

# A MIXED QOS SDRAM CONTROLLER FOR FPGA-BASED HIGH-END IMAGE PROCESSING

Dipl.-Ing. Sven Heithecker, Dipl.-Ing. Amilcar do Carmo Lucas, Prof. Dr.-Ing. Rolf Ernst

Technical University of Braunschweig  
Institute for Computer and Communication Network Engineering  
heithecker,lucas,ernst@ida.ing.tu-bs.de

## ABSTRACT

High-end video and multimedia processing applications today require huge amounts of memory. For cost reasons, the usage of conventional dynamic RAM (SDRAM) is preferred. However, accessing SDRAM is a complex task, especially if multi-stream access, different stream types and realtime capability are an issue. This paper describes a multi-stream SDRAM controller IP that covers different stream types and applies memory scheduling to achieve high bandwidth utilization. Two different architectures are presented and discussed, simulation results with a realistic application configuration demonstrate up to 90% of maximum memory bandwidth utilization. The scheduler IP is suitable for FPGA implementation and is flexible enough to be used in other applications.

## 1. INTRODUCTION

Dynamic RAM memories are important components in multimedia and embedded systems. For example, they are used to store large frames in multimedia and video applications. In high-end applications, such as HDTV or electronic motion pictures, bandwidth is critical. High resolution applications, widely used in motion picture and advertising industries require up to  $2K^1$  resolutions that translate to a data-rate of 2.1 Gbit per second and channel [2], [1].

The very high end of digital cinema (D-Cinema) applications have grown in importance over the last couple of years with a brilliant resolution of 4K per frame [3] and even higher resolutions are to be expected in the future. Real-time processing, such as filtering for up-/and down-scaling, color keying, compression or trick effects at this data rate and precision is beyond the scope of today's workstations and single DSP processors. The market volume for such systems is very small, so ASICs are not economically viable and therefore not an option. Therefore, we follow an FPGA-based system approach as shown in figure 1.

The main problem of DRAM architectures is the long access latency. Current DRAM architectures (DDRAM, DirectRambus, ...) reduce the overhead due to this latency using burst access to several consecutive data words in a memory row. The effect of access latency can be further reduced by exploiting the internal bank structure of a DRAM and switching to another bank while one bank is busy accessing. This technique, called interleaving, increases memory bandwidth, however without reducing the individual memory access latency. If there are enough banks, then

the burst length and, hence, the individual access latency can be reduced without significant loss in memory bandwidth. This technique [7] has already been utilized in the FPGA based memory controller of a commercial HDTV mixer application reaching on average 70% of the SDRAM peak performance with a burst length of only 4 words [8].

With even higher throughputs and more complex access patterns, the latency grows again due to larger intermediate buffers and a higher scheduler complexity. In this paper, we propose a two stage controller architecture consisting of a bank scheduler that is adapted to the memory structure and a request scheduler that reflects the access patterns. We investigate two different buffering and scheduling strategies that are simple enough to run in real-time.

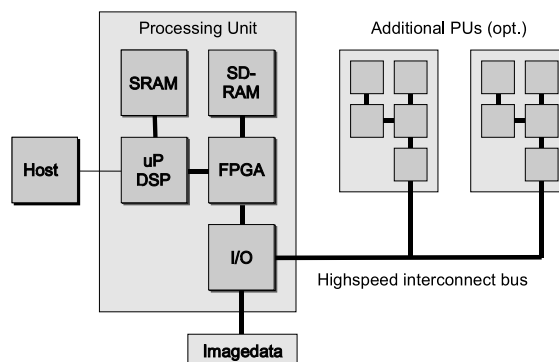


Fig. 1. Flexible image processing platform

After a concise introduction to requirements and related work the two variants of the 2-stage memory controller architecture are described in section 3. Section 4 describes the *SystemC* based evaluation environment, the experimental setup and the results. Section 5 draws a conclusion.

## 2. REQUIREMENTS AND RELATED WORK

A closer look reveals that the following distinct types of accesses have to be served by the SDRAM controller:

- Hard vs. soft realtime. While for some access sequences a maximum latency must be guaranteed, e.g. image streams which must not lose any pixel (hard real-time), other access sequences like user-interface CPU accesses might not suffer

<sup>1</sup>This resolution means 2048x1556 pixels per frame at 30 bit/pixel and 24 pictures/s resulting in 11.4 Mbytes memory requirements per frame

from a short stall (soft real-time). However, it has to be ensured that no access sequence gets completely blocked.

- Fixed address patterns, like image input/output streams or data streams generated by FPGA-based processing accelerators. To increase throughput and to compensate the burst-oriented SDRAM access, data prefetching and large FIFOs which are needed for higher throughput can be used. Since these are usually hard real-time streams, a *maximum latency* and *minimum throughput* must be guaranteed.
- Random address patterns, like DSP memory access in case of cache miss. Since these address patterns are non-deterministic, data prefetching cannot be used. Instead, to minimize stall times, these accesses should be served with *smallest possible latency*, therefore small buffers and short burst lengths should be used. Depending on the CPU task, these access sequences can be hard or soft real-time.

This results in two types of QOS: *guaranteed minimum throughput at guaranteed maximum latency* and *smallest possible latency*.

The Imagine processor [12] uses a configurable memory scheduler [4][11], optimized for the application algorithms that run on the processor. The scheduler is adapted to a specific application and does not distinguish different stream types (hard-real time vs. soft-real-time). The Prophid architecture [6] by Meerbergen et al. describes a dynamic RAM scheduler for the Prophid DSP platform that is focused on streams using large FIFO buffers and round-robin scheduling. Prabhat Mishra et al. [10] provide optimization heuristics for known memory access patterns of a single processor. None of these schedulers support vastly different access types at close to peak SDRAM bandwidth.

Closest to our work is a memory scheduler IP offered by Sonics [14] that also handles different access patterns and service levels at high average memory bandwidth. That bandwidth is similar to that of the scheduler presented here and is close to the maximum possible bandwidth as determined by the memory protocol. While high bandwidth can always be reached by a sufficiently long pipeline, latency is the key problem to reduce cache stalling. The Sonics IP is a complex, 7 stage architecture that has an inherently longer latency than the lean 2-stage architecture which we present here. Furthermore, the Sonics IP pipeline starts with a bank filter, while our findings suggest to put the bank scheduling last to reduce latency. Even though the authors of [14] highlight the importance of low latency access to avoid cache delays, only bandwidth and no latencies are published. Finally, the Sonics IP is a 100MHz ASIC implementation while the scheduler presented here reaches 140 MHz clock speed as an FPGA implementation. Again that high speed is due to the simple yet powerful architecture.

### 3. ARCHITECTURE

Our main design goal was to reach maximum memory bandwidth for multiple access streams with different access sequence types running in parallel. As explained in the introduction, the two objectives (maximum bandwidth and minimum latency) require contradictory buffer types and memory burst lengths.

Fig.2 shows the overall memory controller structure with the two scheduler variants, one with different access paths and distinct FIFOs for different priorities, and one, slightly less complicated architecture with merged access paths and FIFOs. In this paper,

we focus on the two blocks, request scheduler and bank scheduler, their respective buffers and their scheduling strategies.

#### 3.1. 1st Variant: Request Priority Scheduling

##### 3.1.1. Request Scheduler, Bank Buffer

Requests to the SDRAM controller are always made at full burst length. The request scheduler forwards requests of several inputs with different request types, one request per clock cycle, to the bank buffer. The one-per-bank bank buffer FIFOs are small and store the requests sorted by bank.

By applying a priority-based round-robin arbitration similar to [13], a minimum access service level is guaranteed. Requests from high priority inputs are scheduled out of order before requests from standard inputs. To prevent starvation, one slot is always reserved for standard priority requests. Without loss of generality, we use 2 priority levels.

##### 3.1.2. Bank Scheduler

The bank scheduler is responsible for selecting requests from the bank buffer and forwarding them to the tightly coupled access controller for execution. In order to increase bandwidth utilization, two aspects have to be considered: *bank interleaving* and *request bundling*.

*Bank interleaving* is used to hide the bank access latencies. In a common SDRAM setup, an access to a bank in auto-precharge mode and with burstlength 4 takes 4 active cycles in which data is transferred, followed by 4 (read) to 6 (write) passive cycles during which the bank cannot be accessed. However, during these passive cycles, another bank can be accessed. Thus, with at least 3 banks available, it can be guaranteed that each memory bank can be accessed once in a round without stalling. The scheduler maintains bank busy flags which describe the current state of the bank. By only selecting banks which are not busy and by applying a priority-based round-robin arbitration similar to [13], it is guaranteed that successively each bank gets scheduled (assuming a request is available).

Request bundling is used to minimize bus direction switches. On every bus direction switch, tristate cycles have to be inserted (1 for a read-write change, 1-2 for a write-read change, depending on the type of SDRAM) which can lead to a bandwidth decrease of up to 20..27% for alternating read-write accesses. Bundling requests to consecutive read or write sequence alleviates this. To prevent deadlocks, only one request of each bank is allowed in one read or write bundle.

#### 3.2. 2nd Variant: Bank Priority Scheduling

##### 3.2.1. Request Scheduler, Bank Buffer

In this variant, independent paths are provided for different request priorities. Each priority has got its own request scheduler and their own bank buffer FIFOs. Since the same arbitration scheme (without request priorities) is used as in the 1st variant, we still have a deterministic maximum latency and per-input throughput. As explained above, different bank buffer FIFOs for different priorities exist. A later multiplexor selects high priority requests if available, otherwise standard priority requests. To prevent normal requests from being deadlocked by continuous high priority

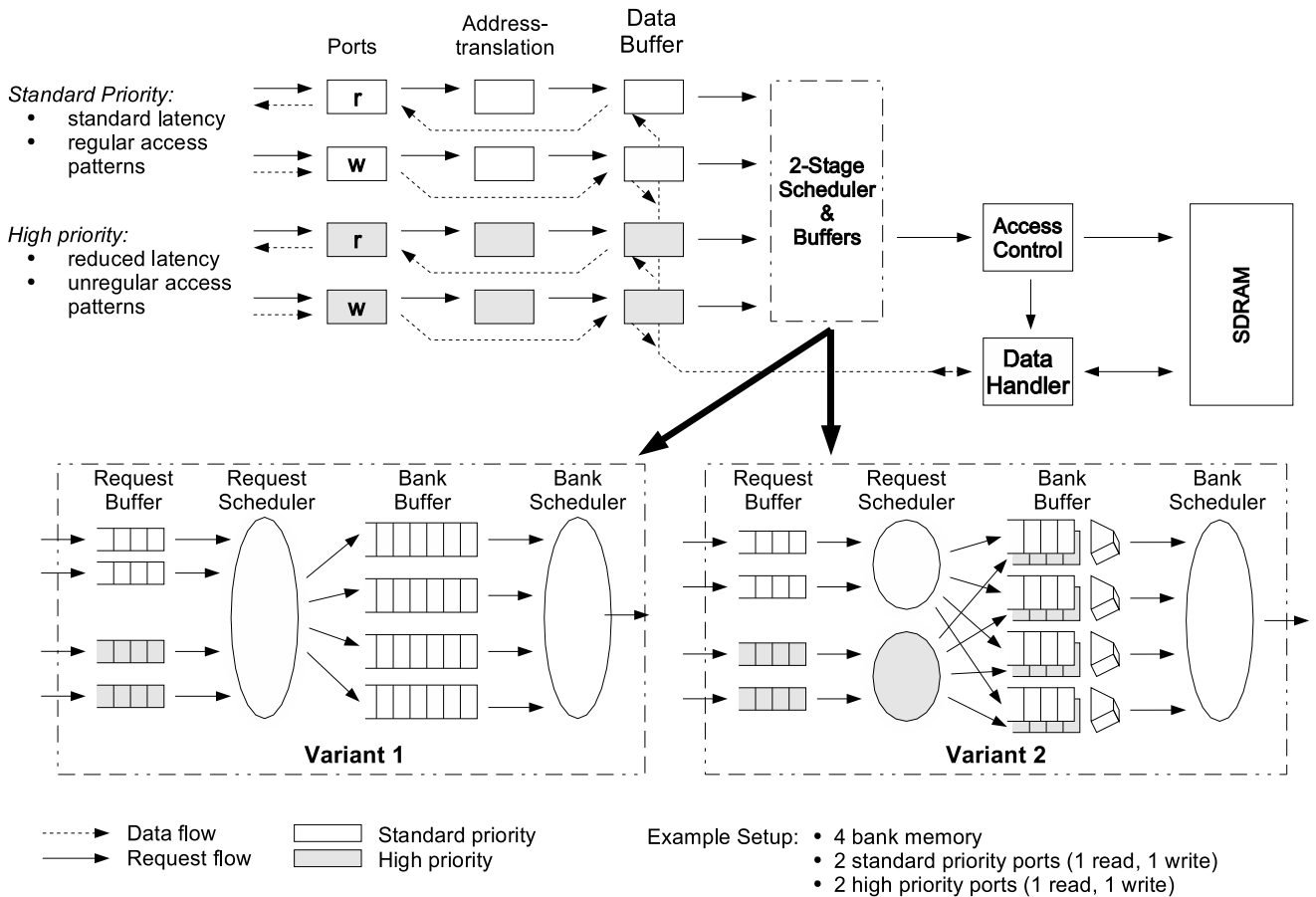


Fig. 2. SDRAM controller architecture

requests, after a sequence of high priority requests one normal request will be selected if available. This is controlled by the Bank Scheduler.

The additional complexity for the scheduler is moderate since the complexity of each scheduler is less compared to the 1st variant (less inputs and no priority selection logic). The bank buffer complexity increases noticeable, especially if more than two priorities are used. However, since it might be possible to reduce each FIFO size, the added complexity will be moderate.

### 3.2.2. Bank Scheduler

The bank scheduler works similar to the 1st variant regarding bank interleaving, request bundling and the round-robin arbitration scheme. In addition, the priority handling has moved from the request scheduler to the bank scheduler. To accomplish this, high priority requests are given precedence over standard priority requests. Only if no high priority requests can be scheduled, a standard priority request is taken. To prevent starvation of standard priority requests, after a consecutive scheduling  $n_{banks}$  high priority requests one standard priority request gets scheduled if available (in this case, the bank buffer FIFO multiplexors select a standard priority request regardless of available high priority requests as stated above).

### 3.3. Additional Modules

Since the rest of the SDRAM controller is the same in both variants, they will be explained only once.

#### 3.3.1. Access Controller

The Access Controller is responsible for SDRAM command issuing, including initialization and refresh cycles. As stated above, and contrary to most other SDRAM controllers, the SDRAM is accessed using auto-precharge mode, which means that a bank is automatically precharged after an access. This results in deterministic bank throughput and access latency which is needed for a global real-time analysis. Also, the complexity of the controller decreases due to a simpler state machine and lack of row address comparators, which makes it more suitable for FPGAs. The Access Controller works with a burst length of 4 words (8 words on DDR-SDRAM).

#### 3.3.2. Address Translation

Due to the auto precharge mode, there is no advantage in mapping access sequences to the same bank and row to exploit row buffer hits as with precharge-based SDRAM controllers [15]. In contrast, the latency increases if two consecutive requests hit the same bank.

To avoid this, request sequences are spread over all banks by permutating and *xor*-ing original address bits analogous to [16]. For many deterministic address patterns, e.g. linear address incrementation during video input, the optimum bank access sequence can be achieved by adapting the algorithm.

### 3.3.3. Data Buffer

The write-request data gets stored inside data buffers. Once the request has been scheduled, the data is transferred to the SDRAM. For read requests, an empty place in the according data buffer is assigned to the request. If the read is executed, the data from the SDRAM is transferred to that place. The bank buffer access is done via tags which are assigned to the requests. This technique allows an easy adaption of the off-chip SDRAM databus width and the usually narrower internal datapathes. It also alleviates the need to store the data inside the buffer FIFOs and the buffer-to-buffer transfers.

## 3.4. Real-Time Aspects

Since all schedulers use round-robin arbitration, and deadlocks due to continuous high-priority requests are prevented by temporarily forcing the scheduling of low-priority requests, we have a deterministic maximum latency and minimum troughput.

## 3.5. Parameters

The controller is parameterizable to serve different applications. The most important parameters are the SDRAM layout (rows, columns, banks) and timing, number of inputs including their priority and the FIFO sizes. Those parameters directly influence the latency and troughput. Increasing the number of ports increases the maximum per-port latency and decreases the per-port maximum troughput (decrease of performance). Increasing the number of banks also increases the per-bank latency and decreases per bank troughput which in turn also affects the worst case performance. However, more banks also increase the possibility that two requests get mapped to different banks which leads to a better average performance. However, if the bank access sequence can be predicted as stated in 3.3.2, then also the worst-case performance is improved.

## 3.6. Implementation

For simulation, we used *SystemC*, because, besides the high simulation speed, it allowed us to model at several abstraction levels in one language and to easily reuse available C models.

To get accurate speed data for implementation in an *Xilinx VirtexII* FPGA, we implemented critical modules (access controller, bank scheduler and data paths) in VHDL, followed by synthesis, placement and routing using Synopsys FPGA compiler 2 and Xilinx tools. The synthesis and routing results indicate a clock speed of 140 MHz.

## 4. EXPERIMENTS

To evaluate the concept, we created a system model consisting of the SDRAM, SDRAM controller as well as several clients. For the experiments, the SDRAM controller and SDRAM model were implemented cycle-accurately at transaction level, except for the

interface between controller and SDRAM which is modeled at RT level.

For the SDRAM we used the same parameter set (table 1) throughout all experiments.

parameter	value
technology	SDR
clock	142 MHz
size	64 MByte
banks	8
databus width	32 bits
burstlength	4
$t_{rcd}, t_{cl}, t_{rp}, t_{wr}$	2
$t_{rc}$	8
refresh	8192 auto-refresh per 64ms

Table 1. SDRAM parameter

## 4.1. Experiments 1a,1b

To test the impact of the bank buffer size regarding maximum troughput and latency, we created one random access, 50% read/write data stream and did several runs with different loads and bank buffer sizes. Access streams are not prioritized for this experiment. We recorded the troughput and latency. Table 2 shows the setup and the results. We made the following observations:

- the maximum reachable SDRAM troughput is 0.9 words/cycle, which means 90% bandwidth utilization.
- increasing the buffer sizes has no major effect on the troughput, but it increases the latency significantly on high loads due to FIFO fill up. On the other side, with 8 banks available, there are enough banks left to achieve close to maximum troughput even if one bank buffer is empty due to its small size. Therefore, a buffer size of 1 is sufficient.

Given the small bankbuffer FIFO size of 1, it is interesting to see if architecture variant 2 with its distinct paths for high and low priority requests provides better QOS than architecture variant 1. Therefore, we created a  $2^{nd}$  access stream to which we gave a higher priority. We tested several load combinations (which together did not exceed the maximum reachable troughput) with both architecture variants and recorded troughput and latency (table 3).

- We clearly see that effect of the priorities in architecture variant 2 is bigger than in architecture variant 1, that means that the latency of the high priority stream is lower while the latency of the low priority stream is higher. That means that even with small buffers, it is better to have different request paths and merge the requests as late as possible. Since the bank buffer FIFO size is only 1, the added complexity is moderate and thus acceptable.
- The priority architecture has no impact on the maximum troughput. Both streams reach their requested troughput.

## 4.2. Experiment 2

While the synthetic experiment 1 gives a good impression of the memory controller performance, we are also interested in the performance under practical conditions. For this purpose, we defined

applied load [w/c]	BB size 1		BB size 2		BB size 5		BB size 10	
	tp [w/c]	lat [cycles]	tp [w/c]	lat [cycles]	tp [w/c]	lat [cycles]	tp [w/c]	lat [cycles]
0.8	0.80	21.3	0.80	21.5	0.80	21.5	0.80	21.5
0.9	0.90	39.9	0.90	42.7	0.90	42.9	0.90	41.0
0.95	0.92	49.0	0.92	64.1	0.92	97.7	0.93	149.6
1.0	0.90	51.7	0.92	64.1	0.92	111.1	0.91	186.3

BB size - bank buffer FIFO size  
w/c - words per cycle  
tp - troughput  
lat - latency

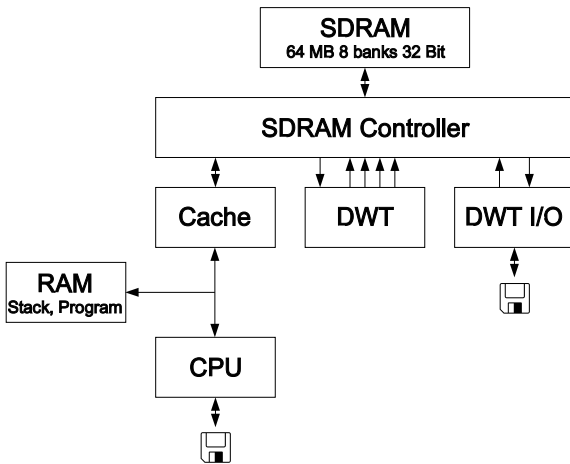
**Table 2.** Experiment 1a: Simulation setup and results

applied load lo / hi [w/c]	architecture variant 1		architecture variant 2	
	troughput lo / hi [w/c]	latency lo / hi [cycles]	troughput lo / hi [w/c]	latency lo / hi [cycles]
0.2 / 0.5	0.2 / 0.5	19.5 / 19.5	0.2 / 0.5	21 / 19
0.2 / 0.7	0.2 / 0.5	40 / 40.5	0.2 / 0.5	55 / 31.5
0.5 / 0.2	0.5 / 0.2	19.5 / 19.5	0.5 / 0.2	20 / 19
0.7 / 0.2	0.7 / 0.2	39.5 / 32.5	0.7 / 0.2	45.5 / 25.5

w/c - words per cycle  
lo - low priority stream  
hi - high priority stream

**Table 3.** Experiment 1b: Simulation setup and results

an example architecture which resembles the setup in [8] and implemented an example application. The experimental setup is shown in figure 3. It consists of a hardware implementation of a discrete wavelet transformation (DWT) algorithm and a CPU with caches.



**Fig. 3.** Experimental setup

#### 4.2.1. Discrete Wavelet Transformation

The DWT implementation is based on [9] and was implemented in *SystemC* at the RT-Level to obtain a cycle true simulation model and a basis for a precise cost and performance estimation via logic synthesis.

After implementation for a *Xilinx VirtexII* FPGA, the CAD tools reported a maximum operating frequency of 100 MHz. This frequency is not a integer fraction of the expected 140 MHz SDRAM memory interface clock. So in order to achieve maximum throughput we used asynchronous FIFOs to transfer data between the two different clock domains. The DWT input consumes one data item (2 pixels) per clock cycle, while the four outputs are idle for 32 cycles, followed by 8 cycles of data output through 4 parallel write ports.

The period of the output streams depends on the image width and on the current DWT level [9]. Because this algorithm, as well as most other image transform or filtering algorithms, has a known memory access pattern, we are able to apply pre-fetching to all input data. This is aided by the fact that the available FIFOs on the FPGA are relatively large. This way, the fixed access pattern can be served by burst memory access.

The data I/O for the DWT uses a separate port, DWT I/O. This port reads and writes a periodic pixel stream.

#### 4.2.2. CPU

For CPU simulation, we used the SimpleScalar SimSafe DLX simulator together and an adapted version of the Dinero cache simulator. The CPU has its own fast local RAM for program code and local variables (stack) and uses the SDRAM for larger data storage.

In the final system, the CPU will be replaced by a DSP, but this architecture is sufficient to generate meaningful burst cache access patterns for these experiments.

#### 4.2.3. Simulation Results

The DWT was setup for a three level 512 x 512 grayscale image, and the CPU performed a compression of a 64 x 64 color image using jpeg from the media testbench [5]. The SDRAM controller was clocked at 140 MHz, the DWT at 100 MHz. The CPU and the cache (2Kb, 4-way associative, block size equal to SDRAM burst length) were clocked at 1 GHz to compensate for the weak DLX performance (as said in Subsection 4.2.2 ideally we would use a

Nr.	DWT	CPU	bank buffer FIFO size	architecture variant 1		architecture variant 2	
				DWT Kcycles	CPU Mcycles / CPI	DWT Kcycles	CPU Mcycles / CPI
1,2,3	no	yes	1, 6/2x6, 6/2x3		5.75 (1.89)		5.75 (1.89)
4,5,6	yes	no	1, 6/2x6, 6/2x3	504		504	
7	yes	yes	1	531	6.75 (2.22)	531	6.61 (2.17)
8	yes	yes	6 / 2x6	532	6.76 (2.22)	538	6.63 (2.18)
9	yes	yes	6 / 2x3			538	6.63 (2.18)

**Table 4.** Experiment 2: Simulation setup and results

powerful DSP instead). We measured the execution time of the DWT and the CPU for several runs with the following parameters:

- Running only CPU, only DWT and both CPU and DWT.
- Bank buffer FIFO size of 1 vs. 6. For architecture 2, we set a) both high and standard priority FIFOs to 6 (2x6) and b) both high and standard priority FIFOs to 3 (2x3). The later (2x3) sizes have about the same complexity as a single FIFO of size 6.

The results are shown in table 4. We made the following observations:

- Architecture 2: setting the bank buffer FIFO size to 2x3 shows no difference compared to the 2x6 setting (nr. 2,5,8 vs. 3,6,9)
- The mutual influence of CPU activity and DWT is small, despite the relatively large load and despite the very different access patterns. This is a result of the flexible 2 stage scheduling scheme with short access bursts and small round-robin time slots leading to low latencies.

## 5. CONCLUSION

Based on earlier experience with a fixed architecture for a reconfigurable HDTV system, we presented two variants of a dynamic RAM scheduler IP that supports several concurrent access sequence types with different requirements including hard real-time periodic sequences and cache accesses with a minimum latency objective. It consists of a 2-stage scheduler, a Bank Scheduler for memory efficiency optimization and a Request Scheduler to arbitrate the access streams. We explained the two variants, the chosen IP parameters and their impact on design performance. The IP has been evaluated in a simulation environment consisting of a processor with cache, an application specific data path for wavelet coding and video I/O, all modeled in *SystemC*. In the evaluation, we demonstrated the high flexibility and efficiency of the 2-stage approach. The simulation data shows a 90% memory bandwidth utilization and adherence to access requirements for a wide range of load scenarios. The IP can easily be adapted to DDR-RAM or RAMBUS DRAM by simply exchanging the memory interface and the bank model.

## 6. REFERENCES

- [1] <http://www.discreet.com>.
- [2] <http://www.quantel.com>.
- [3] <http://www.thomsonbroadcast.com>.
- [4] BRUCE K. KHAILANY, WILLIAM J. DALLY, S. R. Imagine: Media processing with streams. *IEEE Micro* (March/April 2001), 35–46.
- [5] CHUNHO LEE, MIODRAG POTKONJAK, W. H. M.-S. Mediabench: A tool for evaluating and synthesizing multimedia and communications systems. In *International Symposium on Microarchitecture* (1997), pp. 330–335.
- [6] JEROEN A. J. LEITJEN, JEF L. VAN MEERBERGEN, A. H. T. Prophid: A heterogeneous multi-processor architecture for multimedia. In *International Conference on Computer Design* (October 1997), pp. 164–169.
- [7] KERSTEN HENRISS, ROLF ERNST, P. R. Arrangement for processing digital video signals in realtime. *International Patent # WO 01/80549 A1* (October 2001).
- [8] KERSTEN HENRISS, PETER RUEFFER, R. E. A reconfigurable hardware platform for digital real-time signal processing in television studios. In *Proceedings of IEEE Symposium on Field-Programmable Custom Computing Machines* (April 2000).
- [9] PO-CHENG WU, L.-G. C. An efficient architecture for two-dimensional discrete wavelet transform. *IEEE Transactions on circuits and systems for video technology* 11, 4 (April 2001).
- [10] PRABHAT MISHRA, PETER GRUN, N. D. Processor-memory co-exploration driven by a memory-aware architecture description language. In *14th International Conference on VLSI Design* (Jan 2001).
- [11] SCOTT RIXNER, WILLIAM J. DALLY, U. J. K. Memory access scheduling. In *International Symposium on Computer Architecture* (2000), pp. 128–138.
- [12] UJVAL J. KAPASI, WILLIAM J. DALLY, S. R. The imagine stream processor. In *International Conference on Computer Design* (2002).
- [13] WEBER, M. Arbiters: Design ideas and coding styles. In *SNUG, Boston* (2001).
- [14] WEBER, W.-D. Sonics memmax memory scheduler.
- [15] WEI-FEN LIN, STEVEN K. REINHARDT, D. B. Designing a modern memory hierarchy with hardware prefetching. *IEEE Transactions on Computers* 50, 11 (November 2001), 1202–1218.
- [16] ZHAO ZHANG, ZHICHUN ZHU, X. Z. A permutation-based page interleaving scheme to reduce row-buffer conflicts and exploit data locality. In *International Symposium on Microarchitecture* (2000), pp. 32–41.