

FPGA Based High Performance Double-Precision Matrix Multiplication

Vinay B. Y. Kumar · Siddharth Joshi ·
Sachin B. Patkar · H. Narayanan

Received: 15 July 2009 / Accepted: 31 January 2010 / Published online: 18 February 2010
© Springer Science+Business Media, LLC 2010

Abstract We present two designs (I and II) for IEEE 754 double precision floating point matrix multiplication, optimized for implementation on high-end FPGAs. It forms the kernel in many important tile-based BLAS algorithms, making an excellent candidate for acceleration. The designs, both based on the rank-1 update scheme, can handle arbitrary matrix sizes, and are able to sustain their peak performance except during an initial latency period. Through these designs, the trade-offs involved in terms of local-memory and bandwidth for an FPGA implementation are demonstrated and an analysis is presented for the optimal choice of design parameters. The designs, implemented on a Virtex-5 SX240T FPGA, scale gracefully from 1 to 40 processing elements (PEs) with a less than 1% degradation in the design frequency of 373 MHz. With 40 PEs and a design speed of 373 MHz, a sustained performance of 29.8 GFLOPS is possible with a bandwidth requirement of 750 MB/s for design-II and 5.9 GB/s for design-I. This compares favourably with both related art and general purpose CPU implementations.

Keywords High performance computing · Matrix multiplication · Rank-1 scheme · FPGA implementation · Memory-bandwidth trade-off · Scalability

V. B. Y. Kumar (✉) · S. Joshi · S. B. Patkar · H. Narayanan
Department of Electrical Engineering, Indian Institute of Technology, Bombay, Mumbai 400076, India
e-mail: vinayby@iitb.ac.in

S. Joshi
e-mail: siddharth_joshi@daiict.ac.in

S. B. Patkar
e-mail: patkar@ee.iitb.ac.in

H. Narayanan
e-mail: hn@ee.iitb.ac.in

1 Introduction

Field Programmable Gate Arrays (FPGAs) are a form of programmable logic based on configurable routing between logic resources, unlike Application Specific Integrated Circuits (ASICs) like general purpose processors. Designs on FPGAs consist of mapping functions onto *logic slices* and routing them appropriately. They are increasingly being seen as a promising avenue for High Performance Computing (HPC), especially with the introduction of high-end FPGAs like Xilinx Virtex-4/5/6 and the Altera Stratix series. These FPGAs are a very attractive choice due to their abundant local memory, high-speed embedded resources like DSP blocks, PCI-e endpoints etc., apart from the more obvious reasons such as reconfigurability and lower power consumption compared to general purpose hardware. A line of special purpose development boards based on such FPGAs are also available (e.g. Nallatech and Alpha data) which are particularly suitable for HPC; some of these come with tools which help faster algorithm-to-hardware realization with high-level C-like constructs; Maxwell [1, 2], a FPGA parallel computer, is a good example of a system built around such hardware/software ecosystem.

Efficient implementation of matrix-multiplication is an important goal for scientific computing. MEMOCODE [2007] chose acceleration of (complex integer) matrix-multiply as its first HW/SW codesign challenge. Underwood [3] chose matrix multiplication as one of the three main routines for FPGA acceleration in order for HPC. An inspection of Level-3 BLAS routines shows that matrix multiplication (*dgemm*) and triangular equation solution (*dtrsm*) form the building blocks for many important linear algebra algorithms, in fact, the *dtrsm* itself can be expressed in terms of *dgemm*. Matrix-Multiplication, therefore, presents as an important and useful candidate for hardware acceleration. The following resources can serve to provide more perspective in this context: [4, 5].

In this paper, we present designs for double precision floating point matrix multiplication [6], based on the rank-1 update algorithm, targeted at the Virtex-5 SX240T, a high-end Xilinx FPGA. As compared to others this algorithm enables better re-use of data from the input matrices. The processing elements (PEs) use off-the-shelf floating point operators from Xilinx Coregen, resulting in advantages such as the choice of custom-precision, short-design time, portability across Xilinx device generations, better IEEE 754 compliance, etc. The PEs are designed so as to scale linearly in terms of resources with negligible (<1%) degradation in speed. Some of the recent work [7] reports 35% speed reduction associated with scaling, which is typically due to increased routing complexity. The proposed design (II) works well with burst-like input and thus with a high-bandwidth I/O bus like PCI-e—allowing it to scale seamlessly across multiple FPGAs.

The underutilisation of device primitives and the over-dependence on the distributed memory available in FPGAs results in lower performance with respect to scaling, as with some of the discussed related work. The designs presented in this work have evolved by careful use of the high-performing resources on modern FPGAs. Care has been taken to address issues related to scaling for large FPGAs, setting this work apart from related art. The main thrust of our effort has been in reporting the design and

implementation, of this processing element, targeted at FPGAs and the results on the Virtex-5 SX240T.

The following sections are organised as follows—Sect. 2 discusses related work with a quick background on FPGAs, Sect. 3 discusses the underlying algorithm, Sect. 4 elaborately discusses both the designs, Sect. 5 presents an evaluation of the design, Sect. 6 presents an analysis on design parameters, Sect. 7 critically compares our design with the best among the related work, and finally Sect. 9 concludes the paper.

2 Background

The rank-1 update algorithm used in this paper is an elementary idea, variations of which have also been applied to cache-aware computing on general purpose processors [8], though not as aptly. This was chosen to be implemented on an FPGA since it is particularly suitable for the task as verified by both Dou and Prasanna.

Much of the related work are designs targeted and optimised for Virtex II Pro, which is an entry level device for HPC that made floating point computation feasible for the first time on FPGAs.

2.1 Related Work

The two most recent significant designs are those by Dou [9] and Prasanna [7]. They propose linear array based processing elements which are able to sustain their performance by overlapping IO and computation using a technique called memory switching.

Dou has proposed the design of a matrix multiplier highly optimised for the Virtex II Pro series of FPGAs. The design included an optimized custom 12-stage pipelined floating point multiply-accumulator (MAC), but with a few limitations like no support for zero and other denormal numbers. Correcting this requires additional resources and results in a decrease in the design frequency. This design also required the subblock dimensions to be powers of two. The bandwidth requirement was low at 400 MBps with 12.5 Mb of local memory utilization and they report a PE design with a synthesis frequency of 200 MHz accommodating 39 PEs on a Virtex II pro XC2V125, a large hypothetical device, and therefore estimated a 15.6 GFLOPS performance. These being only synthesis results, the real frequency after placing and routing 39 PEs could be less.

Zhou and Prasanna have reported an improved version [7] of their design reported in [10]. The later one reports 2.1 GFLOPS for 8 PEs running at 130 MHz on a cray XD1 with XC2VP50 FPGAs. About 35% speed degradation was observed when scaled from 2 to 20 PEs. In the earlier paper they presented a design with a peak performance of 8.3 GFLOPS for the Virtex II Pro XC2V125, where the clock degradation was 15% when the number of PEs increases from 2 to 24. A recent [4] report, which also discussing other linear algebra operations, shows a similar behaviour with respect to scaling for matrix multiplication.

2.2 FPGAs

For accelerator designs to be more than just an academic exercise the following are important considerations. The design time should be low and allow for extensive testing. The design should be modular in nature and scale with available resources. For integration within an existing system form-factor limitations should be considered as should power and memory. Most HPC systems are based on Infiniband like interconnects between nodes, with nodes having PCI-e for communication with peripherals. The PCI-e connects via the southbridge to the host memory, sharing bandwidth with the host processor. DMA is used in order to transfer data efficiently, and thus accelerators should be compatible with DMA and burst transfer. The overall system should also be oblivious to the presence of the accelerator, requiring minimal modifications to be done to accommodate it. These aspects make FPGA based designs very attractive. Their form factors allow multiple FPGA to fit on existing boards, that can communicate via PCI-e. The cost of being reconfigurable doesn't allow FPGAs to run at clocks as high as modern general purpose processors. However, it does let designs exploit the very low power consumption and their high parallelisability.

FPGAs have a reconfigurable fabric consisting of flip-flops and look-up tables (LUTs) grouped into Configurable Logic Blocks (CLB). The difference between FPGAs results from different arrangements within the CLB and the interconnects between them. The fixed function logic blocks, such as multipliers, and embedded block RAM are 'systematically' interspersed between these. Care should be taken that designs should not be complex from the view of routing between elements. FPGAs have limited resources that facilitate long routing, and can adversely affect the maximum achievable clock frequency if not utilized well. In this work, care has been paid to minimise communication between PEs, keeping the routing complexity low.

The targeted device is from the Xilinx Virtex-5 family [11], based on a 65 nm process, it provides four 6-input LUTs, four flip-flops, multiplexers and carry chains, within a slice, with two slices making a CLB. For our context, we briefly introduce the key FPGA primitives used in this design: the Block RAM (BRAM), the FIFO and the DSP48 slice based Multipliers and Adders. These hard primitives embedded in the Virtex-5 fabric are individually able to clock at speeds greater than 500 MHz while operating at relatively low power.

Block RAM

The BRAMs are 36 bit wide 1 K deep true dual-port SRAMs, true dual-port meaning being able to independently read/write from both ports. They can be used in a variety of width-depth configurations and cascaded if required. Two adjacent BRAMs can be treated as 64 bit wide memories with no additional user logic. They can also be configured as FIFOs with relevant flags available for use.

DSP48 Slices

DSP48E blocks consists of cascadeable, 25×18 bit multipliers and 48-bit adder/subtractr/accumulator. They also allow for functions like shifting, comparisons and others to be implemented. Their ability to be cascaded allows for floating point implementations.

Since designs consist of memory elements feeding computational blocks, the relative placement between the BRAMs and the DSP48 slices is important. The targeted FPGA have these blocks arranged close to each other in special lanes within the fabric. The DSP block and BRAM proximity also true for Altera Stratix series FPGAs.

3 Algorithm

The rank-1 update scheme for matrix multiplication, illustrated in the Fig. 1, has been described here for the convenience of the reader. A rank-1 update is of the form $C \leftarrow \alpha u^T v + C$, where $\alpha = 1$ in this case. The paper partly follows the notation introduced by Dou [9] as both are variations of the same algorithm. Consider A, B and C of dimensions $M \times N$, $N \times R$ and $M \times R$, respectively. The objective is to compute $C = AB$. When a $S_i \times N$ panel of A (say, PA) and a $N \times S_j$ panel of B (say PB) are multiplied, the result is a subblock of the matrix C with dimensions $S_i \times S_j$. The outer product of a k th column vector (u_k) from PA and the k th row vector (v_k) from PB is an intermediate result, the matrix $C_{k;(S_i \times S_j)}$ and accumulation of such results with k ranging over the panel length (from 1 to N) is the required subblock of C.

For an outer product, each element of vector u_k multiplies all elements of vector v_k . Thus, $n(v_k)$ or S_j elements are re-used $n(u_k)$ or S_i times with S_j multiplications each time. That means, if one element of u_k and all the elements of v_k are available to the system then each of the products can be carried out independently. This results in the design proposed, where, broadly, an element from u_k is broadcast to all PEs, and each PE is assigned one element from v_k .

Algorithm 1 Illustration of one $PA \times PB$ computation

```

1: // For simplicity, Let number of PEs be the same as  $S_j$ , let multiplier and adder latency be 1
2: for  $q := 1$  to  $S_j$  do
3:   push  $v_1(q) \Rightarrow$  IBUF
4: end for
5: for  $p := 1$  to  $N$  do
6:   load WBUF  $\Leftarrow$  IBUF
7:   broadcast  $u_p(r) \Rightarrow$  PE $_q$ ,  $q : 1 \rightarrow S_j$ 
8:   for  $r := 1$  to  $S_i$  do
9:     for  $q := 1$  to  $S_j$  do
10:      push  $v_{p+1 \% N}(q) \Rightarrow$  IBUF
11:      // on overflow, cont. pushing next PB block's  $v(\cdot)$ 
12:      @ProcessingElement(q)
13:       $C_p(r, q) \Leftarrow u_p(r) \times v_p(q)$ 
14:      if  $p = 1$  then
15:         $C(r, q) \Leftarrow 0 + C_p(r, q)$ 
16:        push  $C(r, q) \Rightarrow C\_OUT$  // the previous  $PA \times PB$  result
17:      else
18:         $C(r, q) \Leftarrow C(r, q) + C_p(r, q)$ 
19:      end if
20:      end @ProcessingElement(q)
21:    end for
22:  end for
23: end for

```

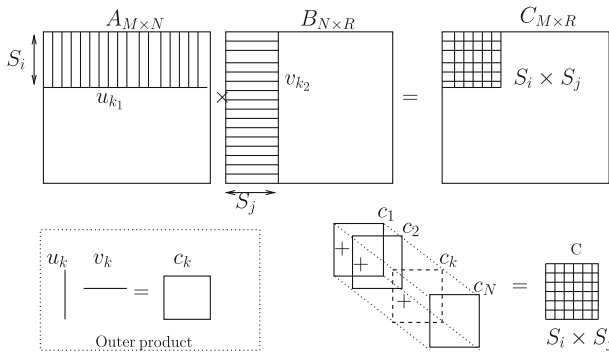


Fig. 1 The rank-1 update scheme

Given this scheme, for any candidate hardware design implementing this elementary idea, the focus now shifts to good PE design, data flow, and effective utilization of FPGA resources. In this case, the data-flow scheme and resource-aware data-path is what essentially sets our results and design apart from related work [7,9].

4 Implementation

This section describes two designs, I and II. The goals for the first design were maximising the parallelism possible, full utilisation of PEs and overlapping I/O and computation and thereby sustaining peak performance. The basic idea of re-using one element for several computations in parallel naturally leads to the broadcasting scheme in the design. The goals of design-I were met at the cost of a high I/O bandwidth and sub-optimal use of available resources. The second design addresses these limitations of the first design by a better utilization of the already existing elements. The section on design II describes its evolution in terms of more effective use of the previous data-path and resources by re-using data more efficiently.

Broadly, broadcasting elements of PA to all PEs and the streaming in of elements of PB to the prefetch registers is central to both the schemes and the relative rate at which they are streamed in, and the manner of their re-use is what essentially differentiates the two. The trade-off between the bandwidth and the local storage enables the use of this scheme in multi-FPGA accelerators.

4.1 Design-1

This design assumes $P = S_j$ and $S_i = S_j$, where P is the number of PEs. The $S_i > S_j$ case is also acceptable. Figure 2 gives an overview of this design. The following enumerated list will describe all the major labeled components, shown in Fig. 3, of the PE. It will be clear shortly that this design requires 2 words per design clock.

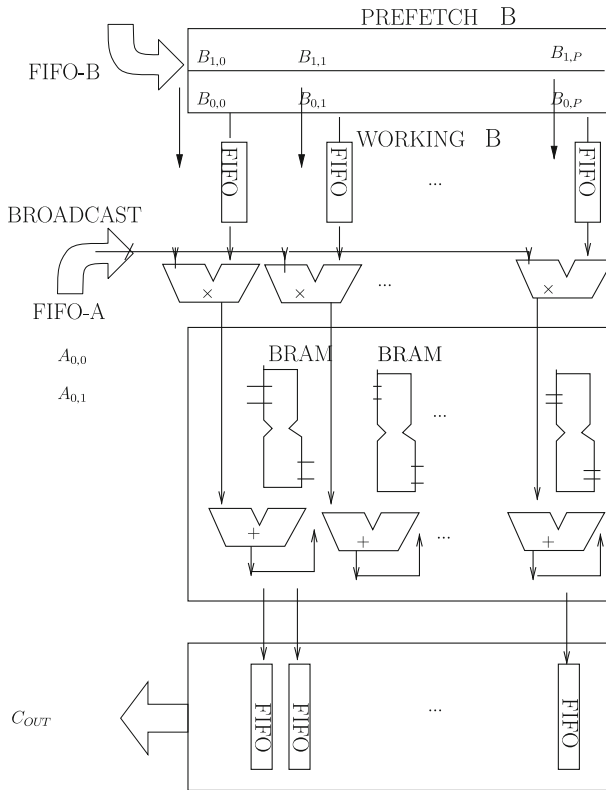


Fig. 2 Overview of Design-I

Component Description

- 1) **B Prefetch Unit:** This unit is used to prefetch S_j elements of the next row of PB while the current row is used. The input to the first of such units is a stream of elements from the matrix B, in a row major fashion. Each unit has one data input and two data outputs: a serial-shift-forward, which happens every clock cycle and a parallel-load-down which happens every P shifts (or S_i clocks if $S_i > S_j$). These P words are available at the output for at most S_i clocks, which exactly how long we need them for all the S_i multiplications.
- 2) **FIFO1:** When the B Prefetch units were connected directly to the multipliers, a severe and unexpected drop in the design frequency was observed. This drop was inferred to be due to the routing overhead in bringing the data lines from a 64 bit register to the 13 DSP48 blocks which make up a double precision multiplier. A FIFO built out of BRAM was placed in the path in order to reduce the length of the routing path thus ensuring the expected design frequency. Due to the physical proximity of the BRAMs to the DSP blocks on Xilinx FPGAs, complicated routing is avoided and also, now, the 64 bits have to route to one BRAM instead of

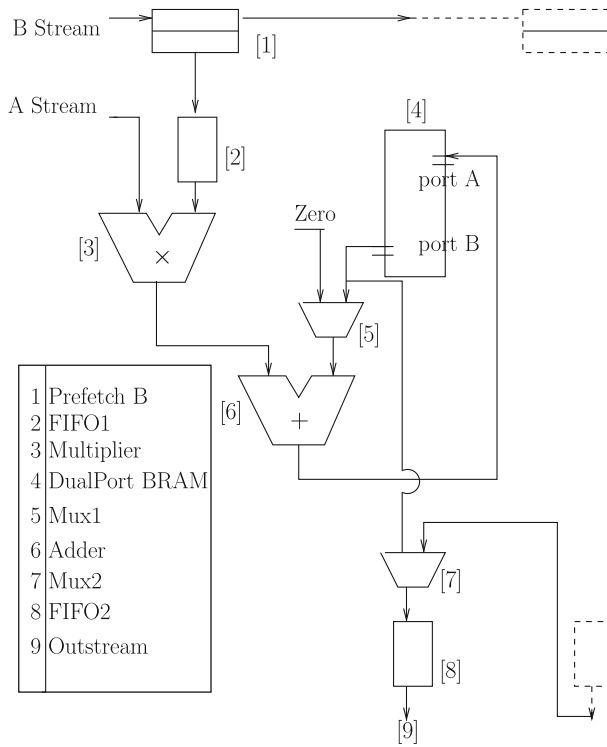


Fig. 3 PE: Design-I

13 DSP48s. This does also increase the latency, but we will address this issue in design-II.

- 3) **Multiplier:** A standard double precision floating point multiplier IP(version v3.0) from Xilinx is used for this block. A latency of 19 cycles gives it a maximum clock speed of about 431 MHz using 13 DSP48 units. One input to the multiplier comes from the prefetch unit via the FIFO1 and the other input is the element from matrix A which is broadcast to all multipliers. The output of this multiplier is one of the inputs to the adder. The recent floating-point v4.0 is superior in terms of area resource(DSP48) usage and latency, especially for Virtex-5 series, but reduces the speed, hence is not used in the design.
- 4) **Dual-Port BRAM:** This dual port blockram is used as the storage space for the accumulation step of the algorithm. The *adder* writes back to the RAM using port A and reads from the RAM using Port B. The output of port B is duplicated as the input to Mux2 as well.
- 5) **Mux1:** When a new panel-panel multiplication starts we not only need to backup the previous panel-panel multiplication result (C) but also ‘reset’ C for the fresh $C \leftarrow u_0^T v_0 + C$ operation. This mux ensures ‘0’ is added to the incoming product stream for S_i cycles, for that fresh update of C , while the previous result C in the BRAM is copied into FIFO2’s, which also needs precisely S_i clock cycles.

- 6) **Adder**: A double precision floating point adder IP (version v3.0) from Xilinx is used for this block. It receives two inputs, one from the multiplier and one from the Mux1. It writes back to blockram at the appropriate location considering its own latency.
- 7) **Mux2**: The mux is used to switch connections between the BlockRAM (Result-backup mode), and the other instances of FIFO2 (Serial-dataout mode).
- 8) **FIFO2**: In order to ensure that there is no stalling in the pipeline the result needs to be backed up. Since both the ports of the result BRAM are busy, a separate memory unit is used for the back-up, in the form of FIFO2. The final updated data of the result sub-matrix 'C' (one column of 'C' when we talk about 1 PE) will be loaded into the corresponding FIFO2s (Result-backup mode) of the PEs. When the result has been copied into the FIFO2, input of the FIFO2 gets connected to the output of FIFO2 of the previous PE, thus allowing us to take the output in a streaming fashion (Serial-dataout mode).

Data Flow-Design 1: The inputs, elements from PA and PB are streamed in column major and row major order respectively. First, one of the rows from PB (v_k) shifts into the B prefetch unit. Once a complete row is shifted-in ($S_j = P$), the 'prefetch-line' registers are full and this data is loaded down to the 'working-line' registers. In the meanwhile the prefetch-line continues to shift-in the next row from PB (v_{k+1}). At this point, when working-line is available and connected to one of the inputs of the multiplier, elements from the corresponding column (u_k) from A are broadcast to the second input of all the multipliers. After a latency period of the multiplier (say, L_m) the first result of multiplication is available at the output along with the appropriate handshaking signals which are used to trigger accumulation of the outer products at the storage area.

Once the pipeline of the multiplier has been filled it shall not be stalled since no data dependencies exist between subsequent multiplications. This allows continuous feed of the data at the maximum design frequency. The result of the addition, available after a latency period (say, L_a), is stored in the BRAM. We will see later that the pipeline may stall in one case. Zero is accumulated with the product stream during the first outer product of each new $PA \times PB$, after which the accumulation happens with the appropriate location in the BRAMs. The FIFOs responsible for the output are loaded with the values parallelly from the BRAMs once the final value of the result subblock $C_{S_i \times S_j}$ is ready. After the loading/backup is completed, these FIFOs switch modes allowing us to stream the data out in a serial fashion.

Merits and Summary: The design described above requires that $S_i \geq S_j = P$, implicitly assuming a bandwidth of 2-words per design clock cycle. PCI-e is capable of such high bandwidths, and is the norm for today's large FPGAs. The merits and demerits of this design have been summarised in the following list, details about the performance and analysis are presented in a later section.

- 1) Overlaps I/O and computation completely. Therefore, except for the initial latency period, all the processing elements (both the floating point operators) are in use all the time, thus sustaining peak performance.
- 2) The design scales seamlessly (<1% reduction in speed) as seen from Table 1.

Table 1 Timing information (post PAR)

No. PEs	SX240T(-2) [MHz]	SX95T(-3) [MHz]
1 PE	374	377
4 PEs	373	374
8 PEs	344	373
16 PE	–	373
19 PEs	–	373
20 PEs	372.8	201
40 PEs	371.7	–

- 3) Uses off-the-shelf Coregen floating point adders and multipliers, allowing for portability across technologies and generations, custom-precision option, better IEEE compliance etc.

Drawbacks:

- 1) The design requires a sustained bandwidth of 2-words per cycle corresponding to the design frequency. For a design speed of, say, 350 MHz this translates to a 5.6 GBps bandwidth requirement, which can well be provided by today's standards, but is high nonetheless.
- 2) Little to almost no flexibility in the choice of S_i and S_j which might affect the overall runtime even with the sustained peak performance.
- 3) S_i cannot be less than the latency of the floating point adder as that would result in a data dependency. This may prove to be a problem for a small number of PEs, practically though this is not a problem due to the large sizes of matrices under question.

4.2 Design-II

Design-I assumed $S_i \geq S_j = P$ and so $S_j = P$ elements of B were being re-used S_i times. The design-II ensures more re-use by allowing S_i and S_j to be greater than P , (i.e., $S_i, S_j \geq P$), however S_j needs to be a multiple of P .

As shown in the Fig. 4 the data-path is similar to that of design-I. The design actually evolved from design-I in an attempt to find a better use of the existing components, especially the dummy FIFO1 used earlier. The following enumerated list describes the major modifications as

- 1) **B Prefetch Unit:** This is the same as described earlier in design-I, however two sets of registers are not necessary. 'BRAM Cache' is made part of the working-line by using it for storage.
- 2) **B Cache:** The design-I employed a dummy FIFO1 (BRAM based) to prevent the design frequency from falling drastically. This component now assumes a central role in design-II. Each BRAM now stores S_j/P consecutive elements of a row from the chosen panel of the matrix B, hence renamed B Cache. Configured in dual port mode, this BRAM can easily implement the working-line required for the design.

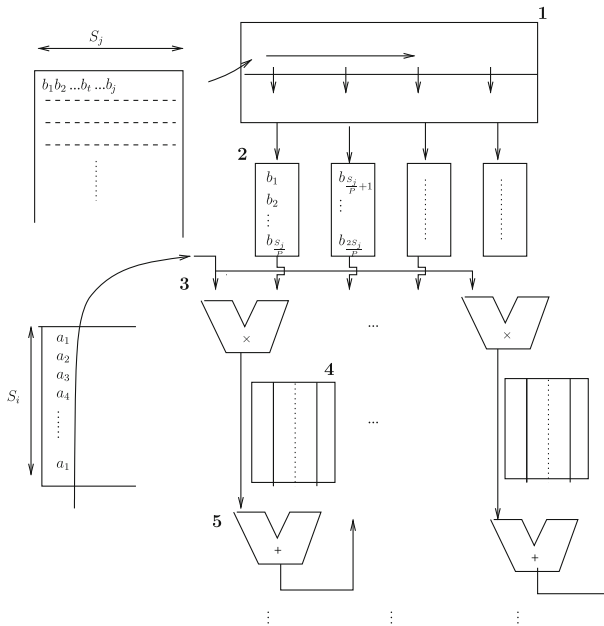


Fig. 4 PE: Design-II

3) **Simple Dual-Port RAM:** To accommodate the larger sizes of S_i and S_j , the storage area has been increased in size and logically segmented into $S_j/n(PEs)$ zones each storing results corresponding to one S_j . Overall, the storage area accounts for the storage of $S_i \times S_j$ elements of the result block. Writing to the appropriate segments is handled by address generation and control.

Design Merits:

- 1) Inherits all the merits from design I—as enumerated earlier.
- 2) Addresses all the identified drawbacks of design I, viz drastically reduces the bandwidth requirement, more flexibility in the choice of S_i and S_j and relaxes the constraint on S_i w.r.t the latency—the details of which are described in the following sections.

Data Flow: Data stream from a row of PB is fed to the prefetch unit as before, but the sequence of data is such that the i th consecutive S_j/P elements from a row with S_j elements, are loaded into the B-Cache storage corresponding to the i th PE. Thus the following sequence is observed, assuming $b_1, b_2, b_3 \dots$ are the contents of PB in row major fashion:

$$b_1, b_{S_j/P+1}, b_{2S_j/P+1}, \dots, b_2, b_{S_j/P+2}, \dots$$

One element from A is used for S_j/P cycles, where it multiplies all S_j elements of a row. During the first outer product computation, the multiplier result is accumulated

with 0 and stored in the BRAM. Thus, one outer product computation takes $\frac{S_i \times S_j}{P}$ cycles to complete, after which the elements from the next column of A are required. As a result of this, the restriction of $S_i \geq L_a$ is relaxed to $\frac{S_i \times S_j}{P} \geq L_a$. But the most important consequence of the new design is that the bandwidth requirement is considerably reduced as a result of a trade-off with local memory usage/data re-use.

Illustrative Example: Consider the product of two square matrices A and B each with dimensions 800×800 . With a design speed of 350 MHz, we consider the following two cases.

Case I: $S_i = S_j = P$; $P = 50$

In this case, the bandwidth required is 2 words per cycle which with 350 MHz means 5.6 GB/s ($= 2 \times 8 \times 350$). One outer product computation in this case takes S_i cycles and therefore one $S_i \times S_j$ subblock computation of the result takes $S_i \times 800$ cycles. For the entire matrix multiplication of $A \times B$, there are 16×16 such subblocks. Therefore, the total number of cycles for $A \times B$ computation is $S_j + S_i \times 800 \times 16 \times 16 = 10240050$.

Case II: $S_i = S_j = 400$; $P = 50$

As $\frac{S_j}{P} = 8$, we see that one word of A is required every 8 clock cycles. So, a bandwidth which gives us 2 words for every 8 cycles, or 0.25 words per cycle or 700 MB/s ($= .25 \times 8 \times 350$) will be sufficient. As for the total computation time, one can see that an $S_i \times S_j$ result subblock computation requires $\frac{S_i \times S_j}{P}$ and there are 4 such blocks here. Therefore, the total number of cycles for $A \times B$ computation is $S_j + \frac{S_i \times S_j}{P} \times 800 \times 4 = 10240400$.

Thus, design II solves the problem using significantly lower bandwidth than the first design. The increase in the cycles required for computations is because of the increased setup time.

5 Design Evaluation

Xilinx ISE 10.1sp1 and ModelSim 6.2e was used for implementation and simulation of the design, respectively. The floating point cores used were the 3.0 versions generated directly from coregen.

The most significant aspect of the design, from the Table 1, appears to be the negligible variation of the speed despite scaling up to 40 PEs, an explanation for which is offered in the comparison section. The design utilises a large number of registers in order to deeply pipeline the architecture, at 20 PEs on an SX95T there is significant decrease in the number of available registers. The drastic reduction in speed, from 373 to 201 MHz, on SX95T is attributed to the expected poor routing when the resource utilization reached $>95\%$ and makes this a corner case.

As shown in Table 2, due to abundance of resources their liberal use is justified. Appropriate pipelining, not shown in the figures, has been done in order to break the critical paths. It can also be seen from the resource usage at 40 PEs that a few more PEs can be accommodated in the SX240T.

Table 2 Resource utilization for SX95T and SX240T devices (post PAR)

No. PEs	DSP48E	FIFO	BRAM	Slice Reg	Slice LUT
1 PE	16	1	2	2511	1374
4 PE	64	4	8	10377	5451
8 PE	128	8	16	20865	10886
16 PE	256	16	32	41841	21750
20 PE	320	20	40	52329	27176
40 PE(sx240)	640	40	80	69%	36%
<i>Resources per device</i>					
Device					
SX240T	1056	516	516	149760	149760
SX95T	640	244	244	58880	58880

Table 3 Resource utilization for Virtex II Pro XC2VP100 (post PAR)

	Tot _{xc2vp100}	U _{15PE}	U _{20PE}
MULT18 × 18s	444	240	304
RAMB16s	444	90	114
Slices	44096	30218 (68%)	37023 (83%)
Speed		133.94 MHz	133.79 MHz

The design was ported to the Virtex 2 Pro XC2VP100 for the sake of comparison and as shown in Table 3, about 20 PEs can be fit with a frequency of about 134 MHz as opposed to 31 PEs and 200 MHz (synthesis), respectively by [9] (In a later usage of the same PE by one of the co-authors of [9], the actual implementation frequency was about 100 MHz [12]).

6 Performance Analysis

We present an analysis of the design parameters listed in Table 4 studying their effect on performance and the constraints they impose. All the analysis is with respect to design-II.

Table 4 List of parameters

Parameters	Meaning
β	Bandwidth in terms of the no. of words per design clock
x_a, x_b	Such that $x_a + x_b \leq \beta$
m	Total amount of local memory
n	Num. of columns of a (or rows of b)

Each element of A is used $\frac{S_j}{P}$ times, in an outer product and therefore the entire computation of the outer product takes $\frac{S_i \times S_j}{P}$ cycles. In order to overlap I/O and computation, the algorithm requires that we prefetch S_j elements of B for the next outer product. We have therefore

$$S_i + S_j \leq \frac{S_i \times S_j}{P} \times \beta \tag{1}$$

We also see that the $S_i \times S_j$ needs to be maximized here. The constraint on memory gives us Eq. 2 which on approximation gives Eq. 3

$$2S_i \times S_j + 2S_j = 2(S_i + 1) \times S_j \leq m \tag{2}$$

$$2(S_i) \times S_j \leq m \tag{3}$$

To maximize $f(S_i, S_j) = S_i \times S_j$, under the constraints Eqs. 1 and 3 we use the Lagrangian constrained optimization method

$$L(S_i, S_j, \lambda, \mu) = S_i \times S_j + \lambda \left(\beta \frac{S_i \times S_j}{P} - (S_i + S_j) \right) + \mu (m/2 - S_i \times S_j)$$

$$\frac{\partial L}{\partial S_i} = S_j - \lambda + \lambda \beta \frac{S_j}{P} - \mu S_j = 0 \tag{4}$$

$$\frac{\partial L}{\partial S_j} = S_i - \lambda + \lambda \beta \frac{S_i}{P} - \mu S_i = 0 \tag{5}$$

Equations 4 and 5 suggest $S_i = S_j$. If we substitute $S_i = S_j = S$, we get

$$\text{Maximize } S \tag{6}$$

$$S \geq \frac{2P}{\beta} \tag{7}$$

$$S \leq \sqrt{\frac{m}{2}} \tag{8}$$

The following analysis for the minimum required bandwidth demonstrates the burst-friendly nature of the design. We know that S_i words of A are required for S_j words of B within $\frac{S_j \times S_i}{P}$ cycles. Thus we get the values for min (x_a) and max (x_b) for a constant bandwidth of β . Thus we get the values for min (x_a) and max (x_b) for a constant bandwidth of β .

$$\min(x_a) = \frac{P}{S_j} \tag{9}$$

$$\max(x_b) = \beta - \frac{P}{S_j} \tag{10}$$

For the case where $S_i = S_j$ equal distribution of bandwidth is the best approach, for other cases a similar analysis results in the appropriate distribution. The availability

of more than the minimum amount of bandwidth means that the excess bandwidth can be used to transfer as much A as required in one go—further trading bandwidth with local storage. This caching also creates time during which the bandwidth can be used for other I/O, allowing for sharing the same bandwidth across multiple FPGA boards.

The design presented here is also analysed with the parameters presented in by prasanna [10]. The local storage utilised is of the order of $O(n)$ where n is the number of PE, thus achieving the optimal latency. The other important parameter mentioned is the IO-bandwidth, which by overlapping IO and computation is optimally used in this design. Due to the use of FPGA primitives we are unable to provide an accurate estimate of the throughput per area parameter.

7 Comparison

The following compares a few aspects of our designs with the recent related work. In particular we compare with Dou [9] and Prasanna [7, 10], the former of which was identified superior to other related work by Craven-2007 [5].

- **Scaling:** As reported previously [7, 10] frequency falls by about 35 and 15%, respectively by scaling to 20 PEs. Our designs show negligible (<1%) degradation in frequency up to 40 PEs. Further, the low-bandwidth requirements and burst-friendly behaviour allows design-II to scale well across multiple FPGAs due to low bandwidth requirement per FPGA.
- **Flexibility:** Dou's design [9] requires matrix subblock dimensions to be powers of 2. Prasanna supports square matrices of limited size in [10] and arbitrary size in [7]. Our designs support arbitrary matrix sizes without placing extra constraints on S_i , S_j .
- **PE/MAC:** Dou's custom MAC [9] is highly optimized for Virtex-2 Pro and may not scale across families of FPGAs. The MAC doesn't support zero and denormal numbers. Our design uses floating-point units from core-generator making the design more flexible (portable, scalable, customizable) along with better IEEE compliance. It is to be noted that these are optimized for Xilinx FPGAs, and specifically the Virtex-5.

We were able to place and route only 20 PEs on Virtex-2 Pro XC2VP100 as opposed to 31 PEs (synthesis-only) by [9] which was possible because of the custom designed MAC which use only 9 18x18 multipliers as opposed to 16 by core generator. But such custom MAC may not be appropriate in the context of, say, Virtex-5 SX240T where there are about 1200 DSP48s and the coregen floating point units are highly optimized to use them effectively.

- **I/O-Computation Overlap:** The designs use a variant of 'pipelining' or buffering for the purpose of overlapping I/O and computations as opposed to memory switching used in related works. This may be a factor in the better scaling of our designs as explained below. Memory switching requires two memory-banks to alternately feed the processing elements. This places constraints on the placement of the memory banks with respect to the processing units. In this implementation, one memory unit feeds another, except for those connected directly to the

processing units. This takes advantage of their physical proximity on the device and the better routing between BRAMs and DSP48 blocks.

8 Extending this Work

8.1 Triangular Matrix Operations

As discussed by Goto [8] the triangular matrix multiplication algorithms as well as the triangular solvers can both be implemented via what they refer to as the General purpose Block times Panel multiply (GEBP). The Triangular Matrix Multiply is cast in terms of panel-panel multiplications, the solver is cast in terms of panel-panel multiplications and a smaller triangular solver. Our implementation of *gemm* is based on panel-panel multiplication and thus can easily replace the GEBP kernel. The host system can offload all such multiplication tasks to the FPGA, leaving it free for other tasks. In the case of TRSM, the host can implement the smaller triangular solver along with required scaling, offloading the panel-multiplication task to the FPGA. Since the FPGA does not face the same limitations of the CPU cache, we do not mind the overhead of expressing TRSM in terms of the GEBP kernel.

8.2 Multi FPGA System

Practical accelerators require that the performance scales with the addition of more FPGAs and across nodes. The bandwidth available to each node however would does not scale accordingly. Here we describe possible architectures of multi-FPGA systems. The simplest system would consist of each FPGA operating independently, receiving the required bandwidth, this however will not scale to more than a few FPGAs. For further scaling, bandwidth between FPGAs can be shared for either elements from A or elements from B , one operand being broadcast to all FGPA, and one remaining independent. The FPGA can also be linked up as a linear array of PE, with the input linear shift registers, feeding the shift registers of neighbouring FPGA. This approach is limited by the increase in the latency associated with filling up the pipeline, and the asynchronous delay of the broadcast. It is estimated that this approach can scale to 2–3 FPGAs sharing the same bandwidth. A combination of the different methods described above can be used for larger systems.

8.3 Virtex-6 FPGAs

This section provides estimates/projections based on the product specifications of the Virtex-6 FPGAs. The two largest FPGA in the SX series are the SX315T and the SX475T. The SX315T can accommodate approximately 50 PEs while the SX475T can easily accommodate more than 65. Assuming a 5% degradation of clock speed, and the same frequency as the Virtex-5 designs, a performance of 35.4 and 46.0 GFLOPS can be expected, respectively.

9 Conclusion

In this paper two designs for matrix multiplication are presented which vividly demonstrate the trade-off between memory and bandwidth. The simplicity of the designs and the use of off-the-shelf floating point units from Xilinx Coregen offer easy reproduction of the design, portability across FPGA families and maintainability along with better IEEE compliance and options such as custom precision. The designs are able to sustain the peak performance, like a few other related work, achieved by use of a technique alternative to memory switching, which also has a favourable impact on routing. Our designs scale well with <1% degradation in speed and design-II further enables scaling across multiple FPGAs. For about 40 PEs, with a design frequency of 373 MHz on Virtex-5 SX240T FPGA, a sustained performance of 29.8 GFLOPS is possible with a bandwidth requirement of 750 MB/s for design-II and 5.9 GB/s for design-I.

The design can be made available upon request. Future work includes it for use with the CRL-India's supercomputer EKA.

Acknowledgments The authors acknowledge Sunil Puranik and others from CRL-India for their insights on HPC; Rahul Badghare and Pragma Sharma for their help in the timing analysis.

References

1. Baxter, R., Booth, S., Bull, M., Cawood, G., Perry, J., Parsons, M., Simpson, A., Trew, A., McCormick, A., Smart, G., Smart, R., Cantle, A., Chamberlain, R., Genest, G.: Maxwell—a 64 fpga supercomputer. In: AHS '07: Proceedings of the Second NASA/ESA Conference on Adaptive Hardware and Systems, pp. 287–294. IEEE Computer Society, Washington, DC, USA (2007)
2. Baxter, R., Booth, S., Bull, M., Cawood, G., Perry, J., Parsons, M., Simpson, A., Trew, A., McCormick, A., Smart, G., Smart, R., Cantle, A., Chamberlain, R., Genest, G.: The fpga high-performance computing alliance parallel toolkit. In: AHS '07: Proceedings of the Second NASA/ESA Conference on Adaptive Hardware and Systems, pp. 301–310. IEEE Computer Society, Washington, DC, USA (2007)
3. Underwood, K.D., Hemmert, K.S.: Closing the gap: Cpu and fpga trends in sustainable floating-point blas performance. In: FCCM, pp. 219–228. IEEE Computer Society (2004)
4. Zhuo, L., Prasanna, V.K.: High-performance designs for linear algebra operations on reconfigurable hardware. *IEEE Trans. Comput.* **57**(8), 1057–1071 (2008)
5. Craven, S., Athanas, P.: Examining the viability of fpga supercomputing. *EURASIP J. Embed. Syst.* **2007**(1), 13–13 (2007)
6. Kumar, V.B.Y., Joshi, S., Patkar, S.B., Narayanan, H.: Fpga based high performance double-precision matrix multiplication. In: VLSID '09: Proceedings of the 2009 22nd International Conference on VLSI Design, pp. 341–346. IEEE Computer Society, Washington, DC, USA (2009)
7. Zhuo, L., Prasanna, V.K.: Scalable and modular algorithms for floating-point matrix multiplication on reconfigurable computing systems. *IEEE Trans. Parallel Distrib. Syst.* **18**(4), 433–448 (2007)
8. Goto, K., van de Geijn, R.: High performance implementation of the level-3 BLAS, accepted 28 Oct 2007
9. Dou, Y., Vassiliadis, S., Kuzmanov, G.K., Gaydadjiev, G.N.: 64-bit floating-point fpga matrix multiplication. In: FPGA '05: Proceedings of the 2005 ACM/SIGDA 13th International Symposium on Field-Programmable Gate Arrays, pp. 86–95. ACM, New York, USA (2005)
10. Zhuo, L., Prasanna, V.K.: Scalable and modular algorithms for floating-point matrix multiplication on fpgas. *IPDPS* **01**, 92 (2004)
11. Xilinx Virtex-5 family User Guide
12. Kuzmanov, G., van Oijen, W.: Floating-point matrix multiplication in a polymorphic processor. In: International Conference on Field Programmable Technology (ICFPT), Dec 2007, pp. 249–252