

# EMBEDDED VERTEX SHADER IN FPGA

Lars Middendorf, Felix Mühlbauer<sup>1</sup>, Georg Umlauf<sup>2</sup>, Christophe Bobda<sup>1</sup>

<sup>1</sup>*Self-Organizing Embedded Systems Group, Department of Computer Science  
University of Kaiserslautern*

<sup>2</sup>*Geometric Algorithms Group, Department of Computer Science  
University of Kaiserslautern*

lmid@gmx.de, (muehlbauer,umlauf,bobda)@informatik.uni-kl.de

**Abstract:** Real-time 3D visualization of objects or information becomes increasingly important in everyday life e.g. in cellular phones or mobile systems. Care should be taken in the design and implementation of 3D rendering in such embedded devices like handhelds devices in order to meet the performance requirement, while maintaining power consumption low. In this work, the design and implementation of a vertex shader on a reconfigurable hardware is presented. The main focus is placed on the efficient hardware/software partitioning of the vertex shader computation, in order to maximize the performance while maintaining a high flexibility. The resulting solution must be compatible to existing vertex shaders in order to allow the large amount of existing program to be easily ported to our platform. A prototype consisting of a PowerPC, peripherals and some custom hardware modules is realized on an FPGA-board. The implementation of a point rendering shows considerable speed up compared to a pure software solution.

**Keywords:** 3D graphics, vertex shader, Direct3D, embedded systems, FPGA

## 1. INTRODUCTION

Rendering of three-dimensional objects in real-time requires much arithmetic performance. This is a problem for embedded systems that are running at low clock speed and often lacks dedicated hardware processing modules like a floating point unit (FPU). In desktop computers, the expensive arithmetic computations related to the rendering of 3D objects are done by specialized stream processing hardware in video cards. Those cards are programmable using small programs called *shaders*. The execution of shaders is the main difference to the CPU. A new instance of the program is invoked for every primitive, vertex or pixel. There are three slightly different types of shaders

for these elements. Each instance can be executed independently of the others because there is no communication possible between instances of the same type. This is advantageous when designing the hardware, because it allows the execution of an arbitrary number of instances in parallel, in order to gain the maximum computation speed. Also, pipeline technique can be used to allow some threads to feed parameters to other threads waiting for them in the pipeline. As a result the available hardware can be used more efficiently.

We developed a hardware accelerator for executing vertex shaders that is in particular useful for embedded systems, because it uses very few hardware resources, in this case FPGA slices. It is a kind of coprocessor that is directly connected to the CPU by a fast bus. The main program running on the CPU loads the shader code and all inputs into this coprocessor. While the coprocessor is running the shader, the main program accomplishes further computations in parallel until the results can be read back.

It is important to minimize the resource usage of the hardware, because the number of available slices in a FPGA is very limited. The clock speed is also very low and we have to maximize the utilization of sub-components in all cycles. The scheduling of the threads is therefore pre-calculated and stored as part of the shader code. Hence, the control-logic consists only of the program counter and a few multiplexers, that route the data flow, thus enough space is left on the device to implement floating point calculations. The multiplexer configuration is stored in a table with one row for each cycle. In this work a shader converter that generates the control table from a Direct3D9[4] vertex shader was developed. The shader converter performs the scheduling of all operation on the generated hardware unit, analyzes and optimizes the data flow and maps the calculations to operations of our ALU. Currently we can execute four threads on the ALU in parallel. This allows a speed-up of factor four compare to the software implementation of the shader.

The rest of the work is organized as follows: Section 2 provides the basics of vertex shader while section 3 introduces some work related to custom implementation of vertex shaders. Section 4 explain our implementation. A naive co-design approach is first explained, followed by a more efficient one. Also the design decisions for the hardware software partitioning are explained. The results obtained on a prototype implemented on a Xilinx Virtex 4 evaluation platform are given in 5. Finally section 6 concludes the work and provides some indication on the future directions.

## 2. VERTEX SHADER

In a rendering process, each 3D-point, also called *vertex* must traversed a set of computing stations, the *render pipeline* until the final step wher it can be drawn on the display. The stations consist of a coordinate transformation

(object  $\rightarrow$  world, world  $\rightarrow$  camera) stage, an illumination, a clipping, a projection, a scaling to screen resolution step, and finally the step to approximate the float values to integer values is performed.

To simplify spatial calculations in computer graphics *homogeneous coordinates*  $(x, y, z, w)$  are used. Transformations like *translations, rotations, scalings, shearings, projections*, etc. can be mapped to  $4 \times 4$  matrices and can be combined to only one matrix by multiplying the corresponding matrices. Thus, a transformation of a vertex by a certain list of transformations can be realized by one matrix-vector-multiplication.

For illumination calculations the dot product (scalar product) is very important, because the light intensity depends on the angle between surface normals and light sources. Normals can be transformed similar to vertices which is advantageous when filling the surface normals together with the vertices of the scene into the render pipeline to speed up processing.

In conclusion, each stage of the render pipeline executes mainly matrix and vector operations using all values which are involved like coordinates, surface normales, surface attributes, lightning parameters, etc.

There are several vertex shader versions for different hardware. We focus on implementing a subset of the smallest version 1.1 [4]. All versions use a RISC instruction set. Each instruction can read from up to three registers and write to one result register. Almost every register is 128 bits wide and stores four 32 bit floating point numbers. Hence most of the commands operate on vectors with four components. The individual components of a vector can be reordered and duplicated while reading from a register and there is a write mask for every component of the result register. This improves flexibility and allows optimizing calculations. It is for example possible to get a cross product with two instructions. Because our hardware is scalar-based, all write and swizzle-masks are free and should be used to improve performance. The shader converter analyzes the data flow for each individual component and if a result is not used, the calculation is removed on a per-component basis.

There are global and local registers. Each instance of the shader has its own set of local registers consisting of temporary and output registers. It is not allowed to write global registers which makes parallelizing possible, because there is no synchronization required and no operation in one thread depends on results calculated in another thread.

Vertex Shader 1.1 does not support jumps or subroutines. A detailed description of the instruction set can be found in the DirectX SDK [4].

### 3. RELATED WORK

Lots of work has been done already in the domain of accelerating graphics applications utilizing FPGAs in general. Some of them are listed in [8, 6]. Often

ten a combination of a desktop computer and a FPGA builds the computing unit. The need of 3D graphics visualization in embedded systems is still growing with the increasing spreading of mobile multimedia systems in everyday life like cellular phones and PDAs. Even the MPEG H.264 standard which is the video coding for next-generation multimedia involves rendering of 2D and 3D deformable mesh geometry [5].

The still continuing miniaturization has led to highly integrated chips and finally to so-called SoCs (system on chip). Here all components and peripherals are placed on a single chip like processors, hardware accelerators, bus and peripheral controllers and allow a PLB (printed circuit board) independent redesign or update of applications which is an important advantage.

Sohn et al. introduced a multimedia co-processor for mobile applications using an ARM-10 processor and fixed-point arithmetic [1]. The company Bitboys developed an vector graphics processor targeting for high-end multimedia cellular phones which is available as IP core for SoCs integration and can process SVG and OpenVG object data [7].

We are particular interested in a system in which custom hardware can co-habit with software. Also, the system should provide enough flexibility to ease the redesign and also allow a run-time adaptation, while maintaining the performance high and the power consumption low. The next sections explain our solution to this problem.

## 4. IMPLEMENTATION

Our target platform in this project was a Xilinx Virtex 4 evaluation board featuring a Virtex4-FX12 FPGA. This FPGA contains an embedded PowerPC 405 processor, on-chip memory (BlockRAM) and miscellaneous DSP functions. We use the external DDR-RAM as video frame buffer to store 3D object data. A simple system on chip with DDR-RAM controller, VGA out module and system bus needs already half of the available slices of the FPGA. Because, floating point hardware modules are expensive, we tried to avoid or reuse them as much as possible. Thus, an efficient design considering speed and chip area has to be found.

### Basic Design

In a first design a field of 32 registers combined with an adder and a multiplier unit and an instruction memory was drawn up. This *co-processor* is directly connected to the main processor via the FCM bus, which allows to extend the native PowerPC instruction set with custom instructions that are executed by a user-defined configurable hardware accelerator.

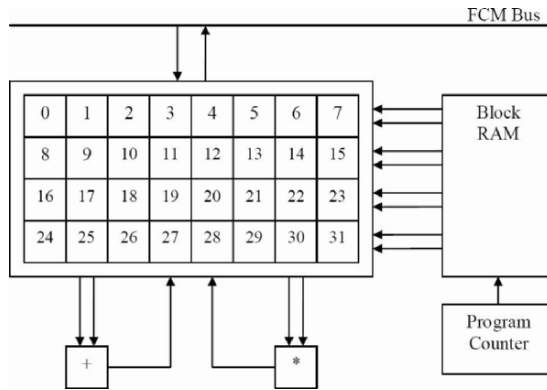


Figure 1. First design: Field of registers

The data words read from the BlockRAM (see Figure 1) specify which registers supply the input values for the arithmetic units and to which register each result should be written back.

The implementation of this design is very straight forward and also expandable for further operations like division or square root. So the two(or more) operations are executed simultaneously. Unfortunately, the design needed huge multiplexers and address decoders leading to very high resources consumption. The complete chip area was filled by this first version of the design.

### Final Design

In order to improve the first design, the idea was to exchange the expensive registers, previously realize using the on-chip available logic (LUTs) to on chip memory, namely dual ported BlockRAM. These can hold up to 512 values each (compared to 32 for the register field) but only two read respectively write accesses are possible simultaneously. Because the dot product needs 8 input values and since only one value can be provided by a BlockRAM, we duplicated data to 8 BlockRAMs in order to be able to read eight values simultaneously. In order to keep the consistency in all 8 BlockRAMs all write request are dispatched to all 8 BlockRAMs (Figure 2). In the following this BlockRAM unit is called *register array*. This new design consumes very few slices and also provides much space for provisional results. Compared to registers a memory read access takes one clock cycle and could cause additional delays in the computation. However, due to the saving of slices, an efficient design of the ALU will compensate the lost in the BlockRAM usage (Figure 3).

We next explain the components of the final design (FCM Controller and ALU) in more detail.

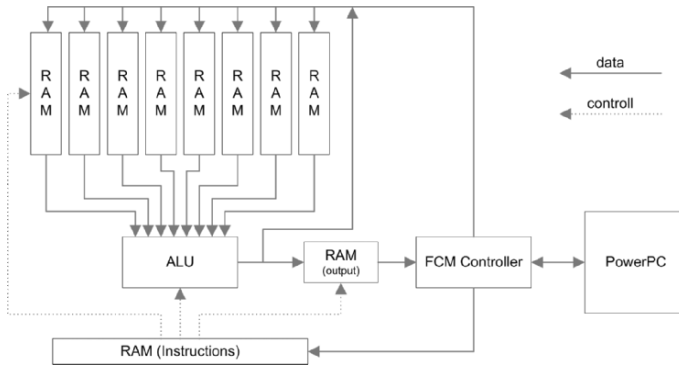


Figure 2. Final design

**FCM Controller.** The control module implements the interface to the FCM bus. The CPU is able to write to the register array, which can hold up to 128 4D vectors, and to the instruction memory with a maximum of 512 opcodes. The FCM controller is also able to read back the results from the output memory. Because the FCM instructions to handle double words provide only 5 address bits an additional 2 bit register is used to access all 512 possible memory locations.

The final result of the shader is stored in a special output RAM that is written by the ALU and read by the CPU and the FCM control module. The output RAM is addressed independently from the register array which allows copying the shader results to an arbitrary position. The additional RAM also saves one multiplexer that, otherwise, would be needed at the address lines of the register array to switch between the ALU and the control module.

**ALU.** The ALU has eight floating-point input variables, one input port that is used to select the equation, and one output port which can be used to get the result as shown in Figure 3.

In every cycle there are nine different outputs available. One of them is selected by the output multiplexer and controlled by the current instruction code. Some results are intermediate result of longer calculations and have therefore a shorter latency. Table 1 lists the instruction set of the ALU. The arithmetic units use pipelining to save hardware resources and cause delays which are also shown in the table. When generating instructions for the ALU this behaviour has to be taken into account. Still, the ALU can accept one set of values per clock cycle.

Every input variable can be pre-multiplied by -1 before it is read into the ALU. This is implemented as an exclusive-or between the sign bit of the floating-point number and the corresponding instruction bit. Because of de-

lays the `slt` instruction that is used for comparisons, minimum, maximum and absolute value, the parameters  $e, f, g, h$  must be provided twice. The `rsq` command returns a rough approximation for the inverse of square root which is much more likely to be used as the square root itself, e.g. for normalization of vectors. Usually for a more precise result one step of the Newton iteration (formula:  $x_{n+1} = \frac{1}{2}x_n(3 - x_0x_n^2)$ ) is sufficient [3].

Table 1. ALU commands

<i>command</i>	<i>result</i>	<i>delay</i>	<i>notes</i>
dot4	$a \cdot b + c \cdot d + e \cdot f + g \cdot h$	19	4D dot product
dot2	$a \cdot b + c \cdot d$	14	2D dot product
mult4	$a \cdot b \cdot c \cdot d$	18	multiplication
mult2	$a \cdot b$	9	multiplication
div	$a/b$	27	division
rsq	0x5F3759DF $-(a \gg 1)$	2	start value for newton iteration of $\frac{1}{\sqrt{a}}$
slt	if $(a \cdot b + c \cdot d < 0)$ then $(e \cdot f)$ else $(g \cdot h)$	14	input values are needed after 5 clock ticks again
int2float	$float(a)$	6	converts integer to float
float2int	$int(a)$	6	converts float to integer

**Instruction Format.** All instructions have a fixed length of 128 Bit, because of the eight input registers. The whole instruction can be divided into four words with the layout shown in Table 2. The input values are read from memory at position `src*` and inverted according to `inv*`. The ALU result is stored at index `dst`, if `we` (write enable) is set and in the extra output RAM if `oe` (output enable) is set. To avoid an extra function de-multiplexer for each ALU command a selection bit was arranged (see remaining entries in Table 2).

Table 2. Instruction format

<i>Word</i>	<i>0:8</i>	<i>9:17</i>	<i>18:26</i>	<i>27</i>	<i>28</i>	<i>29</i>	<i>30</i>	<i>31</i>
cmd0	src0	src1	dst	we	inv0	inv1	inv2	inv3
cmd1	src2	src3	out	oe	inv4	inv5	inv6	inv7
cmd2	src4	src5	-	div	rsq	slt	mult2	dot2
cmd3	src6	src7	-	f2i	i2f	mult4	dot4	-

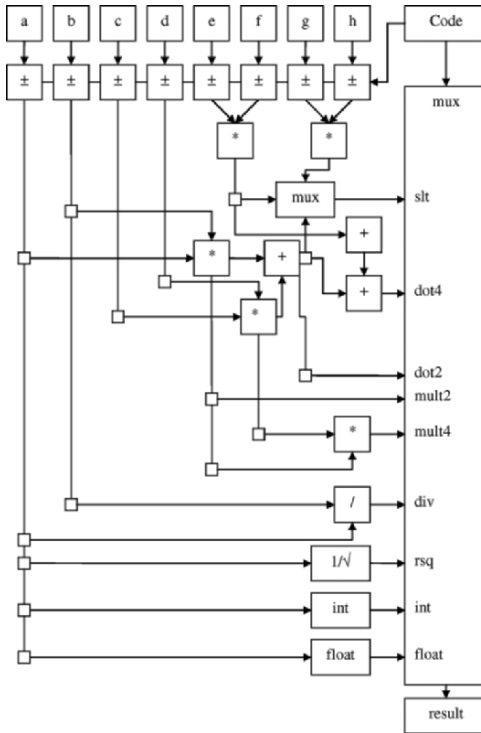


Figure 3. ALU

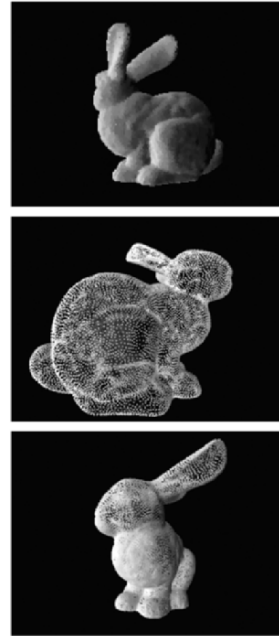


Figure 4. This bunny consist of approximately 20.000 vertices [2]

## Vertex Shader Converter

To generate the ALU opcodes a given vertex shader program is compiled with DirectX SDK[4]. Using the syntax analysis for the resulting code a data flow graph is build up, which points out the dependencies between input, provisional and output values. Now, vector operations are mapped to scalar (ALU) operations and long processing chains are move to the program start, while considering the delays caused by the arithmetic sub-units. Multiplications by -1 can be handled directly by the ALU input stage. Diversions which are not available in *Vertex Shader 1.1* and therefore are realized by multiplication with the inverse, can be processed directly by the ALU. Sometimes algebraic conversions can help to map calculations to the optimized dot product (e.g.  $(a + b)c \rightarrow ac + bc$ ). Also the usage of the `slt` command is more practical.



## 5. RESULTS

The most important disadvantage of this implementation is the limitation to one result per cycle. This means that a matrix-vector multiplication takes at least four cycles. The high latency of certain operations is not really a problem, but it is different for almost every instruction, so that it can be difficult or impossible to fully load the ALU. Because of the strict requirements, not all commands could be directly implemented in hardware. For example a single `mul` instruction that multiplies two vectors component wise can take up to four cycles. But usually the instruction is part of a more complex calculation and the shader converter can merge the previous and following calculations so that the whole block may be mapped to four larger instructions that also take four cycles.

On the other hand the ALU can calculate a 4D dot product every cycle. It has been chosen to be specially optimized, because it is a very important and often used operation. Even the most simple but useful shader does a vector-matrix multiplication that can be calculated using four 4D dot products. A large number of other instructions using only multiplications and additions can be reordered and mapped to dot products. But the most important reason for the dot product is the fact that it has only one scalar result and fits perfectly to the limited register array that can only write one result value. It is also slightly cheaper than a parallel componentwise multiplication and addition because it only needs three addition modules.

There is the possibility to output directly the intermediate 2D dot product and to skip the last addition for a lower latency. This can be useful when interpolating between two vectors. The additional multiplier outside of the dot product gives the ability to multiply four floating-point numbers in one cycle. This is important for calculating multi-linear functions that could otherwise only be achieved by a large number of cumbersome repeated high latency dot products.

This design cannot be enhanced any further. Adding another instruction type extends the multiplexer at the output of the ALU and leads to increased complexity. The timing constraints will not be met and the required clock speed of 100MHz cannot be achieved.

The new hardware component has been tested with a mesh viewer. The viewer is running on the PowerPC CPU, but the vertex shader can be calculated either in software or hardware to compare the performance. The triangles are not filled and the mesh is rendered as a point model (bunny model from [2]). We want to measure the speed of the vertex calculations and in a real application the expensive triangle filling would also be done in hardware. For each configuration 100 frames have been rendered several times with different point counts. The time spans are very precisely measured directly on the board

with a special 64 bit register that counts the CPU cycles. The vertex shader consists of six instructions that calculate the coordinates and the lighting from a directional light source.

Comparing the results both for software and hardware it is obvious that the hardware accelerated version is much faster, see Table 3. The last column of Table 3 contains the ratio between software and hardware performance. The ratio is higher when rendering more vertices, because there is a fixed overhead per frame for clearing the color and depth buffers.

Table 3. Performance results for 100 frames [sec].

<i>Vertex Count</i>	<i>Hardware</i>	<i>Software</i>	<i>Software/Hardware</i>
5000	3.607s	13.32s	3.693
10000	4.832s	24.25s	5.019
15000	6.104s	35.22s	5.770
20000	7.409s	46.26s	6.244

## 6. CONCLUSION

We have introduced a hardware accelerator for a vertex shader. Our design consumes few resources (slices) on FPGA, while supporting almost all functions of the common language for such data processing *Vertex Shader 1.1*. Compared to a software only version a significant speed advantage could be achieved. This application is suitable for the domain of embedded systems.

## REFERENCES

- [1] Jerald Yoo Ju-Ho Sohn, Jeong-Ho Woo and Hoi-Jun Yoo. Design and test of fixed-point multimedia co-processor for mobile applications. In *DATE 2006*, 2006.
- [2] Leif Kobbelt. Special effects SS05, 2005. [http://www-i8.informatik.rwth-aachen.de/old-site/teaching/ss05/praktikum\\_sfx/](http://www-i8.informatik.rwth-aachen.de/old-site/teaching/ss05/praktikum_sfx/) [date: 09/12/2006].
- [3] Chris Lomont. Fast inverse square root, 2003. <http://www.math.purdue.edu/~clomont/Math/Papers/2003/InvSqrt.pdf> [date: 2003].
- [4] Microsoft Corporation. DirectX sdk, 2006. <http://www.microsoft.com/directx> [date: 09/12/2006].
- [5] Iain Richardson. *H.264 and MPEG-4 - video compression*. Wiley, 2003.
- [6] Henry Styles and Wayne Luk. Customising graphics applications: Techniques and programming interface. In *IEEE Symposium on Field-Programmable Custom Computing Machines 2000*, 2000.
- [7] symbian.com. Bitboys introduces vector graphics processor for mobile devices at game developers conference. [www.symbian.com](http://www.symbian.com), 2005.
- [8] David Thomas and Wayne Luk. *Implementing Graphics Shaders Using FPGAs*, page 1173. Lecture Notes in Computer Science. Springer Berlin / Heidelberg, 2004.