

# A Novel Machine Paradigm to Accelerate Scientific Computing

Reiner W. Hartenstein, Jürgen Becker, Rainer Kress, Helmut Reinig

University of Kaiserslautern

Erwin-Schrödinger-Straße, D-67663 Kaiserslautern, Germany

Fax: ++49 631 205 2640, email: abakus@informatik.uni-kl.de

## Abstract

*A large potential to accelerate scientific algorithms lies in a hardware structure, which is a direct match to the structure of the algorithm. The MoM-3 is a universal accelerating co-processor, which can be configured to implement performance critical parts of an algorithm directly in hardware. The MoM-3 architecture supports fine grain parallelism by use of a highly reconfigurable ALU (rALU), and the increase of the effective memory bandwidth by avoiding addressing overhead. This is achieved by novel hardware concepts based on a new machine paradigm, which uses (one or more) data sequencers instead of an instruction sequencer. The operation principles of the MoM-3 are illustrated by a computationally intensive FIR filtering example. Compared a high-end SPARC station, actual speed-up factors between 12 and 62 (corresponding to technologically normalized speed-ups between 24 and 189) have been obtained for FIR filter and JPEG image compression algorithms.*

## Keywords

- novel machine paradigm, universal accelerator, problem solving environment
- data address generators, fine grain parallelism, numerical algorithms
- FIR filters, JPEG algorithm

## 1. Introduction

Scientific computing provides the greatest challenges to modern workstations and even supercomputers. A lot of different computer architectures have been presented, which take into account characteristics, that are common to many scientific algorithms. Vector processors [3] speed up operations on large arrays of data by the use of pipelining techniques. Parallel multi-processor architectures [11] benefit from the fact, that many operations on large amounts of data are independent from each other. This allows to distribute these operations onto different processors (or processing elements) and execute them in parallel. But all of these architectures basically still follow the von Neumann machine paradigm with a fixed instruction set, where the sequence of instructions triggers the accesses to data in memory and the data manipulations.

The Map-oriented Machine 3 (MoM-3) is an architecture based on the Xputer machine paradigm [2]. Instead of a hardwired ALU with a fixed instruction set, an Xputer has a reconfigurable ALU based on field-programmable devices. All data manipulations, which are performed in the loop bodies of an algorithm, are combined to a set of compound operators. Each compound operator matches a single loop body and takes several data words as input to produce a number of resulting data words. The compound operators are configured into the



field-programmable devices. After configuration, an Xputer's "instruction set" consists only of the compound operators as they are required by the algorithm actually running on the Xputer. The combination of several operations of a high level language description to one compound operator allows to introduce pipelining and fine grain parallelism to a larger extend, as can be done in fixed instruction set processors. E.g. intermediate results can be passed along in the pipeline, instead of writing them back to the register file after every instruction. Since many scientific algorithms compute array indices in several nested loops, the sequence of data addresses in a program trace shows a regular pattern. This leads to the idea to have complex address generators compute such address sequences from a small parameter set, which describes the address pattern. And instead of an instruction sequencer as a centralized control to trigger the operations in the reconfigurable ALU, the address generators themselves serve as a decentralized control. They automatically activate the appropriate compound operator, each time a new set of input data is fetched from memory and the previous results have been written back. This so-called data sequencing mechanism directly matches the loop structure of the algorithm, where the index computations serve as a means to provide the right data to ever the same operations.

Field-programmable devices available commercially today are not well suited to implement arithmetic operations on wordlengths of 32 bits or more. Therefore, we developed our own architecture, the so-called rDPA (reconfigurable datapath architecture). Compared to other SRAM-based architectures like XILINX 4000 series [8], AT&T ORCA [7], or Altera's FLEX family [6], the rDPA is quite coarse grain. It provides a small array of so-called datapath units (DPU), where each can be configured to implement any operator from C programming language, as well as some others, that are stored in an extensible library. The wordlength of a single DPU is 32 bits. The configuration code for the rDPA and the address generators is derived from C programming language, without requiring further user interaction or compiler directives to obtain satisfactory results.

The following section introduces the hardware structure of the MoM-3, including the address generators and the rDPA. Afterwards, the features of the C compiler for the MoM-3 are outlined. The fourth section explains the way algorithms are executed on the MoM-3 by means of an example. The final sections provide a performance comparison and conclude the paper.

## 2. The MoM-3 Hardware

The MoM-3 architecture is based on the Xputer machine paradigm [2]. MoM is an acronym for Map oriented Machine, because the data memory is organized in two dimensions like a map. Instead of the von Neumann-like tight coupling of the instruction set to the data manipulations performed, an Xputer shows only a loose coupling between the sequencing mechanism and the ALU. That's why an Xputer efficiently supports a reconfigurable ALU (rALU). The rALU contains compound operators which produce a number of results from a number of input data (figure 1). All input and output data to the compound operators is stored in so-called scan windows (SW). A scan window is a programming model of a sliding window, which moves across data memory under control of a data address generator. All data in the scan windows can be accessed in parallel by the rALU operator. The rALU operators are activated every time a new set of input data is available in the scan window. This so-called data sequencing mechanism is deterministic, because the input data is addressed by Generic Address Generators (GAGs). They compute a deterministic sequence of data addresses from a set of algorithm-dependent parameters. An Xputer Data Sequencer contains several Generic Address Generators running in parallel, to be able to efficiently cope with multiple data sources and destinations for one set of compound operators.



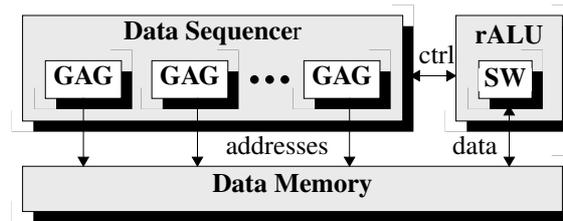


Figure 1. Block Diagram of an Xputer

In the MoM-3, the Data Sequencer is distributed across several computational modules (C-Modules), as can be seen in figure 2. The MoM-3 includes up to seven C-Modules. Each C-Module consists of a Generic Address Generator (GAG), an rALU subnet and at least two megabytes of local memory. All C-Modules can operate in parallel when each Generic Address Generator accesses data in its local memory only. Apart from the local memory access, two means of global communication are available. First, the rALU subnets can exchange internal results with their neighbours without disturbing parallel activity. Second, the Generic Address Generators can access data memory and rALU subnets on other C-Modules using the global MoMbus. This can be done only sequentially, of course. The global MoMbus is used by the MoM-3 controller (M3C) as well, to reconfigure the address generators and the rALU whenever necessary. The MoM-3 controller is the coordinating instance of the Data Sequencer, which ensures a well-defined parallel activity of the Generic Address Generators by selecting the appropriate parameter sets for configuration. Via the MoMbus-VMEbus interface, the host CPU has access to all memory areas of the MoM-3 and vice versa. The Generic Address Generators have DMA capability to the host's main memory, reducing the time to download a part of an application for execution on the MoM-3. The host CPU is responsible for all disk I/O, user interaction and memory allocation, so that the MoM-3 completely uses the functionality of the host's operating system.

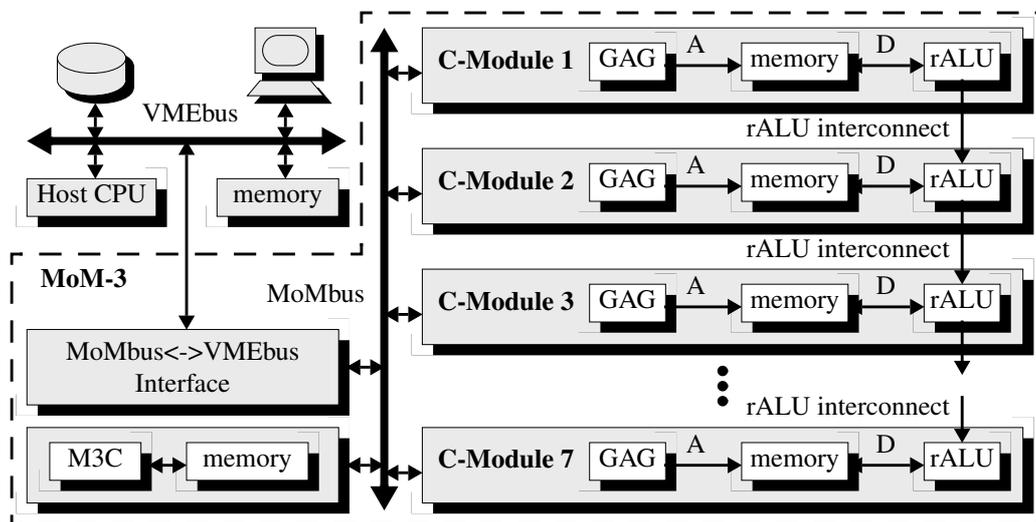
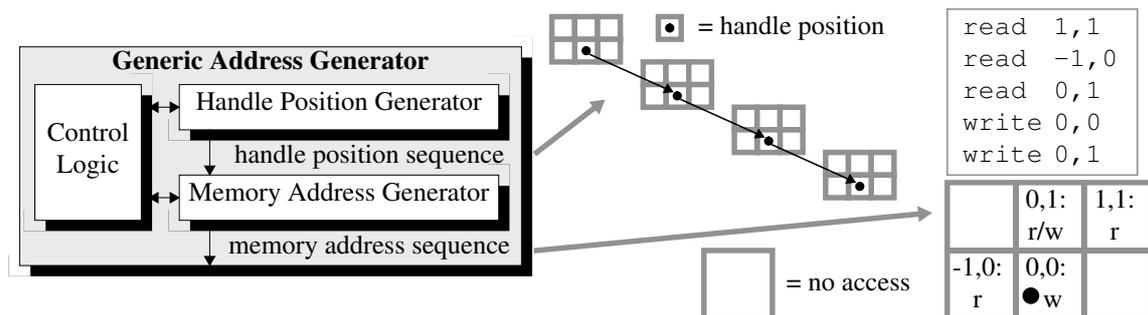


Figure 2. The MoM-3 Hardware

The printed circuit boards with the MoM-3 controller and the C-Modules reside in their own 19 inch rack with the MoMbus backplane. The only connection to the host computer is through the bus interface. That way, the MoM-3 could be interfaced to a different host with minimal effort. The part of the bus interface that is injected into the host's backplane would have to be redesigned according to the host's bus specifications. All other parts could remain the same.

## 2.1 The Data Sequencer

The MoM-3 Data Sequencer consists of up to seven Generic Address Generators and the MoM-3 controller. Each Generic Address Generator controls one scan window. It operates in a two-stage pipeline. The first stage computes handle positions for the scan windows (Handle Position Generator). A handle position consists of two 16-bit values for the two dimensions of the data memory. The sequence of handle positions describes how the corresponding scan window is moved across the data memory (figure 3). Such a sequence of handle positions is called scan pattern. A Handle Position Generator can produce a scan pattern corresponding to four nested loops at the most. It is programmed by specifying a set of parameters, such like starting position, increment value, and end position of a loop, each both for the x and y dimension of the data memory.



**Figure 3. Generic Address Generator**

The second pipeline stage computes a sequence of offsets to the handle positions, to obtain the effective memory addresses for the data transfers between the scan window / rALU and the data memory. Therefore this stage is called Memory Address Generator. The range of offsets may be  $-32$  to  $+31$  in both dimensions of the data memory. The sequence of offsets may be programmed arbitrarily, but at most 250 references to the data memory can be made from one handle position. The offset sequence may be varied at the beginning and at the end of the loops of the Handle Address Generator. This allows to perform extra memory accesses to fill a pipeline in the rALU at the beginning of a loop, as well as additional write operations to flush a pipeline at the end. The address parts of the two dimensions may be combined to a real linear memory address in four ways: one row of two-dimensional memory may consume an address space of 10, 12, 14, or 16 bits. This allows to adjust the “size” of the data memory to the size of the processed data, to reduce wastage of address space. The software environment contains a memory management module, which makes use of the unused area, from the end of a data row to the next appropriate power-of-two boundary. If the algorithm uses another array of data, which fits into the remaining space, such two data arrays are placed in memory, one after the other. After the concatenation of the two address parts, a 32-bit base address is added, to obtain the effective memory address. The base address typically is the starting address of the data array referenced by this Generic Address Generator, as determined by the memory management software. The Memory Address Generator itself is a three stage pipeline. The first stage performs the offset calculations and the combination of address parts to a linear address. The second stage adds the 32-bit base address, and the third stage handles the bus protocol for data transfers.

The Generic Address Generators run in parallel. They synchronize their scan patterns through the activations of the rALU. All scan patterns proceed until either they detect an explicit synchronization point in the offset sequence of the Memory Address Generator, or they

are blocked by a write operation to memory, waiting for the rALU to compute the desired result.

The hardware controlled generation of long address sequences allows to access data in memory every 120 ns, using 70 ns memory devices and a conventional, non-interleaved memory organization. A further speed-up of memory accesses could be obtained by interleaved memory banks and by the introduction of static RAM caches, like in conventional computers. With caches, it would make sense to have a programmable cache controller in the MoM-3. Since the hardware generated address sequences are deterministic, a compiler could compute a cache update strategy at compile time, which would result in the minimum number of cache misses. This should provide better performance than the probabilistic cache update strategies usually applied.

## 2.2 Reconfigurable ALU

One rALU-subnet of the MoM-3 is shown in figure 4. It contains an rDPA array (reconfigurable datapath architecture) made of eight rDPA chips, arranged in a two by four array. Each rDPA chip contains an array of four by four datapath units (DPU). A datapath unit can imple-

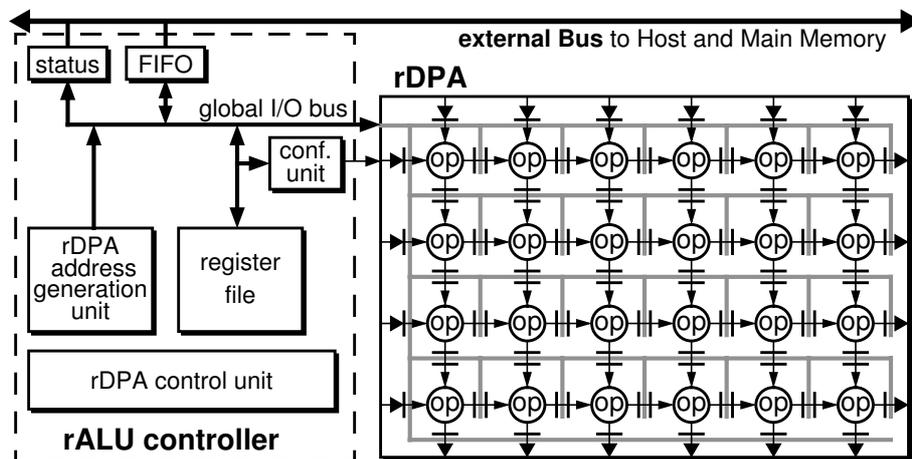
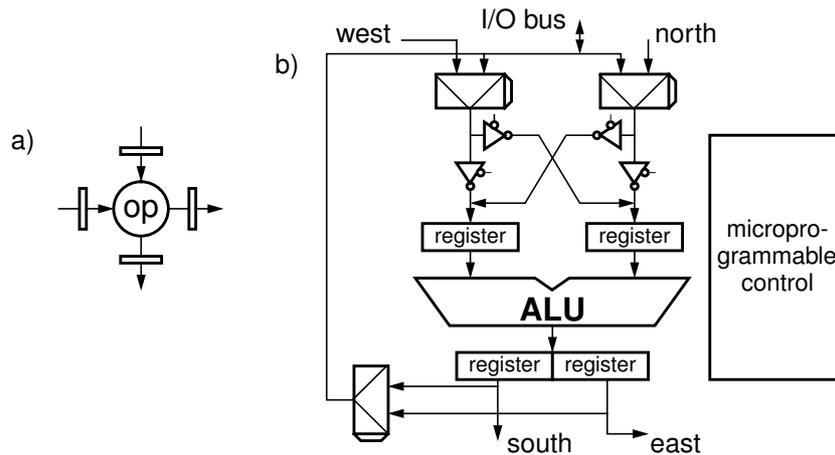


Figure 4. One subnet of the reconfigurable data-driven ALU

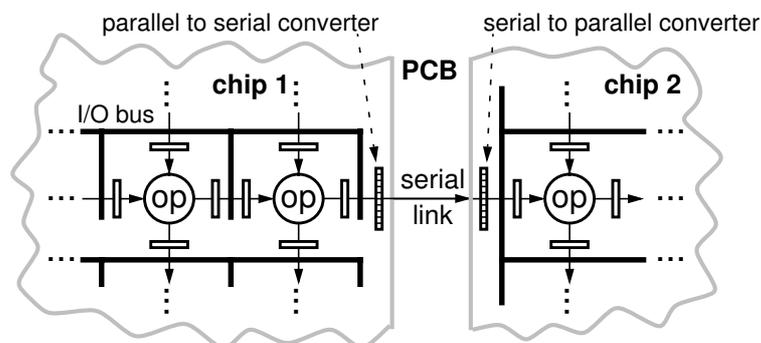
ment any unary or binary operator from C language on integer or fixed-point data types up to 32 bits length. Multiplication and division are performed by means of a microprogram, whereas all other operators can be executed directly (figure 5). Additionally, operators like multiplication with accumulation are available in the library of configurable operators. Floating-point arithmetic cannot be done in a single datapath unit. But it is possible to build a pipelined floating-point operator using several adjacent datapath units, e.g. two DPUs for normalization and denormalization and at least one further DPU for the operation. This requires a sequential transfer of multiple intermediate values along the DPU interconnect, which are used as operands to a pipelined floating-point operation. This can be done, because each datapath unit consists of a conventional ALU, a microprogram memory and a sequencer for horizontal microcode (figure 5b). The microprogram for a floating-point operator simply contains multiple data transfers from the preceding datapath unit before the operation starts. Each datapath unit has two input registers, one at the north and one at the west, and two output registers, one at the south and one at the east. All registers can be used for intermediate results throughout the computation. The registers in the datapath units in fact are a distributed implementation of the scan window of the Xputer paradigm. The datapath units can read 32-bit intermediate results from their neighbours to the west and to the north and pass 32-bit results to

their neighbours in south and/or east direction. Along with the data, two condition bits are transferred. Conditional expressions can be mapped onto the rDPA array by computing both results, for the true and the false case of the condition, in different datapath units. Additional datapath units are required to evaluate the condition and to set the condition bits accordingly. Only the result, which corresponds to the evaluation of the condition, as specified in the condition bits, is passed on to the following DPUs as a result of the conditional expression. A global I/O bus allows input and output data to be written directly to a datapath unit without passing them from neighbour to neighbour. The datapath units operate data-driven. The operator is applied each time, new input data is available either from the bus or from a neighbouring datapath unit. This decentralized control simplifies pipelining of operations, when each takes a different time to execute.



**Figure 5. Datapath unit (DPU): a) symbolic representation; b) block diagram**

The array of datapath units expands across rDPA chip boundaries in a way that is completely transparent to the compiler software. To overcome pinout restrictions, the neighbour to neighbour connections are reduced to serial links with appropriate converters at the chip boundaries (figure 6). A pair of converters behaves like a special-purpose datapath unit, restricted to routing operations in the appropriate direction. To the programming software, the serial links combined with the preceding DPU appear to be a DPU with an increased latency for data transfers. Although pipelining drastically reduces the penalty of the conversion from 32 bits to 2 bits, this still may turn out a bottleneck in the current rDPA prototype with some algorithms. With state of the art IC technology and packages larger than 144 pins, we could build an rDPA with 66 MHz serial links (instead of 33 MHz) and four bits wide communications channels. This four-fold speed-up would overcome the shortfalls of the current prototype.



**Figure 6. rDPA array expansion across chip boundaries**

In addition to the rDPA array, an rALU controller circuit interfaces to the MoMbus. It provides a register file as a kind of cache for frequently used data, and controls the data transfers on the global I/O bus. The rDPA, in combination with the rALU controller, supports pipelining on operator level (e.g. floating point operations are implemented as a pipeline across several datapath units), pipeline chaining [3], and pipelining of loop bodies, as shown in the example in section 4.2. That way, the compound operators of subsequent loop iterations are computed in parallel in a pipeline. Each of the loop iterations is finished to a different degree, depending on its progress in the pipeline.

### 3. MoM-3 C Compiler

The C compiler for the MoM-3 takes an almost complete subset of ANSI C as input. Only constructs, which would require a dynamic memory management to be run on the MoM-3 are excluded. These are pointers, operating system calls and recursive functions. Since the host's operating system takes care of memory management and I/O, the software parts written for execution on the MoM-3 do not need such constructs. Especially for scientific computing, the restrictions of the C subset are not that important, since FORTRAN 77 lacks the same features and is most popular in scientific computing. There are no extensions to C language or compiler directives required to produce configuration code for the MoM-3. The compiler computes the parameter sets for the Generic Address Generators, the configuration code for the rDPA arrays, and the reconfiguration instructions for the MoM-3 controller, without further user interaction. First, the compiler performs a data and control flow analysis. The data structure obtained allows restructurations to perform parallelizations like those done by compilers for supercomputers. These include vectorization, loop unrolling, and pipelining on expression level. The next step performs a re-ordering of data accesses to obtain access sequences, which can be mapped well to the parameters of the Generic Address Generators. Therefore, the compiler generates a so-called data map, which describes the way the input data has to be distributed in the data memory to obtain optimized hardware generated address sequences. After a final automatic partitioning, data manipulations are translated to a rALU description, and the control structures of the algorithm are transformed to Data Sequencer code. An assembler for the Data Sequencer translates the Data Sequencer code to parameter sets for the Generic Address Generators and a reconfiguration scheme for the MoM-3 controller.

The rALU description is parsed by the ALE-X assembler (arithmetic and logic expression language for Xputers). It generates a mapping of operators to datapath units, merges DPU operators where possible, and computes a conflict-free I/O schedule, which matches the operators' speed, to keep the datapath units as busy as possible. The separation of the rALU code generation from the rest of the compiler allows to use other devices than the rDPA to build a MoM-3 rALU, without having to rewrite the C compiler with all its optimizations. For a new rALU, only an assembler generating rALU configuration code from ALE-X specifications has to be written. This is not too difficult, because ALE-X describes only arithmetic and logic expressions on the scan windows' contents.

A more detailed description of the MoM-3 programming environment can be found in [10]. The most important benefit of the MoM-3 C compiler is the fully automatic code generation. The programmer neither has to be a hardware expert, to guide hardware synthesis with appropriate constraints or compiler directives, nor has he to interfere with different stages of the compilation to get reasonable results. All parallelizations are done from a C source without requiring special constructs like vector operations [9] or parallel flow control [4] to detect parallelism, in contrast to many parallelizing compilers for supercomputers. This allows to compile the same algorithm specification for execution on a conventional computer and on the MoM-3.



## 4. Example: two dimensional FIR filter

The operation of the MoM-3 and the way an algorithm is adapted to the non-von Neumann architecture of the MoM-3 is illustrated with a two-dimensional FIR filter. Although this is not one of the grand challenge problems, like fluid dynamics or geometric modelling, it has many characteristics in common with many typical problems of scientific computing. A large amount of data has to be processed with ever the same sequence of operations, and the algorithm is organized in nested loops. Furthermore it is easy to comprehend, to allow to concentrate on the operation principles of the MoM-3, and how they correspond to the structure of the more important grand challenge problems. The two dimensional FIR (finite impulse response) filter is one whose impulse response processes only a finite number of nonzero samples. The equation for a general two dimensional FIR filter is

$$y_{nm} = \sum_i \sum_j k_{ij} \cdot x_{n-i, m-j} \quad (1)$$

### 4.1 Straightforward Implementation

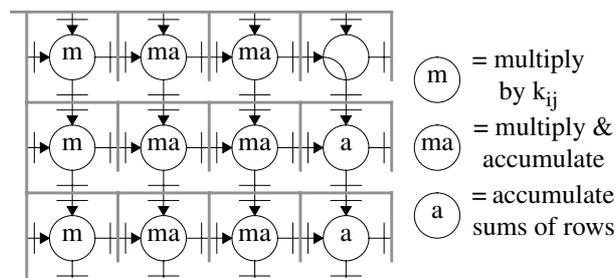
In this example a two dimensional filtering of second order is shown, that is indices  $i$  and  $j$  have a range of zero to two. Given the order of a filter and its weights  $k_{ij}$ , the most efficient implementation on a von Neumann computer is to unroll the two loops of the filter kernel, to create a large expression computing the new value of  $y_{nm}$  (see figure 7). The weights  $k_{ij}$  now are constants, so that the compiler can optimize the multiplications by replacing them with a matching sequence of additions and shift operations.

```

For (n = 0; n < maxn; n++) {
  For (m = 0; m < maxm; m++) {
    y[n][m] = k00*x[n][m]    + k01*x[n][m+1]    + k02*x[n][m+2]
              + k10*x[n+1][m] + k11*x[n+1][m+1] + k12*x[n+1][m+2]
              + k20*x[n+2][m] + k21*x[n+2][m+1] + k22*x[n+2][m+2];
  }
}
    
```

**Figure 7. Source code of a 3 \* 3 two-dimensional FIR filter**

The same source code produces an efficient implementation on the MoM-3 as well. The C compiler vectorizes the expression and performs loop unrolling to the extend of the capacity of the rDPA array. With multiplication and accumulation being a valid operator, the expression to compute  $new[y][x]$  takes three by four datapath units. The multiplications and summation of the rows are done in a three by three array of multipliers and multiplier-accumulators. A column of three routers/adders combines the sums of rows to the complete result (figure 8).



**Figure 8. rALU operator for a single 3 \* 3 2-D FIR filter operation**

Taking the chip boundaries into account (with the costly parallel-to-serial converters) eight of these expressions can be computed in parallel, each in its own rDPA chip. The configuration for the whole rDPA array simply is a concatenation of two by four of these operators, leaving a spare row to adjust to the chip boundaries. The scan window made of the DPU registers now consists of ten words in x direction and three words in y direction, instead of the three by three words of a single filter operator.

During vectorization, the compiler performs a partial unrolling of the inner loop, so that eight consecutive loop iterations are executed in parallel, and the loop index x is incremented by eight. The number of vectorizable loop iterations is taken from the capacity of the rDPA array, which can implement eight  $3 * 3$  FIR filter operators in parallel. Furthermore, the C compiler knows from the hardware configuration file the number of available C-Modules and splits the input array into a corresponding number of stripes. The required overlap of two can be computed from the data dependency analysis. That way all C-Modules may operate in parallel to speed-up computation by a factor of seven, at the most. The ALE-X assembler detects, that most of the input values to the eight parallel expressions are used several times in different multiplications, and stores these values in the register file for quick access. Furthermore, the two input values in x direction, that overlap with the next loop iteration are read from the register file as long as the pipeline is full throughout the inner loop. The compiler generates a Memory Address Generator configuration that first fills the pipeline by reading ten columns of three input values. Within the inner loop, only eight columns have to be read from memory. The overlapping two columns are taken from the register file.

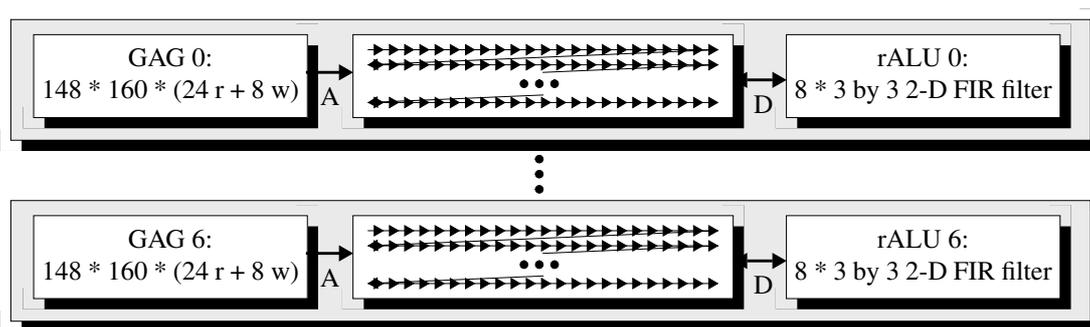


Figure 10. Compiler generated scan pattern for a  $3 * 3$  2-D FIR filter

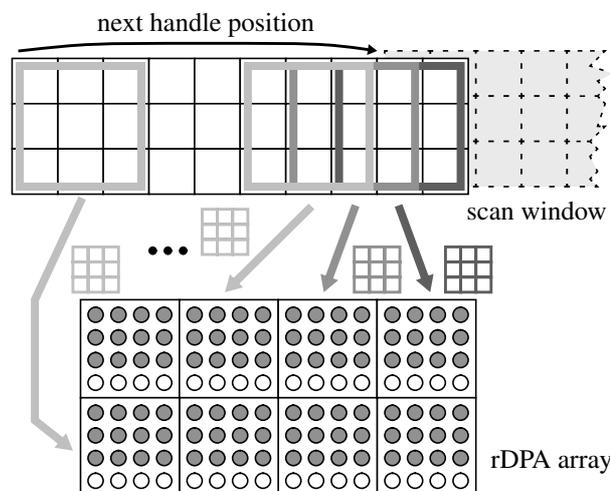


Figure 9. Scan window and rALU configuration of a vectorized  $3 * 3$  2-D FIR filter

The only rearrangements to the data are the distribution of the stripes onto the C-Modules. Therefore the data map is quite straightforward. The input values are stored row by row, each strip of rows in a separate C-Module memory. All Generic Address Generators compute the same address sequence. Scanning the memory with a three by ten window row by row, at each scan window position ten input values are read from memory and eight values are written as results. All values required at overlapping scan window positions are fetched from the register file for subsequent accesses. The row by row scan pattern for the scan windows exactly matches the two nested for-loops of the algorithm in figure 7, with the loop index of the outer loop being incremented by eight.

Knowing this structure of the compiler output, a performance estimation is quite simple. Because eight filter operations are performed in parallel, the computation time for the multiplication is not the bottleneck, but the memory and register file access times to provide the large number of input data. For each group of eight result values,  $3 * 8 = 24$  inputs have to be fetched from memory, each memory access taking 120 ns. Two columns of inputs are reused in the next loop iteration and can be fetched from the register file in 60 ns each. Two of the ten input columns are used in the computation of two results, requiring  $2 * 3 = 6$  additional register file accesses. Six input columns are used, each in three result computations, totalling in another  $6 * 3 * 2 = 36$  register file accesses. The total time to compute eight result values and store them in memory is  $24 * 120 \text{ ns} + 8 * 120 \text{ ns} + 6 * 60 \text{ ns} + 6 * 60 \text{ ns} + 36 * 60 \text{ ns} = 6720 \text{ ns}$ . A single row includes  $1280 - 2 = 1278$  positions. Eight positions are computed in parallel, resulting in  $1278 / 8 = 160$  iterations of the inner loop. One stripe takes  $(1024 + 6 * 2) / 7 = 148$  rows, which is equal to the number of iterations of the outer loop. Using seven C-Modules in parallel, it takes  $148 * 160 * 6720 \text{ ns} = 0.159 \text{ s}$  to filter a 1280 by 1024 pixel image. The filter weights may be chosen arbitrarily, because the execution time is bound by memory I/O speed. On the MoM-3, the result would be computed in the same time, regardless whether a pixel is 8 bits wide or 32 bits, because even 32-bit multiplications would be faster than memory access time in this example. The corresponding performance figures for von Neumann computers in table 1 were based on the same C source. All variables were declared as registers, the ones with the highest reference count first. The source code was compiled using GNU gcc compiler, with all optimizations turned on, including the exact specification of the microprocessor type.

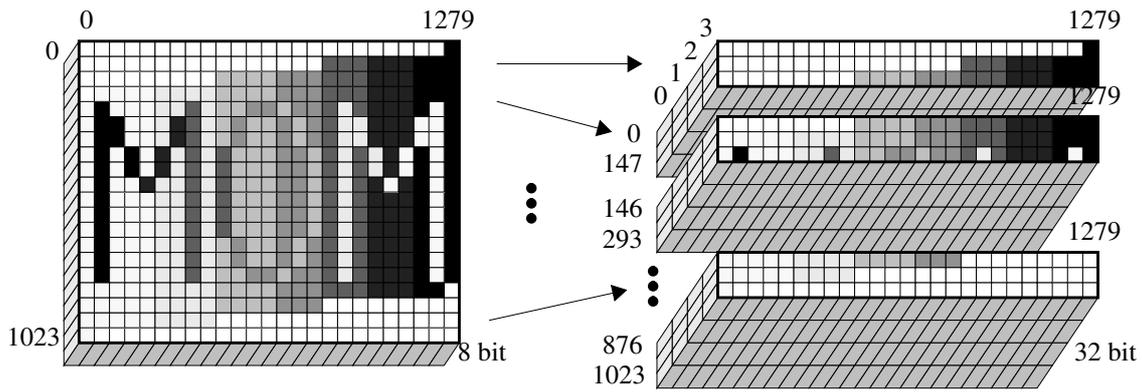
## 4.2 Manually Optimized Implementation

From the algorithm description in figure 7, the MoM-3 C compiler cannot make use of all optimization possibilities. If the input data is only 8 bits wide (a typical case for most grayscale images), four values could be transferred in a single 32 bit memory access, reducing I/O time by a factor of four. The data map for a packed representation of a 1280 by 1024 pixel image, using eight bits per pixel, can be seen in figure 11.

Instead of using the register file, values which are required in subsequent loop iterations could be passed to the next DPU using the neighbour to neighbour interconnect. Since these transfers can be done in parallel, they are much faster than register file accesses. Furthermore, if the eight FIR filter operators would compute the results for eight rows in parallel, instead of eight results in a single row, the neighbour to neighbour connections could be utilized to a far larger extend. Instead of six overlapping positions in the next loop iteration, twenty positions would overlap, so that fewer inputs would have to be read from memory and more inputs would be passed on from DPU to DPU.

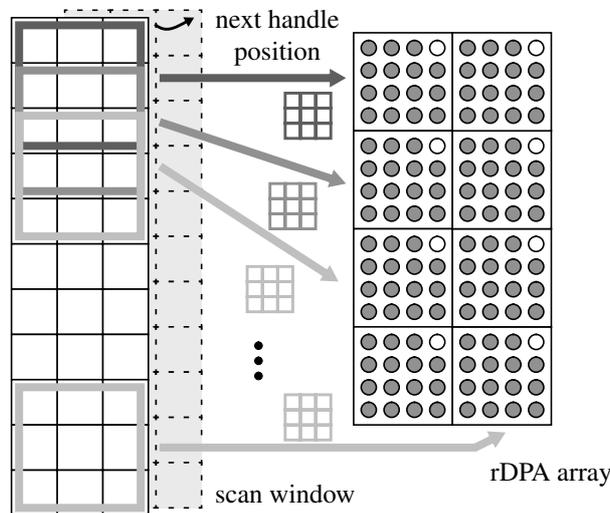
The straightforward compiler implementation doesn't make use of the fact, that the next row of results requires to read again most of the inputs, which have been read during the computation of the previous row of results. Computing eight rows of results in parallel, the reused inputs can be fetched from the register file instead of the memory. The rALU configuration and





**Figure 11. Data map for an optimized 2-D FIR filter**

the scan window corresponding to these manual optimizations of the filtering algorithm can be found in figure 12.

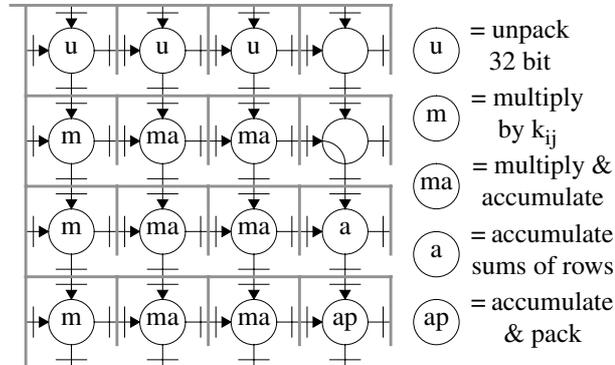


**Figure 12. Manually optimized scan window and rALU configuration**

In order to achieve such optimizations, manual changes to the input source are required. The passing of the input values to the next DPU for the following loop iteration has to be specified explicitly. Since the compiler performs a vectorization, proceeding from the inner loops to the outer loops, the parallelization of the row computations has to be done explicitly. Packed transfers of four 8-bit inputs in one 32-bit memory access cannot be specified solely in the C source. The unpacking of 32 bit values can be done in a single DPU, which accepts a 32 bit value and four times passes an eight bit value to the next DPU for computation. But this has to be done in the ALE-X assembler source. A description of that operation sequence in C would not pass unchanged through the compiler's transformation and optimizations steps, so that the ALE-X assembler would not recognize an unpack operator in the resulting description and merge it to a single DPU.

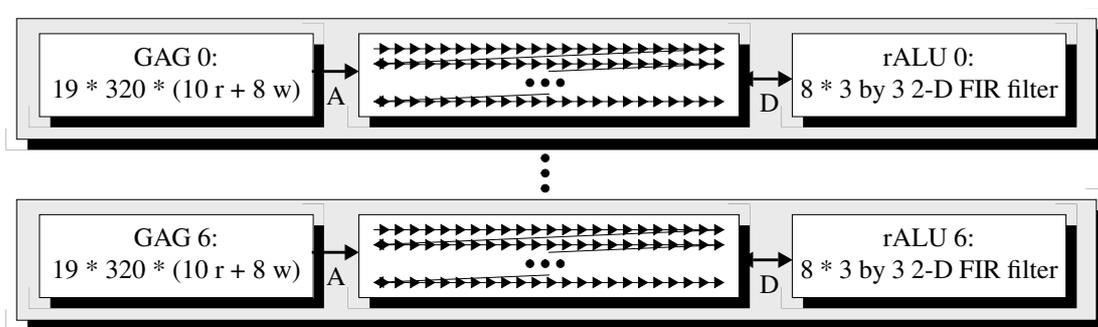
An rALU operator (figure 13) to compute a single result value consists of three unpacking DPUs ("u"), a three by three array of multiply-accumulate DPUs ("m" and "ma"), and three DPUs to accumulate the sums of rows ("a"). The last of these DPUs performs the packing of four result values to a 32 bit word as well ("ap"). The partial sums of rows are transferred from west to east to the next DPU. The operation is pipelined on expression level, because each mul-

tiply-accumulate DPU accepts the next input operand (of the subsequent loop iteration) after it has calculated its partial sum and passed the result to the next DPU to the east. That way, the whole rDPA array can be kept busy, with filter operators of consecutive loop iterations finished to different degrees in the pipeline. One filter operator fits into a single rDPA chip, so that the parallel-to-serial converters do not impose a delay on the operation. Eight of these operators fit into the rDPA array of a single C-Module.



**Figure 13. rALU operator for a 3 by 3 filter operation (manually optimized)**

The scan window moves one position to the right in the inner loop, but eight positions downwards in the outer loop, because eight rows are processed in parallel (figure 14). That way the overlapping inputs in each row can be transferred in parallel to the next DPU to the south, where they are required as inputs for the following iteration. The overlapping inputs of the adjacent rows of the eight operators are fetched from memory only once and then taken from the register file in half of the time.



**Figure 14. Scan Pattern for a 2-D FIR Filter**

One column of eight result values requires the following time for I/O: ten values have to be read from memory, and eight written back. Furthermore, two input values are fetched a second time from the register file and six values are fetched twice from the register file for reuse in another filter operator. The total I/O time is  $3000 / 4 = 750$  ns. The division by four stems from the fact, that every memory or register access transfers the values for four iterations. One pipeline stage of the multiplication and accumulation takes 420 ns in the average case, because the multiplications with constants are resolved to shift and add operations. The parallel-to-serial converters, taking 600 ns, do not account, because each filter operator is local to one rDPA chip. Although optimized, this implementation is still I/O bound. The optimized algorithm takes  $148 / 8 = 19$  iterations of the outer loop. The inner loop contains  $1280 - 2 = 1278$  iterations, each producing a column of eight results. The total time to filter a 1280 by 1024 8-bit grayscale image is  $19 * 1278 * 750$  ns = 0,018 s.

```

For (n = 0; n < maxn; n++) {
    x00 = x[n][0]; x01 = x[n][1];
    x10 = x[n+1][0]; x11 = x[n+1][1];
    x20 = x[n+2][0]; x21 = x[n+2][1];
    For (m = 0; m < maxm; m += 3) {
        x02 = x[n][m+2]; x12 = x[n+1][m+2]; x22 = x[n+2][m+2];
        y[n][m] = k00*x00 + k01*x01 + k02*x02
                + k10*x10 + k11*x11 + k12*x12
                + k20*x20 + k21*x21 + k22*x22;
        x00 = x[n][m+3]; x10 = x[n+1][m+3]; x20 = x[n+2][m+3];
        y[n][m+1] = k00*x01 + k01*x02 + k02*x00
                  + k10*x11 + k11*x12 + k12*x10
                  + k20*x21 + k21*x22 + k22*x20;
        x01 = x[n][m+4]; x11 = x[n+1][m+4]; x21 = x[n+2][m+4];
        y[n][m+2] = k00*x02 + k01*x00 + k02*x01
                  + k10*x12 + k11*x10 + k12*x11
                  + k20*x22 + k21*x20 + k22*x21;
    }
}

```

**Figure 15. Manually Optimized C Source Code**

For a fair comparison, the C source for the von-Neumann computers is manually optimized as well. Instead of copying the input values to the next position, three iterations of the inner loop are unrolled and a cyclic buffering scheme is applied. This eliminates the copy time for the microprocessors, because they cannot rely on a register to register interconnect to do the copy operations in parallel. The resulting source code can be seen in figure 15. The “in” array is declared as “unsigned character” of course, to allow the compiler to do the same optimizations on multiplications as on the MoM-3. The “inXY” variables and the loop counters are declared as register variables, where the sequence of declarations takes into account, how often each variable is referenced. In this source code, there is no explicit unrolling of the outer loop like in the optimized MoM-3 source, because a standard microprocessor cannot process multiple FIR filter operators in parallel. Table 1 reveals, that the manual optimizations improved performance on the conventional computers by more than a factor of two.

## 5. Performance Evaluation

The performance figures for some example algorithms are given in table 1. The first column lists the algorithms. The two-dimensional FIR filter algorithms are measured for different kernel sizes. For all implementations, the weights of the kernel are constants compiled into the code, allowing the compilers to replace multiplications with optimized shift and add sequences. Input to the filter is a 1280 by 1024 pixel grayscale image using 8 bits per pixel. Only the time to filter the image in memory is counted, excluding disk I/O operations. The first group of performance figures are measured on the output of the compilers, with all optimizations turned on. The second group of figures are measured on manually optimized code, that takes into account that the same values are multiplied to different weights in succeeding steps of the filter algorithm (see section 4.2). This allows to reduce memory cycles by storing these values internally. Multiplications and accumulations are done in the same DPU. The additional cost of the addition is weighed out far by the reduction of the required rDPA array size. Since all DPUs fit into a single chip in the 3 \* 3 case, the comparably high cost (in execution time) of the serial links can be avoided. The MoM-3 uses all seven C-Modules in parallel on overlapping strips of the input image to speed up the filter operation. The JPEG compression algo-



rithm is based on the IJG program “cjpeg” [5], inserting time measuring code around the compression kernel, after the input image is read. To exclude disk output, output is redirected to “/dev/null”. A 704 by 576 RGB colour image using 24 bits per pixel is used as input for the time measurements. The JPEG compression algorithm is adapted to the MoM-3 so that seven C-Modules operate in parallel. A more detailed description of the MoM-3 implementation of this algorithm can be found in [1].

Algorithms	CPU Time (in seconds)			Speed-up Factors		
	68020, 16 MHz [seconds]	Sparc 10, 50 MHz [seconds]	MoM-3, 33 MHz [seconds]	MoM-3 vs. 68020, 16 MHz	MoM-3 vs. Sparc 10/51, 50 MHz	MoM-3 NT (new technology ) vs. Sparc <sup>a</sup>
3x3 2-D FIR <sup>b</sup>	365.13	2.65	0.159	2296	16.7	33.3 <sup>c</sup>
5x5 2-D FIR <sup>b</sup>	1088.20	4.73	0.368	2957	12.9	25.7 <sup>c</sup>
7x7 2-D FIR <sup>b</sup>	1784.30	7.96	0.674	2647	11.8	23.6 <sup>c</sup>
3x3 2-D FIR <sup>d</sup>	167.14	0.71	0.018	9286	39.4	78.9 <sup>c</sup>
5x5 2-D FIR <sup>d</sup>	451.38	1.88	0.038	11878	49.5	145 <sup>e</sup>
7x7 2-D FIR <sup>d</sup>	743.38	3.60	0.058	12817	62.0	189 <sup>e</sup>
JPEG <sup>f</sup>	74.50	1.51	0.128	582	11.8	23.6 <sup>c</sup>

**Table 1. Performance evaluation: MoM-3 vs. Sparc 10 and current host computer**

- a. Sparc 10/51 versus hypothetical MoM-3 based on Sparc’s technology
- b. 1280 \* 1024 grayscale image, 8 bits per pixel, in memory, compiler optimizations only
- c. critical factor is solely memory access time: additional speed-up of two
- d. 1280 \* 1024 grayscale image, 8 bits per pixel, in memory, manually optimized
- e. critical factor is speed of serial links: 4 bit vs 2 bit with larger packages, and 66 MHz vs 33 MHz serial link clock, but only an additional speed-up of approximately three, because with new technology critical factor is memory access time
- f. IJG cjpeg compression, 704 \* 576 RGB image, 24 bits per pixel, excluding disk I/O

The second column of table 1 gives the CPU times in seconds for the ELTEC host, which runs an MC68020 at 16 MHz. The third column describes the performance of a SPARCstation 10/51 running at 50 MHz. The CPU times of the conventional computers are all based on compilations (GNU gcc) with all optimizations turned on, producing output specially adapted to exactly the kind of processor inside the computer. The fourth column indicates the time to run the same algorithm on the MoM-3 prototype (33 MHz version).

The three rightmost columns of table 1 illustrate the speed-up factors obtained compared to the hosting ELTEC workstation and a modern SUN Sparc 10/51. The last column takes into account, that our prototype is built with an inferior technology compared to a modern Sparcstation. The notes at the bottom of the table explain how the extra speed-up factors would be achieved.



## 6. Future Work

Especially to fulfil the high requirements of scientific computation, some performance improvements will have to be made to the hardware. To improve memory access time, multiple interleaved memory banks can be used, with an appropriate change to the compiler, to guarantee a conflict-free access scheme by a proper distribution of data in the data map. The floating-point performance of the rDPA array can be improved further, by adding hardware support for the specific requirements of floating-point pipeline stages in each datapath unit. For all situations, where an equal amount of I/O data has to be processed in each iteration of the inner loop of an algorithm, high-speed I/O channels (e.g. CCD cameras, or disk arrays) can be integrated easily into the MoM-3. By interfacing them as if they were a combination of a data memory and a Generic Address Generator, I/O channels could offer a number of data words as if they were fetched from a high-speed memory.

On the software side, an automatic partitioning compiler is under development, which takes full-featured ANSI C as input and creates two C source files as output. The first C output is an ANSI C description of the operations to be performed on the host computer, including the interface calls to activate the MoM-3 for appropriate parts of the algorithm. The second C output follows the syntax of the C subset and describes those parts of the algorithm to be executed on the MoM-3.

## 7. Conclusions

The MoM-3, a configurable accelerator has been presented, which provides acceleration factors in three orders of magnitude for its outdated host computer, and speed-ups in the range of 12 to 62 when compared to a state of the art workstation. High acceleration factors can still be maintained, when working from a high level input language. The compilation is done from an almost complete subset of standard C language. Good results are achieved, without requiring user interaction or special compiler directives to generate code for field-programmable devices and the controlling hardware of the MoM-3. The custom designed rDPA circuit provides a coarse grain field-programmable hardware, especially suited for 32-bit datapaths for pipelined arithmetic. A new sequencing paradigm supports multiple address generators and a loose coupling between sequencing mechanism and ALU. The address generation under hardware control leaves all memory accesses free for data transfers, providing de facto a higher memory bandwidth from the same memory hardware. The loose coupling of the data sequencing paradigm allows to fully exploit the benefits of reconfigurable computing devices. The combination of hardwired address generators and configurable devices is the key to speed up both data manipulations and data accesses - and still maintain a programming environment familiar to a conventional software developer.

The MoM-3 prototype is currently being built as a co-processor to a VMEbus-based ELTEC workstation. The address generators, the MoM-3 controller and the rALU controller have just returned from fabrication and are now under test. All three are integrated circuits based on 1.0  $\mu$ m CMOS standard cells, a technology made available by the EUROCHIP project of the CEC. The rDPA circuits will be submitted to fabrication in a 0.7  $\mu$ m CMOS standard cell process soon. All MoM-3 performance figures were obtained from simulations of the completed circuit designs, which were integrated into a simulation environment, describing the whole MoM-3 at functional level. The C compiler and its development environment are implemented in C programming language, running under SunOS on SPARCstations.



## References

- [1] R. W. Hartenstein, J. Becker, R. Kress, H. Reinig, K. Schmidt: A Reconfigurable Machine for Applications in Image and Video Compression; European Symposium on Advanced Services and Networks / Conference on Compression Techniques and Standards for Image and Video Communications, Amsterdam, March 1995
- [2] R. W. Hartenstein, A. G. Hirschbiel, K. Schmidt, M. Weber: A Novel Paradigm of Parallel Computation and its Use to Implement Simple High-Performance Hardware; Future Generation Systems 7 (1991/92), p. 181-198, Elsevier Science Publ., North-Holland, 1992
- [3] Jaap Hollenberg: The CRAY-2 Computer System; Supercomputer 8/9, pp. 17–22, July/September 1985
- [4] K. Hwang: Advanced Computer Architecture: Parallelism, Scalability, Programmability; McGraw-Hill, 1993
- [5] T. G. Lane: cjpeg software, release 5; Independent JPEG Group (IJG), Sept. 1994
- [6] N.N.: FLEX 8000 Data Book; Altera, Inc., 1993
- [7] N.N.: ORCA Preliminary Data Sheet, AT&T Microelectronics, 1994
- [8] N. N.: The XC4000 Data Book; Xilinx, Inc., 1994
- [9] John Reid: Fortran shapes up to the future; Scientific Computing World, January 1995
- [10] K. Schmidt: A Program Partitioning, Restructuring, and Mapping Method for Xputers; Ph. D. Thesis, University of Kaiserslautern, 1994
- [11] S.A. Zenios, R.A. Lasken: The Connection Machines CM-1 and CM-2: Solving Nonlinear Network Problems; Int'l. Conf. on Supercomputing, pp. 648–658, ACM Press, 1988

