-in-the-middle approach. Their method requires 2^{16} chosen plaintexts and has a time complexity of $2^{44.5}$ encryptions which results in about two days of computation employing 50 dual core CPUs at the cost of approximately €10000. The widely adoption of KEELOQ in practice, paved the way to side-channel analysis as a further viable option for attacking chips that implement it. In [17] the first successful DPA and DEMA attacks on KEELOQ implementations applied to both Identify Friend or Foe (IFF) and code hopping devices, are presented. The attack is prevented if a 60-bit seed value, with good random properties, is employed for the key derivation. Nevertheless, considering the other commonly implemented options of the cipher, the authors reported how to reveal a manufacturer key from a receiver using a few 1000 power traces, and how to recover the device key of a remote control with as few as 10 traces. In [13] the authors apply algebraic techniques to cryptanalyze the cipher. This attack employs the entire codebook, 2^{27} encryptions and has an estimated success probability of 44%. The results of a brute-force attack, implemented on the FPGA-based code-breaker COPACOBANA, are reported in [11]. The authors claim the secret key recovery of a remote control in less than 0.5 seconds if a 32-bit seed is used and in less than 6 hours in case of a 48-bit seed. The case of a 60-bit seed needs in the worst case about 1011 days at the cost of approximately \$10000. However, the technical effort needed to build an FPGA-based code-breaker, and even more one the size of the COPACOBANA, is much greater than that needed to carry out a GPU-based attack. Moreover, while FPGA-based code-breakers require specialized hardware, the GPU-based attack can benefit of a large installed base of CUDA-enabled devices, allowing distributed attacks to be carried out by groups of users, or by botnets.

VII. CONCLUSIONS

In this paper, we report our experience with *bit-level parallelism* in GPGPU programming, using as a case study the brute force attack on the KEELOQ cipher. We proposed a full redesign of the computation strategy from the original hardware implementation-oriented algorithm to reach high performance in parallel software, by exploiting SIMD techniques down to the bit level. We report a speedup of $\times 40$ speedup in the computation time with respect to a CPU brute force attack, even though only consumer-grade hardware is used.

REFERENCES

- [1] "C KEELOQ Implementation," [On line] <http://cryptolib.com/ciphers/keeloq/KeeLoq.c>, December 2011.
- [2] Andrea Di Biagio and Alessandro Barenghi and Giovanni Agosta and Gerardo Pelosi, "Design of a parallel AES for graphics hardware using the CUDA framework," in *IPDPS*. IEEE, 2009, pp. 1–8.
- [3] Andrey Bogdanov, "Linear slide attacks on the KEELOQ block cipher," in *Inscrypt*, ser. Lecture Notes in Computer Science, Dingyi Pei and Moti Yung and Dongdai Lin and Chuankun Wu, Ed., vol. 4990. Springer, 2007, pp. 66–80.
- [4] Eli Biham, "A fast new des implementation in software," in *FSE*, ser. Lecture Notes in Computer Science, Eli Biham, Ed., vol. 1267. Springer, 1997, pp. 260–272.
- [5] Erik Lindholm, John Nickolls, Stuart Oberman, and John Montrym, "NVIDIA Tesla: A Unified Graphics and Computing Architecture," *Micro, IEEE*, vol. 28, no. 2, pp. 39–55, 2008.
- [6] Giovanni Agosta and Alessandro Barenghi and Fabrizio De Santis and Andrea Di Biagio and Gerardo Pelosi, "Fast Disk Encryption through GPGPU Acceleration," in *PDCAT*. IEEE Computer Society, 2009, pp. 102–109.
- [7] Giovanni Agosta and Alessandro Barenghi and Fabrizio De Santis and Gerardo Pelosi, "Record Setting Software Implementation of DES Using CUDA," in *ITNG*, Shahram Latifi, Ed. IEEE Computer Society, 2010, pp. 748–755.
- [8] John D. Owens, Mike Houston, David Luebke, Simon Green, John E. Stone, and James C. Phillips, "GPU Computing," *Proceedings of the IEEE*, vol. 96, no. 5, pp. 879–899, May 2008.
- [9] John Nickolls and Ian Buck and Michael Garland and Kevin Skadron, "Scalable parallel programming with cuda," *ACM Queue*, vol. 6, no. 2, pp. 40–53, Mar. 2008.
- [10] Khronos OpenCL Working Group, "OpenCL - The Open Standard for Parallel Programming of Heterogeneous Systems," [On line] <http://www.khronos.org/opencl/>, January 2011.
- [11] Martin Novotny and Timo Kasper, "Cryptanalysis of KeeLoq with CO-PACOBANA," in *Workshop on Special Purpose Hardware for Attacking Cryptographic Systems (SHARCS'09)*, 2009.
- [12] Microchip Technology Inc., "Security and Authentication Design Center-KEELOQ® 3 Development Kit," [On line] http://www.microchip.com/stellent/idcplg?IdcService=SS_GET_PAGE&nodeId=2074, December 2011.
- [13] Nicolas Courtois and Gregory V. Bard and David Wagner, "Algebraic and slide attacks on KEELOQ," in *FSE*, ser. Lecture Notes in Computer Science, Kaisa Nyberg, Ed., vol. 5086. Springer, 2008, pp. 97–115.
- [14] NVIDIA Corporation, "CUDA Technology," [On line] <http://www.nvidia.com/CUDA>, Sep. 2008.
- [15] Owen Harrison and John Waldron, "AES Encryption Implementation and Analysis on Commodity Graphics Processing Units," in *CHES*, ser. Lecture Notes in Computer Science, Pascal Paillier and Ingrid Verbauwhede, Ed., vol. 4727. Springer, 2007, pp. 209–226.
- [16] Sebastiaan Indestege and Nathan Keller and Orr Dunkelman and Eli Biham and Bart Preneel, "A practical attack on KEELOQ," in *EURO-CRYPT*, ser. Lecture Notes in Computer Science, Nigel P. Smart, Ed., vol. 4965. Springer, 2008, pp. 1–18.
- [17] Thomas Eisenbarth and Timo Kasper and Amir Moradi and Christof Paar and Mahmoud Salmasizadeh and Mohammad T. Manzuri Shalmani, "On the power of power analysis in the real world: A complete break of the KEELOQ code hopping scheme," vol. 5157, pp. 203–220, 2008.