

Efficient execution of the MD5 algorithm in password recovery

Sam Mitchell and Nathaniel Weidler
Department of Electrical and Computer Engineering
Utah Stat University
Logan, Utah 84322

e-mail: samuel.alan.mitchell@gmail.com, NWeidler@gmail.com

Abstract—Summarize project and results (executive summary).

I. INTRODUCTION

Describe the problem and what you're doing to address it (i.e. building architecture-informed md5 password cracker for software and hardware; how password cracking operates).

A. Related work

What is needed to understand your optimisations and/or implementations. Do you draw on existing techniques? If so, describe and cite them.

B. Structure of paper

Which sections discuss which aspect of the work.

II. SOFTWARE IMPLEMENTATION

Overview and objectives MD5 is a computationally simple operation with minimal execution time. An inefficient implementation of MD5 is sufficient for everyday hashing needs, however password recovery applications require an efficient method to test large numbers of passwords.

Parallel computing is an effective method to quickly analyze large sets of data. OpenMP is one parallel computing implementation which has been in use for C/C++ since 1998 [1]. The authors utilized this parallel computing method with 16 threads, separating the set of passwords into 16 equal portions.

Further optimization methods were selected or developed. These methods are: Single Instruction Multiple Data, loop unrolling, string precalculation, and first block precalculation.

A. SIMD

Single Instruction Multiple Data (SIMD) is a functionality of standard processors which allows a single operation to operate on multiple integers in parallel [2]. This is effective due to the current register size (128-bits) of modern cpus.

The MD5 hashing algorithm operations are performed on 32-bit integers. This allows for four (4) simultaneous operations on each core, which led us to believe that the algorithm would see a 300% speedup.

Our implementation of SIMD initially slowed down the hash — a 37% decrease in performance when compared to the nominal case. Reasons for this are not confirmed, however

we theorize that the stack operations resulting from context switches caused this slow-down.

In the following optimizations, the SIMD implementation is shown to outperform the Single Instruction Single Data case.

B. Loop unrolling

Loop unrolling is an elementary optimization technique utilized to remove the overhead from loop operations [3]. The loop was unrolled to a specific length to assist in string precalculation. In the nominal (no SIMD) case, loop unrolling to 26 iterations gave slight improvements (2%). In the SIMD case, loop unrolling to 169 iterations resulted in a substantial speed increase. Nominal to SIMD unrolled resulted in a 103% speedup, while SIMD to SIMD unrolled resulted in a 220% speedup.

C. String precalculation

One significant delay in MD5 calculation is string creation. We strategically combined loop unrolling with precalculation methods. The loop was originally unrolled to have 26 iterations (SIMD had $\frac{26*26}{4} = 169$ iterations), which allowed precalculation of the first characters of the string.

MD5 is calculated on 16 32-bit numbers, and on a length seven (7) string there are 13 32-bit numbers that are always 0. The last 32-bit number is $7*8=56$. This is consistent on every length seven (7) string, so these values were directly coded into the MD5 hash.

When inserted to the already unrolled system, string precalculation resulted in a 36% speedup in the SIMD case.

D. First block precalculation

The first block of the MD5 calculation can be precomputed due to the repetitive nature of the password cracker — the first block only accepts the first four characters of the password. The first block was computed in function `G_MD5()` and stored for the successive 168 hashes. While this saved operations, the anticipated speedup was a maximum of $\leq \frac{1}{64}$ or 1.5% due to the nature of MD5. Testing this optimization revealed a 3% speedup, which was unexpected and as of yet unexplained.

E. Results

As expected, the combination of SIMD and loop unrolling was the optimization with the greatest speedup. When combined, these methods result in fewer comparisons and branches.

Table I: Benchmarking performed on a computer with two (2) Intel Core i7-4770K CPUs. Each CPU has 4 cores, 3.5 GHz, released June 2013.

Implementation	Hash rate (million hashes per second)
md5pc_base	40
md5pc_O0	20
md5pc_O1	155
md5pc_O2	157
md5pc_O3	158

String precalculation was an unexpectedly effective optimization. This was not an improvement in MD5 calculation efficiency, but a removal of excessive overhead operations.

Benchmarking was performed on all seven (7) character combinations of a lowercase, alphabetic character set. Comparisons were made between optimizations as well as best-case to the nominal OpenMP optimization. The comparisons of each optimization are noted in the respective section. The best-case system is compared to the nominal using the -O1, -O2, -O3, and -fomit-frame-pointer compiler flags. Results are found in Table I.

Compiler flags are very important in optimizing execution time of a process. The 690% speedup between md5pc_O0 and md5pc_O3 hash rates shows that performing operations in an out-of-order manner is often more effective than readable code.

F. Password recovery

A hash of an unknown password was provided (aebc994aa5b00a0308c9fd257bf63ebd). In order to ensure that no passwords would be missed, password testing was performed in discrete groups: (1) 1-6 character strings, (2) 7 character strings, (3) 8 character strings, and so on. These tests were performed using the same processor used in the benchmarking tests in Table I.

The 1-6 character strings were tested without optimizations because it was executed in an acceptably short time. No results were found. 5.91 second execution time.

The 7 character string was tested using the program used in benchmarking. No results were found. 51.15 execution time.

The 8 character string was tested using a modified version of the benchmarking program. The resulting string "tyygsuef" was found in 6:18.85 (minutes:seconds).

In terms of total execution time, the testing was completed in 7:15.91. This does not represent testing every string of lengths 1-8, because the program ceased executing once the password was found.

III. HARDWARE IMPLEMENTATION

Overview and objectives; details of implementation. Include architecture diagram. Discuss benchmark procedure. Ways to improve performance and/or discussion of how higher performance could be achieved. Performance is dependent on implementation (GPU vs. FPGA). FPGA or GPU details (cores, gates, clock, etc). Use a table, if need be.

IV. CONCLUSION

Summarize results and lessons learned.

REFERENCES

- [1] OpenMP Architecture Review Board, "The OpenMP API specification for parallel programming," <http://openmp.org/wp/>, 2015, [Online; accessed 02-Feb-2015].
- [2] B. Hubert. (2004) Simd and other techniques for fast numerical programming with gcc. (Accessed: 4 February 2015). [Online]. Available: ds9a.nl/gcc-simd/
- [3] J. J. Dongarra and A. Hinds, "Unrolling loops in fortran," *Software: Practice and Experience*, vol. 9, no. 3, pp. 219–226, 1979.