



| | |
|--------------------|--|
| Title | A model for peak matrix performance on FPGAs |
| Author(s) | Lin, CY; So, HKH; Leong, PHW |
| Citation | The 19th IEEE International Symposium on Field-Programmable Custom Computing Machines (FCCM 2011), Salt Lake City, UT., 1-3 May 2011. In Conference Proceedings, 2011, p. 251-251 |
| Issued Date | 2011 |
| URL | http://hdl.handle.net/10722/158706 |
| Rights | Annual IEEE Symposium on Field-Programmable Custom Computing Machines (FCCM). Copyright © IEEE Computer Society. |

A Model for Peak Matrix Performance on FPGAs

Colin Y. Lin and Hayden K.-H. So
 Dept. of Electrical & Electronic Engineering
 The University of Hong Kong
 Hong Kong
 Email: {linyu, hso}@eee.hku.hk

Philip H.W. Leong
 School of Electrical and Information Engineering
 The University of Sydney
 Sydney, Australia
 Email: philip.leong@sydney.edu.au

Abstract—Computations involving matrices form the kernel of a large spectrum of computationally demanding applications for which FPGAs have actively been utilized as accelerators. The performances of such matrix operations on FPGAs are related to underlying architectural parameters such as computational resources, memory and I/O bandwidth. A model that gives bounds on the peak performance of matrix-vector and matrix-matrix multiplication operations on FPGAs based on these parameters is presented. The architecture and efficiency of existing implementations are compared against the model. Future trends in matrix performance on FPGA devices are estimated based on the performance model and system parameters from the past decade.

I. MATRIX OPERATIONS ON FPGAS

For the purpose of this work, we assume that initially all input data are stored off-chip and must be loaded onto the FPGA fabric during the operation. Similarly, the result of computation must be stored off-chip. We make no assumption about the mechanism of data I/O as that is irrelevant to the derivation of the performance model. Furthermore, we assume the input matrices and vectors are large.

The following architectural parameters were considered in this work: on-chip logic/DSP resources (r), on-chip memory (m), I/O bandwidth (b), number of MACCs (k), and clock frequency (f).

In a matrix-vector multiplication, $y = Ax$, the entire input vector x must be reused in the calculation of each of the element in y while each row of data in A is used only once. Therefore, x is stored using on-chip memory while A is streamed through the FPGA to generate results. Blocking is used if the on-chip memory is limited.

The performance bound is found to be determined by the values kf and b . If $kf < b$, the performance is computation resource bounded and the maximum achievable operation per second equals to $\text{OPS}_{max} = \left\lfloor \frac{R}{r_{macc}} \right\rfloor \cdot 2f$. If $b < kf$, then the performance is bounded by I/O bandwidth. In this case, $\text{OPS}_{max} = 2B$.

In a matrix-matrix multiplication, $C = AB$, every data element from both matrix A and B are reused more than once. Therefore, the amount of on-chip memory m that may be used as temporary storage is an important design parameter. Optimized I/O scheme must therefore be used to

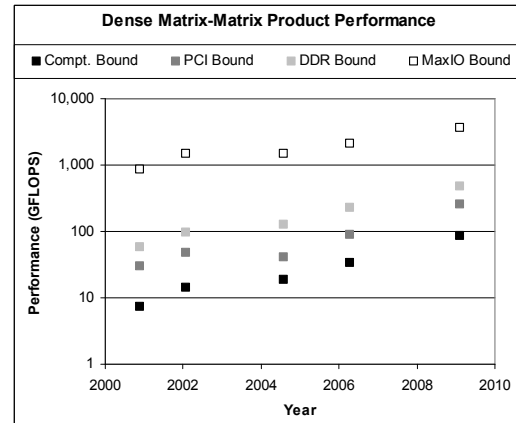


Figure 1. Matrix-matrix multiplication performance bounds for Xilinx Virtex series FPGAs.

utilize on-chip memory efficiently for the input matrices A and B .

In this work, the on-chip memory is mainly used to store the partial results of matrix C . The values are updated in-place as A and B are streamed on and off chip for computation.

With such assumptions, the performance of matrix-matrix multiplication have been found to depend on k , f , m and b . If $2kf < \sqrt{mb}$, the performance is bounded by compute resource availability. In this case, the maximum performance is $\text{OPS}_{max} = \left\lfloor \frac{R}{r_{macc}} \right\rfloor \cdot 2f$. If $\sqrt{mb} < 2kf$, then the performance is bounded by on-chip memory and I/O bandwidth. In this case, $\text{OPS}_{max} = \sqrt{MB}$.

II. TECHNOLOGY TRENDS & CONCLUSION

Figure 1 shows the memory, I/O and computation bound on matrix-matrix multiplication for Xilinx Virtex series FPGAs over the past decade. Computational resources have been the limiting factor in the past. However, with the amount of logic resources increasing according to Moore's Law, it is possible that I/O and on-chip memory will become the dominant bound in the near future.