

In fact, the interaction field of a many-body system is usually expressed as empirical potential among the atoms (particles) of the many-body system under study and it should include energy terms that represent bonded and non-bonded (Van der Waals and Coulombic forces) interactions. By limiting our attention to the simulation of a biological system, a typical force field used in classical MD simulations appears to be like follows

$$V = \sum_{\text{bonds}} K_b [b - b_0]^2 + \sum_{\text{angles}} K_\theta [\theta - \theta_0]^2 + \sum_{\text{dihedrals}} K_\phi [1 + \cos(n\phi - \delta)] + \sum_{\text{improper}} K_\xi [\xi - \xi_0]^2 \quad (10a)$$

$$+ \sum_i \sum_j \epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] + \sum_i \sum_j \frac{q_i q_j}{4\pi\epsilon r_{ij}}. \quad (10b)$$

The first four terms (Eq. (10a)) represent the bonded potential. The last one (Eq. (10b)) gives the non-bonded energy contribution to the total potential and it is the most time-consuming task ( $\approx 85\text{--}90\%$  of the total CPU time).

Each pair interaction  $ij$  in Eq. (10b) has to be calculated to obtain the energy due to the non-bonded forces. The number of pairs to be calculated is halved using Newton's third law  $F_i = -F_j$ .

The number of pair interactions can be further reduced including only atom pairs within a cutoff distance. Therefore, at each step a generic atom  $i$  has a set of  $j$  different atoms to interact with (the non-bonded pair list). This list can be updated every  $n$  steps, with  $n \approx 10$ .

Parallelization can be achieved with a particle decomposition scheme, i.e. a number of  $i$  particles is assigned to each processor (home particles) [7]. Then, each processor will calculate independently the pair interactions of the  $i$  home particles with the  $j$  particles.

It is worth noting that in a shared memory machine the replicated data model is automatically obtained, i.e. data about every atom in the system is available to each processor, while in a distributed memory architecture position, and undergone force of each particle, have to be communicated between the machine processors.

Once obtained the interaction forces on each atom, the Newton's equations of motion can be numerically integrated using, for example, the leap-frog algorithm [7].

The integration is not a computationally demanding task but it can be nevertheless performed in parallel on each home particle. A longer integration time step permits to simulate longer times with the same number of steps. If the bond-stretching vibrations have to be taken into account, the time step cannot be longer than 0.25 fs. The application of distance constraints between bonded atoms, with algorithms like SHAKE [32], permits to increase the time step to 2 fs.

After the non bonded interaction calculations, SHAKE is the most computational-intensive part of the classical MD calculation. Modern MD simulations of biological systems include a great number of solvent molecules on which it is possible the application of independent distance constraints by different processors. Constraints application to protein or DNA, on the contrary, is better performed by one processor. A good load balancing can be obtained with a functional parallelism, i.e. assigning different tasks (solvent distance constraints, solute distance constraints, bonded calculations, etc.) to different processors.

### 3.3.1. Test case

GROMACS version 1.6 [4], from Groningen BIOSON Research Institute, is an efficient parallel MD code, particularly suited to simulate biological systems. GROMACS has been parallelized using a message-passing model. Both PVM [33] (Parallel Virtual Machines) and MPI [34] (Message Passing Interface) have been utilized within the GROMACS software.

As biological test case we chose to simulate the Major Histocompatibility Complex macromolecule class II (MHC II) complexed with an antigenic peptide and water molecules.

Table 6  
MHC II simulated on a Compaq ES40 cluster

SMP mach.	Nodes	E.T. (s)	Speedup	% parallel code
1	1	932	1	N.A.
	2	467	2.0	99.8
	3	346	2.7	94.3
	4	281	3.3	93.1
2	2	470	2.0	99.1
	4	281	3.3	93.1
	6	215	4.3	92.3
	8	168	5.5	93.7
3	3	549	2.7	93.8
	6	216	4.3	92.2
	9	146	6.4	94.9
	12	114	8.2	95.8
4	4	284	3.3	92.7
	8	168	5.5	93.7
	12	113	8.2	95.9
	16	101	9.2	95.1

Table 7  
MHC II simulated on a IBM 8x SP3 cluster

SMP mach.	Nodes	E.T. (s)	Speedup	% parallel code
1	1	1594	1	N.A.
	8	313	5.1	91.8
2	16	191	8.3	93.9

MHC II are immunoglobulin-like proteins present in variety of cells, which bind antigens from endocytosed plasma membrane and extracellular proteins. Antigen-MHC II complexes are recognized at the surface of the cell by helper T-cells promoting activation of an immune response.

MHC II is composed by 7776 atoms, the antigenic peptide by 280 atoms and there are 17 296 water molecules giving a total of 59 863 atoms. The system was simulated for 100 steps with a time step of 2 fs; a non-bonded cutoff radius of 1.8 nm was used updating the non-bonded pair list every 10 steps.

Table 6 (cluster of Compaq ES40, 4x ev6@500 MHz) and Table 7 (cluster of IBM SP3, 8x PWIII@222 MHz) show the test case results. In column one the number of SMP machines involved in the test run is shown; in column two the total number of processors and in the remaining columns the elapsed time, the speedups of the MPI and the portion of parallelized code as extrapolated by the Amdahl law, are reported. The last data were calculated using the following equation:  $\alpha = \frac{S \times p - p}{S \times p - 1}$  where  $S$  is the speedup and  $p$  the number of processors.

It is worth noting the variability of the extrapolated parallel portion of the code that varies from a maximum of 99.8% (one Compaq ES40 machine, 2 processors) to a minimum of 91.8% (one IBM SP3 machine, 8 processors).

This behaviour is due to the simultaneous use of different parallelization schemes: particle decomposition for the non-bonded force calculation that represents the 90.2% of serial calculation time; functional parallelism for the remaining portion of the code.

For this reason the parallel efficiency of the code strongly depends on the simulated system and on the condition of the simulation. Nevertheless, the weight of the computational core (the non-bonded force calculation) can be considered as the minimum portion of parallel code obtainable in a real simulation.

### 3.4. Atomic and Molecular Physics: The electron–molecule scattering treated with the Single Center Expansion (SCE) method

The application of the central field model in quantum mechanics yields to the factorization of the wavefunction in its radial and angular part and it has been and is nowadays, applied to the quantum description of bound and continuum electronic states of atoms [35] while several attempts have been made in the past to extend it to bound molecular systems [36].

One of the most successful application of this model has been the treatment of the scattering processes which occur in the collision of electrons with polyatomic non-linear targets, known as the Single Center Expansion (SCE) method. A general discussion on the computational aspects of this model has been given before [37] and specific results have also been analyzed elsewhere [38,39]. We will therefore only outline here the basics of the generation of the SCE numerical wavefunction and its related molecular properties. The full details on the mathematical and numerical aspects of SCE procedure are given elsewhere. We refer the interested reader to a recent paper which describes the SCeLib code a computational procedure for the SCE pre-processing phase, combined in a single set of programs [11].

#### 3.4.1. The SCE wavefunction and related molecular properties

In the present computational approach one needs a general-purpose quantum chemistry code that is employed to generate the Single Determinant description, (near to the Hartree–Fock — HF — limit) of the target electronic wavefunction, and an interface with a numerical procedure that can give us all the necessary quantities as being referred to the molecular Centre Of Mass (C.O.M.) as expansion centre [37].

In most of the numerical methods employed to solve the scattering equations one then converts the CC equations into a set of coupled radial equations by making first a single centre expansion (SCE) of the bound and continuum functions and then by integrating over all the angular coordinates. Limiting our attention to the ground state wavefunction, one therefore writes down the bound orbitals as expansions around the centre of mass, from which the body-fixed (BF) frame of reference originates:

$$\phi_i(\mathbf{r}) = r^{-1} \sum_h u_{hi}^i(r) X_{hi}^{l, \mu}(\theta, \phi). \quad (11)$$

Here  $i$  labels a specific, multicenter molecular orbital (MO), which contributes to the density of the bound electrons in the nonlinear target, while the indices  $|p\mu\rangle$  for the continuum functions label one of the relevant Irreducible Representations (IR) and one of its components, respectively. The index  $h$  labels a specific basis, at a given angular momentum  $l$ , for the  $p$ th IR that one is considering [40].

In order to perform expansions (11), one needs, therefore, to start from the multicenter wavefunction which describes the target molecule and then generate by quadrature each of the  $u_{hi}^i(r)$  coefficients; they were obtained for the first time by numerical quadrature of the multicenter GTO's given as Cartesian Gaussian functions [41]. The  $X$  angular functions used in (11) are symmetry-adapted, generalized harmonics build as linear combinations of spherical harmonics  $Y_l^m(\theta, \phi)$  which, for a given  $l$ , form a basis of the  $(2l + 1)$ -dimensional IR of the full rotation group [40,42].

After one obtains the radial coefficients  $u_{hl}^k(r)$ , each bound one-electron wavefunction is also expanded about the C.O.M. in terms of the  $X$  functions

$$u_i(\mathbf{r}) = r^{-1} \sum_{hl} u_{hl}^k(r) X_{hl}^{pk}(\theta, \phi) \quad (12)$$

then, the one-electron density function can be written as usual

$$\rho(\mathbf{r}) = \int |\psi(x_1, x_2, \dots, x_N)|^2 dx_2 \dots dx_N = 2 \times \sum_i |u_i(\mathbf{r})|^2, \quad (13)$$

where the factor 2 is due to the sum over spin, and the  $i$  sum is over each doubly occupied orbital. Once the quantity  $\rho(\mathbf{r})$  is obtained from the bound-state wavefunction  $u_i(\mathbf{r})$  of Eq. (13), it can be expanded in terms of symmetry-adapted functions belonging to the  $A_1$  irreducible representation as

$$\rho(\mathbf{r}) = r^{-1} \sum_{hlm} \rho_{hlm}(r) X_{hlm}^{A_1}(\theta, \phi), \quad (14)$$

where

$$\rho_{hlm}(r) = 2 \times \sum_i \int_0^\pi \sin(\theta) d\theta \int_0^{2\pi} u_i(\mathbf{r}) \cdot u_i(\mathbf{r}) d\phi. \quad (15)$$

Once the SCE one-electron density has been computed from Eqs. (14) and (15) one is able to derive from it all the molecular quantities which depend on the electronic distribution, either by integration, like the electronic Static Potential  $\int \rho(\mathbf{r}) d\mathbf{r}$ , or by differentiation like the density gradient ( $\nabla\rho$ ) and Laplacian ( $\nabla^2\rho$ ).

By using the SCFLib [43] API (Application Program Interface) one has access to a large set of these molecular properties being calculated with respect to the molecular COM and for the first time we coded in it the gradient and Laplacian of the electronic density and static potential as explained in details in the next section.

### 3.4.2. Numerical implementation details

A feature common to all SCE molecular properties is that they can be expressed as

$$f(r, \theta, \phi) = \sum_{hlm} F_{hlm}(r) X_{hlm}^{A_1}(\theta, \phi), \quad (16)$$

where the  $f$  represent a given molecular property, function of the spherical coordinates as dot product of a radial component  $F$  times an angular analytic function  $X$ .

In a similar manner are expressed the gradient and Laplacian of the  $f$ :

$$\begin{aligned} \nabla f(r, \theta, \phi) &= \sum_{lm} \nabla [F_{lm}(r) X_{lm}(\theta, \phi)] \\ &= \sum_{lm} \left[ \frac{dF_{lm}}{dr} X_{lm} \hat{\mathbf{e}}_r + \frac{F_{lm}}{r} \left( \frac{\partial X_{lm}}{\partial \theta} \hat{\mathbf{e}}_\theta + \frac{m}{\sin\theta} X_{lm} \hat{\mathbf{e}}_\phi \right) \right] \end{aligned} \quad (17)$$

and

$$\begin{aligned} \nabla^2 f(r, \theta, \phi) &= \sum_{lm} \nabla^2 [F_{lm}(r) X_{lm}(\theta, \phi)] \\ &= \sum_{lm} \left[ \frac{d^2 F_{lm}}{dr^2} + \frac{2}{r} \frac{dF_{lm}}{dr} - \frac{l(l+1)}{r^2} \right] X_{lm} \end{aligned} \quad (18)$$

but we let the interested reader to the specific literature [44] for the necessary mathematical backgrounds and numerical implementation details.

On Eqs. (16)–(18) above, the angular part is analytic and the radial part is numerical as derived from the SCE procedure. This last fact suggest that when one want to calculate a given property over a given  $(r, \theta, \phi)$  grid, the most time consuming section will be the calculation of the  $F_{lm}(r)$ ,  $F'_{lm}(r)$  and  $F''_{lm}(r)$ . In order to improve the calculation of the  $F_{lm}(r)$  and of its  $r$  derivatives, a natural way to proceed is to post-fit it after the SCE procedure, with a suitable fitting function.

In a first approach, we used a cubic spline fitting [45] of the  $F_{lm}(r)$  function. The spline fitting has the property that the fitted function can be very efficiently evaluated at each  $r$  point but more important, that one can obtains its first and second derivatives with respect to  $r$  from the same fitting parameters of the  $F_{lm}(r)$ . Thus, when one has fitted the  $F_{lm}(r)$  function with a cubic spline, by calculating the four parameters needed at each  $r$  point, a single call to an evaluation routine can efficiently return the  $F_{lm}(r)$ ,  $F'_{lm}(r)$  and  $F''_{lm}(r)$ .

For what concerns the angular part of our SCE functions, a natural limit arise on the floating point representation of the  $X_{lm}$  (i.e. the  $Y_{lm}$ ); in fact, it is well known that over a certain value of  $l$  one reaches the limit of the double precision floating point arithmetic. To overcome this limit, one can either use a quadrupole precision floating point arithmetic (currently available on 64 bits computers) or use specifically designed routines for multiple precision fit/evaluation of the  $Y_{lm}$ , and even available on Internet [46].

### 3.4.3. Numerical results

In this section we will report some calculation examples by limiting the resulting data to the pre-processing part of the whole electron–molecule SCE scattering procedure. These results refer to those recently reported in the literature [11] but are enriched with more informations on the inner profiling of the SCELlib code used for the test. This could offer a way to the reader to better analyze the relative importance of the various code sections with the aim to extract the numerical intensive kernels and evaluate them in detail. However, we should warn again the reader that we are referring to the pre-processing part of the whole electron–molecule computational procedure, and that the successive scattering calculation is likewise CPU and I/O demanding as stated in a recently published paper [47], where a new method of integration of the scattering equations has been presented.

An even more interesting aspect to discuss of the pre-processing part, would be the fact that through an interface code — the SCELlib API — one is able to use the pre-processing data to map a given molecular property (static potential, electron density, etc.) over a generic Cartesian grid. This of course, could be of more general use in many related scientific areas, but we should leave the interested reader to some preliminary results obtained using a prototype version of the SCELlib API set of programs [43].

As a test case to discuss, we chose a  $C_{2v}$  molecule, the  $SO_2$  system, with a number of electrons able to produce a timing long enough to avoid any random measurement error in the parallel runs. For the  $SO_2$  molecule, we started from its HF/D95\*\* optimized geometry ( $R(SO) = 1.423 \text{ \AA}$ ,  $\widehat{OSO} = 118.1^\circ$ ) and using this electronic solution to derive the SCE wavefunction of Eq. (12) with the SCELlib package. Furthermore, we used a  $R, \theta, \phi$  of  $300 \times 96 \times 84$  for a total of 2 419 200 points,  $L_{MAX}$  was set to 10 and the Centre Of Mass (COM) chosen as expansion centre. This required a scratch memory usage of about 300 MB, which is small enough to be easily managed by the server class machines under test. In the followings, we decided to show only the results coming from the most time-consuming part of the SCELlib package, the *sphint* function, which performs the Single Centre Expansion of the GTO wavefunction of a SCF run. In this program, the grid evaluation of the GTO wavefunction and the subsequent angular integration are carried out and to show its performances the Table 8 reports the overall elapsed time of the *sphint* code section together with the corresponding speed-ups on the IBM, Compaq and SUN SMP parallel machines of the  $p$ -threaded SCELlib version using from 1 to 14 nodes.

It is evident from the reported data that all the machines show a fairly linear speed-up and in the case of the ES4500 this trend is maintained up to  $n = 14$ . We should note however, that in this latter case we observe a super-linear speed-up between 4 and 12 nodes which is quite unusual for this kind of parallel architecture where for this type of floating-point intensive codes one usually find a performance bottleneck beyond 8 CPUs. This particular behaviour has been described and explained in detail in Ref. [11] and it will not repeated here.

Table 8  
 SO<sub>2</sub> SCELlib *sphint* parallel runs on IBM, SUN and Compaq SMP architectures. Elapsed Time (E.T.) in seconds and speedup vs. the number of nodes

Nodes	IBM 43P-260		Compaq ES40		SUN ES4500	
	E.T. (s)	Speedup	E.T. (s)	Speedup	E.T. (s)	Speedup
1	208.9	1.0	104.1	1.0	393.3	1.0
2	104.1	2.0	52.6	2.0	183.3	2.1
4			25.9	4.0	89.6	4.4
8					45.1	8.7
12					32.0	12.3
14					28.5	13.8

More interesting to discuss here, is the analysis of the relative performances of the leading computational parts of the SCELlib code. To this end, we report in the followings, the profiling informations of the test performed on the same molecular system cited above but over a smaller grid (100R points) and serially ran over a single CPU Alpha ev67@667 MHz CPU, with the standard *prof* Unix profiler under the Tru64 V5.0A Operating System.

%time	seconds	cum %	cum sec	procedure (file)
51.6	17.4873	51.6	17.49	CalcMO (<scelib>)
47.0	15.9287	98.6	33.42	sphint (<scelib>)
0.7	0.2529	99.4	33.67	pow_di (<scelib>)
0.3	0.0967	99.7	33.77	SubPointC (<scelib>)
0.1	0.0488	99.8	33.81	gangf (<scelib>)
0.1	0.0283	99.9	33.84	MulPointC (<scelib>)
0.1	0.0215	99.9	33.86	Plm (<scelib>)
0.0	0.0088	100.0	33.87	setinp (<scelib>)
0.0	0.0029	100.0	33.88	norm (<scelib>)
0.0	0.0029	100.0	33.88	rotgto (<scelib>)
0.0	0.0020	100.0	33.88	monorm (<scelib>)
0.0	0.0010	100.0	33.88	mkgrid (<scelib>)

It is evident from the data reported that 99.9% of the computing time refers to the part of the code running in parallel that is, the *sphint* function. In fact, apart from the preparation routines (from *setinp* to *mkgrid*) the rest of the functions are called from *sphint* and among these, the most time consuming one is *CalcMO*. This behaviour confirm the linear speedup found over the three parallel SMP architectures used for the tests where we used a *p*-threaded version of the *sphint* function, so that the whole computational kernel was eligible to be run in parallel.

It is interesting to point out the fact that two routines, *CalcMO* and *sphint*, account for the majority of the CPU time spent for the calculation. In the former function the Gaussian wavefunction is calculated over a convenient  $(\theta, \phi)$  grid and then used by the latter (the "caller") to perform a Gauss-Legendre (Gauss-Chebyshev) angular integration. These two steps are necessary to evaluate the  $u_{hl}^i(r)$  terms of Eq. (11) and once carried out, the one-electron density can be calculated and from it, all the necessary molecular properties referred to a Single Centre of reference. The computing time of the evaluation/integration steps is strictly bound to the low level routines used by the two functions cited above. While the evaluation step depends from the *exp* function calculation of the Gaussian

basis set (see Eq. (3)), the integration phase is totally bound to the evaluation of the angular spherical harmonics  $Y_{lm}$  functions (see Eq. (11)). This behaviour is quite interesting because it shows explicitly the dependence of the two major stages of the pre-processing calculation phase performed in *sphint* from the basic low-level routines *exp* and  $Y_{lm}$ .

#### 4. The commonly shared numerical behaviour

The first element of similarity among the codes we analyzed so far, is that we have a *code implementation dependency*. By this we mean that the user is not guaranteed that the same binary code is executed on different architectures, due to the fact, for example, that different algorithms can be used by low-level vendor routines at run-time.

A second element of similarity can be found by looking at the pre-processing behaviour of these codes (a situation that we refer to as *job type dependency*). By this we mean that the user is not guaranteed that the same code sections will perform exactly at the same computational rates when applied to different molecular systems. In examples, if you change your molecule geometry or, even if you change some run-time parameters in the input, you are not assured that the same code section will perform at a constant computational rate.

This is a common behavior of many codes in this and other scientific areas, which have years of development behind, containing several thousands of FORTRAN/C source lines: it is well known to the users (end-users or developers) of the codes mentioned above, that even the traditional (serial) profiling stage can be a cumbersome activity due to the hidden complexity of the source files.

Nonetheless, we have seen as they share specific features which make possible to derive some common conclusions on the best numerical environment for those applications. We would point out in fact as stated in Eq. (3), the central role played by the *exp* function in quantum chemistry and in molecular physics. At the same time, the BLAS low-level routines are of paramount importance in materials science as well as in quantum chemistry as reported in Eqs. (1) and (9). Furthermore, the evaluation of the most time-consuming part of the potential used in classical MD of biomolecules and reported in Eq. (10b) closely relates to the expression of any SCE molecular property and its gradient/laplacian of Eqs. (11)–(18), where the dependency on  $1/r$  at various ranks is clearly evident.

These last comments suggested us to try a sort of summary of the numeric needs of these closely related computational areas. We will sketch these requirements in a short report (Table 9) where we would focus the attention on those low-level routines which are the computational cores of the above cited codes. These computational kernels are generally small, cleaned sections of code that we would expect to be ported on highly optimized silicon chips and dedicated to the high performance scientific computing community.

Table 9  
A tentative attempt to summarize the most time consuming low-level routines needed in (bio)chemical-physics simulations

Function class	Function type
Basic mathematical functions	<i>exp</i> , <i>r</i> , ...
Linear algebra basic functions	BLAS L1-3, ...
Polynomial functions	Legendre, Hermite, Chebyshev, ...
Simple fitting functions	Splines, <i>N</i> rank polynomials, ...
Series expansions	Fast Fourier transform, ...

## 5. Conclusions

We have seen in the previous sections how many aspects join apparently different scientific fields, in the computational area of (bio)chemical-physics. From the computational side these similarities become even more evident and a common set of functions and routines are frequently used in the majority of the codes adopted in this area of computing.

This suggests a reasonable request to electronic engineers: why do not exploit the porting of these computational kernels on optimized silicon chips? We are really confident that this could improve the performances of the codes used in these area by several order of magnitude. By the way, we are not underestimating the inner complexity in the design and manufacture of these *special-purpose* (SP) devices, but we would conversely focus attention of the designers on the possible benefits of these SP chips to this (not so) specific field of application.

First of all, let imagine this SP chip integrated as a co-processor in any serial or parallel architecture by means of the PCI (Peripheral Component Interconnect) connector of its board. A part from the key technical issues regarding the data bandwidth toward the *General Purpose* (GP) CPU, this component could co-operate with it to solve specific numeric problem, by leaving to the main GP processor the remaining of the workload (I/O, devices management, etc.). This hypothetic system with SP and GP processors coupled together could be a viable High Performance Computing (HPC) solution for this and many other computational areas, providing one is able to easily administer the SP processors within the GP system.

In fact the SP boards could be designed to solve many different specific numerical problems (not necessary only in the area of (bio)chemical-physics) but, in our opinion, the whole system should conform to the following requirements:

- The SP chip must have an Application Program Interface (API) so that it can be easily used as the corresponding software functions they refer to.
- The SP API library should be as conformant as possible to the corresponding software counterparts and in a second step, the SP API should be able to conform to the underlying hardware present (see, e.g., the role played by the GL/OpenGL library in the graphic area of computing).
- The best hardware solution for this computational area would be a combination of GP/SP processors with the maximum configuration flexibility (e.g., you change your GP system and maintain the SP boards on the new system).

This combination of GP/SP chips could operate at impressive computational rates with minimal increasing of the hardware complexity as the PC market clearly shows. In fact the idea of merging GP and SP devices is not new: Intel Co., has introduced three years ago the MMX SP (SIMD) chip into its series of Pentium processor with significant performance gains in executing multimedia X86 instructions.

This last point open the discussion on the possible benefits of a mixed GP/SP architecture for this area of computing, a topic which is outside the scope of this paper. We would in any case point out that, differently from other computational areas like High Energy Physics or Astrophysics,<sup>6</sup> the codes mentioned above require CPU as well as memory and storage I/O top performances. Hence, the mixed GP/SP solution to the computing environment could surely represent the best viable way to High Performance Computing in the short-medium terms.

Furthermore, the expected performances of the SP chips dedicated to this High Performance area of Computing are at least of one order of magnitude better than the corresponding GP RISC processors (the SP chip will eventually perform with greater performances over the GP counterpart and stated by the published vendor benchmarks). This, together with the expected feasibility of the SP board over a standard PCI bus, can improve the large diffusion of this kind of HPC solution with immediate economic benefits for the producer.

<sup>6</sup> The interested reader can refer to the articles on the parallel SP machines, APE and GRAPE, presented in this issue for any specific computational details.



## Acknowledgements

We wish to thank the kind help of Prof. F.A. Gianturco and Prof. L. Colombo for many helpful discussions on the topics covered by this paper, Dr. A. Pieretti for the Gaussian98 results and the CASPUR computational centre for providing the computer architectures used in this work.

## References

- [1] E.g. see, E. Clementi, G. Corongiu (Eds.), *METECC-95*, STEF, Italy, 1995.
- [2] M.J. Frish, G.W. Trucks, H.B. Schlegel, G.E. Scuseria, M.A. Robb, J.R. Cheesman, V.G. Zakrzewsky, J.A. Montgomery, R.E. Stratmann, J.C. Burant, S. Dapprich, J.M. Millam, A.D. Daniels, K.N. Kudin, M.C. Strain, O. Farkas, J. Tomasi, V. Barone, M. Cossi, R. Cammi, B. Mennucci, C. Pomelli, C. Adamo, S. Clifford, J. Ochtersky, G.A. Petersson, P.Y. Ayala, Q. Cui, K. Morokuma, D.K. Malick, A.D. Rabuck, K. Raghavachari, J.B. Foresman, J. Cioslowski, J.V. Ortiz, B.B. Stefanov, G. Liu, A. Liashenko, P. Piskorz, I. Komaromi, R. Gomperts, R.L. Martin, D.J. Fox, T. Keith, M.A. Al-Laham, C.Y. Peng, A. Nanayakkara, C. Gonzalez, M. Challacombe, P.M.W. Gill, B.G. Johnson, W. Chen, M.W. Wong, J.L. Andres, M. Head-Gordon, E.S. Replogle, J.A. Pople, *Gaussian98* (Revision A.7), Gaussian Inc., Pittsburgh, PA, 1998.
- [3] See, e.g., W.F. van Gunsteren, X. Daura, A.E. Mark, GROMOS force field, in: N.L. Allinger, T. Clark, J. Gasteiger, P.A. Kollman, H.F. Schaefer III, P. Schreiner (Eds.), *Encyclopedia of Computational Chemistry*, Vol. 2, Wiley and Sons, London, 1998, pp. 1211–1216.
- [4] S. Achterop, R.v. Drunen, D.v.d. Spoel, A. Sijbers, H. Keegstra, B. Reitsma, M.K.R. Renardus, *Gormacs: A Parallel Computer for Molecular Dynamics Simulations*, Proceedings of Physics Computing'92, Vol. 1, World Scientific, Singapore, 1993.
- [5] R. Car, M. Parinello, *Phys. Rev. Lett.* 55 (1985) 2471.
- [6] C.Z. Wang, K.M. Ho, in: I. Prigogine, A.A. Rice (Eds.), *Advances in Chemical Physics*, Vol. LXXXIX, p. 651.
- [7] H.J.C. Berendsen et al., Selected parallel applications and practical elements, in: *Aspects of Computational Science*, NCF, 1995, Chapter 6.
- [8] D.G. Green, K.E. Meacham, M. Surridge, F. van Hoessel, H.J.C. Berendsen, Parallelization of Molecular Dynamics Code. GROMOS87 parallelization for distributed memory architecture, in: *Methods and Techniques in Computational Chemistry: METECC-95*, STEF, 1995.
- [9] D. Roccatano, R. Bizzarri, G. Chillemi, N. Sanna, A. Di Nola, *J. Comput. Chem.* 19 (1998) 685.
- [10] E.g. see, <http://www.casput.it>.
- [11] N. Sanna, F.A. Gianturco, SCELlib: A parallel computational library of molecular properties in the single center approach, *Comput. Phys. Comm.* 128 (2000) 139.
- [12] Digital Extended Mathematical Library (DXML), <http://www.digital.com/hpc/software/dxm1.html>.
- [13] Engineering and Scientific Subroutine Library (ESSL), [http://www.research.ibm.com/acte/Opt\\_Lib/Topic\\_OptLibraries.html](http://www.research.ibm.com/acte/Opt_Lib/Topic_OptLibraries.html).
- [14] See, e.g., a recent unpublished work by K. Goto on some BLAS subroutines implemented on the Alpha ev56/ev6 microprocessors, by directly contacting the author at the e-mail [goto@statlab.rim.or.jp](mailto:goto@statlab.rim.or.jp).
- [15] R.G. Parr, W. Yang, *Density Functional Theory of Atoms and Molecules*, Oxford University Press, New York, 1989.
- [16] L. Colombo, M. Rosati, Parallel tight-binding molecular dynamics simulations on symmetric multiprocessing platforms, *Comput. Phys. Comm.* 128 (2000) 108.
- [17] G. Chillemi, A. Pieretti, M. Rosati, N. Sanna, The Performance of chemical-physics codes on CASPUR parallel architectures, *CASPUR Tech. Rep.* #10, in preparation.
- [18] A.J. van der Steen (Ed.), *Aspects of Computational Science*, NCF Publ., Den Haag, The Netherlands, 1995.
- [19] C.C. Roothaan, *Rev. Mod. Phys.* 23 (1951) 69.
- [20] G.G. Hall, *Proc. Roy. Soc. (London) A* 205 (1951) 541.
- [21] E. Clementi (Ed.), Appendices 7C–7E in *MOTECC-90, Modern Techniques in Computational Chemistry*, ESCOM, Leiden, The Netherlands, 1990.
- [22] R.C. Raffanetti, *Chem. Phys. Lett.* 20 (1973) 335.
- [23] For a discussion of SCF convergence and stability see, H.B. Schlegel, J. McDouall, in: C. Ogretir, I.G. Csizmadia (Eds.), *Computational Advances in Organic Chemistry*, Kluwer, The Netherlands, 1991.
- [24] W.J. Hehre, L. Radom, P.v.R. Schleyer, J.A. Pople, *Ab Initio Molecular Orbital Theory*, Wiley and Sons, New York, 1986.
- [25] M.P. Allen, D.J. Tildesley, *Computer Simulations of Liquids*, Oxford University Press, Oxford, 1991.
- [26] L. Colombo, in: D. Stauffer (Ed.), *Annual Review of Computational Physics IV*, World Scientific, Singapore, 1996, p. 147.
- [27] I. Kwon, R. Biswas, C.Z. Wang, K.M. Ho, C.M. Soukoulis, *Phys. Rev. B* 49 (1994) 7242.
- [28] <http://www.netlib.org/lapack/>.
- [29] G.H. Golub, C.F. Van Loan, *Matrix Computations*, 3rd edn., Johns Hopkins University Press, Baltimore, 1996.
- [30] J.J. Dongarra, D.C. Sorensen, A fully parallel algorithm for the symmetric eigenvalue problem, *SIAM J. Sci. Stat. Comp.* 8 (1987) 175–186.

- [30] <http://www.nag.com/numeric/FS/FS.html/>.
- [31] <http://www.netlib.org/blas/>.
- [32] J.-P. Ryckaert, G. Cicotti, H.J.C. Berendsen, *J. Comput. Phys.* 23 (1977) 327.
- [33] A. Beguelin, J. Dongarra, A. Geist, R. Manjchek, V. Sunderam, Oak Ridge Natl. Laboratory Tech Report TM-11826 1, 1991.
- [34] W. Gropp, E. Lusk, A. Skjellum, *Using MPI Portable Parallel Programming with the Message Passing*, MIT Press, 1994.
- [35] B.H. Bransden, C.J. Joachain, *Physics of Atoms and Molecules*, J. Wiley and Sons, New York, 1987, p. 290.
- [36] See, e.g., R. Moccia, *J. Chem. Phys.* 40 (1964) 2164, 2176, 2186.
- [37] See, e.g., F.A. Gianturco, D.G. Thompson, A.K. Jain, in: W.M. Huo, F.A. Gianturco (Eds.), *Computational Methods for Electron Molecule Collisions*, Plenum, New York, 1994.
- [38] See, e.g., F.A. Gianturco, R.R. Lucchese, N. Sanna, *J. Chem. Phys.* 104 (17) (1996) 6482.
- [39] See, e.g., F.A. Gianturco, R.R. Lucchese, N. Sanna, *J. Phys. B: At. Mol. Opt. Phys.* 32 (1999) 2181.
- [40] S.L. Altmann, The Cubic group, *Proc. Camb. Phil. Soc.* 53 (1957) 343.
- [41] F.A. Gianturco, R.R. Lucchese, N. Sanna, *J. Chem. Phys.* 100 (9) (1994) 1.
- [42] H.H.H. Homeier, E.O. Steinborn, *J. Mol. Struct. (Theochem)* 368 (1996) 31.
- [43] N. Sanna, F.A. Gianturco, SCELib API: an interface to the parallel computation of molecular properties, in preparation.
- [44] N. Sanna, Vector spherical harmonics: Concepts and applications to the single centre expansion method, *Comput. Phys. Comm.* 132 (2000) 66.
- [45] W.H. Press et al., *Numerical Recipes*, 2nd edn., Cambridge University Press, 1992.
- [46] See, e.g., J.M. Smith, F.W.J. Olver, D.W. Lozier, Extended-range arithmetic and normalized legendre polynomials, in: *ACM Trans. Math. Software*, March 1981, pp. 93–105.
- [47] R. Curik, F.A. Gianturco, N. Sanna, *J. Phys. B* (2000) to be published.

# Attachment 5A

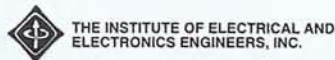


# Splash 2

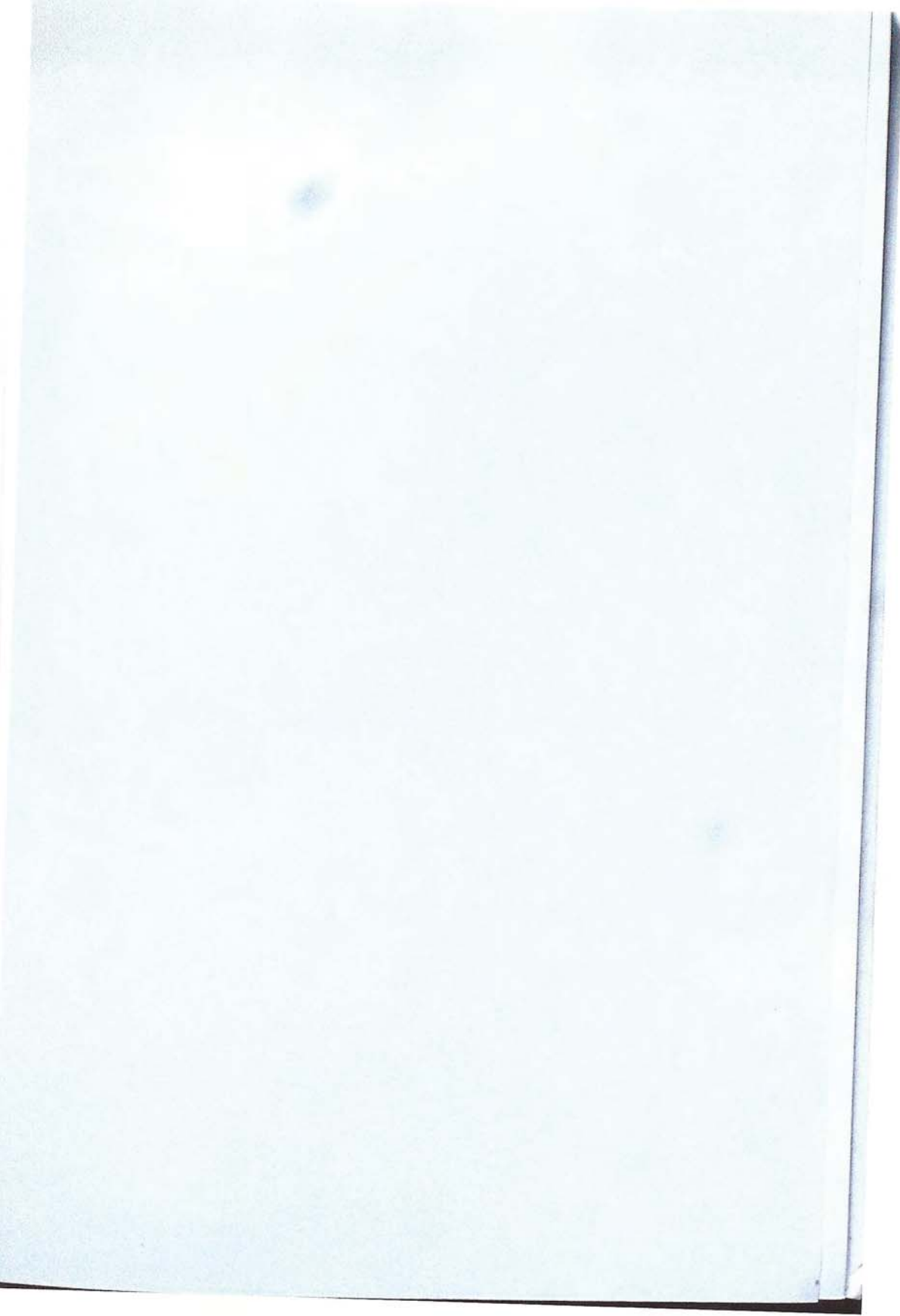
FPGAs in a Custom Computing Machine

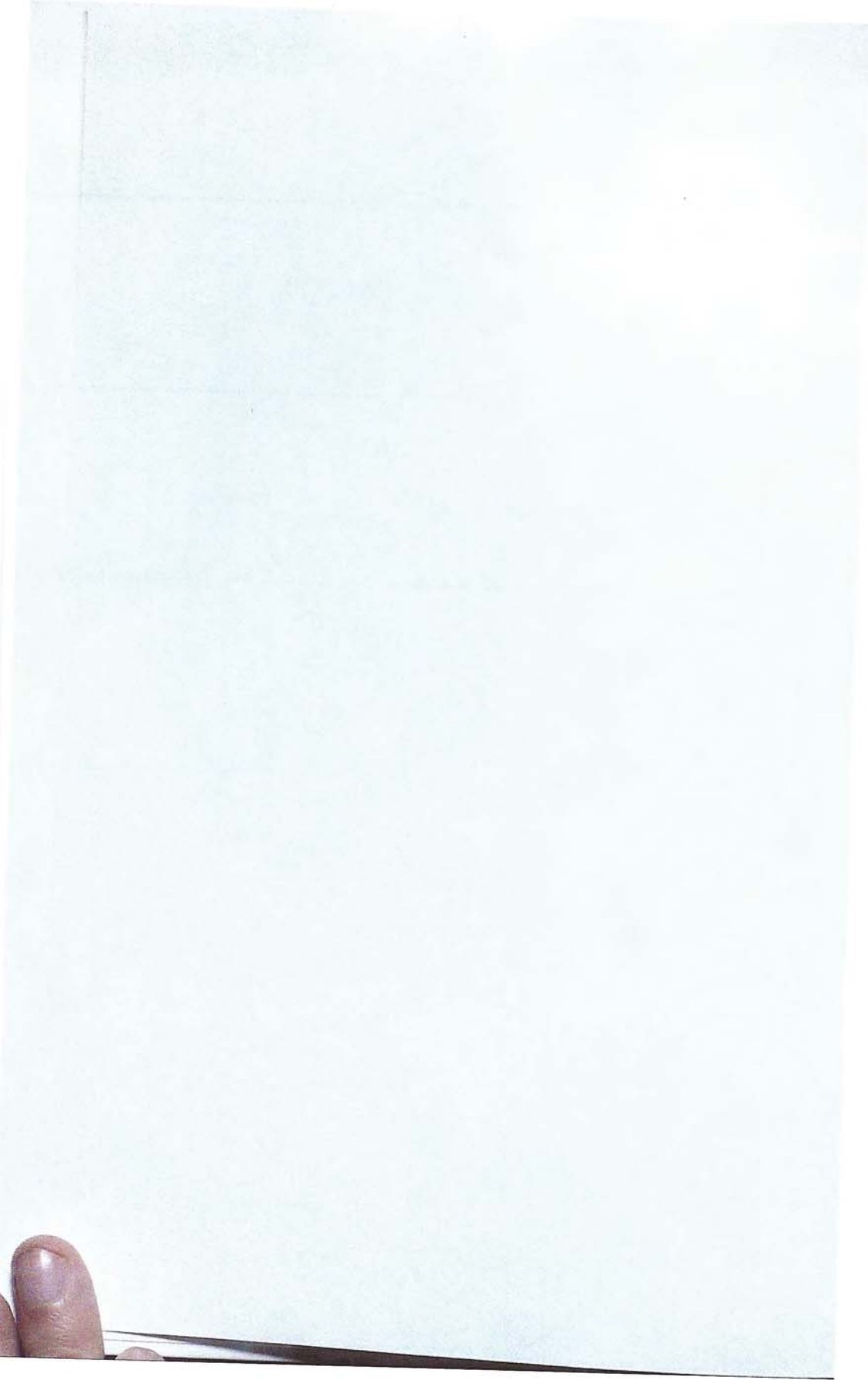


**Duncan A. Buell**  
**Jeffrey M. Arnold**  
**Walter J. Kleinfelder**









# **Splash 2**

## **FPGAs in a Custom Computing Machine**



Xilinx, the Xilinx logo, XC3090, XC4010, XBLOX, XACT, LCA, and Configurable Logic Cell are trademarks of Xilinx, Inc.  
CM-2 and Paris are trademarks of Thinking Machines Corporation.  
VMEbus is a trademark of Motorola Corporation.  
SPARC and SPARCstation are trademarks of SPARC International, Inc. Products bearing a SPARC trademark are based on an architecture developed by Sun Microsystems, Inc. SPARCstation is licensed exclusively to Sun Microsystems, Inc.  
UNIX is a trademark of UNIX System Laboratories.  
Sun, Sun Workstation, SunOS, and SBus are trademarks of Sun Microsystems, Inc.  
Design Compiler and FPGA Compiler are trademarks of Synopsys, Inc.  
DEC is a trademark of Digital Equipment Corporation.  
Verilog is a trademark of Cadence Design Systems, Inc.

# Splash 2

## FPGAs in a Custom Computing Machine

Duncan A. Buell  
Jeffrey M. Arnold  
Walter J. Kleinfelder  
*Editors*  
*Center for Computing Sciences*  
*Bowie, Maryland*



**IEEE Computer Society Press**  
Los Alamitos, California

Washington • Brussels • Tokyo

---

Library of Congress Cataloging-in-Publication Data

Buell, Duncan A.  
Splash 2: FPGAs in a custom computing machine / Duncan A. Buell,  
Jeffrey M. Arnold, Walter J. Kleinfelder.  
p. cm.  
Includes bibliographical references and index.  
ISBN 0-8186-7413-X  
1. Spash 2 (Computer) 2. Electronic digital computers—Design  
and construction. I. Arnold, Jeffrey M. II. Kleinfelder, Walter J.  
III. Title.  
QA76.8.S65B84 1996  
004.2 '2—dc20

95-47397  
CIP



IEEE Computer Society Press  
10662 Los Vaqueros Circle  
P.O. Box 3014  
Los Alamitos, CA 90720-1264

Copyright © 1996 by The Institute of Electrical and Electronics Engineers, Inc. All rights reserved.

Copyright and Reprint Permissions: Abstracting is permitted with credit to the source. Libraries are permitted to photocopy isolated pages beyond the limits of US copyright law, for private use of their patrons. Other copying, reprint, or republication requests should be addressed to: IEEE Copyrights Manager, IEEE Service Center, 445 Hoes Lane, P.O. Box 1331, Piscataway, NJ 08855-1331.

IEEE Computer Society Press Order Number BP07413  
Library of Congress Number 95-47397  
ISBN 0-8186-7413-X

Additional copies may be ordered from:

IEEE Computer Society Press  
Customer Service Center  
10662 Los Vaqueros Circle  
P.O. Box 3014  
Los Alamitos, CA 90720-1264  
Tel: +1-714-821-8380  
Fax: +1-714-821-4641  
Email: cs.books@computer.org

IEEE Service Center  
445 Hoes Lane  
P.O. Box 1331  
Piscataway, NJ 08855-1331  
Tel: +1-908-981-1393  
Fax: +1-908-981-9667  
mis.custserv@computer.org

IEEE Computer Society  
13, Avenue de l'Aquilon  
B-1200 Brussels  
BELGIUM  
Tel: +32-2-770-2198  
Fax: +32-2-770-8505  
euro.ofc@computer.org

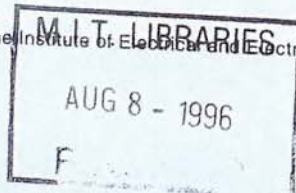
IEEE Computer Society  
Ooshima Building  
2-19-1 Minami-Aoyama  
Minato-ku, Tokyo 107  
JAPAN  
Tel: +81-3-3408-3118  
Fax: +81-3-3408-3553  
tokyo.ofc@computer.org

Assistant Publisher: Matt Loeb  
Technical Editor: Dharma P. Agrawal  
Acquisitions Assistant: Cheryl Smith  
Advertising/Promotions: Tom Fink  
Production Editor: Lisa O'Conner  
Cover Image: Dan Kopetzky, Center for Computing Sciences

Printed in the United States of America



The Institute of Electrical and Electronics Engineers, Inc



QA76.8  
.S65  
.B84  
1996

# Contents

---

<i>PREFACE</i>	<i>xi</i>
<b>1 CUSTOM COMPUTING MACHINES: AN INTRODUCTION</b>	<b>1</b>
1.1 Introduction	1
1.2 The Context for Splash 2	4
1.2.1 FPGAs	4
1.2.2 Architecture	5
1.2.3 Programming	6
<b>2 THE ARCHITECTURE OF SPLASH 2</b>	<b>10</b>
2.1 Introduction	10
2.2 The Building Blocks	11
2.3 The System Architecture	12
2.4 Data Paths	13
2.5 The Splash 2 Array Board	16
2.5.1 The Linear Array	16
2.5.2 The Splash 2 Crossbar	16
2.5.3 Xilinx Chip X0 and Broadcast Mode	17
2.6 The Interface Board and Control Features	17

<b>3</b>	<b>HARDWARE IMPLEMENTATION</b>	<b>19</b>
3.1	Introduction	19
3.2	Development Board Design	21
3.3	Interface Board Design	21
3.3.1	DMA Channel	23
3.3.2	XL and XR	23
3.3.3	Interrupts	24
3.3.4	Clock	24
3.3.5	Programming and Readback	24
3.3.6	Miscellaneous Registers	25
3.4	Array Board Design	25
3.4.1	Processing Element	26
3.4.2	Control Element	28
3.4.3	External Memory Access	28
3.4.4	Crossbar	28
3.4.5	Programming and Readback	29
3.4.6	Miscellaneous Registers	29
<b>4</b>	<b>SPLASH 2: THE EVOLUTION OF A NEW ARCHITECTURE</b>	<b>31</b>
4.1	Splash 1	31
4.2	Splash 2: Thoughts on a Redesign	34
4.3	Programming Language	36
4.4	Choice of FPGAs	37
4.5	Choice of Host and Bus	38
4.6	Chip-to-Chip Interconnections	39
4.7	Multitasking	42
4.8	Chip X0 and Broadcast	43
4.9	Other Design Decisions	43
<b>5</b>	<b>SOFTWARE ARCHITECTURE</b>	<b>46</b>
5.1	Introduction	46
5.2	Background	47
5.3	VHDL as a Programming Language	49
5.3.1	History and Purpose of VHDL	50
5.3.2	VHDL Language Features	50
5.3.3	Problems with VHDL	51
5.4	Software Environment	51

5.5	Programmer's View of Splash 2	55
5.5.1	Programming Process,	55
5.5.2	Processing Element View,	56
5.5.3	Interface Board View,	57
5.5.4	Host View,	57
<b>6</b>	<b>SOFTWARE IMPLEMENTATION</b>	<b>60</b>
6.1	Introduction	60
6.2	VHDL Environment	60
6.2.1	Splash 2 VHDL Library,	61
6.2.2	Standard Entity Declarations,	61
6.2.3	Programming Style,	64
6.3	Splash 2 Simulator	66
6.3.1	Structure,	66
6.3.2	Configuring the Simulator,	67
6.3.3	Input and Output,	68
6.3.4	Crossbar and Memory Models,	68
6.3.5	Hardware Constraints,	70
6.4	Compilation	70
6.4.1	Logic Synthesis,	70
6.4.2	Physical Mapping,	71
6.4.3	Debugging Support,	71
6.5	Runtime System	72
6.5.1	T2: A Symbolic Debugger,	72
6.5.2	Runtime Library,	73
6.5.3	Device Driver,	74
6.6	Diagnostics	75
<b>7</b>	<b>A DATA PARALLEL PROGRAMMING MODEL</b>	<b>77</b>
7.1	Introduction	78
7.2	Data-parallel Bit C	80
7.2.1	dbC Overview,	80
7.2.2	dbC Example,	81
7.3	Compiling from dbC to Splash 2	82
7.3.1	Creating a Specialized SIMD Engine,	83
7.3.2	Generic SIMD Code,	84
7.3.3	Generating VHDL,	84
7.4	Global Operations	88
7.4.1	Nearest-Neighbor Communication,	88

7.4.2	Reduction Operations, 89	
7.4.3	Host/Processor Communication, 91	
7.5	Optimization: Macro Instructions	92
7.5.1	Creating a Macro Instruction, 93	
7.5.2	Discussion, 94	
7.6	Evaluation: Genetic Database Search	94
7.7	Conclusions and Future Work	95
<b>8</b>	<b>SEARCHING GENETIC DATABASES ON SPLASH 2</b>	<b>97</b>
8.1	Introduction	97
8.1.1	Edit Distance, 98	
8.1.2	Dynamic Programming Algorithm, 98	
8.2	Systolic Sequence Comparison	100
8.2.1	Bidirectional Array, 100	
8.2.2	Unidirectional Array, 103	
8.3	Implementation	104
8.3.1	Modular Encoding, 105	
8.3.2	Configurable Parameters, 106	
8.3.3	Bidirectional Array, 107	
8.3.4	Unidirectional Array, 107	
8.4	Benchmarks	107
8.5	Discussion	108
8.6	Conclusions	108
<b>9</b>	<b>TEXT SEARCHING ON SPLASH 2</b>	<b>110</b>
9.1	Introduction	110
9.2	The Text Searching Algorithm	111
9.3	Description of the Single-Byte Splash Program	113
9.4	Timings, Discussion	114
9.5	Outline of the 16-bit Approach	115
9.6	Conclusions	116
<b>10</b>	<b>FINGERPRINT MATCHING ON SPLASH 2</b>	<b>117</b>
10.1	Introduction	117
10.2	Background	120

10.2.1	Pattern Recognition Systems, 121	
10.2.2	Terminology, 122	
10.2.3	Stages in AFIS, 123	
10.3	Splash 2 Architecture and Programming Models	125
10.4	Fingerprint Matching Algorithm	125
10.4.1	Minutia Matching, 126	
10.4.2	Matching Algorithm, 127	
10.5	Parallel Matching Algorithm	128
10.5.1	Preprocessing on the Host, 131	
10.5.2	Computations on Splash, 132	
10.5.3	VHDL Specification for X0, 133	
10.6	Simulation and Synthesis Results	134
10.7	Execution on Splash 2	137
10.7.1	User Interface, 137	
10.7.2	Performance Analysis, 137	
10.8	Conclusions	139
<b>11 HIGH-SPEED IMAGE PROCESSING WITH SPLASH 2</b>		<b>141</b>
11.1	Introduction	141
11.2	The VTSplash System	142
11.3	Image Processing Terminology and Architectural Issues	143
11.4	Case Study: Median Filtering	150
11.5	Case Study: Image Pyramid Generation	153
11.5.1	Gaussian Pyramid, 154	
11.5.2	Two Implementations for Gaussian Pyramid on Splash 2, 155	
11.5.3	The Hybrid Pipeline Gaussian Pyramid Structure, 157	
11.5.4	The Laplacian Pyramid, 157	
11.5.5	Implementation of the Laplacian Pyramid on Splash 2, 159	
11.6	Performance	159
11.7	Summary	163
<b>12 THE PROMISE AND THE PROBLEMS</b>		<b>166</b>
12.1	Some Bottom-Line Conclusions	166
12.1.1	High Bandwidth I/O Is a Must, 166	
12.1.2	Memory Is a Must, 167	
12.1.3	Programming Is Possible, and Becoming More So, 168	
12.1.4	The Programming Environment Is Crucial, 168	
12.2	To Where from Here?	169



12.3	If Not Splash 3, Then What?	171
12.3.1	Architectures,	172
12.3.2	Custom Processors,	173
12.3.3	Languages,	174
12.4	The “Killer” Applications	177
12.5	Final Words	178
<b>A</b>	<b>SPLASH 2 DEVELOPMENT—THE PROJECT MANAGER’S SUMMARY</b>	<b>179</b>
<b>B</b>	<b>AN EXAMPLE APPLICATION</b>	<b>186</b>
	<b>REFERENCES</b>	<b>190</b>

## Splash 2

### FPGAs in a Custom Computing Machine

*edited by Duncan A. Buell, Jeffrey M. Arnold, and Walter J. Kleinfelder*

Details the complete Splash 2 project—the hardware and software systems, their architecture and implementation, and the design process by which the architecture evolved from an earlier version machine. In addition to the description of the machine, this book explains why Splash 2 was engineered. It illustrates several applications in detail, allowing you to gain an understanding of the capabilities and the limitations of this kind of computing device.

The Splash 2 program is significant for two reasons. First, it is part of a complete computer system that achieves supercomputer like performance on a number of different applications. The second significant aspect is that this large system is capable of performing real computations on real problems. In order to understand what happens when the application programmer designs the processor architecture of the machine that executes his programs, it is necessary to see the system as a whole. This book looks in-depth at one of the handful of data points in the design space of this new kind of machine.

#### Contents:

- Custom Computing Machines: An Introduction
- The Architecture of Splash 2
- Hardware Implementation
- Splash 2: The Evolution of a New Architecture
- Software Architecture
- Software Implementation
- A Data Parallel Programming Model
- Searching Genetic Databases on Splash 2
- Text Searching on Splash 2
- Fingerprint Matching on Splash 2
- High-Speed Image Processing with Splash 2
- The Promise and the Problems
- An Example Application



Published by the IEEE Computer Society Press  
10662 Los Vaqueros Circle  
P.O. Box 3014  
Los Alamitos, CA 90720-1314

IEEE Computer Society Press Order Number BP07413  
Library of Congress Number 95-47397  
ISBN 0-8186-7413-X

ISBN 0-8186-7413-X



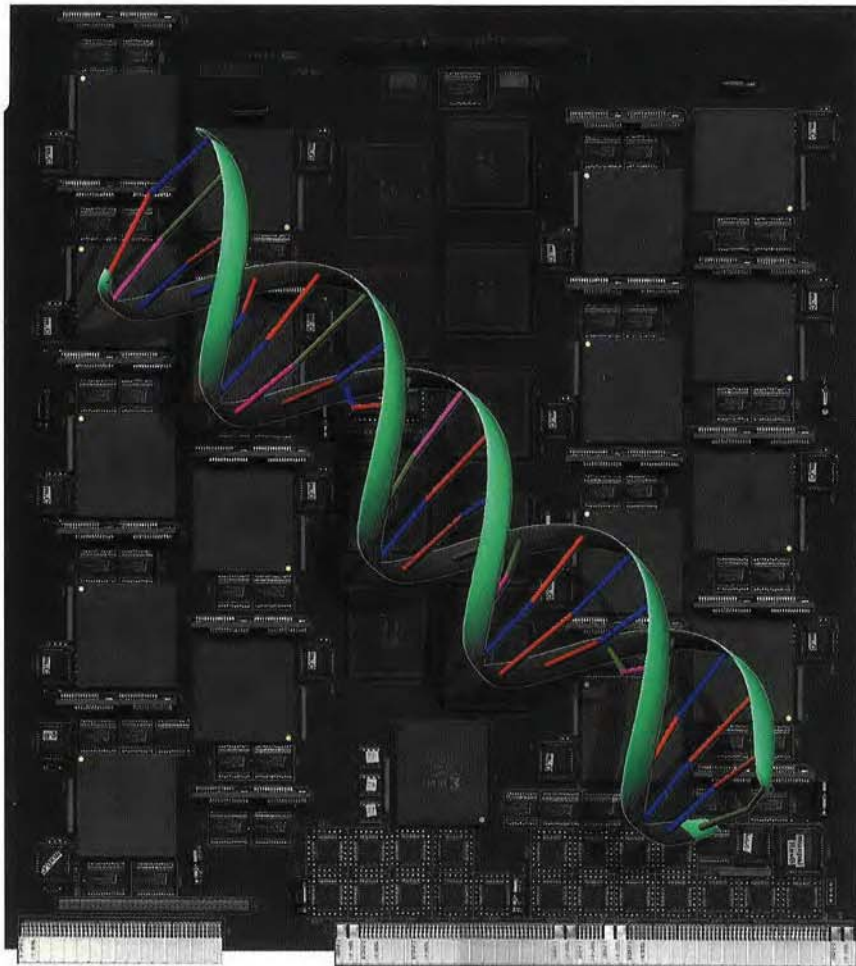
# Attachment 5B

QA 76  
.8  
.S68 S68  
1996

# Splash 2

**FPGAs in a Custom Computing Machine**

**LANDOVER  
GenColl**



**Duncan A. Buell  
Jeffrey M. Arnold  
Walter J. Kleinfelder**



# **Splash 2**

## **FPGAs in a Custom Computing Machine**

Xilinx, the Xilinx logo, XC3090, XC4010, XBLOX, XACT, LCA, and Configurable Logic Cell are trademarks of Xilinx, Inc.

CM-2 and Paris are trademarks of Thinking Machines Corporation.

VMEbus is a trademark of Motorola Corporation.

SPARC and SPARCstation are trademarks of SPARC International, Inc. Products bearing a SPARC trademark are based on an architecture developed by Sun Microsystems, Inc. SPARCstation is licensed exclusively to Sun Microsystems, Inc.

UNIX is a trademark of UNIX System Laboratories.

Sun, Sun Workstation, SunOS, and SBus are trademarks of Sun Microsystems, Inc.

Design Compiler and FPGA Compiler are trademarks of Synopsys, Inc.

DEC is a trademark of Digital Equipment Corporation.

Verilog is a trademark of Cadence Design Systems, Inc.

# Splash 2

## FPGAs in a Custom Computing Machine

Duncan A. Buell  
Jeffrey M. Arnold  
Walter J. Kleinfelder  
*Editors*  
*Center for Computing Sciences*  
*Bowie, Maryland*



**IEEE Computer Society Press**  
Los Alamitos, California

Washington • Brussels • Tokyo

---



**Library of Congress Cataloging-in-Publication Data**

Buell, Duncan A.

Splash 2: FPGAs in a custom computing machine / Duncan A. Buell,  
Jeffrey M. Arnold, Walter J. Kleinfelder.

p. cm.

Includes bibliographical references and index.

ISBN 0-8186-7413-X

1. Spash 2 (Computer) 2. Electronic digital computers—Design  
and construction. I. Arnold, Jeffrey M. II. Kleinfelder, Walter J.  
III. Title.

QA76.8.S65B84 1996

004.2 '2—dc20

95-47397  
CIP



IEEE Computer Society Press  
10662 Los Vaqueros Circle  
P.O. Box 3014  
Los Alamitos, CA 90720-1264

Copyright © 1996 by The Institute of Electrical and Electronics Engineers, Inc. All rights reserved.

*Copyright and Reprint Permissions:* Abstracting is permitted with credit to the source. Libraries are permitted to photocopy isolated pages beyond the limits of US copyright law, for private use of their patrons. Other copying, reprint, or republication requests should be addressed to: IEEE Copyrights Manager, IEEE Service Center, 445 Hoes Lane, P.O. Box 1331, Piscataway, NJ 08855-1331.

IEEE Computer Society Press Order Number BP07413  
Library of Congress Number 95-47397  
ISBN 0-8186-7413-X

*Additional copies may be ordered from:*

IEEE Computer Society Press  
Customer Service Center  
10662 Los Vaqueros Circle  
P.O. Box 3014  
Los Alamitos, CA 90720-1264  
Tel: +1-714-821-8380  
Fax: +1-714-821-4641  
Email: cs.books@computer.org

IEEE Service Center  
445 Hoes Lane  
P.O. Box 1331  
Piscataway, NJ 08855-1331  
Tel: +1-908-981-1393  
Fax: +1-908-981-9667  
mis.custserv@computer.org

IEEE Computer Society  
13, Avenue de l'Aquilon  
B-1200 Brussels  
BELGIUM  
Tel: +32-2-770-2198  
Fax: +32-2-770-8505  
euro.ofc@computer.org

IEEE Computer Society  
Ooshima Building  
2-19-1 Minami-Aoyama  
Minato-ku, Tokyo 107  
JAPAN  
Tel: +81-3-3408-3118  
Fax: +81-3-3408-3553  
tokyo.ofc@computer.org

Assistant Publisher: Matt Loeb  
Technical Editor: Dharma P. Agrawal  
Acquisitions Assistant: Cheryl Smith  
Advertising/Promotions: Tom Fink  
Production Editor: Lisa O'Conner  
Cover Image: Dan Kopetzky, Center for Computing Sciences

Printed in the United States of America



The Institute of Electrical and Electronics Engineers, Inc

QA76  
.8  
.S68 S68  
1996

# Contents

---

<b>PREFACE</b>	<b>xi</b>
<b>1 CUSTOM COMPUTING MACHINES: AN INTRODUCTION</b>	<b>1</b>
1.1 Introduction	1
1.2 The Context for Splash 2	4
1.2.1 FPGAs,	4
1.2.2 Architecture,	5
1.2.3 Programming,	6
<b>2 THE ARCHITECTURE OF SPLASH 2</b>	<b>10</b>
2.1 Introduction	10
2.2 The Building Blocks	11
2.3 The System Architecture	12
2.4 Data Paths	13
2.5 The Splash 2 Array Board	16
2.5.1 The Linear Array,	16
2.5.2 The Splash 2 Crossbar,	16
2.5.3 Xilinx Chip X0 and Broadcast Mode,	17
2.6 The Interface Board and Control Features	17

v

<b>3</b>	<b><i>HARDWARE IMPLEMENTATION</i></b>	<b>19</b>
3.1	Introduction	19
3.2	Development Board Design	21
3.3	Interface Board Design	21
3.3.1	DMA Channel,	23
3.3.2	XL and XR,	23
3.3.3	Interrupts,	24
3.3.4	Clock,	24
3.3.5	Programming and Readback,	24
3.3.6	Miscellaneous Registers,	25
3.4	Array Board Design	25
3.4.1	Processing Element,	26
3.4.2	Control Element,	28
3.4.3	External Memory Access,	28
3.4.4	Crossbar,	28
3.4.5	Programming and Readback,	29
3.4.6	Miscellaneous Registers,	29
<b>4</b>	<b><i>SPLASH 2: THE EVOLUTION OF A NEW ARCHITECTURE</i></b>	<b>31</b>
4.1	Splash 1	31
4.2	Splash 2: Thoughts on a Redesign	34
4.3	Programming Language	36
4.4	Choice of FPGAs	37
4.5	Choice of Host and Bus	38
4.6	Chip-to-Chip Interconnections	39
4.7	Multitasking	42
4.8	Chip X0 and Broadcast	43
4.9	Other Design Decisions	43
<b>5</b>	<b><i>SOFTWARE ARCHITECTURE</i></b>	<b>46</b>
5.1	Introduction	46
5.2	Background	47
5.3	VHDL as a Programming Language	49
5.3.1	History and Purpose of VHDL,	50
5.3.2	VHDL Language Features,	50
5.3.3	Problems with VHDL,	51
5.4	Software Environment	51

5.5	Programmer's View of Splash 2	55
5.5.1	Programming Process,	55
5.5.2	Processing Element View,	56
5.5.3	Interface Board View,	57
5.5.4	Host View,	57
<b>6</b>	<b>SOFTWARE IMPLEMENTATION</b>	<b>60</b>
6.1	Introduction	60
6.2	VHDL Environment	60
6.2.1	Splash 2 VHDL Library,	61
6.2.2	Standard Entity Declarations,	61
6.2.3	Programming Style,	64
6.3	Splash 2 Simulator	66
6.3.1	Structure,	66
6.3.2	Configuring the Simulator,	67
6.3.3	Input and Output,	68
6.3.4	Crossbar and Memory Models,	68
6.3.5	Hardware Constraints,	70
6.4	Compilation	70
6.4.1	Logic Synthesis,	70
6.4.2	Physical Mapping,	71
6.4.3	Debugging Support,	71
6.5	Runtime System	72
6.5.1	T2: A Symbolic Debugger,	72
6.5.2	Runtime Library,	73
6.5.3	Device Driver,	74
6.6	Diagnostics	75
<b>7</b>	<b>A DATA PARALLEL PROGRAMMING MODEL</b>	<b>77</b>
7.1	Introduction	78
7.2	Data-parallel Bit C	80
7.2.1	dbC Overview,	80
7.2.2	dbC Example,	81
7.3	Compiling from dbC to Splash 2	82
7.3.1	Creating a Specialized SIMD Engine,	83
7.3.2	Generic SIMD Code,	84
7.3.3	Generating VHDL,	84
7.4	Global Operations	88
7.4.1	Nearest-Neighbor Communication,	88

7.4.2	Reduction Operations, 89	
7.4.3	Host/Processor Communication, 91	
7.5	Optimization: Macro Instructions	92
7.5.1	Creating a Macro Instruction, 93	
7.5.2	Discussion, 94	
7.6	Evaluation: Genetic Database Search	94
7.7	Conclusions and Future Work	95
<b>8</b>	<b>SEARCHING GENETIC DATABASES ON SPLASH 2</b>	<b>97</b>
8.1	Introduction	97
8.1.1	Edit Distance, 98	
8.1.2	Dynamic Programming Algorithm, 98	
8.2	Systolic Sequence Comparison	100
8.2.1	Bidirectional Array, 100	
8.2.2	Unidirectional Array, 103	
8.3	Implementation	104
8.3.1	Modular Encoding, 105	
8.3.2	Configurable Parameters, 106	
8.3.3	Bidirectional Array, 107	
8.3.4	Unidirectional Array, 107	
8.4	Benchmarks	107
8.5	Discussion	108
8.6	Conclusions	108
<b>9</b>	<b>TEXT SEARCHING ON SPLASH 2</b>	<b>110</b>
9.1	Introduction	110
9.2	The Text Searching Algorithm	111
9.3	Description of the Single-Byte Splash Program	113
9.4	Timings, Discussion	114
9.5	Outline of the 16-bit Approach	115
9.6	Conclusions	116
<b>10</b>	<b>FINGERPRINT MATCHING ON SPLASH 2</b>	<b>117</b>
10.1	Introduction	117
10.2	Background	120

10.2.1	Pattern Recognition Systems, 121	
10.2.2	Terminology, 122	
10.2.3	Stages in AFIS, 123	
10.3	Splash 2 Architecture and Programming Models	125
10.4	Fingerprint Matching Algorithm	125
10.4.1	Minutia Matching, 126	
10.4.2	Matching Algorithm, 127	
10.5	Parallel Matching Algorithm	128
10.5.1	Preprocessing on the Host, 131	
10.5.2	Computations on Splash, 132	
10.5.3	VHDL Specification for X0, 133	
10.6	Simulation and Synthesis Results	134
10.7	Execution on Splash 2	137
10.7.1	User Interface, 137	
10.7.2	Performance Analysis, 137	
10.8	Conclusions	139
<b>11 HIGH-SPEED IMAGE PROCESSING WITH SPLASH 2</b>		<b>141</b>
11.1	Introduction	141
11.2	The VTSplash System	142
11.3	Image Processing Terminology and Architectural Issues	143
11.4	Case Study: Median Filtering	150
11.5	Case Study: Image Pyramid Generation	153
11.5.1	Gaussian Pyramid, 154	
11.5.2	Two Implementations for Gaussian Pyramid on Splash 2, 155	
11.5.3	The Hybrid Pipeline Gaussian Pyramid Structure, 157	
11.5.4	The Laplacian Pyramid, 157	
11.5.5	Implementation of the Laplacian Pyramid on Splash 2, 159	
11.6	Performance	159
11.7	Summary	163
<b>12 THE PROMISE AND THE PROBLEMS</b>		<b>166</b>
12.1	Some Bottom-Line Conclusions	166
12.1.1	High Bandwidth I/O Is a Must, 166	
12.1.2	Memory Is a Must, 167	
12.1.3	Programming Is Possible, and Becoming More So, 168	
12.1.4	The Programming Environment Is Crucial, 168	
12.2	To Where from Here?	169

12.3	If Not Splash 3, Then What?	171
12.3.1	Architectures,	172
12.3.2	Custom Processors,	173
12.3.3	Languages,	174
12.4	The “Killer” Applications	177
12.5	Final Words	178
<b>A</b>	<b>SPLASH 2 DEVELOPMENT—THE PROJECT MANAGER’S SUMMARY</b>	<b>179</b>
<b>B</b>	<b>AN EXAMPLE APPLICATION</b>	<b>186</b>
	<b>REFERENCES</b>	<b>190</b>

# Preface

---

The Splash 2 project began at the Supercomputing Research Center<sup>1</sup> in September of 1991 and ended, with success, in the spring of 1994. Splash 2 is an attached processor system using Xilinx XC4010 FPGAs as its processing elements. As such, it is a *custom computing machine*. That is to say that much of what would be the instruction set architecture of the processing elements is not specified except in the details of the program developed by the application programmer. Although a higher-level block diagram of processing elements, memories, interconnect, and dataflow exists in the hardware structure of Splash 2, the details of the instruction set architecture level of the machine will vary from one application to the next.

The Splash 2 project is significant for two reasons. First, Splash 2 is part of a complete computer system that has achieved supercomputer-like performance on a number of different applications. By "complete computer system" we mean that SRC created or caused to be created an extant hardware system (replicated a dozen times), a complete programming and runtime environment, and a collection of application programs that exercised the unique hardware.

The second significant aspect of Splash 2 is that we were fortunate enough to be able to build a large system, capable of performing real computations on real problems. One common complaint about performance results on novel computing machines or environments is that results on small problems cannot be accurately extrapolated to large problems. The Splash 2 system that was designed and built is a full-sized machine and does not suffer from this defect.

To get to the point: why a book?

<sup>1</sup>Renamed the Center for Computing Sciences in May 1995, but referred to throughout this book as SRC.



This is a novel computing machine. In order to understand what happens when the application programmer is permitted, indeed required, to design the processor architecture of the machine that will execute his program, it is necessary to see the system as a whole. Programmability and problems to be run on this machine both had major influences on its architecture, just as its architecture and its unique nature influence the kinds of problems one would expect to program for this machine and the nature of that programming. And standing between the user and the machine, as the old joke goes, is a new kind of programming environment and an evolving understanding of how this environment must allow the use of the hardware, without forcing every programmer to be a hardware-design engineer.

At the first IEEE Workshop on FPGAs for Custom Computing Machines, one of the industrial attendees remarked that, although nearly everyone would agree, as part of a coffee-room discussion or the like, that it would be *interesting* to think about building a "computer" using FPGAs, no one in management (except perhaps at SRC and DEC PRL) had put up the commitment necessary in time and money to do a real test of the idea. It was then remarked that, given the nature of the marketplace and of engineering management, these first attempts had to be successful in order to open the door for future work. We feel we have been successful, and we offer in this book an in-depth look at one of the handful of data points in the design space of this new kind of machine.

We would hope that this book has a broad appeal and is readable with understanding by nearly all computer scientists and computer engineers. To the hardware designer, perhaps we can offer a new look at programming applications on a moderately general FPGA-based computing machine instead of designing circuits for a specific board incorporating FPGAs. The engineering world has viewed FPGAs, to a great extent, as the next logical step in a continuum of electronic devices; we offer, we feel, the option of viewing them much more broadly than that. To the computer architect we offer a variant hardware platform and an indication of how that general platform can be used. Much of computer architecture is a compromise between the functionality desired and the limits of what can be built given existing technology; we offer the use of a new technology that can offer, to a limited extent now, and could offer much more generally later, greatly increased functionality. For some of those who have hard problems in computation, we offer much of the power of special-purpose hardware without the inflexibility and uncertain delivery times of hardware. The long-term task is not to map a high-level language to a particular architecture or range of architectures, but in some sense to create for each application program a suitable architecture to which the high-level language will be mapped. And to the language designers and compiler writers we offer a world to conquer. We have presented one imperfect but usable approach to programming such a computing machine, and we trust that others interested in the critical problem of making these machines programmable can learn both from what we did right and what we did wrong.

Chapter 1 discusses the general concept of Custom Computing Machines, of which Splash 2 is one example. Chapters 2 and 3 describe at a high level and then in some detail the hardware architecture of Splash 2. Chapter 4 covers the design considerations and decisions in arriving at the second-generation Splash 2 architecture. We present this chapter at the end of the section on hardware, on the basis that it is easier to understand variations in a design when those variations are compared against something concrete.

Chapters 5 and 6 describe, also first at a higher and then at a lower level, the software architecture of Splash 2. All the application programs in the latter chapters were done using VHDL as an applications programming language and these tools in support. The main goal of the Splash 2 project was to show that software, as described in Chapters 5 and 6, could make a computer using FPGAs as its processing elements into something that reasonable people would call “programmable,” and, in that sense, the heart of the Splash 2 project is in Chapters 5 and 6. Throughout the life of Splash 2, however, there has been an alternative view of programming. This view is reflected in Chapter 7 on the Splash 2 version of the programming language dbC. The approach taken in dbC is to permit the programmer to use a version of C as the programming language. It is the compiler which then becomes responsible for, in essence, recognizing the instruction set architecture necessary to execute the program and then creating in the FPGA the requisite registers, logic units, and control.

Chapters 8 through 11 then describe four different applications programmed to conclusion on Splash 2. The first of these—the sequence comparison problem—was the driving application, in the sense that funding for Splash 2 was based on its perceived ability to perform this computation. This and the text processing application were done at SRC.

The Splash 2 project team was fortunate in that SRC’s parent company, the Institute for Defense Analyses, issued two contracts, to Virginia Polytechnic Institute and to Michigan State University, for applications work on Splash 2 in image processing and fingerprint identification. Both applications seemed good matches with the Splash 2 architecture but lay outside the normal realm of SRC’s research program. The faculty members involved have each prepared a chapter on these applications.

We close in Chapter 12 with some opinions and speculations about the future. In an appendix, the project manager presents a chronology of the entire Splash 2 project.

It is incumbent on us, and a genuine pleasure, to thank the Center for Computing Sciences of the Institute for Defense Analyses and the CCS Director, Francis Sullivan, for supporting us in our writing and editing of this book. All royalties will be donated to the Center for Excellence in Education, formerly known as the Rickover Foundation, in McLean, Virginia. The Center for Excellence in Education supports science and engineering education through its sponsorship of the Research Science Institute each summer for high school seniors, its Role Models and Leaders Project in Washington, D.C., Los Angeles, and Chicago for promising women and minority high school students intending to study science and engineering, and its support and preparation of the United States Informatics Olympiad team each year.

The Splash 2 project was a success in large part due to the ability of those who were involved nearly full-time, but it might not have taken the course it did had the hard-core Splash 2 players not had the chance to get advice and occasional help from a much larger group of experts both at SRC and elsewhere.

We acknowledge, therefore, the help and advice of Nate Bronson, Dan Burns, Bill Carlson, Neil Coletti, Maripat Corr, Steve Cuccaro, Hillory Dean, Chuck Fiducia, Brad Fross, Charles Goedeke, Maya Gokhale, Frank Hady, Dzung Hoang, Bill Holmes, Amy Johnston, Elaine (Davis) Keith, Dan Kopetzky, Andy Kopser, Steve Kratzer, Jim Kuehn, Sara Lucas, Michael Mascagni, Marge McGarry, John McHenry,

Ron Minnich, Lindy Moran, Fred More, Mark Norder, Lou Podrazik, Dan Pryor, Craig Reese, Paul Schneck, Brian Schott, Nabeel Shirazi, Doug Sweely, Dennis Sysko, Mark Thistle, Bob Thompson, Ken Wallgren, Alice Yuen, Neal Ziring, and Jennifer Zito.

*Duncan A. Buell*  
*Jeffrey M. Arnold*  
*Walter J. Kleinfelder*  
Bowie, Maryland  
March 1996

# CHAPTER 1

---

## Custom Computing Machines: An Introduction

*Duncan A. Buell<sup>1</sup>*

### 1.1 INTRODUCTION

It is a basic observation about computing that generality and efficiency are in some sense inversely related to one another; the more general-purpose an object is and thus the greater the number of tasks it can perform, the less efficient it will be in performing any of those specific tasks. Design decisions are therefore almost always compromises; designers identify key features or applications for which competitive efficiency is a must and then expand the range as far as is practicable without unduly damaging performance on the main targets.

This thesis has certainly been true in processor architecture of computers aimed at computationally intense problems. Vector processors and vector supercomputers have targeted computationally intense, array-oriented floating point problems, usually in the hard sciences and engineering, but have not sacrificed the necessary speed on their core applications in order to be all things to all people. Thus, on computationally intense problems that do not fit well on traditional supercomputers, perhaps due to such things as integer arithmetic or scalar code, fast workstations can often outperform supercomputers.

To counter the problem of computationally intense problems for which general-purpose machines cannot achieve the necessary performance, special-purpose processors, attached processors, and coprocessors have been built for many years, especially

<sup>1</sup>A version of this chapter appeared as Buell and Pocek [11] and is used with permission.

in such areas as image or signal processing (for which many of the computational tasks can be very well defined). The problem with such machines is that they are special-purpose; as problems change or new ideas and techniques develop, their lack of flexibility makes them problematic as long-term solutions.

Enter the FPGA. Field Programmable Gate Arrays, first introduced in 1986 by Xilinx [34], were seen rather immediately by a few people to offer a totally new avenue to explore in the world of processor engineering. The great strength of the computer as a tool is in its ability to be adapted, via programming, to a multitude of computational tasks. The possibility now existed for an FPGA-based computing device not only to be configured to act like special-purpose hardware for a particular problem, but to be reconfigured for different problems and for this reconfiguration to be a programming process. By being more than single-purpose, such a machine would have the advantage of being flexible with at least a limited range of different applications. By being programmable, such a machine would open up "design of high-performance hardware" to individuals who can "design hardware" in an abstract sense but not a concrete sense. Finally, by being designed to operate as if they were hardware, the applications for these machines can achieve the hardware-like performance one gets from having explicitly parallel computations, from not having instructions and data fetched and decoded, and from having the ability to design processing units that reflect precisely the processing being done.

It is no exaggeration to say that machines using FPGAs as their processing elements have demonstrated that very high performance on an absolute scale, and extraordinary performance when measured against price, is possible with this technology. The PerLe machines built at DEC's Paris Research Laboratory have been programmed on a number of applications with impressive results [7, 8, 9, 31]: An implemented multiplier can compute a 50-coefficient, 16-bit polynomial convolution FIR filter at 16 times audio real time. An implementation of RSA decryption executes at 10 times the bit rate of an AT&T ASIC. A Hough Transform implementation for an application in high-energy physics achieves a compute rate that a standard 64-bit computer could not equal without a 1.2 GHz clock rate.

As can be seen from the later chapters of this book, some of the applications programmed on Splash 2 have achieved similarly promising results. It was a general observation made by those involved in the Splash 2 project that, on applications that fit the machine, one Splash 2 Array Board could deliver approximately the compute power of one Cray YMP processor. One of the commercial licensees of the Splash 2 technology sells its system for about \$40,000. Of course, not all applications fit well, and most that do not fit well actually fit very badly indeed, but this is nonetheless a performance-to-price ratio substantial enough to warrant continued investigation and experimentation.

The idea of reconfiguring a computer is certainly not new. The Burroughs B1700 had multiple instruction sets with different targets (Fortran, COBOL, and such) implemented with different microcode. In another way, the Paris functions of the Thinking Machines Corp. CM-2 differed from one version to the next. In the former case, standard views of hardware instructions were implemented. In the latter case, with a novel machine and a new architecture, we presume that an effort was made to implement function calls that users were seen to need and to delete unneeded functions when the instruction store ran out.

A certain amount of disagreement exists over what label to give to these machines and how to refer to them. The group at DEC's Paris research lab refers to their machine as a Programmable Active Memory (PAM) [7, 9, 29, 31]. Another commercial entity uses the term "virtual computer" [12]. From Brown University we have Processor Reconfiguration through Instruction Set Metamorphosis (PRISM) [4, 5, 6, 32]. We have already used the term "FPGA-based computing device" here. That none of these are truly satisfactory was evidenced in the spring of 1994 when the `comp.arch.fpga` newsgroup was discussed and established; the most contentious point was over the name. Both the new newsgroup and the IEEE workshops we have organized use the term FPGA, thus being in some sense bound in terminology to a particular technology (unless, of course, one can convince the developers of the next technology that its name should allow FPGA as its acronym). We have chosen to use the term Custom Computing Machine (CCM). None of these terms is perfect, but we believe that this one is no worse than any of the others.

The work on CCMs also differs from what is considered reconfigurable computer architecture, in that the term "reconfigurable" usually refers to a much higher level of the system. In the case of CCMs, that which is reconfigurable and significant by virtue of its reconfigurability is the "processor architecture" itself. A reconfigurable computer, by contrast, is likely to be either a multiprocessor in which the interconnections among the processors can change, or a heterogeneous machine with processors of different kinds that a user can choose to include or exclude in the view he/she has of "the computer."

It is to be emphasized that this is not a mature computing technology and that CCMs are not a panacea for all problems in high-performance computing. Among the many issues and problems are the following:

1. Are FPGAs large enough, or will they become large enough, so that a significant unit of computation can be put on a single chip so as to avoid both the loss of efficiency in going off-chip and the problems in partitioning a computation across multiple chips?
2. With current technology, even in the best of circumstances, the user must be exposed to the hardware itself. What is the level of understanding about chip architecture, signal routing delays, and so forth, that a "programmer" must know in order to use a CCM? How much more must be known in order to obtain the performance that would warrant using a CCM instead of a general-purpose computer?
3. If these machines are to be viewed as "computers," then they must be capable of being programmed. How will this be done? What sort of programming language is appropriate? How do we "compile" when we have eliminated the underlying machine model?
4. Granting the point that these are limited-purpose, but not special-purpose, machines, what is the range of architectures needed to cover the spectrum of applications for which these machines make sense?
5. What are the cost/performance figures necessary to make this a viable approach for getting a computing task done? General-purpose machines are cheaper and easier to use but can be slow. ASICs and special-purpose devices are faster but more expensive in small quantities, take longer to develop, and are hard to modify. Where might CCMs fit between these two?

One problem faced by those involved in Splash 1 and Splash 2 at the Supercomputing Research Center<sup>2</sup> has been a stubborn refusal from some quarters to believe that achieving high performance on a CCM is possible without a design or programming agony so great as to be offputting to all but the most dedicated of "application designers/programmers." Even in the face of our evidence to the contrary, the case has been difficult to make to some critics.

The case is especially hard because what is needed is to build a complete hardware system, to create or cause to be created a programming environment worthy of being called a programming environment, and then to develop a variety of different applications so that the proof of concept is complete. Further, since the goal is to demonstrate a competitive performance with more expensive and sophisticated machines, the CCM must be big enough to do real work and to be part of complete computational processing environments; it cannot be just a toy machine suitable only for doing kernels of problems. To our knowledge, only two such machines have been built that meet these criteria—Splash 2 and the larger of the DEC PAM systems.

The goal in this book is to present the Splash 2 experience as a whole. Splash 2 was designed and developed in an iterative process from top to bottom to top and back again.

## 1.2 THE CONTEXT FOR SPLASH 2

### 1.2.1 FPGAs

FPGAs in general have a wide spectrum of characteristics, but the FPGAs used for CCMs have been of two distinct types. The Xilinx XC4010, a typical example of one type, is a chip containing a  $20 \times 20$  square array of Configurable Logic Blocks (CLBs) [34]. Each CLB can serve one of three functions, either as two flipflops, or as Boolean functions of nine inputs, or as 32 bits of RAM. The function use has two 4-input functions, each producing an output; these two bits combine with a ninth input in a 3-input Boolean function. The RAM usage simply takes advantage of the fact that the 4-input functions are done with lookup tables to allow the input bits to be viewed as addresses.

Connecting the CLBs to one another and to special Input Output Blocks (IOBs) on the periphery of the chip are routing resources running from CLB to CLB, skipping one CLB, or running the full length of the chip. Configuration of the FPGA is done by loading a bit file onto on-chip RAM.

In contrast to the relatively coarse granularity of the Xilinx chips, the second type of FPGA, by Algotronix, Ltd., and by Concurrent Logic, Inc. [1, 13], is fine-grained. The Algotronix chip is a  $32 \times 32$  array of 2-input, 1-output Boolean function logic cells, with the signal lines running only point to point from one cell to its neighbors in each rectangular direction.<sup>3</sup>

<sup>2</sup>The Supercomputing Research Center was renamed the Center for Computing Sciences in May 1995, but will be referred to throughout this book as SRC.

<sup>3</sup>Algotronix is now a part of Xilinx.

To a first-order approximation, the chips first marketed by Concurrent Logic and now by Atmel<sup>4</sup> resemble the Algotronix chip. Interestingly enough, the unscientific best-guess estimates at the SRC in the fall of 1991, when Splash 2 was being designed, suggested that the then-high-end XC4010 and Concurrent Logic CLi6000 chips had roughly the same “compute power” in spite of the radically different architectures.

### 1.2.2 Architecture

There have been several architectures proposed and built for CCMs. Although any taxonomy runs the risk of pigeonholing some particular machine into a category distasteful to its designer, the following is a reasonable categorization.

**Special-Purpose Devices.** The first and most obvious use of FPGAs for CCMs is in special-purpose machines built to perform a particular computation or kind of computation and not intended to be very flexible, except perhaps from one instance of the problem to the next. There have been several machines built for neural network computations (Ganglion, for example [15]). Here, the computation is clearly parallel, the individual compute nodes are neither standard nor very large, and one feature of neural nets is that a moderately high degree of connectivity is desired among the compute nodes; but the precise connectivity and multiplier constants at each compute node vary from application to application. Other applications for which special-purpose devices have been built include statistical physics, embedded control, and network control [14, 19, 25, 33].

Somewhat more general than a special-purpose device, but still very much in a narrow band of applications, is the use of an FPGA-based computer for rapid prototyping, not just of ASICs or of single circuit boards, but of an entire system. A CCM can be a complete system—processors, memory, data path, and so on—at the block diagram level, and the characteristics and details needed can be programmed into functioning hardware. Similarly, in an appropriate niche market, a CCM could be used in low-volume applications, cheaper in development cost than special-purpose hardware but faster than what one could obtain from a programmed microprocessor.

**Coprocessors.** One of the most tantalizing possible uses for FPGAs as compute resources is as coprocessors tightly coupled to the main processor of a computer. The development of RISC processors has meant that some instructions that used to be part of a processor’s repertoire are no longer present; these functions must now be performed in software routines that are inherently slower. Some computations have natural kernels that have never been part of the instruction set architecture of any processor. Much of the PRISM [4, 5, 6, 32] work has focused on two points: 1) the language, compiler, and system issues involved in determining that a particular core computation occurs frequently enough that it warrants being put onto the coprocessor, and 2) arranging the computation so that “hardware” exists in the FPGA coprocessor when it is needed and that data can be transferred to and from the coprocessor at speeds great enough to make use of the coprocessor worthwhile. Several such machines are described in [17, 18, 21, 22, 23, 24, 30]. In a similar vein, the SRC worked with Thinking Machines Corp. on the production of the “CM-2X,”

<sup>4</sup>National Semiconductor also had rights to and sold a version of this chip.



described by Cuccaro and Reese [16], a 512-node CM-2 in which Xilinx XC3090 chips replaced the floating-point chips as coprocessors.

**Attached Processors.** As we have said, there have been two notable examples of FPGA-based attached processors—the PeRLe-0 and PeRLe-1 machines built by DEC's Paris research lab [8, 9, 29, 31], and the Splash 1 and Splash 2 machines built at SRC [2, 3, 20, 26]. The PeRLe-1 board featured 23 Xilinx XC3090 chips, with the core computational unit being a  $4 \times 4$  grid, connected by a TURBOchannel to a DEC workstation. Splash 2, in contrast, is used primarily either as a linear array of 16 XC4010 chips per board or with data being broadcast to the 16 chips simultaneously. Both have achieved supercomputer performance on a range of applications including image processing, computational science benchmarks, data compression and encryption benchmarks, and molecular biology.

Some machines, which have for one reason or another been built with a particular purpose in mind, are general enough that they would find wider application. The CHAMP machine, described by Box [10], designed for image processing, is certainly among these. Other examples are described in Quenot et al. [27] and Raimbault et al. [28].

### 1.2.3 Programming

Notwithstanding the tremendous effort necessary to engineer the hardware of a CCM, the fundamental test of these machines is, and no doubt will continue to be, a software problem. Regardless of the architecture or the potential peak performance of the machine, if the effort to achieve that peak requires either an extraordinarily important problem or a fanatically dedicated user, the machine cannot be termed truly successful.

By this criterion we believe that CCMs have not yet silenced all their critics, but that we have turned important corners and have achieved a genuine understanding of the needed directions for research and development.

One problem in developing software for a CCM is that the programming process is far more vertically complex than for a standard computer. At the highest level are all the usual problems encountered when looking for performance from a computer—the user must be generally aware of the architecture of the CCM and program accordingly. But even from this top level working down, issues from deep within the FPGA must be dealt with. Must one partition the computation onto distinct chips in advance, or will an automatic partitioner be able to obtain sufficient utilization and speed? Even after partitioning, will a given chip be so densely packed with logic that routing delays will reduce the speed below a minimally acceptable level?

Apart from issues such as these, there are other factors to be considered. Logic synthesis and placement and routing on Xilinx chips presently takes several minutes to an hour for a chip with any substantial fraction of the logic used in a given application. How long will users be willing to wait, in this era of interactive computing, between iterations of this process? To what extent will they be willing to program in a language that is not Fortran or C? At what level can they or should they get involved in the performance-improving details of logic synthesis and/or placement and routing in order to gain the necessary speed improvements of a given application?

These problems actually differ from one kind of CCM to another. In an attached processor, the entire computation (or a definable portion of it) is taking place on the CCM. In a coprocessor system, the CCM portion must be extracted from the existing code (possibly with the help of compiler directives or annotations). The CCM code is similarly different from one sort of machine to another. On a coprocessor, one can assume that some effort might be expended by a user to optimize a particular "instruction," and the key issues would lie in recognizing its applicability and arranging for data to be delivered to and retrieved from the coprocessor. In an attached processor system, the CCM code could normally be much larger, allowing for more optimization (and thus more of the structure of the source to be obscured).

A final issue in programming is worthy of mention. The Algotronix chips have a feature that the Xilinx chips do not; part of the logic of the chip could be reconfigured without having the rest of the chip affected by the change. The reason for this lies in the different routing resources. On the Xilinx chip, no block of the chip can be assumed to be free of signals routed to or from some other block of the chip. On the Algotronix chip, however, the possibility of swapping hardware designs in and out like programs in and out of virtual memory was incorporated in the design from the beginning. The potential of such a feature for a CCM coprocessor is obvious.

#### REFERENCES

- [1] Algotronix Ltd., *The Configurable Logic Data Book*, Algotronix Ltd., Edinburgh, Scotland, UK, 1990.
- [2] J.M. Arnold, D.A. Buell, and E.G. Davis, "Splash 2," *ACM Symp. Parallel Algorithms and Architectures*, ACM Press, New York, 1992, pp. 316–322.
- [3] J.M. Arnold et al., "The Splash 2 Processor and Applications," *Proc. Int'l Conf. Computer Design*, CS Press, Los Alamitos, Calif., 1993, pp. 482–485.
- [4] P.M. Athanas, "Functional Reconfigurable Architecture and Compiler for Adaptive Computing," *Proc. 1993 Int'l Phoenix Computer and Comm. Conf.*, CS Press, Los Alamitos, Calif., 1993, pp. 49–55.
- [5] P.M. Athanas and H.F. Silverman, "An Adaptive Hardware Machine Architecture for Dynamic Processor Reconfiguration," *Proc. Int'l Conf. Computer Design*, CS Press, Los Alamitos, Calif., 1991, pp. 397–400.
- [6] P.M. Athanas and H.F. Silverman, "Processor Reconfiguration through Instruction Set Metamorphosis: Architecture and Compiler," *Computer*, Vol. 26, No. 3, Mar. 1993, pp. 11–18.
- [7] P. Bertin, *Mémoires Actives Programmables: Conception, Réalisation et Programmation*, PhD thesis, Université Paris 7, 1993.
- [8] P. Bertin, D. Roncin, and J. Vuillemin, "Programmable Active Memories: A Performance Assessment," in G. Borriello and C. Ebeling, eds., *Research on Integrated Systems*, MIT Press, Cambridge, Mass., 1993, pp. 88–102.
- [9] P. Bertin and H. Touati, "PAM Programming Environments: Practice and Experience," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1994, pp. 133–139.
- [10] B. Box, "Field Programmable Gate Array Based Reconfigurable Preprocessor," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1994, pp. 40–49.

- [11] D.A. Buell and K.L. Pocek, "Custom Computing Machines: An Introduction," *J. of Supercomputing*, Vol. 9, 1995, pp. 219–230.
- [12] S. Casselman, "Virtual Computing and the Virtual Computer," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 43–49.
- [13] Concurrent Logic Inc., *Cli6000 Series Field-Programmable Gate Arrays*, Concurrent Logic Inc., Sunnyvale, Calif., 1992.
- [14] C.P. Cowen and S. Monaghan, "A Reconfigurable Monte-Carlo Clustering Processor (MCCP)," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1994, pp. 59–66.
- [15] C.E. Cox and W. Ekkehard Blanz, "Ganglion—a Fast Hardware Implementation of a Connectionist Classifier," IBM Research Report RJ8290, *Proc. 1991 IEEE Custom Integrated Circuits Conf.*, IEEE Press, Piscataway, N.J., 1991, pp. 6.5.1–6.5.4.
- [16] S.A. Cuccaro and C.F. Reese, "The CM-2X: A Hybrid CM-2/Xilinx Prototype," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 121–131.
- [17] B.U. Heeb, *Debora: A System for the Development of Field-Programmable Hardware, and Its Application to a Reconfigurable Computer*, PhD thesis, VDF, Informatik Dissertationen 45, ETH Zürich, Zürich, Switzerland, 1993.
- [18] B.U. Heeb and C. Pfister, "Chameleon, a Workstation of a Different Color," in H. Grünbacher and R.W. Hartenstein, eds., *Field Programmable Gate Arrays: Architectures and Tools for Rapid Prototyping*, Springer-Verlag, Berlin, 1993, pp. 152–161.
- [19] H.-J. Herpel et al., "A Reconfigurable Computer for Embedded Control Applications," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 111–121.
- [20] D.T. Hoang, "Searching Genetic Databases on Splash 2," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 185–192.
- [21] C. Iseli and E. Sanchez, "Spyder: A Reconfigurable VLIW Processor Using FPGAs," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 17–25.
- [22] C. Iseli and E. Sanchez, "A C++ Compiler for FPGA Custom Execution Units Synthesis," *Proc. IEEE Symp. FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1995, pp. 173–179.
- [23] C. Iseli and E. Sanchez, "Spyder: A SURE, SUPerscalar and REconfigurable, Processor," *J. of Supercomputing*, Vol. 9, 1995, pp. 231–252.
- [24] P. Marchal and E. Sanchez, "CAFCA (Compact Accelerator for Cellular Automata): The Metamorphosable Machine," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1994, pp. 66–72.
- [25] S. Monaghan and C.P. Cowen, "Reconfigurable Multi-Bit Processor for DSP Applications in Statistical Physics," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 103–111.
- [26] D.V. Pryor, M.R. Thistle, and N. Shirazi, "Text Searching on Splash 2," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 172–178.
- [27] G.M. Quénot et al., "A Reconfigurable Compute Engine for Real-Time Vision Automata Prototyping," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1994, pp. 91–101.

- [28] F. Raimbault et al., "Fine Grain Parallelism on a MIMD Machine Using FPGAs," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 2-9.
- [29] M. Shand, P. Bertin, and J. Vuillemin, "Hardware Speedups for Long Integer Multiplication," *Proc. ACM Symp. Parallel Algorithms and Architectures*, ACM Press, New York, 1990, pp. 138-145.
- [30] S. Singh and P. Bellec, "Virtual Hardware for Graphics Applications Using FPGAs," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1994, pp. 49-59.
- [31] J. Vuillemin et al., "Programmable Active Memories: Reconfigurable Systems Come of Age," *IEEE Trans. VLSI Systems*, to be published in Mar. 1996.
- [32] M. Wazlowski et al., "PRISM II: Compiler and Architecture," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 9-17.
- [33] L.F. Wood, "High Performance Analysis and Control of Complex Systems Using Dynamically Reconfigurable Silicon and Optical Fiber Memory," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 132-142.
- [34] Xilinx, Inc., *The Programmable Gate Array Data Book*, Xilinx, Inc., San Jose, Calif., 1993.

## CHAPTER 2

---

# The Architecture of Splash 2

*Duncan A. Buell and Jeffrey M. Arnold<sup>1</sup>*

### 2.1 INTRODUCTION

In this chapter we present the higher-level architecture of the Splash 2 system. This architecture is what an application programmer would normally be expected to see. Although the current admirable trend in general-purpose computing is to allow the programmer to perform computations without being required to understand or even be aware of hardware structures, it has always been the case in high-performance computing that knowledge of architectural features is necessary. Programmers on vector machines learn how to vectorize their algorithms and how to write code from which compilers can extract vector computations. Programmers on massively parallel machines must study I/O and data layout patterns. Similarly, programmers of Splash 2 must understand the architecture of the machine in order to make effective use of it. More correctly, they must understand the architecture in order to make *any* use of it. Unlike more common machines with a longer history, we have not yet reached the point at which custom computing machines can be used without paying reasonably close attention to the hardware.

Splash 2, as can be seen in the following discussion, has a substantial generality in its structure. Although generality in an architecture can be a very desirable feature, such generality can be anathema to effective use of the machine if all possible details must be considered for routine use. Consequently, every effort was made

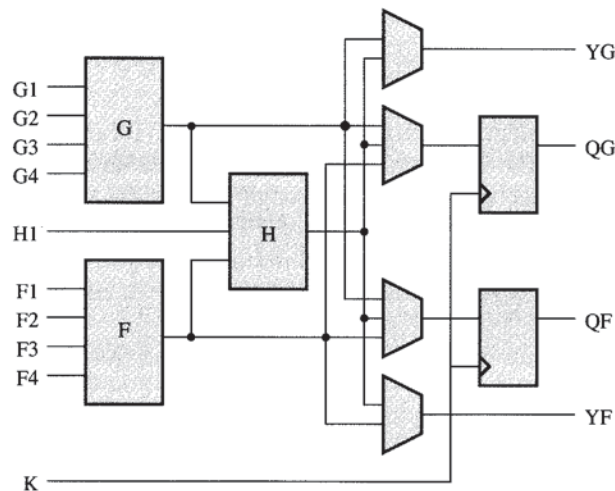
<sup>1</sup>A version of this chapter appeared as Arnold et al. [1] and is used with permission.

to provide standard avenues for a programmer to use features of the architecture in standard ways. The system was designed with the strong intent that most applications would have data streaming linearly past the FPGA processing elements or have data broadcast to them in SIMD fashion. These capabilities were thus supported both in hardware and software and in example computations and programs, and a reading of the architectural description should be done with a focus on how the architecture supports these two models of computation.

## 2.2 THE BUILDING BLOCKS

The basic building block from which Splash 2 is made is the Xilinx XC4010 FPGA [3, 4]. As mentioned in Chapter 1, the XC4010 contains a  $20 \times 20$  array of Configurable Logic Blocks (CLBs). The XC4010 CLB (shown in Figure 2.1) contains three lookup tables and two flip-flops. Two tables, labelled F and G, can each implement any Boolean function of up to four inputs. The outputs of the F and G functions can also be combined with a ninth input, H1, to form a single Boolean function of nine inputs. The YF output of the CLB can be taken from the output of either the F table or the H table. Similarly, the YG output can come from either the G or the H table. The F and G tables can also be configured to appear as individual  $16 \times 1$ -bit RAMs or a single  $32 \times 1$  RAM. Not shown in Figure 2.1 is additional fast carry logic, which allows a single CLB to implement a two-bit full adder. An additional wire allows the carry-out of one CLB to be connected directly to the carry-in of an adjacent CLB.

Figure 2.2 illustrates the routing structure of the XC4000 series FPGA. Connecting the CLBs are three types of signal routing resources including a single-length interconnect between adjacent switch boxes ("S" in Figure 2.2), a double-length interconnect between alternate switch boxes, and a set of long lines that span the width and height of the chip. The switch boxes contain programmable switches that allow each segment to connect to three others. Configuration of the FPGA is done by loading a bit file into on-chip RAM; the hardware to do this in Splash 2 is implicit in our description in this chapter of the general architecture and is discussed in greater length in Chapter 6.



**FIGURE 2.1** Xilinx XC4000 CLB Architecture

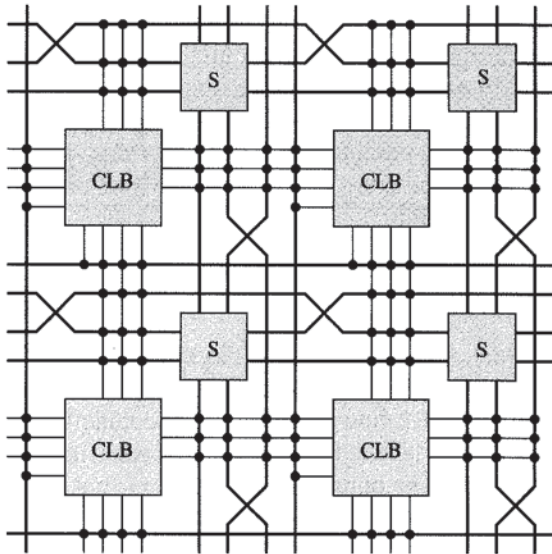


FIGURE 2.2 Xilinx XC4000 Routing Structure

The 400 CLBs can be viewed as 800 flip-flops, which can in turn be thought of as a maximum of twenty-five 32-bit “registers,” where by “register” we mean to include registers, adders, comparators, multiplexors, and similar basic structures. For example, a 16-bit object requires eight CLBs. Adders, subtractors, and comparators are implemented by “rippling” the fast carry output from one CLB to the next. In order to reduce the signal propagation time for the carries, one would normally want to have the CLBs physically adjacent to one another in the final design. The Xilinx-supplied tools attempt to do this, and the “Hard Macros” supplied by Xilinx can be used to guarantee that a logic object is placed into contiguous CLBs.

In addition to their use as registers, the Boolean function use of the CLBs is necessary to implement the rather more random control logic that will exist in any program, so the available number of “registers” is certainly always less than the maximum.

### 2.3 THE SYSTEM ARCHITECTURE

Splash 2 is an attached processor system. Although it was not designed for the purpose of being an attachment specifically for Sun workstations, the system as designed uses a SPARCstation 2 as a host and attaches to the host through the SBus.

The overall system architecture is pictured in Figure 2.3. An SBus adapter card is placed in the host and connects via a cable to the Interface Board of Splash 2. The Interface Board and the Splash 2 Array Boards reside in a separate cabinet on a Futurebus+ backplane.

Splash 2 is designed to execute either synchronously with the host or asynchronously as an attached processor. Programs for Splash 2 are loaded on the system by the host through the SBus connection. In some applications, the processing on Splash 2 is then driven by a clock on the Interface Board, and data is delivered to and taken from Splash 2 by DMA channels on the Interface Board. The Interface Board

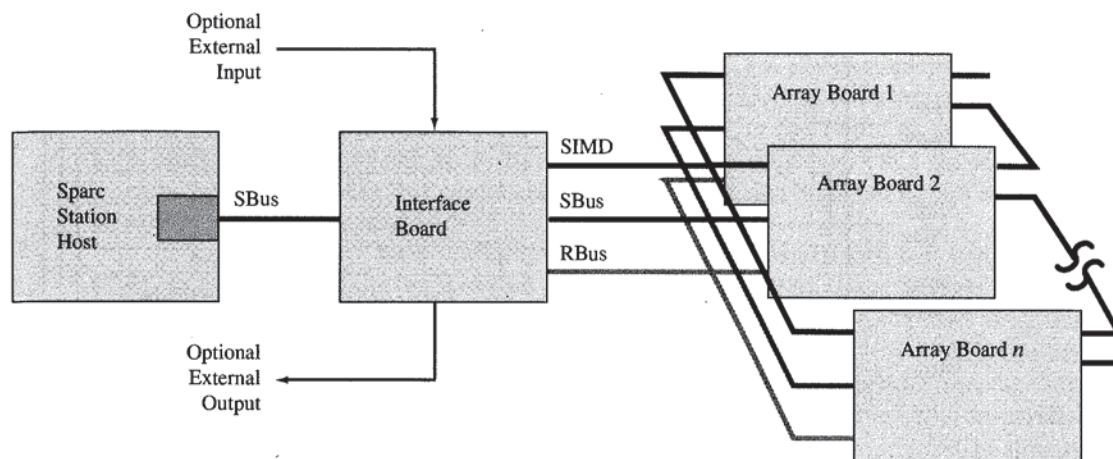


FIGURE 2.3 Splash 2 System Architecture

can also be configured to accept clock and data from an external source; in these applications, the role of the host is only to load programs and provide high-level control. Also supported are synchronous execution of Splash 2 via software clocking and slave data transfers.

The system architecture is designed to accommodate up to 16 Array Boards. A system containing eight Array Boards was built and functioned correctly, but most systems were built with two or four Array Boards in a 5-slot chassis that is 15" wide, 28" deep, and 24" high. The largest possible system was never built because it would have required an expensive special version of the cabinet. There seemed to be no reason based on the eight-board system's operating characteristics to expect any problems with the larger system, and assembling a larger number of smaller systems from the same number of Array Boards allowed a greater breadth of applications to be tested.

## 2.4 DATA PATHS

Each Splash 2 Array Board contains 17 Xilinx XC4010 FPGA chips [4] as its processing elements (see Figure 2.4). Sixteen of these are connected in a linear array to create a linear data path and the seventeenth provides a broadcast capability to the other 16 chips. To each of these 17 chips is attached 512 Kbytes of memory. Reflecting this basic design, there are three different paths by which data can be delivered to or taken from the Array Boards.

The primary models of computation that were intended to be supported by the Splash 2 architecture were a SIMD or broadcast-of-data model and a linear (but not restricted to "systolic") model. The programmer viewing Splash 2 as a SIMD machine sees, among other things, a 36-bit-wide data path from the Interface Board down the SIMD Bus to each Array Board simultaneously. Xilinx chip X0 on each Array Board can then broadcast the SIMD Bus data to the other FPGAs on its Array



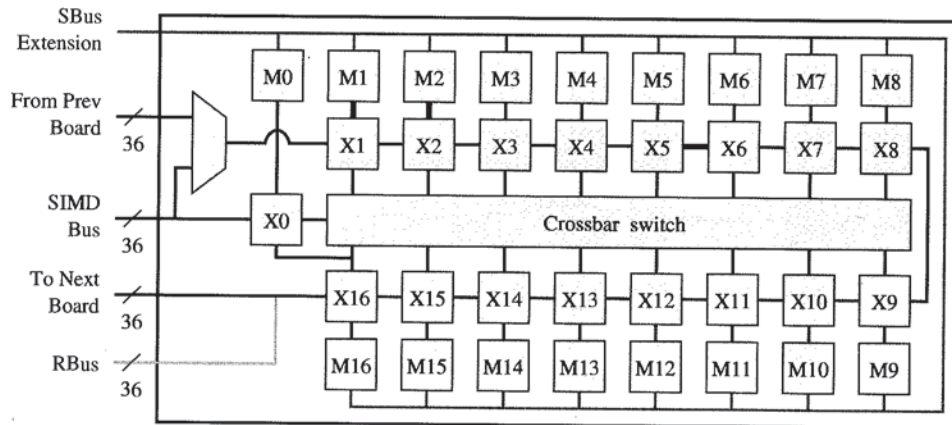


FIGURE 2.4 Array Board Architecture

Board. This mode of transferring data to the Splash 2 system was used, for example, in the text matching computation described later in this book.

Viewed as a machine with a linear data path, the SIMD Bus can be used to transmit data from the Interface Board to the first FPGA on the first Array Board. The data can then be moved through the linear data path on that board, then to the first FPGA on the second Array Board, and so on. Data from the last FPGA on the last Array Board returns to the Interface Board via the RBus. The linear path is also 36 bits wide, and is bidirectional (except for the initial segment along the SIMD Bus), so that data can be streamed in both directions for correlation computations; this was done for some versions of the DNA sequence comparison program detailed in Chapter 8. The definition of “last Array Board” is based upon the contents of a register on the Interface Board. This register can be changed during a program’s execution, so the length of the processor array can be changed dynamically.

Data coming from the SPARCstation 2 host is 32 bits wide. The 36-bit-wide data path in Splash 2 arises naturally from this and from the need and desire to have tag bits on the data. Although data coming from an external signal could genuinely be 36 bits wide, in most applications the tag bits are set and read by Xilinx FPGAs on the Interface Board. Since the Splash 2 system executes asynchronously with the host, it is routine for Splash 2 to be able to begin executing before data can be delivered from the host. One use for the tag bits, therefore, is to serve as a “valid data” signal. In linear mode, the Xilinx chips on the Array Boards would pass “data” down the linear data path immediately upon startup, but would not actually begin processing that data until a valid data tag appeared. Similarly, the Xilinx FPGA on the Interface Board that was handling output would discard any “result” it received until a valid data tag appeared on the output path.

Another use for the tag bits arises when the machine is used as if it were a SIMD machine. It is possible to broadcast a 32-bit word to the broadcasting X0 chip as well as a 4-bit instruction opcode. (Actually, of course, there is no structure to the 36 bits being broadcast, so any bits not needed for data could be used for an opcode.) This opcode could be used by X0 to process the data or control the broadcast just as with any other hierarchical SIMD machine.

These paths are not necessarily mutually exclusive, and these configurations are not hardware-controlled by something like a mode bit. For example, it is possible to use the SIMD Bus for broadcast of data to all FPGAs simultaneously, but then to use the linear data path as a sort of “back door” for the return of results or for necessary neighbor-to-neighbor communications. The only restriction on these data paths is that the SIMD Bus can be used only for delivering data and not for returning results to the Interface Board.

In either of the above modes, data from the Interface Board can come either from the host or from an external signal (see Figure 2.5). Each Xilinx XC4010 FPGA on the Interface Board is programmable by the user (some standard programs for common applications also exist). Xilinx chip XL (for “Xilinx Left”) handles incoming data for delivery to the SIMD Bus; in addition to setting the tag bits, if necessary, it handles the DMA transfers from the host, possible serial-to-parallel or parallel-to-serial data conversions, or similar massaging of the input data stream. Xilinx chip XR similarly handles data on the RBus, which would normally be the output path from the Splash 2 system. In some applications, especially in circumstances when postprocessing of the results was necessary, XR actually performed that part of the computation. This was true, for example, in the DNA sequence comparison computation.

The third means by which data can be provided to or retrieved from Splash 2 is directly through the 8.5 Mbytes of memory on each Array Board. This memory can be read by or written from the host as part of its normal address space via an

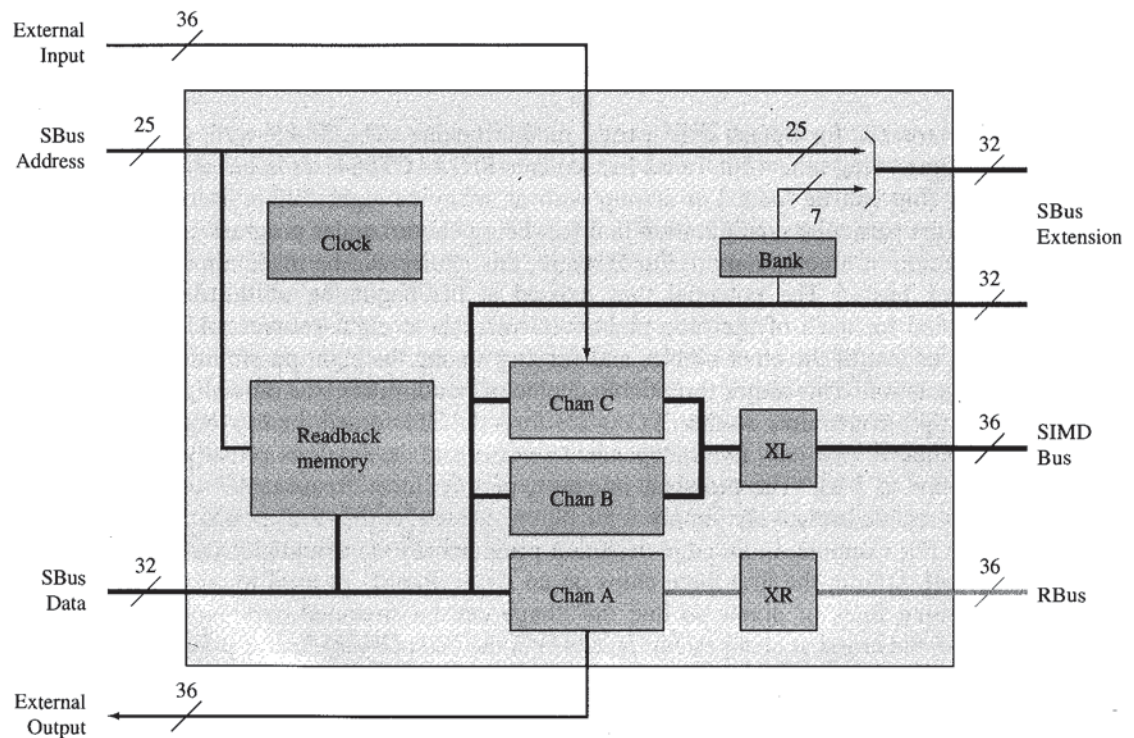


FIGURE 2.5 Interface Board Architecture

SBus extension independent of the FPGAs and their linear data path. The memory, however, is not dual-ported; during such memory read/write operations, the FPGAs are prevented from executing (and thus possibly accessing memory themselves). Thus, this mode of data transfer is not intended to be suitable for highly interleaved accesses of memory by Splash 2 and its host, but rather for bulk transfers before or after large phases of a given computation.

## 2.5 THE SPLASH 2 ARRAY BOARD

### 2.5.1 The Linear Array

The Splash 2 Array Board is detailed in Figure 2.4. Each Array Board contains 17 Xilinx XC4010 FPGA chips as processing elements. Sixteen of these, X1 through X16, form the processing array and are connected with a 36-bit-wide data path linearly and via a crossbar. To each FPGA is connected 512 Kbytes of memory. Throughout the Splash 2 system, the normal data object has been assumed to be 32 bits, augmented where possible and sensible with four tag bits. Here, in the connection from FPGA to memory, we find the one instance in which this design has been compromised. Three 36-bit-wide data paths, 18 bits for a memory address, and 32 bits for memory data would have left far too few of the 160 total pins for controlling each FPGA. The compromise was to reduce the memory data width to 16 bits.

### 2.5.2 The Splash 2 Crossbar

The crossbar for Splash 2 is a truly unique feature. The 36-bit-wide path is made by aggregating nine 4-bit Texas Instruments SN74ACT8841 crossbar chips [2]. Each such chip can be loaded at startup with as many as eight different configurations, with the particular configuration in effect being chosen under program control during execution of a computation. Furthermore, this choice can be made almost on a tick-by-tick basis.<sup>2</sup> The potential thus existed at the beginning of the design of this machine for each of the nine nibbles to have up to eight sources and destinations independent of the other nibbles and varying among the eight possibilities during the computation. This rather formidable choice of possibilities was only slightly reduced when pin constraints on the FPGAs X1 through X16 forced the low-eight nibbles to be paired so that only five independent sources and destinations actually exist on the machine as built. The crossbar, however, permits most “reasonable” configurations to be realized relatively simply.

For example, in an edge-detection program written essentially just for practice by Jeff Arnold, the first three chips on an Array Board are used to circularly buffer incoming lines of pixels so that the image can be streamed continuously into the board; the crossbar changes configuration at the end of every line of pixels to produce the effect of a circular buffer of three input lines on which a  $3 \times 3$  filter can be applied.

<sup>2</sup>This is “almost” tick-by-tick only because one cannot reverse source and destination in one clock period.

In other applications, the presence of the crossbar permitted the programmer to get beyond the rigid structure of a linear data path by “jumping ahead” in time/space in order to maintain a tightly pipelined, systolic-like computation. Although the specific chips chosen for the crossbar were the source of a later problem (more is said about this in Appendix A), we have been unable to find other examples of machines in which processor-to-processor communication can be changed as rapidly or with as much variety as in Splash 2.

### 2.5.3 Xilinx Chip X0 and Broadcast Mode

The seventeenth Xilinx chip, labelled X0, performs several functions that provide much of the flexibility of the Splash 2 architecture. Three bits from X0, controlled by a program that must be loaded into X0, select which of the eight configurations of the crossbar are in effect at any given point in time. For static configurations lasting throughout a given phase of a computation, this “program” controlling the crossbar is invisible to the programmer; if a varying crossbar is desired, however, the program to control the crossbar must be written by the programmer as part of the complete application.

The other major function of X0 is to broadcast data received on the SIMD Bus to the other 16 Xilinx FPGAs. This is possible because X0 and X16 share wires into the crossbar. Clearly, of course, both FPGAs must not be permitted to drive signals simultaneously on these wires, but this does not normally limit the range of usage; in situations in which X0 needs to serve as a broadcast chip, X16 normally does not need access into the crossbar. When X0 is broadcasting, X16 is receiving just like any other of the FPGAs. Since the crossbar is bidirectional, X0 can in fact receive data from the crossbar as well, adding to its ability to control execution on its Array Board.

Chip X0, like the other FPGAs, has a  $256K \times 16$ -bit memory attached to it with a 16-bit data path. In most instances where custom computing machines such as Splash 2 have been built, the use of memory for lookup tables has been important in achieving high performance. The rather limited compute resources on an FPGA requires the use of such memory to reduce the need for processing logic. This is especially true of chip X0.

Most massively parallel SIMD machines have had a front-end processor; chip X0 can, to the limit of its own capability, serve that function in Splash 2. As is described in Chapter 9 on the text processing application, X0 can perform some general computations and data preparation. It is also possible to use the four tag bits (or, for that matter, any other of the 36 bits of the data path) as instruction opcodes to chip X0. In such a situation, the memory would be used to store a “microcode subroutine,” which would be executed by a small finite-state or other machine implemented on X0.

## 2.6 THE INTERFACE BOARD AND CONTROL FEATURES

A detailed description of the Interface Board is presented in Chapter 3 on hardware and implementation details, but now we discuss several aspects of the Interface Board that have to be considered as part of the higher-level architecture of the Splash 2

system. To permit control of parallel programs running searches until some particular event occurs, global `or` and global `valid` bits run to FPGA X0 from each of the 16 FPGAs X1 through X16 on an Array Board. Global `or` and global `valid` bits from X0 of each Array Board are then wire-ORed on the backplane to a register on the Interface Board and appear as inputs to FPGAs XL and XR on the Interface Board. On each Array Board, the `or` and `valid` bits are bidirectional, allowing further control by X0 of computations on the Array Board.

The clock that drives Splash 2 as an asynchronous attached processor resides on the Interface Board. Because various programs could be expected to run (and in fact do run) at widely differing speeds, a clock module was chosen whose frequency could be tuned to the speed at which the synthesized Xilinx chip program could run. Chip XL has control of the clock-regulated computation on the Array Boards; the Array Boards can be single-stepped, *n*-stepped, or allowed to run freely. To reduce the programming overhead for routine computations, several default programs for XL were written at an early stage in the development of the system and can be selected by a programmer from the library. This is also true of programs for the FPGAs on the Array Boards.

The alert reader will already have noted that the address bits available from the host on the SBus are insufficient to address all the memory potentially available on a full-sized Splash 2 system. The address extension is done on the Interface Board; the particulars of this process appear in Chapter 3 on hardware implementation details.

A final feature of the Splash 2 system is the ability to load or store a configuration state into the Xilinx chips. Readout of the state is invaluable for debugging, program optimization, and monitoring program behavior.

## REFERENCES

- [1] J.M. Arnold, D.A. Buell, and E.G. Davis, "Splash 2," *Proc. ACM Symp. Parallel Algorithms and Architectures*, ACM Press, New York, 1992, pp. 316–322.
- [2] Texas Instruments Inc., *The SN74ACT8800 Family Data Manual (SCSS006A)*, Texas Instruments Inc., Dallas, 1988.
- [3] S.M. Trimberger, ed., *Field Programmable Gate Array Technology*, Kluwer Academic Publishers, Boston, 1994.
- [4] Xilinx, Inc., *The Programmable Gate Array Data Book*, Xilinx, Inc., San Jose, Calif., 1993.

# CHAPTER 3

---

## Hardware Implementation

*Walter J. Kleinfelder and Jeffrey M. Arnold*

### 3.1 INTRODUCTION

Figure 3.1 illustrates the system architecture of Splash 2. The system consists of a Futurebus+ backplane enclosure containing one Interface Board and up to 13 Splash 2 Array Boards,<sup>1</sup> and a SPARCstation 2 host computer with Adapter Board. The Adapter Board plugs into the host computer's internal SBus and extends the address and data bus to permit host-resident programs to directly address the memory and control registers in the Splash 2 system. The Adapter Board also provides the interface logic required to permit Splash 2 to perform direct memory access (DMA) transfers to and from the host memory and to generate SBus interrupts. The Splash 2 enclosure is connected to the host system via a cable between the Adapter Board and the Interface Board. To complete the linear data path, each Array Board is also connected to its two neighboring boards through a separate custom backplane in addition to the Futurebus+ backplane.

During the course of the Splash 2 project, two different interface boards were developed to provide the connection between the host computer and the Splash 2 Array Boards. The Development Board was designed with minimal functionality but with an extensible wire-wrap core to allow prototyping of various features. The final Interface Board was built in printed circuit technology and incorporates a number of higher-level functions. Both interfaces extend the host address and data buses to the

<sup>1</sup>This is the largest enclosure fabricated for the Splash 2 system. The architecture supports up to 16 Array Boards.

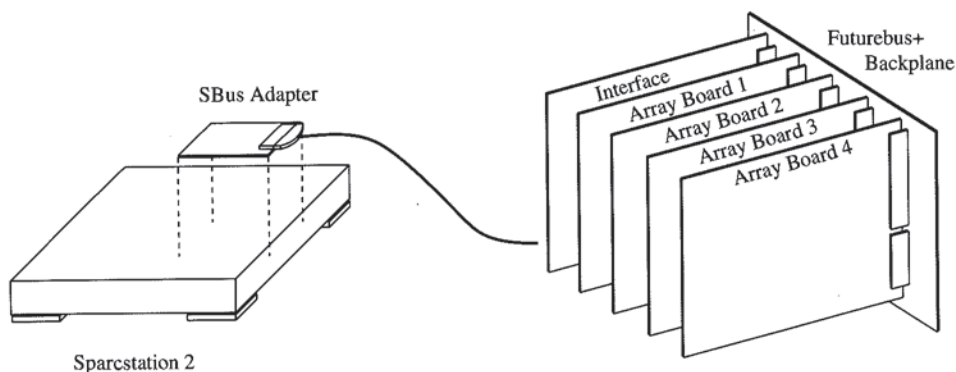


FIGURE 3.1 Splash 2 System Architecture

backplane memory bus, permitting the host to read and write memory and memory-mapped control registers on the Array Boards. All data transfers between the host and the interface are 32 bits wide, the word size of the SBus. Because the 25 bits of physical address space for a single SBus slot is insufficient to address the entire Splash 2 memory space, the address is extended to 32 bits by a 7-bit bank register on the Interface Board.

The Splash 2 system may transfer data to and from the host system memory via DMA. The Interface Board contains up to three independent DMA channels implemented as optional daughter boards that may be plugged onto the Interface Board. In addition to supporting DMA, this daughter board arrangement allows high-speed external input and output devices to be connected directly to Splash 2, bypassing the host SBus. For example, an external video input may be brought directly into Splash 2 by replacing one DMA channel with a specially designed daughter board.

The linear data path extends from the Interface Board along the SIMD Bus, through the set of Array Boards in daisy-chain fashion, and back to the Interface Board along the RBus. The SIMD Bus is a 36-bit unidirectional bus driven by the Interface Board and connected to each Array Board in the system. The Array Board daisy chain and the RBus are 36-bit-wide bidirectional data paths. The linear data path can therefore be used to pass data in either the "left-to-right" direction or the "right-to-left" direction. Left-to-right is defined to be from the SIMD Bus through the Array Boards and back on the RBus. Right-to-left is defined to be from the Interface Board down the RBus to the last, or rightmost, Array Board through the daisy chain and terminating at the first, or leftmost, Array Board. Two data streams can pass simultaneously through the array in opposite directions, with one stream following the left-to-right direction on a subset of the bits of the linear data path and the second stream moving right-to-left on the remaining bits.

Each Splash 2 Array Board contains 16 Processing Elements, X1–X16, with direct 36-bit connections between adjacent elements. Each element can also communicate with all other elements on the same board through a programmable crossbar. A seventeenth Control Element, X0, controls the crossbar and provides support logic. The sixteen Processing Elements and the Control Element each consist of a Xilinx XC4010 FPGA and a 256K × 16 static RAM.

### 3.2 DEVELOPMENT BOARD DESIGN

The Development Board was designed to serve a variety of purposes during the early development of Splash 2. The original goal of the Development Board was to support the initial debugging of the Array Board design and the system software while the design of the Interface Board proceeded in parallel. The flexibility of the Development Board also made it a convenient vehicle on which to prototype various components of the final Interface Board design. The Development Board eventually became a critical tool for the instrumentation and debugging of the DMA transfer protocol. A modified version of the Development Board was also used as a test fixture for the DMA and Clock daughter boards.

The design philosophy behind the Development Board was to keep the hardware simple by moving as much control as possible to the host software. This was accomplished by placing virtually every signal on the backplane under the direct control of the host computer. Readable and writable registers are connected to the SIMD Bus and RBus data (32 bits each) and tags (4 bits each), and the RBus size and direction controls. Read-only registers provide access to the interrupt request and global OR signals. The system clock mechanism is a register that, when written by the host, generates a single pulse of 100 nsec duration. To generate successive clock pulses the host must write repeatedly to the clock register. Using this mechanism the SPARCstation host is able to achieve a maximum clock rate of 4 MHz. To improve the performance of stream-based applications an additional address is decoded that, when written, loads the SBus data into the SIMD register and then generates a single clock pulse.

The physical design of the Development Board consists of three main components: the SBus interface, the backplane interface, and the wire-wrap core. The data path portion of both the SBus and the backplane interfaces are implemented with surface-mount printed circuit technology along two edges of the board. The center of the board contains a large grid of holes for wire-wrap socket pins. This prototyping area is used to implement the control state machines and to experiment with the DMA and Clock circuits.

### 3.3 INTERFACE BOARD DESIGN

The Splash 2 system buses are implemented on the P896.2 Futurebus+ profile A backplane with a 128-bit extension. The Splash 2 Interface Board plugs into Slot 1 of the Futurebus+ backplane and accepts the cable from the SBus Adapter Board in the host computer. The Interface Board is responsible for generating all signals required in the backplane and is structured to drive up to 16 Splash 2 Array Boards. The principal functions of the Interface Board include SBus control and data transfer, system clock generation, and pre- and postprocessing of data to and from the Array Boards.

Figure 3.2 illustrates the Splash 2 Interface Board architecture. The host data bus (SD[0:31]) is gated and buffered to drive the backplane memory data bus for memory-mapped reads and writes. The host address bus (SA[2:24]) is buffered and decoded locally for accesses to the Interface Board. The Bank Register is loaded by the host with a 7-bit value. The SBus address and the Bank Register together form



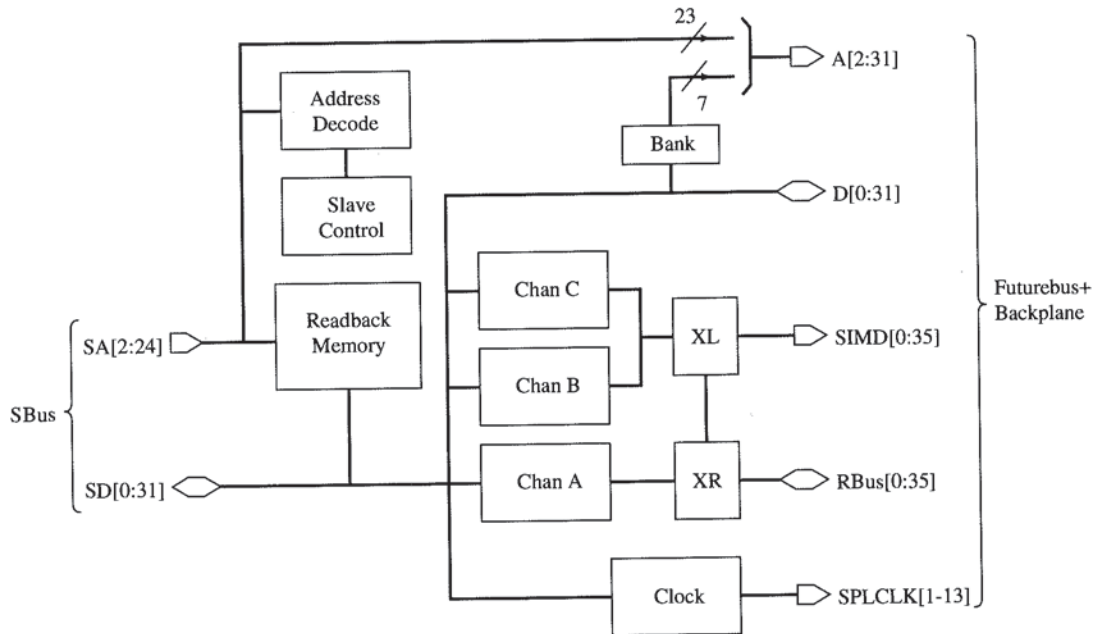


FIGURE 3.2 Interface Board Architecture

the Array Board memory address,  $A[2:31]$ . Since only 32-bit accesses are supported, the two least significant bits of the address ( $A[0:1]$ ) are always zero.

The slave control state machine receives the decoded address and generates internal read and write timing signals in response to SBus slave cycles. The slave machine also controls the timing of the SBus acknowledge signal, ensuring sufficient access time for the various registers and memories. SBus read operations to the Array Board are acknowledged in 9 SBus clock cycles, while write operations are acknowledged in 8 cycles. In accesses to the facilities on the Interface Board, writes are acknowledged in 3 cycles and reads in 4 cycles.

The clock circuit generates the system clock signal for the Splash 2 Array Boards and can be programmed by the host to various frequencies. During execution of Splash 2 programs this signal clocks the user-defined circuitry. To aid debugging of applications, the clock generator can be programmed to stop, single-step, or step a fixed number of times. To prevent DMA data overruns or underruns the clock generator can also be stopped and restarted by user-defined logic in XL or XR. To minimize the clock skew across the system, separate clock signals (SPLCLK) are driven to each Array Board.

XL and XR are user-programmable XC4010 FPGAs that provide the interface between the DMA channels and the backplane bus. XL controls Channels B and C and drives the 36-bit SIMD bus. XR controls Channel A and can receive data from or drive data to the 36-bit RBus. The direction of the RBus is determined by the RDIR output of XR. Data may also be passed between XL and XR through a separate 36-bit bidirectional path.

### 3.3.1 DMA Channel

The DMA channels perform SBus-compatible burst transfers between the host memory and a 256-word FIFO. Each DMA channel contains an address register, a transfer count register, and a control register. The address register contains the address of the host memory buffer. The transfer count register contains the number of burst transfers to perform. The control register contains an enable bit and a direction bit. All DMA transfers are performed in 16-word bursts, the largest supported by the SBus. The FIFOs contain programmable high- and low-watermark registers, which permit the DMA channel to determine when to request a transfer. When the channel is enabled and the direction bit is set to read from the host memory, if the FIFO has space for at least 16 words, then an SBus READ operation is requested. Similarly, when the direction bit is set to write to the host, memory, and if at least 16 words of data are available in the FIFO, then an SBus WRITE operation is requested.

After each burst transfer completes the address and transfer count registers are updated. The address register is incremented by 64 to point to the next block of 16 words (64 bytes), while the transfer count is decremented by 1. When the transfer count reaches zero an SBus interrupt is requested.

The SBus side of the DMA channel is 32 bits wide, but the FIFO and the Splash 2 data paths are 36 bits wide. When transferring data from the host, the word is extended to 36 bits by concatenating the contents of a 4-bit tag register to the data. When transferring data to the host, the 4 tag bits are saved in a register.

The DMA channel allows direct loading and unloading of the FIFO data from the host by mapping the input and output data registers of the FIFO into the host's address space. This feature allows the host operating system software to handle the boundary conditions of transfers that are not aligned on 64-byte boundaries. See Chapter 6 for more details on DMA data alignment.

### 3.3.2 XL and XR

The primary function of the two user-programmable FPGAs, XL and XR, is to perform pre- and postprocessing on the input and output data streams. DMA Channel A is controlled by XR, while Channels B and C are controlled by XL. Both XL and XR receive the Splash 2 system clock and a free-running clock that is synchronous with the system clock. Either XL or XR may stop and restart the system clock when the FIFOs, in their respective DMA channels, are empty or full. When XL or XR are used to stop the system clock the free-running clock may be used to drive the controlling state machine, allowing it to restart the system clock when the condition has cleared.

DMA Channels B and C share a common 36-bit-wide data bus with XL, which may select the channel from which to receive data. XL is also responsible for driving the SIMD bus, typically with the input data from one or both of Channels B and C. XR sits between DMA Channel A and the RBus and is typically used to postprocess result data from the RBus before sending it back to the host through Channel A. A separate 36-bit bidirectional bus connects XL and XR. This bus may be used to coordinate clock control, to close the loop through the linear array by passing data from the RBus back to XL and the SIMD bus, or to pass input data from XL through XR to the linear data path in the right-to-left direction. For example, an application can receive

input data from both Channels B and C simultaneously, sending the Channel B data along the left-to-right direction and the Channel C data moving from right to left.

### 3.3.3 Interrupts

The Interface Board receives up to 16 individual interrupt requests, one from each Splash 2 Array Board. Interrupt requests can also be generated by XL, XR, the DMA channels, and the clock. The interrupt circuit logically ANDs each request with a corresponding bit in a mask register and then ORs these results together to form a single SBus interrupt request which, when enabled by a bit in the control register, is passed to the SBus. When handling an interrupt, the host can read the contents of the interrupt register to determine which of the possible sources made the request. If the request came from one of the Array Boards, the host can then interrogate a similar register on the requesting Array Board to determine which PE initiated the request.

### 3.3.4 Clock

The system clock can be selected from two possible sources: the programmable clock generator or a software-generated clock pulse. The clock generator, or "hardware clock," is a daughter card that plugs onto the Interface Board. The heart of the hardware clock is an Analytic Instruments FS-30 programmable frequency synthesizer, which has a frequency range of less than 1 Hz to 30 MHz. The frequency of the system clock is set by the host and is asynchronous with respect to the SBus clock. Both XL and XR have the ability to immediately stop the system clock, typically in response to a DMA channel nearing the full or empty mark. The host computer also has the ability to start and stop the clock generator, and may program the generation of a specific number of clock cycles. Special synchronization circuitry ensures that the first and last clock pulses are not truncated. The output signal has a nominal duty cycle of 50% plus or minus 5 nsec.

The "software clock" is a register very similar to the clock register on the Development Board. In the interest of performance, a bit in the control register determines whether writes to the software clock generate one or two pulses of 100 nsec each. Another bit in the control register selects either the hardware clock or the software clock to drive the system clock.

### 3.3.5 Programming and Readback

The configuration and state readback mechanisms for all of the user-programmable FPGAs are implemented on the Interface Board using a single 256K × 32 memory to store the bitstreams. All of the FPGAs on a single Array Board are programmed simultaneously using the serial configuration mode of the XC4010. Prior to programming, the host merges the 17 individual bitstreams (one each for X0 through X16) into a single 17-bit-wide stream. This operation is known as "corner turning." The corner-turned configuration stream is then loaded into the Interface Board memory, and the programming sequence is begun by an on-board state machine. This state machine reads sequential locations from the memory and performs write operations over the SBus extension to a special address on the selected Array Board. A base

address register contains the location of the first word in the bitstream and configuration stops when the address counter reaches the top of the memory. Seventeen of the 32 bits are used to store the data for the 17 user-programmable FPGAs on the Array Board. Two additional bits may optionally contain the configuration streams for XL and XR. The remaining bits of the Interface Board memory are not used.

The same basic mechanism is used to perform state readback. Another state machine on the Interface Board reads the internal state information from all 17 FPGAs on a given Array Board and stores the data into successive locations in the Interface Board memory. Readback terminates when the address counter reaches the top of memory. Once the FPGA state information is stored in the Interface Board memory the SPARCstation host can retrieve the specific state bits of interest to the programmer.

Since the configuration and readback operations employ the SBus extension between the Interface Board and the Array Boards, during both configuration and readback the Splash 2 system will not respond to the SBus; any host attempt to access the Splash 2 system during one of these operations will result in a bus time-out.

### 3.3.6 Miscellaneous Registers

There are several control and status registers on the Interface Board that are accessible to the host. The configuration and readback mechanism contains an 18-bit base address register and a control register with a direction bit (programming versus readback) and a start bit.

The main control and status register (CSR) contains the "bypass" mode bit which, when set, disables XL and XR and enables the bypass registers. These registers allow the Interface Board to mimic the behavior of the Development Board. A separate bypass register contains the RBus size and direction bits and the broadcast bit for use in bypass mode.

The CSR also provides access to the signals that control the programming and readback of XL and XR, including the PROGRAM, INIT, and DONE pins of the Xilinx chips. Another signal, RBTRIG, is used to initiate the readback operation in XL and XR. The clock source select and the interrupt enable bits are also in the CSR.

There are three levels of reset available on the Interface Board. At the lowest level is the system reset signal, which is connected to the host computer's power-on reset. At the second level is the "panic" bit in the CSR. This signal is used to reset various control state machines on the Array Boards. At the highest level is the "program reset" bit in the CSR, which is connected to the individual Processing Elements' global set/reset (GSR) signal.

The Interface Board also contains a 32-bit read-only identification (ID) register implemented using DIP switches. The ID switches are set to contain board version information and a unique serial number.

## 3.4 ARRAY BOARD DESIGN

Figure 3.3 shows a block diagram of the Array Board architecture. The Array Board contains 16 Processing Elements (PEs) arranged in a linear array. A seventeenth FPGA-memory pair, referred to both as the Control Element or as the seventeenth

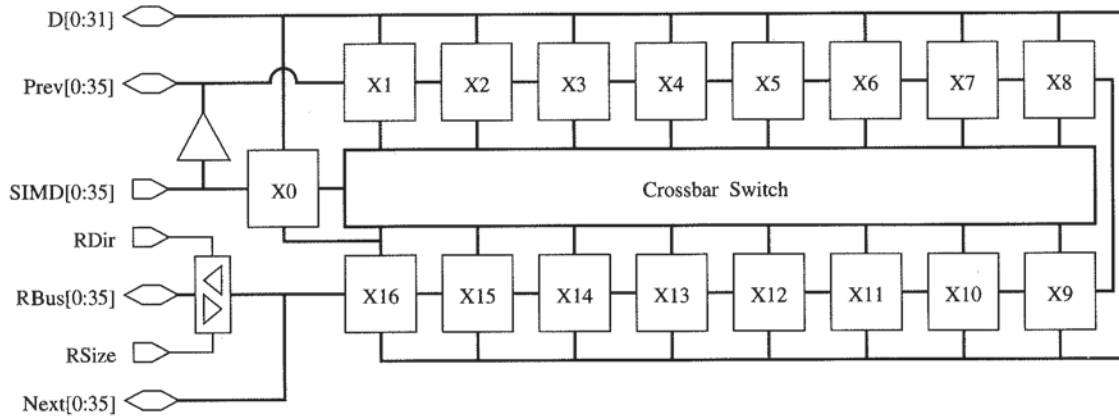


FIGURE 3.3 Array Board Architecture

Processing Element PE0 in some applications, manages the crossbar and can send data to or receive data from the crossbar. The 36-bit linear data path enters the array from the previous board, continues through adjacent PEs, and exits the array to the next board. The first, or leftmost, board in the system detects that it is in slot 2 and receives its input data from the SIMD bus instead of the previous board connection. The last, or rightmost, board is determined by comparing the slot number of each board to the RBus Size value on the backplane. The selected board then either sends data to the RBus or receives data from the RBus, as determined by the state of the RDIR backplane signal. Array Boards that are not at either end of the array simply communicate with adjacent boards via unbussed pins in a custom backplane extension.

Along the front edge of the Array Board are 18 light-emitting diodes (LEDs). One LED (green) is connected directly to power and indicates whether the board is receiving power. A second LED (red) is connected to an output pin of the Control Element (X0). The remaining 16 LEDs (amber) are each connected to an output pin of each PE. These LEDs are available for use by application programs and are used extensively by the diagnostic software.<sup>2</sup>

### 3.4.1 Processing Element

The organization of the Splash 2 Processing Element is shown in Figure 3.4. The PE consists of a Xilinx XC4010 FPGA and 256K × 16 RAM. The RAM is implemented with four 256K × 4 static RAM chips with 20 nsec access time mounted on a ZIP package. The memory control state machine is implemented in a 22V10 programmable logic device (PLD).

The Processing Element FPGA has four principal data paths, corresponding approximately to the four sides of the chip. There is a 36-bit-wide bidirectional data path to each of the two neighboring PEs (to the left and right), a 41-bit interface to the central crossbar, and a 36-bit interface to the local RAM. The crossbar interface consists of a 36-bit-wide bidirectional data path and five output enable control signals.

<sup>2</sup>One programmer wrote a program that scrolled a banner of text across the boards.

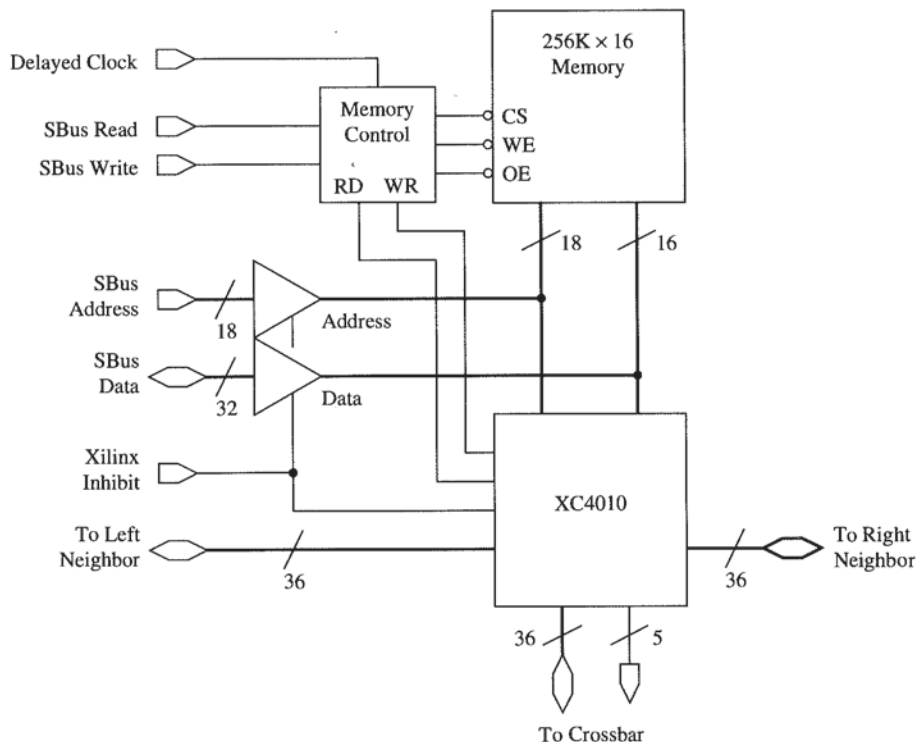


FIGURE 3.4 Splash 2 Processing Element

The 36 bits of data are arranged in four groups of 8 bits and one group of 4 bits, with each group controlled by a separate output enable. The RAM interface consists of a 16-bit data path, an 18-bit address, and separate read and write control signals.

The memory control device is used to present a purely synchronous interface to the programmer. To initiate a memory read operation the FPGA asserts the read control signal and the address at the rising edge of the system clock. Data from the memory is available on the next rising edge of the clock. To initiate a write operation the FPGA asserts the write control signal, the address, and the data on the rising edge. The write enable pulse is generated by the memory control PLD from a delayed version of the system clock. The interface circuitry and the RAM timing guarantee that a write pulse is not applied to the RAM until the address and data have met the setup requirements. The write pulse is released in time to meet the required hold time. To guarantee that these constraints are satisfied, it is necessary to register the memory interface signals in the IOB flip-flops of the FPGA.

There are several other signals available to the Processing Element, including:

- system clock
- broadcast bit from the control element
- program reset signal from the Interface Board
- two handshake register bits

- global OR result and valid bits connected to the control element
- Xilinx Disable bit
- LED control signal

### 3.4.2 Control Element

The organization of the Control Element (X0) is very similar to that of the Processing Element, consisting of a Xilinx XC4010 and a 256K×16 RAM. The memory interface of the Control Element is identical to that of the Processing Element. The 36-bit SIMD bus is connected directly to X0. X0 is responsible for selecting the crossbar configuration in use at any given time through a 3-bit “crossbar select” port. X0 may also read or write the 36-bit crossbar through the port it shares with PE X16. An output port allows X0 to override X16’s crossbar output enables, effectively taking control of the 36-bit data path.

The bidirectional global OR and Valid bits from each of the 16 PEs are connected to X0. X0 in turn may also drive the open collector systemwide global OR and Valid signals on the backplane. These signals are intended to be used to permit X0 to perform Array Board-level synchronization, and then to participate as the board’s representative in systemwide synchronization.

The single-bit broadcast signal from the backplane is an input to X0, which may then drive the board-level broadcast signal to the 16 PEs.

### 3.4.3 External Memory Access

The SBus extension bus permits the host to directly read and write the PE memories. Since the PE memory is not dual-ported, the address and data bus of each memory is shared between the FPGA and the SBus extension. Therefore it is necessary to ensure that the FPGA does not interfere with SBus accesses to the memory, and vice versa. The host accomplishes this by stopping the system clock and asserting the “Xilinx Disable” signal in the Array Board control register prior to any memory access. The system software must guarantee that the Xilinx Disable pin of each PE is wired to the internal Global Tri-State (GTS) signal within the FPGA. Once the memory access is complete, the host must clear the Xilinx Disable bit and restart the clock.

References from the host to the PE memories are multiplexed such that each 32-bit host access is converted to two 16-bit accesses to sequential memory locations. The RAM multiplexor is clocked with the SBus clock, not the programmable system clock, so the memory can be accessed by the host while the system clock is stopped.

### 3.4.4 Crossbar

The crossbar network permits the communication between processing elements via a selection of preprogrammed configurations. The Texas Instruments SN74ACT8841 chip, shown in Figure 3.5, is used for this application. The 8841 chip has sixteen 4-bit bidirectional ports, which may be connected in any desired pattern. Each port’s output may be selected from any of the other ports. The output port selection is

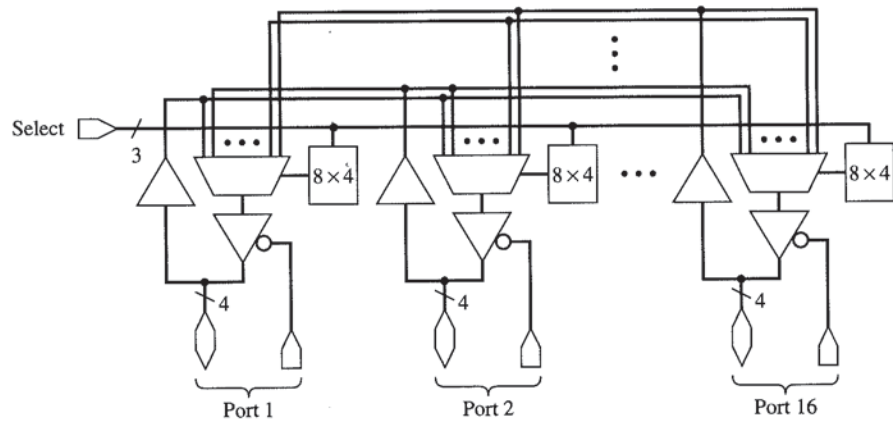


FIGURE 3.5 Architecture of the TI 8841 Crossbar Chip

controlled by an  $8 \times 4$  register file. The register file entry in use during a given clock cycle is chosen by the external 3-bit Select signal. The output of each port is independently controlled by a separate output enable pin.

On the Splash 2 Array Board, nine 8841 chips are coupled to form the central crossbar. The output enable pins are grouped together in pairs to form five controls for each Processing Element. Each FPGA therefore supplies five signals to control its corresponding crossbar connection, arranged as four 8-bit paths and one 4-bit path.

The crossbar array is preloaded by the host with up to eight connection configurations. During execution, three bits from the Control Element X0 select the desired preloaded configuration, which X0 can change from clock cycle to clock cycle. The crossbar chips are individually loaded, since their configurations will not necessarily be identical.

### 3.4.5 Programming and Readback

Configuration and readback of the Array Board FPGAs is accomplished by performing write and read operations over the extended SBus. The PROGRAM bit in the Array Board control register causes all 17 FPGAs (X0–X16) to enter program mode. Subsequent writes to the “configuration” register cause each FPGA to extract one bit from the data word and load it into its internal configuration. The INIT and DONE status signals are available as separate registers to be read by the host.

State readback is performed in a similar manner. When written, the RBTRIG bit in the control register puts all 17 FPGAs in readback mode. Subsequent reads from the configuration register return the internal state information to the SBus. An on-board state machine generates the configuration clock (CCLK) timing for both configuration and readback.

### 3.4.6 Miscellaneous Registers

Most of the registers on the Array Board and the SBus address decoding are actually implemented in an eighteenth Xilinx FPGA. This FPGA is not user-programmable, however, but is configured at power-on from a set of on-board PROMs. The “panic”



signal from the Interface Board will also force a reload of the controller's configuration.

All of the programmable features and mode controls are accessible to the host through the SBus address space. All registers are aligned on 32-bit boundaries. The registers are organized into two separate pages of the address space, one to be available in user mode and the other accessible only in supervisor mode. The user mode registers include the control register, the handshake register, the configuration registers, the version and serial number register, and the crossbar configuration registers. The supervisor mode space contains the interrupt status and mask registers.

The control register contains a number of signals effecting the operation of the Array Board, including the PROGRAM and RBTRIG for loading and unloading the FPGAs. The control register also contains the Xilinx Disable signal used to force the FPGA pins to their high impedance state and the Handshake Direction signal used to set the direction of the handshake register. The configuration registers include addresses from which the host may read the INIT and DONE status signals from each FPGA and the location to which the configuration bitstream is written or the readback stream is read.

The handshake register is an asynchronous communication channel between the host and the Processing Elements. One bit of the handshake register is connected to each of the 17 FPGAs. The bits themselves are bidirectional, but the direction of all 17 bits is determined by the Handshake Direction bit of the control register. The version and serial number register contents are hard-coded in the controller's configuration PROM. The version number changes with each revision of the controller PROM program, and the serial number is unique to each Array Board.

The crossbar configuration information is loaded into a set of 4-bit registers within each of the nine TI 8841 crossbar chips. These registers are mapped into the SBus address space of the Array Board, allowing the host to write directly to the TI chips.

The interrupt status register latches the state of the interrupt bits from each of the 17 FPGAs. The controller combines the latched status with the mask information to form the board-level interrupt signal. The interrupt status register is cleared when read by the host.

# CHAPTER 4

---

## Splash 2: The Evolution of a New Architecture

*Duncan A. Buell*

The preceding two chapters have described the hardware designed and built as Splash 2. It is important, however, to trace the design process that led us to the artifact we have today, otherwise there is the danger of seeing only the extant machine and not the potential variations. In this chapter we examine the decisions that led to the final architecture.

### 4.1 SPLASH 1

The germ of the idea for Splash 1 [2] apparently came from Dick Kunze and Paul Schneck at SRC in late 1986. Discussion took place among Kunze, Schneck, and Dick Lipton from Princeton, in part due to a realization that a Splash-like processor would be a generalization of the special-purpose P-NAC (Princeton Nucleic Acid Comparator) that Lipton was having built. P-NAC was designed to execute the edit-distance (approximate string matching) algorithm used in comparing DNA sequences against each other.

By the spring of 1987, the essential Splash 1 architecture had been laid out; at a meeting among the SRC and Princeton principals held at SRC on February 27, 1987, the linear array of Xilinx chips and memories had already taken shape. As with any such new system, there were several variants that were considered from time to time but never adopted. One early thought was that a 128-board system could be built. This never got beyond the concept stage. Another idea that surfaced again and again and resurfaced briefly in the later design of Splash 2 was the possibility of

including floating-point chips in the array path. This idea was studied for Splash 1 but never adopted, the stumbling block being in part that the floating-point chips operated at a fixed speed whereas Splash 1 designs ran at speeds that were dependent on the programming. This implied that the floating point chips could not be substituted one for one with Xilinx chips, leading to rather complicated control paths. Further, while applications that made use of the Xilinx capability—reconfigurable processor architecture—and applications that made use of the floating-point power could be envisioned, there seemed only a limited benefit from mixing the two. Given that numerous floating-point accelerators exist and that the real point of the experiment was to demonstrate the power of FPGAs for computing, there seemed to be no overwhelming reason to add the floating-point capability to Splash 1. In the later design of Splash 2 the subject came up once again. By then the consensus was that Splash 1 was and Splash 2 would be processors with a niche that lies outside the world of floating-point computation. Given more serious consideration for Splash 2, however, was the idea of including a fast microprocessor on each Array Board to provide more general compute capability close to the Xilinx chips themselves.

The final Splash 1 processor was a single multiwire board that plugs into the VMEbus of a Sun workstation (see Figure 4.1).

Each board contained 32 Xilinx XC3090 FPGA chips X0 through X31 as PEs connected in a linear array by a 32-bit-wide path. Chips X0 and X31 could be similarly connected to form a ring, were it necessary to route data around the ring more than once or to send data in both directions through the FPGAs. Data synchronization on and off the board was handled by a pair of FIFOs controlled by X0 and X31, respectively. Between each pair of interior Xilinx chips was a 128K × 8 RAM with an 8-bit-wide path to the FPGAs.

The Xilinx XC3090 chips in Splash 1 had a maximum clock rate of 32 MHz. To accommodate Splash 1 designs that could not be run at maximum speed, usually due to placement and routing problems or to the inability of the VMEbus to deliver data at a sufficiently high rate, the clock rate could be set in factors of two from 1 MHz to 32 MHz.

A three-Xilinx-chip (one for input from the host, one for output to the host, and one in the middle) Splash 1 board, nicknamed PUDDLE, was wire-wrapped by hand in order to gain an understanding of the hardware and expose unforeseen problems. This board became operational in March of 1988. When it had been thoroughly debugged, a full 32-chip board was wire-wrapped and become operational at the end of 1988. Finally, schematics for the multiwire “production” version were finished in mid-February of 1989, and the final boards were fully tested just in time for SRC’s 1989 summer workshop on Splash 1 applications.

Programming of Splash 1 was originally done with the Xilinx-supplied XACT editor; later tools included the Viewlogic schematic capture package. From the outset, however, there were difficulties in developing application codes; especially for individuals unused to hardware design, programming with XACT was not easy. To make the machine more accessible, the Logic Description Generator (LDG), a higher-level language whose output could be mapped to the Xilinx chips, was designed and implemented at SRC through the fall of 1988 and the winter of 1988–89 by Maya Gokhale [3]. In addition to the software for direct execution, a debugger called *Trigger* was used extensively.

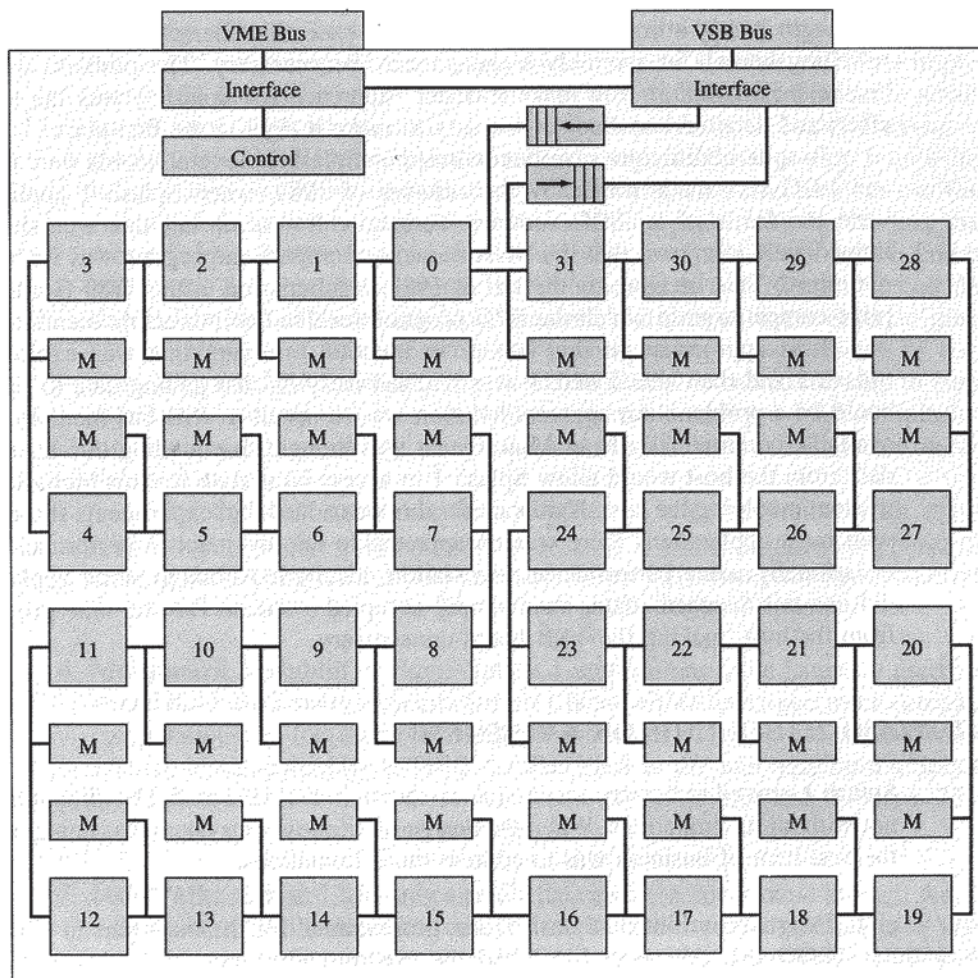


FIGURE 4.1 Splash 1 Architecture

These new tools permitted some escape from the low-level details of hardware design, but programming the Xilinx chips was still a nontrivial task. Problems existed at several levels. Splash 1 programming was still hardware design. Counters and sequencers had to be explicitly constructed and connected to the logical units that they controlled. Timing information about a design was difficult to obtain, and the inability to perform complete state readback and restore hindered the debugging process. One of the hardest problems centered on the *apr* (Automatic Place and Route) software from Xilinx. This program ran slowly and often failed to completely route a chip's design. It was often, if not usually, necessary to provide *apr* with hand-placed designs in order to get fully routed designs with acceptable execution speeds as a result of *apr*'s work. One major focus—which has been successful—in moving from Splash 1 to Splash 2 has been to eliminate the need for such low-level effort in order to get a working design. Programmers have for years been accustomed to the fact that, if speed is the object, the working program (usually in a high-level

language) is only a first step; detailed study of execution characteristics and possibly rewriting kernels in assembly language may be necessary. The problem with this “make it right before you make it better” approach for Splash 1 was the level of effort and detailed knowledge necessary to make it right in the first place.

In spite of difficulties, work continued on Splash 1. Several boards were finished and LDG was made robust by the summer of 1989, when Splash 1 applications were the focus of an SRC summer workshop. It was during this workshop and immediately afterward that the DNA sequence comparison program was written and optimized. This became, by the fall of 1989, a submission to the 1989 Gordon Bell prize competition, in which the SRC program received an honorable mention.

It is appropriate in this history to mention one path that was explored for Splash 1 and then abandoned. It was realized early on that getting data to Splash 1 could be a problem. An apparent solution was to install a VSB bus memory board, available commercially from Motorola. It was thought that loading this board with data from the host would allow Splash 1 to access large data streams multiple times without involving the host. Hardware for this was added, but experiments showed that with most applications there was no appreciable improvement in performance, and occasionally some performance degradation occurred. Although some applications did use this “cached” data, most simply accepted a stream for one-time processing from the host, making the VSB board unnecessary.

## 4.2 SPLASH 2: THOUGHTS ON A REDESIGN

Splash 1 proved to be very successful, as shown in Gokhale et al. [2], although it was not without its limitations. When the design of a follow-on system was contemplated, the first item of business was to address those limitations.

- 1. Programmability:** Splash 1 was programmed for the most part in Gokhale’s LDG [3]. Although LDG had the obvious advantage of having been done expressly for programming Splash 1, it had the disadvantages associated with being an internally developed system. In addition, the Xilinx tools were unequal to the task of supporting code development for Splash 1. Users often had to perform placement themselves in order for the `apf` tool to succeed in routing a design. Further, many of the problems in debugging a design required more detailed knowledge of the Xilinx design than could be obtained from the tools. These problems combined made Splash 1 difficult to program by individuals unused to hardware design. One of the key issues in planning Splash 2, then, was to make it *programmable*, to make it a processing system that would be usable, without (undue) agony, by a wide range of programmers.
- 2. I/O Speed:** Many of the original applications for Splash 1 were strongly I/O-bound. The VMEbus can deliver about 4 Mbytes/sec (in slave mode; VME using DMA would be significantly better) from the host to a Splash 1 board, but an application running at 16 MHz needs a bandwidth of 32 Mbytes/sec in order to consume and produce one character per clock tick. Any follow-on system would have to overcome the I/O bottleneck of Splash 1 in order to be truly successful.

**3. Memory:** The primary uses of the memory chips in Splash 1 were for lookup tables and for storing microprograms to be executed by state machines implemented in the Xilinx chips. In some instances the lookup tables were small, but in some applications, such as integer multiplication, a single reference to memory would replace complicated and slow logic and the speedup from using lookup tables would be limited only by the memory size. The memory use was encumbered, however, by problems in sustaining peak rates, by the fact that memory loads had to be done down the linear data path of the Splash board, and by the fact that the memories were connected to two Xilinx chips on the linear data path. This last problem required that programmers exercise great care to separate in time the access to memory and the transmission of data from Xilinx to Xilinx. On the other hand, the ability of one Xilinx chip to store data into a memory that the next Xilinx chip would read provided on Splash 1 a communication capability that was useful but which was not retained in the Splash 2 design.

The Splash 1 memories were also not as large as desirable for many applications, and the requirement that memory reads and writes take place by passing data through the FPGAs caused unnecessary complications, requiring that a special Xilinx program for memory read/write be written and used.

**4. Multiboard Scalability:** Some Splash 1 applications quite naturally used (or would have used) more than one Splash 1 board either for larger, more complex computations or for multiphase computations. In order to use a multiboard system, it was necessary to bring the data back to the Sun host from a Splash 1 board and then to send it from the host to the next Splash 1 board. This aggravated the I/O bottleneck problem.

**5. Data Path:** Splash 1 had only a single data path—a linear route through the 32 Xilinx chips. While the linear (which is sometimes systolic) paradigm is very powerful and its application to Splash has been very successful, there were many applications that either could not be done or whose efficiency suffered because the linear path was the only data path. Given that high performance seemed to require careful control of the data pipeline past the somewhat limited processing resources, the cost of transmitting data to the appropriate processing element needed to be diminished.

In Splash 1, the data from the host passed directly into the Splash board, so that handling of the FIFOs and any preconditioning of the data (merging of two input streams, for example) had to be done by the Xilinx chips on the Splash board. This often complicated the programming. Strong suggestions had surfaced early that performance might be substantially improved if the data preconditioning were moved out of the general processing array.

Finally, on Splash 1 it was often observed that 32-bit data widths were sufficient for the data but that extra tag bits sent along with the data would have been very useful. Not having the extra bits either complicated the design of programs or required extra clocks (a major negative factor in a highly pipelined program) in order to transmit the necessary control information from FPGA to FPGA. An extension of the input/output data path width from 32 to 36 bits would probably remedy this shortcoming.

- 6. Clock:** The Splash 1 clock had only power-of-two speeds. Designs that came close to running at 32 MHz could only be run at 16 MHz, for example.

The above list addressed the known and specific limitations of the Splash 1 boards as built. Once a redesign of a Splash-like system was contemplated, however, all the earlier design decisions were reviewed. The original design of Splash 2 by Andy Kopser was given in [4] in a preliminary architecture description. On September 12, 1991, however, all these decisions came up for review at the first of a series of architectural design meetings. These meetings were intended to start from first principles to design a Splash-like FPGA-based processor; although Kopser's earlier thoughts were taken into account, none of his conclusions was accepted as given—all were subject to further scrutiny. Among these were the following:

1. Choice of FPGA chips
2. Choice of host and connecting bus
3. The linear array and any other interconnection of the FPGAs
4. Multiprogramming of multiple Array Boards

#### 4.3 PROGRAMMING LANGUAGE

The decision to use VHDL [5] as the language in which to program applications for Splash 2 was actually made quite early. The use of VHDL would clearly be a compromise. In its favor were the facts that it is a defined standard, that it is a programming language (at least in simulation mode), and that it is supported by commercial tools for both simulation and synthesis. The commercial tools also provide a programmer with most of the bells and whistles of a debugging environment that are now expected by users. Finally, the goal of Splash 2 was to demonstrate the ability to program logic into an FPGA-based machine; although a high-quality translation of a high-level language to Xilinx bitfiles would be necessary for performance, we did not feel that we wanted to make it a significant part of the Splash 2 project itself.

Working with an off-the-shelf VHDL system would not, however, address all the issues involved in programming Splash 2. Quite apart from the "religious" issue that VHDL is an Ada derivative while most modern programmers are using C, there would be known problems both "above" and "below" the VHDL level. At the time Splash 2 was begun, it was not clear that it would be possible to drive a VHDL simulation from a general C-language interactive front end. The user's view of the programming environment might necessarily be that of the VHDL vendors' tools—which were designed for use by engineers doing circuit or VLSI design and might seem unduly foreign or even hostile to application programmers. More important, due to the need to have some example applications achieve high performance, it was not clear that the output of the logic synthesis would produce a Xilinx bitfile that would use crucial performance features of the FPGAs. Much of the success of Splash 1 had come when the programmers had specifically controlled from LDG the resources on the XC3090 FPGAs. A serious question was whether an outside vendor whose tools were aimed at a very different target consumer would provide the resource utilization that we would need.

It has thus always been assumed that VHDL is not perfect and that some language more like C should be developed. We realized, however, that we didn't know enough *a priori* about programming Splash 2 to permit development of "the right language." VHDL was therefore viewed as an acceptable middle ground, with the hope that in the process of programming Splash 2 in VHDL for a varied list of initial applications, enough would be learned about the programming model appropriate for Splash 2 to permit language development after the fact. Meanwhile, useful work would be accomplished by those brave pioneers who had coded the original applications in VHDL.

#### 4.4 CHOICE OF FPGAS

SRC had gained extensive experience with and understanding of the Xilinx XC3090 chips, much of which would translate to the new XC4000 series chips, but the question was opened as to whether a different vendor's chip might be more desirable. Prominent among the options was the Concurrent Logic, Inc. FPGA. The two chips appear remarkably alike to a "computer designer," despite the extremes of granularity between the two products. The Xilinx XC4010 has 400 Configurable Logic Blocks (CLBs) in a square array. Each CLB takes two sets of four inputs and produces any Boolean function of each set, then any Boolean function of the two bits of result together with a ninth input signal. With this coarse structure, Xilinx advertises the XC4010 as roughly equivalent to a gate array of 10,000 gates [6].

The Concurrent Logic chip, by contrast, is very fine-grain. The high-end chip at the time was the CLi6005, a  $56 \times 56$  array of cells that in most modes serves to produce one output from three inputs. Concurrent Logic advertised its CLi6005 chip as being roughly equivalent to a gate array of 5,000 gates, the similarity between the Xilinx and the Concurrent Logic figures perhaps saying more about the basic complexity of 1992 silicon technology than about the clear superiority of one vendor over another.

In the final analysis, three factors were decisive:

1. The fact that the Xilinx chips were a known quantity made it necessary to have a very good reason to change.
2. The delivery schedule for the Concurrent Logic chips was some months behind that for the Xilinx chips.
3. Most important, as a technical matter, the Concurrent Logic chips had 108 I/O pins compared to the Xilinx's 160. As it was envisioned at the time the decision was made, even the Xilinx's 160 I/O pins seemed insufficient. This was borne out by later experience.

In the process of deciding on a chip, it was necessary to compare not just the chips but to take into account their features and the processing power per Array Board that could be accommodated. A feature new to the XC4010 chip was a fast carry internal to the CLBs, which makes arithmetic computations faster and requires less programming and fewer CLBs. Further, the number and quality of the interconnection lines had increased, which would help more applications run at higher speeds. Finally,



the new chips allow for the use of CLBs as a 32-bit RAM, configured either as  $32 \times 1$  bit or as  $16 \times 2$  bits.

The major difference between the XC3090 and the XC4010 chips, however, was in the basic size and structure—the XC4010s have 400 Configurable Logic Blocks (CLBs) instead of the previous 320, each CLB has nine input lines instead of five, and the maximum speed is 40 MHz instead of 32 MHz. The improvements in the FPGAs to be used would permit Splash 2 to have 17 Xilinx chips on an Array Board instead of the previous 32. This was both a conscious decision and a necessity. The newer chips were, in the packaging available at the time, physically somewhat larger, and it was not possible to put 32 of them on a single Array Board along with the memories and the crossbar (to be described in Section 6). The hope was that because the newer chips were each more powerful than the old chips, and because it had been more often the case with previous applications that they were I/O-limited rather than processor-limited, the plan to use half as many chips per Array Board, each perhaps somewhat less than twice as powerful, would provide a reasonable processor-to-I/O balance.

As it has come to pass, the decision to use the Xilinx chips has not been without its problems. The VHDL tools used in programming Splash 2 reduced the VHDL program to the gate level in the synthesis step rather than constructing efficient CLB designs, so a natural inefficiency exists in the use of the VHDL language for the Xilinx FPGAs. (To a great extent these issues were addressed in working with Synopsys on their FPGA Design Compiler.)

#### 4.5 CHOICE OF HOST AND BUS

In the design of Splash 2, there were no tacit assumptions, and even the choice of host and bus were open for discussion. The options were narrowed considerably, however, by various practical considerations. Sun workstations continued to be the norm at SRC, and that, coupled with Sun's dominant position in the overall workstation market, made it hard to really consider abandoning Sun as a host. The constant realization existed, however, that the object of study was "the Splash 2 attached processor" and not "attached processors for Sun workstations." By constantly keeping in mind the fact that the choice of hosts contained a nontrivial degree of arbitrariness, we were able to avoid embedding the machine so deeply in the Sun milieu that it could not be re-engineered at modest cost for a different host.

The choice of the particular host was not so arbitrary. The realistic options were the SPARC 1+ and the SPARCstation 2, using the SBus with either. Experiments on SPARC 1+ workstations at SRC showed, however, that the SBus as implemented did not match the SBus as documented by Sun. Since the SBus was the logical choice (for reasons described below), this forced the decision in favor of the SPARCstation 2.

These decisions, together with the decisions on the Futurebus+ backplane, were not made in a vacuum. Another hardware-build was under way at SRC at the same time, also a workstation enhancement a few months ahead of Splash 2 in development, and the decision was made that Splash 2 would use the same bus, backplane, and so on. The real goal of Splash 2 was to demonstrate the capability of FPGA-based computing, and the use of hardware in common with the other project would permit more effort to be directed to the specifics of meeting that goal.

Having decided to use the SPARCstation 2, the decision to use the SBus was rather easy. By far the most limiting hardware feature of Splash 1 was the 4 Mbytes/sec peak data rate of the VMEbus on the Sun host. While the VMEbus protocol is rated as high as 16 Mbytes/sec peak transfer rate, existing implementations of VME buses do not reach that rate. Above all else, we did not want Splash 2 to be, as was Splash 1, I/O bound, and the SBus appeared to provide at least an order of magnitude higher data rate than the Sun VMEbus. Early estimates were 38 Mbytes/sec; tests now show that a CPU-loaded SPARC 2 can sustain about 40 Mbytes/sec through the SBus via DMA and that an unloaded machine can deliver as much as 54 Mbytes/sec. Unfortunately, these are peak transfer rates from the same 16-word DMA buffer of the host. Without further work on the drivers on the host, perhaps to include double-buffering, the data transfer via DMA takes place during only about 40 percent of the time, so transfer rates actually are limited to 20 to 25 Mbytes per second.

Along with the decisions on host and bus, driven by the need to provide Splash 2 with data, the plan for data connections that did not go through the host was an early feature of the design. As is remarked upon by nearly every experimenter with workstation attachments (including Bertin et al. [1]), workstation discs are too slow to produce volumes of data at high speed, memories are too small, and the Ethernet connections simply cannot sustain the load. The external connection could be to a traditional supercomputer functioning as an I/O device (which a supercomputer does quite admirably) or to an array of discs (as one might find in a text search or database search application).

Several different external data connections seemed desirable and potentially usable. An early plan to include several such connections on the Interface Board was dropped in favor of Wally Kleinfelder's idea to put such connections onto a daughterboard. This would allow the DMA to be built and tested on an early-version Interface Board while the final Interface Board was being designed and would also allow future daughterboard designs such as HiPPI.

Finally, another early decision was that the Sun and Splash 2 would use different clocks. It was a simple matter, then, to leave open the option of having the external data also carry the clock to be used by Splash 2.

#### 4.6 CHIP-TO-CHIP INTERCONNECTIONS

One of the major decisions in designing Splash 2 was the choice of chip-to-chip data paths. In this, the 160 I/O pins on the XC4010 chip turned out to be one of the forcing factors.

With 160 I/O pins, one can implement four 36-bit data paths but have only 16 pins left over for control. Even by dropping the tag bits (bits 35-32 of the data path) one cannot get five 32-bit ports (this is exactly 160, leaving nothing for control). We are, therefore, necessarily designing for a four-port array.

The obvious extension from a one-dimensional array would be a two-dimensional array of Xilinx chips. This consumes four ports, so memory connections would have to be shared with the chip-to-chip connectors, as on Splash 1. Two such arrays are given in Figures 4.2 and 4.3.

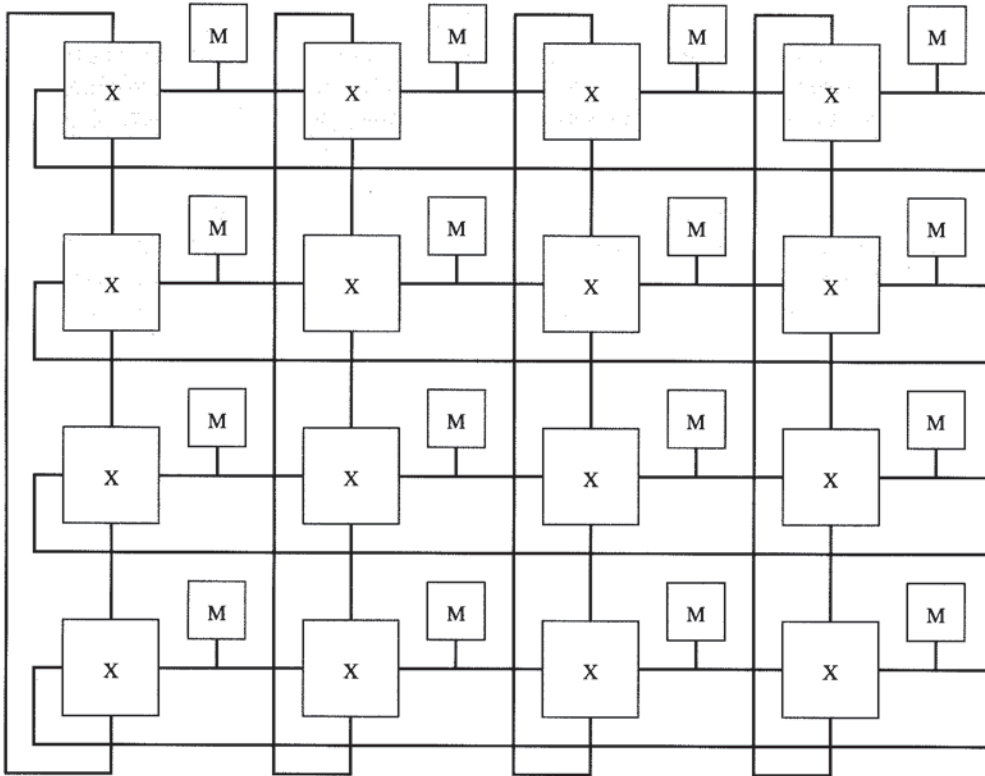


FIGURE 4.2 Two-dimensional Toroidal Mesh

In Figure 4.2 we easily obtain a 2-D mesh, but we cannot easily string the chips together into a 1-D array if the memories are also being used. That problem is fixed in the design of Figure 4.3—at the expense, of course, of the 2-D mesh itself.

The basic problem is simple: If the memories are connected to two processor chips by the same lines that are used for processor-to-processor connection, then a linear array of processors with memories in between uses up two of the four ports per chip. If the memories are active, then with only two ports per processor we can achieve only a linear array.

Similar objections ruled out the use of busses, and a major step was taken in the decision to connect each processor to its own local memory. Some capability was lost with this decision. There were a few Splash 1 programs whose efficiency was due to the ability of the FPGAs to transfer data through the memories; the stored output of one FPGA's work could be accessed in an arbitrary order by the next FPGA. But this would have required a wider data path to accommodate two ports, or the time-multiplexing of access to memory by the FPGAs, and the advantages seemed not worth the hardware investment or the increased complexity of programming the memory use.

With one of the four ports per chip used for memory, three ports become available for chip-to-chip communication. Various three-port configurations were considered, including the chordal ring. A chordal ring (a ring with regularly or irregularly

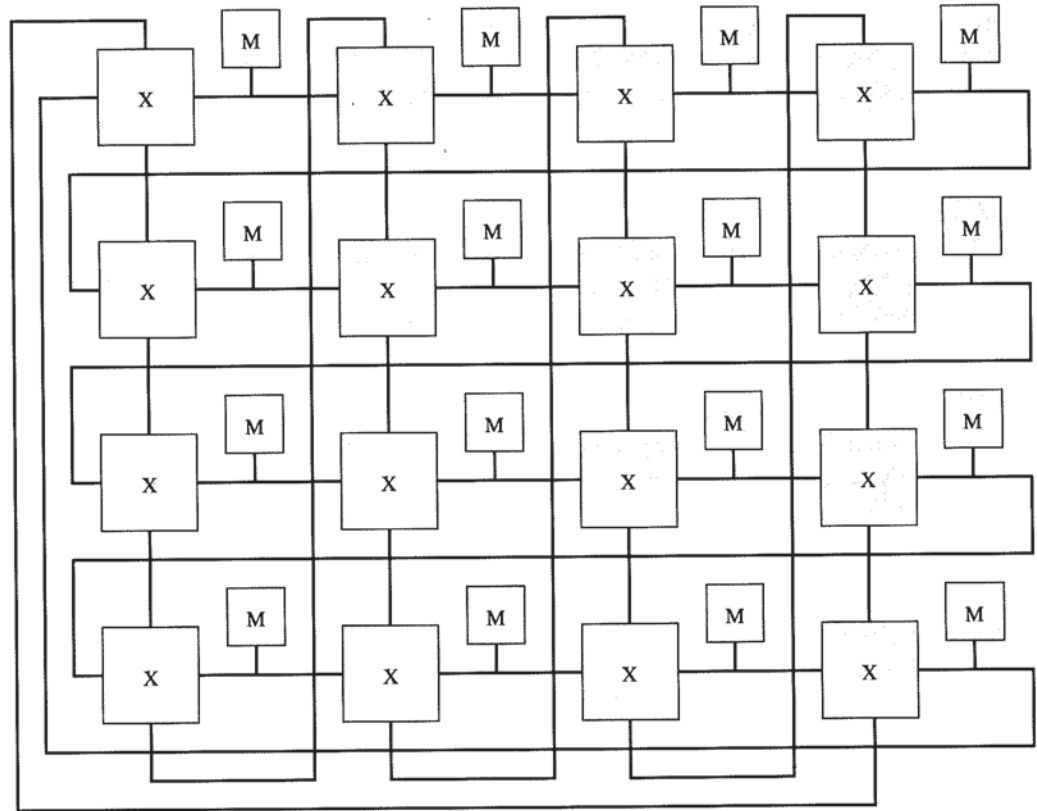


FIGURE 4.3 Two-dimensional Toroidal Mesh with Shift on Wraparound Connections

spaced “chords” as additional connections) can be used as a linear array, or a linear array with shortcuts, and had the advantage of not being tied to “nice” numbers (like 16) of processor chips. A possible drawback, though, is that “normal” programming patterns do not (yet?) include the chordal ring as routine.

For these and similar reasons, the eventual choice was a linear array with two ports per chip and a crossbar connecting the chips with the other port. This choice was made easier when made in conjunction with a choice of the TI reconfigurable crossbar chips (TI SN74ACT8841 were used). Each such chip is a 16-port, 4-bit-wide crossbar and can be programmed with eight different configurations. In this way, data paths in nibble sizes could be programmed (although pin limitations later limited this to byte sizes), and a wide variety of communication patterns could be accommodated.

Although it turned out later to lead to serious problems, the choice of the TI chip seemed an excellent one at the time. In the manner in which we would use the chip, the latency across it was one tick, so the crossbar communications would not differ from communications down the linear array. This would simplify the programming task—slower chips would require the programmer to insert pipeline stages in a program and then to synchronize them carefully. The multiple configurations permitted by the TI chip seemed to provide the logical connectivity needed, and the ability

to switch configurations with at most a one-tick delay was a very attractive option. None of the more sophisticated switch chips available at the time, nor the use of FPGAs to implement the switch, offered these features.

#### 4.7 MULTITASKING

Up to now, we have been discussing the architecture of the Splash 2 Array Board or the data paths in and out. The decisions on these were originally made in Kopsler's Splash 2 design and then reaffirmed by the architectural committee. One of the design decisions that was changed radically was an original thought that each Splash 2 Array Board would have its own input and output FIFOs and that each Array Board could run a separate Splash 2 process multiprogrammed from the host. This would have required extensive control hardware on the Interface Board as well as complete software protocols for the Sun host's control and context switching of processes running independently and asynchronously on the Splash 2 attached processor.

The possibility of allowing Splash 2 Array Boards to work on separate tasks was hotly discussed for several sessions at the weekly architecture meetings. The goal—which, in the final analysis, seemed impossible to achieve—was that arbitrary subsets of Splash 2 Array Boards could be chained together, each running distinct processes, possibly communicating with each other directly, possibly communicating through the host, and possibly not communicating with each other at all.

It was even envisioned at this time that different users might run different programs on Splash 2 concurrently, in addition to the situation in which a single user might have multiple independent Splash 2 processes.

In the end, all such plans were discarded. The final Splash 2 system is an attached processor in which all the Array Boards in a given system form a linear chain; the only variations in configuration are that broadcast to all Array Boards simultaneously is possible and that the physical chain of Array Boards can be logically shortened. Although use of the entire system can be time-shared, no partitioning of the system for concurrent execution of independent processes in different partitions can occur.

The major factors in this decision were: that the algorithmic complexity of controlling the independent processes would require too much hardware support if it were designed to run at the necessary speeds; that the complexity of the "back-end" network controlling the subsets of Splash 2 Array Boards into chains would be too great; and that insufficient real estate existed even on the large Array Boards planned for Splash 2 to allow for the Xilinx chips, memories, and crossbar, as well as FIFOs, DMA controllers, bus arbitration with the Interface Board, and network communication with the other Splash 2 Array Boards in a chain. A final concern was that each subsystem would have to be able to run at a different clock rate so that maximum efficiency of the Splash 2 processes could be obtained. This would clearly have necessitated complex mechanisms to arbitrate bus and DMA access for data movement on and off Splash 2.

Somewhat reluctantly, then, Splash 2 became a system all of which, at any point in time, would be assigned to a single process. The Array Boards were to form a linear array (although broadcast was still possible). The FIFOs and DMA control for each Array Board were consolidated into one pair of input and one pair

of output FIFOs using DMA channels and moved onto the Interface Board. Interrupts and global AND/OR were similarly cascaded from each Xilinx chip to a board-level register and from board level to a system-level register on the Interface Board. The inclusion of Xilinx chips XL and XR on the Interface Board would provide for control of data transfer, clock (even a clock supplied by the external input), and tag bits independent of the Splash 2 Array Boards. In Splash 1, such control had usually been done in the first array chip, leading to asymmetry and crowded designs. With proper programming of XL and XR, the asynchronies of DMA transfer and external input and clock should not be seen by the Splash 2 Array Boards themselves, and the XL and XR programs should function much like a system I/O library.

With the adoption of this more conservative plan, some applications were given up, but the general opinion was summed up by Ron Minnich: "Now, at last, I think we have a real chance that this thing can be built."

#### 4.8 CHIP X0 AND BROADCAST

One of the casualties in moving the I/O off the individual Array Boards onto the Interface Board was that it was no longer as simple to envision broadcast of data from the host or from an external interface to all the processor FPGAs simultaneously; yet this programming model was seen to be equally as important to support as the simpler model in which long streams of data would pass down the linear array. Unfortunately, the basic power-of-two dilemma existed. Earlier Splash 1 programs had occasionally been complicated or suffered decreased efficiency because FPGAs X0 and X31 had handled I/O; this took away somewhat from the power-of-two advantage of the 32-long linear array. When the I/O handling was moved from the beginning and end of the linear array on each Array Board to FPGAs XL and XR on the Interface Board, the power-of-two structure returned to the processing array when viewed as a linear array.

Now, however, something needed to be done for broadcasting data. With some reluctance, FPGA X0 was added to each Array Board. The reluctance came primarily from the realization that 17 is not a very elegant number. Reading and writing the memory on each Array Board would become more complicated (having made a decision to include another FPGA on the Array Board, there was little dispute over making it look as much like the other FPGAs as possible, so it was given the same memory as the other FPGAs), reading and writing the configuration and state of the FPGAs would also be more complicated, and most inelegant of all, X0 would have to share lines into the crossbar with some other FPGA.

Despite the inelegance, the Array Board architecture with the 17 FPGAs has proved to be successful. The complications from not having power-of-two structures have been more than compensated for by the greatly increased ability to move data into the Array Boards and to the PEs.

#### 4.9 OTHER DESIGN DECISIONS

In addition to these "coarse" decisions on the architecture, a number of other changes were made to the Splash 1 design. Data paths were fixed at 36 bits in width. This would accommodate 32-bit words and 4-bit tags carried along with them. An earlier

plan for a 40-bit crossbar was scaled back to 36 bits. No good use could be immediately envisioned for the extra four bits, and the I/O pins were needed elsewhere. The only place in which Splash 2 is not a 36-bit machine is in the Xilinx-to-memory path. This path is 16 bits wide, largely because the 160-pin Xilinx XC4010 chip cannot support three 36-bit data paths (two ports for the linear array and one into the crossbar) and a path to memory at least 32 data bits wide ( $3 \cdot 36 + 32 = 140$ , leaving only 20 pins for Xilinx control and memory address). It was thought to be absolutely necessary that memory be accessible in every clock period, making multiplexing of data and address infeasible.

One change from Splash 1 to Splash 2 was to add a separate memory read/write path that did not require going through the Xilinx chips. The memories could now be directly read/written from the Sun host over the SBus. They are not, however, dual-ported; the FPGAs must be inactive during the read/write operations. This change allows tables to be loaded in bulk and results to be read from the memories without requiring the circuitous path through the Xilinx chips. However, since the memories were 16 bits wide and the "natural" word size of both Splash 2 and its Sun host is 32 bits, an obvious question arose as to where the conversion from 32 to 16 bits would take place. In this case, no answer is perfect. Placing halfword data on word boundaries in the host is easy for the programmer but wasteful of memory space on the host and of I/O bandwidth to Splash 2. Packing halfwords two to one into Sun words is a slight annoyance for the programmer but wastes neither bandwidth nor memory. This latter choice prevailed after it was determined that, in fact, the memories could be double-cycled on the Splash 2 Array Boards fast enough to keep up with accesses from the host; the host would never know that the Splash 2 memories were 16 and not 32 bits wide.

One of the suggestions made and discussed was whether to include on each Splash 2 Array Board a microprocessor to perform tasks not easily or efficiently done by the Xilinx chips. Although this was a suggestion made at a time when independent execution of Splash 2 Array Boards in a multiprogrammed environment was being considered, it was an idea that had a wider context. Many of the uses to which FPGAs for computing have been put have been to augment the instruction set of a microprocessor; one could easily imagine a Splash 2 Array Board being viewed in this way by a user focusing attention on an on-board microprocessor. In the end this idea was not pursued, largely because the control of the Array Boards and programming would be overly complex. The intent, as discussed in the next paragraphs, was to be able to program the machine in a high-level language. It appeared, however, that in order to make effective use of the microprocessor, it would be necessary to control and synchronize processing and data movement on the Array Boards at a very low level. Splash 1 applications had been, and Splash 2 applications were expected to be, highly pipelined, but in order to do this on a Splash 2 Array Board, the output of the microprocessor's compiler would have to mesh closely with the "program" for the Xilinx FPGAs. The software effort for the Xilinx part of the Array Board seemed difficult enough without the complications that the on-board microprocessor would add.

Throughout the design and architecture refinement, the emphasis was on producing a machine that could be programmed at a moderately high level. This fundamental assumption about what it would take to produce a successful computing engine had an effect, as mentioned, on many of the design decisions, especially

the decisions involving the degree of independence of one board from another and the control of board-level entities. For example, the inclusion of a microprocessor on each Splash 2 Array Board would have required programmers to code not only asynchronous Sun/Splash 2 execution but also to coordinate the interaction of the on-board microprocessor and the Xilinx chips. It did not seem clear that this could be done in a high-level language in a way that would be tolerated by programmers, and similarly, it did not seem clear that a compiler could readily be written that would deliver the performance that users could have a right to expect.

The discussion about programming continues to this day. VHDL is, on the one hand, sufficiently high-level and sufficiently modern to be recognized and accepted as a "programming language." On the other hand, it does retain many of the quirks of hardware description, and mastering the methods for getting around these quirks does not render much less arcane the art of VHDL. It is not "C-like" and cannot by itself be made to be "C-like."

At the heart of many of the issues surrounding the programming of Splash 2 is the fact that the architecture at present is completely exposed to the user, who sees, in essence, the memory address and data registers, the specific data paths, and so forth.

#### REFERENCES

- [1] P. Bertin, D. Roncin, and J. Vuillemin, "Programmable Active Memories: A Performance Assessment, in G. Borriello and C. Ebeling, eds., *Research on Integrated Systems*, MIT Press, Cambridge, Mass., 1993, pp. 88-102.
- [2] M. Gokhale et al., "Building and Using a Highly Parallel Programmable Logic Array," *Computer*, Vol. 24, No. 1, Jan. 1991, pp. 81-89.
- [3] M. Gokhale et al., "The Logic Description Generator," Tech. Report SRC-TR-90-011, SRC, Bowie, Md., 1990.
- [4] A. Kopser, "Splash 2: Architectural Motivation," tech. report, SRC, Bowie, Md., 1991.
- [5] D.L. Perry, *VHDL*, McGraw-Hill, New York, 1991.
- [6] Xilinx, Inc., *The Programmable Gate Array Data Book*, Xilinx, Inc., San Jose, Calif., 1993.



# CHAPTER 5

---

## Software Architecture

*Jeffrey M. Arnold*

### 5.1 INTRODUCTION

As we saw in Chapter 4, the Splash 1 system was programmed at the logic gate level with the macro language LDG [4]. This meant that the process of developing applications for Splash 1 was very labor-intensive, requiring a detailed understanding of the internal structure of the Xilinx devices. For this reason, applications programmers with little or no hardware experience found Splash 1 extremely difficult to program. The result was that there were never more than half a dozen proficient Splash 1 “programmers,” and these were people with extensive hardware design backgrounds.

When we set out to design a software environment for Splash 2, our main objective was to improve the ease of programmability of the system, opening it up to a much larger audience of applications developers. The specific design goals of the Splash 2 software environment were to:

- select or develop a procedural language for writing applications
- provide a rich debugging environment that did not require a detailed understanding of the hardware
- provide a smooth and efficient interface between the host computer and Splash 2
- develop a comprehensive set of diagnostic tools for hardware development and maintenance
- leverage commercial off-the-shelf technology wherever possible

With these goals in mind we chose to base the Splash 2 programming environment on the VHSIC<sup>1</sup> Hardware Description Language (VHDL) [6, 10] and modern Computer Aided Design (CAD) tools such as simulation and logic synthesis. Applications for Splash 2 are developed by writing behavioral descriptions of algorithms in VHDL, which are then iteratively refined and debugged within the Splash 2 simulator. During the course of this iteration, the VHDL implementation is manually partitioned by the programmer into a set of individual FPGA programs. Once the partitioned implementation is determined to be functionally correct in simulation, it is compiled and optimized to produce a network of logic gates. This *gate list* is then mapped onto the FPGA architecture by automatic placement and routing tools to form a loadable FPGA object module. Static timing analysis tools are applied to the object module to determine the maximum operating frequency and the set of critical paths. This information is fed back to the user, who may choose to manually optimize the design. The runtime system provides the interface between the host computer and the Splash 2 system and consists of a C language library and an interactive symbolic debugger.

This chapter presents the architecture of the Splash 2 software system. We begin with a background discussion of the underlying CAD technologies that make custom computing possible. We then proceed to justify our choice of VHDL as the programming language of Splash 2. Next is a discussion of the architecture of the programming environment and the system software. Finally, we present the models of the system the programmer sees at each of several levels of abstraction.

## 5.2 BACKGROUND

The success of Splash 2, and of custom computing in general, has been made possible by the confluence of two important technologies: infinitely reprogrammable logic arrays (static RAM-based FPGAs) and high-level CAD software. Over the past few years the CAD industry has made significant advances in automatic generation of hardware design from high level specification. This process may be divided into two steps: *logic synthesis* and *physical mapping*. Logic synthesis is the process by which procedural descriptions of algorithms are mapped into Boolean logic gates, bypassing traditional structural techniques such as schematic capture. The physical mapping process converts the resulting gate list into a specific hardware technology, such as the static RAM- (SRAM-) based FPGAs used in Splash 2. Together, these technologies move the task of application development for custom computing from the realm of hardware design into the realm of software programming.

Figure 5.1 illustrates the flow through the Splash 2 program development process from design entry through hardware configuration. There are two feedback loops in this flow. The inner loop is used to establish the functional correctness of a program by simulating the design and observing the response to a set of test vectors. The outer loop constructs the physical implementation by synthesizing and optimizing the logic and then mapping the result into the FPGA technology. A static timing analyzer

<sup>1</sup>The Very High Speed Integrated Circuits (VHSIC) program was an initiative funded by the U.S. Department of Defense in the late 1970s and early 1980s.

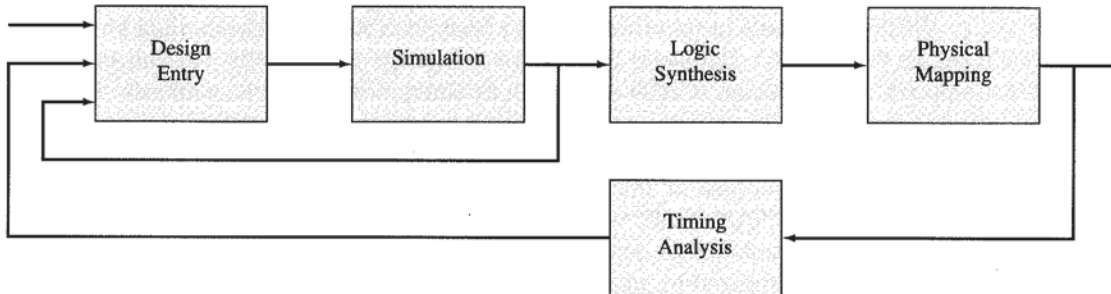


FIGURE 5.1 Splash 2 Program Development Process

is used to predict performance and identify potential bottlenecks. The programmer may use this information to determine overall system performance and possibly guide further optimization.

Logic synthesis [2, 3, 13] is the process of converting a high-level description of an architecture into an optimized logic implementation. The input to the synthesis process is typically in the form of a procedural or mixed procedural and structural description of the intended architecture. The logic synthesis tools extract control and data flow information from this description and produce a set of Boolean equations and module instances that perform the desired function. This internal representation is then optimized to meet user specifications of area and delay. Since the design has not been mapped into the logic blocks of the FPGA technology at this stage in the synthesis process, the optimization must be based upon estimates of logic block packing, logic propagation delays and a fan-out-dependent statistical model of the routing network. Many of the parameters that control the optimization may be set by the user, allowing trade-offs to be made between minimizing area and maximizing performance. The output of the synthesis process is a list of technology-independent logic gates.

The physical mapping [3, 16] process converts the generic gate list produced by logic synthesis into a configuration bitstream for the particular FPGA by *partitioning* the gates into logic blocks, *placing* the logic blocks into the FPGA, and *routing* the signal nets between the blocks. The partitioning phase groups the combinational logic gates into Boolean functions that will fit in the lookup tables of the logic blocks (3 and 4 inputs for the XC4010) and assigns registers to the flip-flops of the logic and I/O blocks. During the partitioning phase it is often possible to trade chip area (gates) for speed by replicating functions that have a high fan-out. Unfortunately, this trade-off requires a close coupling between the synthesis and mapping processes that is not present in today's tools.

The placement step accepts the partitioned design and determines a good placement for the logic blocks in the FPGA array [16]. Most FPGA placement algorithms use a stochastic optimization algorithm such as simulated annealing to minimize a cost function such as total net length. Traditional integrated circuit routing techniques are based on decomposing the area available for wiring into rectangular "channels" that can be routed independently. Unfortunately, this approach does not work well for FPGAs, because the interconnect resources are fixed in place. Therefore, most FPGA routers use a form of maze router that does not decompose the

problem into independent routing channels. User-specified timing requirements can be used to guide the router by building a detailed delay model from the physical interconnect parameters of the FPGA. Finally, the router must handle additional constraints imposed by the FPGA such as special clock networks and carry-chain routing.

Once the detailed routing is complete, the static timing analyzer is able to make an accurate prediction of the maximum operating frequency and determine the critical paths of the design based upon the known logic block and routing resource delays. The programmer may use the critical path information to manually optimize the design or restructure the program for resynthesis. The delay information extracted from the design by the static timing analyzer may also be used to construct a structural simulation model of the design, which in turn can be used to perform detailed timing simulation.

A great deal of research and development in the area of FPGA design tools is taking place in academia and industry, with the result that the quality of the available tools is rapidly improving. We therefore felt it was efficacious to leverage “commercial off-the-shelf” technology for Splash 2 as much as possible, allowing ourselves to concentrate on the system integration issues.

### 5.3 VHDL AS A PROGRAMMING LANGUAGE

One of the most important objectives of the Splash 2 software effort was to move the task of application development from the realm of hardware engineering to the realm of software programming. This desire led to several selection criteria for a “production” programming language for Splash 2. Among these criteria were support for the use of procedural as well as structural specification, and the ability to build higher levels of abstraction through encapsulation of function. To support high-performance applications, we felt that the language should include an escape mechanism to allow the programmer to explicitly specify hardware details. Finally, the language had to be directly executable to allow interactive source-level debugging of application programs.

In the early stages of the Splash 2 effort we explored the option of developing our own language based upon a subset of C. Such a language would have the advantage of familiarity to most users, be directly executable on a wide variety of platforms, and come complete with a rich development environment. However, we felt that the task of compiling a subset of C into hardware would quickly become a major research project in its own right, detracting from the Splash 2 system development effort. Therefore we chose to focus our efforts on system integration, leveraging commercial logic synthesis tools by basing the Splash 2 programming environment on VHDL.

Ultimately, we believe the best programming model for custom computing machines is to develop higher-level programming languages that can be compiled into a form suitable for input to commercial CAD tools. Such a language would synthesize an application-specific architecture, perhaps use VHDL as an intermediate language, and use commercial logic synthesis in the “assembly” process. One such effort, based upon the dbC language [5], is described in Chapter 7.

### 5.3.1 History and Purpose of VHDL

VHDL evolved from an effort to develop a design specification and interchange language common to all of the participants of the VHSIC program [8]. The language traces its roots back to a planning session in 1981, although the initial development effort was not begun until 1983. The importance of this work became clear to the broader engineering community with the first release of the language and simulator in 1985. A standards committee of the IEEE was established to further refine the language, which was released in 1987 as IEEE STD-1076 [6]. IEEE standards are reviewed and renewed every five years, and as part of the 1992 renewal of VHDL the language was extended to include a number of new features, such as a foreign-language interface, impure functions, and shared variables [7].

### 5.3.2 VHDL Language Features

Rather than develop an entirely new language, the designers of VHDL chose to base the syntax and semantics of their language upon an existing well-defined standard, Ada [11]. Many of the high-level programming features of Ada are therefore found in VHDL. Like Ada, VHDL is a strongly typed language with user-definable and -extensible data types. Structured objects such as vectors, arrays, and records are fully supported. Operators, functions and procedures may be overloaded on the data types of arguments and return results. VHDL supports data abstraction through the use of packages, which present a clean interface to objects and operations on objects while insulating the programmer from the details of the object implementation. VHDL explicitly represents concurrency and synchronization through the `Process` and `Wait` constructs and supports the automatic inference of registers and latches through signal assignment within sequential processes. VHDL also supports a wide range of abstraction levels by allowing the mixture of behavioral and structural representations, with `Generate` constructs and `Generic` parameters to control the instantiation of structural components.

VHDL also includes a number of features specifically designed to support simulation. File input and output are supported directly by the language, and the `TEXTIO` package is provided to support formatted ASCII I/O. Dynamic storage allocation is supported through the use of *access* types (that is, pointers), the object allocator `new`, and the implicit `Deallocate` procedure. The `assert` statement may be used to check that a specified condition is true. If the condition is not true, an error at one of several different severity levels may be reported. Although these language features have no direct analog in physical hardware (that is, they are not *synthesizable*), together they greatly facilitate the implementation of a system simulator, as is shown in Chapter 6.

The compilation process for VHDL is separated into an *analysis* phase and an *elaboration* phase, which are roughly analogous to compilation and object module loading in a conventional programming language compiler. VHDL provides the programmer with a great deal of control over the compilation process by deferring the binding of generic parameters and architecture instances until elaboration time. The elaboration time binding is controlled by the `Configuration` statement, which allows the user to specify the architecture to use for each component instance in the design and to override any generic parameter values passed to the architecture. This

in turn allows the user to select component architectures from a library and to control the instantiation of those components without requiring a detailed understanding of the library implementation.

### 5.3.3 Problems with VHDL

VHDL is not a panacea. VHDL is a large language with many features, which often takes a long time to learn. The syntax, although very similar to Ada, is unfamiliar to many programmers, who may find it verbose and cumbersome. The stateless nature of VHDL functions and procedures forces the use of structural representations for complex state machines. Finally, although VHDL has explicit constructs for concurrency and synchronization, many programmers find that coordinating many parallel fine-grain tasks can be difficult and time-consuming.

When we began development of Splash 2 in 1991, some features of the VHDL language were not supported by commercial synthesis tools; in particular, the use of `Generic` parameters, multidimensional arrays, and constant folding for multiply and divide operations were unsupported. The level of compliance of the tools has improved significantly over the last several years, and today there are very few VHDL constructs that synthesis tools cannot handle.

The other leading candidate for the role of Splash 2 programming language was the Verilog [15] hardware description language. Like VHDL, Verilog supports both simulation and synthesis from the same source code, so there was no fundamental impediment to using Verilog for Splash 2. The syntax of Verilog is closer to the C language and thus would be more familiar to many programmers. We felt, however, that Verilog would not be as rich a *programming* language as VHDL, because it did not have many of the language features we were looking for. The built-in data types of Verilog are very closely tied to hardware constructs such as wire-AND logic and high impedance (tristate) drivers, and there is no support for building abstract data types above these. Verilog also does not support the overloading of operators or procedures based upon data type. For these reasons we felt that we would not be able to provide the same level of abstraction with Verilog that we could with VHDL.

## 5.4 SOFTWARE ENVIRONMENT

The VHDL programming environment for Splash 2 consists of a system simulator, a logic synthesis package, a VHDL library that is common to both tools, and a SunOS-based runtime system. The Splash 2 simulator is a hierarchical model of the Splash 2 system comprising a set of VHDL models for each of the components of the system. The simulator provides a framework for the development and debugging of applications. Within the simulator, an application program is able to interact with the system exactly as it would with the physical hardware. The system models also verify that the application program meets various hardware constraints, such as memory sequencing and setup and hold times. The user may also specify crossbar configurations and initial memory contents with separate ASCII files, which are read by both the simulator and the runtime system.

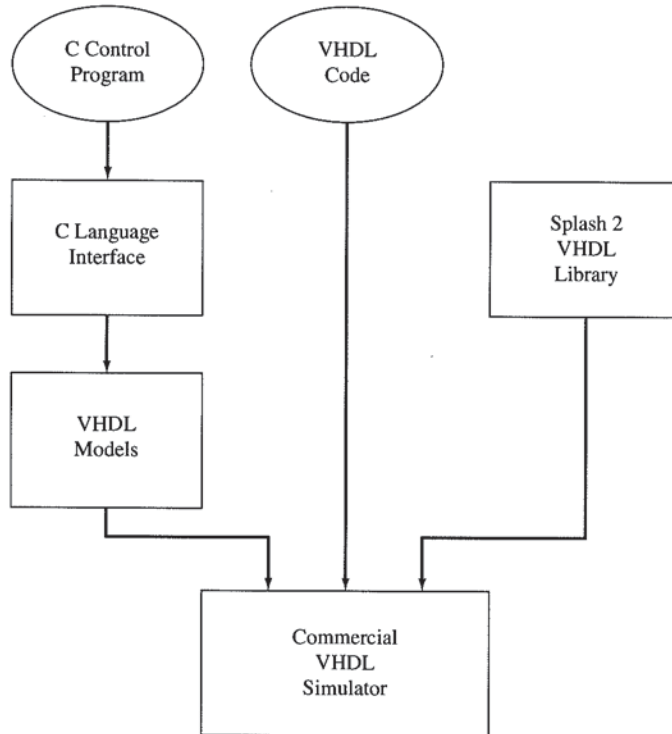
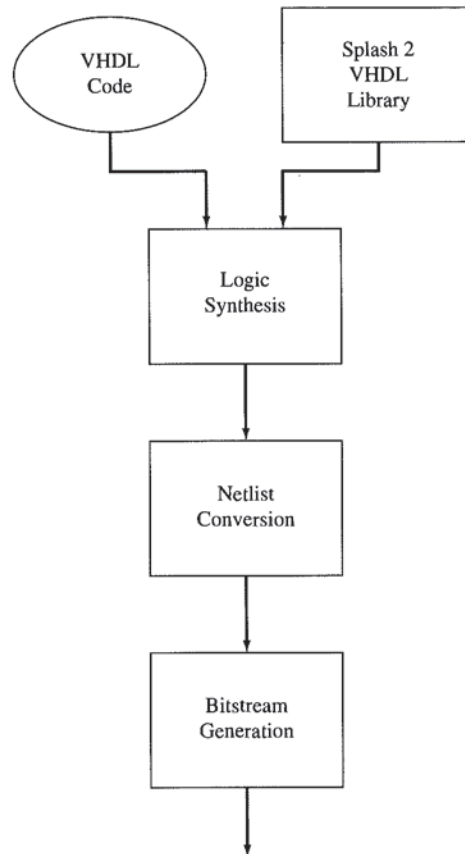


FIGURE 5.2 The Splash 2 Simulation Environment

The components of the simulation environment are shown in Figure 5.2. The ovals represent the two components of the user's application code: the VHDL program(s) for the computing elements and the C control program, which will run on the host computer. The C Language Interface is an optional piece of the environment that allows the simulation to be controlled by the same program that will run on the host. The VHDL Models block is the set of simulation models for the system, including the central crossbar, the external memories, and the Interface Board. The Splash 2 VHDL Library contains a set of data types, constants, procedures, and components designed to facilitate the interface between the application VHDL code and the rest of the system and to provide access to the Xilinx hard macros. Hard macros are predefined components, such as adders and counters, which provide guaranteed performance. Hard macros also provide the only access to special hardware features such as the fast-carry logic. Finally, the Commercial VHDL Simulator provides the simulation engine and the graphical user interface.

The VHDL simulation environment allows Splash 2 applications to be developed in either a top-down or bottom-up fashion. Top-down design is supported by beginning with a single high-level VHDL model for the entire Splash 2 system and iteratively descending through levels of hierarchy corresponding to the structure of the simulator down to the computing element, adding detail at each level. Bottom-up design is supported through the use of a library of default components for all of the pieces of the system except for the element being developed. As each element is completed, the corresponding library component is replaced with the actual design.



**FIGURE 5.3** The Splash 2 Compilation Environment

A mix of logic synthesis and standard compilation techniques are used to compile the VHDL programs into FPGA configurations, as shown in Figure 5.3. The VHDL Code that was developed in the simulation environment (Figure 5.2) is compiled with the same VHDL Library used to produce the Splash 2 object module. The logic synthesis tools from Synopsys Inc. [13, 14] map the VHDL code into a gate list. During the course of the Splash 2 project we used two different generations of Synopsys logic synthesis tools: the version 2.2 Design Compiler [12] and the version 3.0 FPGA Compiler [14].

At the beginning of the project we chose what was then the state-of-the-art Synopsys Design Compiler as the basis of our compiler. This tool was not tailored specifically to the FPGA technology and therefore required some customization to suit our needs. We developed a technology library that allowed the Design Compiler to produce a generic gate list from a reasonable subset of VHDL, and a net list conversion program called `edif2xnf`. `Edif2xnf` parsed the hierarchical EDIF net list, flattened the structure, and produced another file in Xilinx Net list Format (XNF) that was suitable for mapping onto the physical hardware by the Xilinx-provided bitstream generation tools [17]. Along the way it also performed some minor optimizations specific to both Splash 2 and the FPGA architecture.



Our experiences with the Synopsys Design Compiler and with our own `edif2xnf` program were fed back to Synopsys to help direct the development of their FPGA Compiler product. By the middle of 1993 Synopsys released their version 3.0 FPGA Compiler [14], which was able to compile logic directly into Xilinx RAM-based lookup tables and produce XNF net lists. The FPGA Compiler removed the need for our custom technology library and `edif2xnf`, but we found that some minor modification of the net list was still necessary. The program `xnfer` was written to fix XNF net list errors and to automatically insert logic common to all Splash 2 designs, such as the control for the internal Global Tri State (GTS) signal.

The major components of the runtime environment are shown in Figure 5.4. There are two host software interfaces to the Splash 2 system, a C library, which can be linked into an application-specific control program, and an interactive symbolic debugger called T2. Both interfaces are built upon the same underlying runtime system, `libsplash.a`, and provide the same functionality. The runtime system implemented by `libsplash.a` allows the user to open the device, map the Splash 2 memory into the host address space, establish input and output data streams, and control the system clock. The clock may be singly stepped, multiply stepped, or allowed to run free. The user can establish software handlers for interrupts generated by individual Processing Elements. The runtime library and the hardware diagnostic suite rely on the services provided by the Unix device driver, including memory

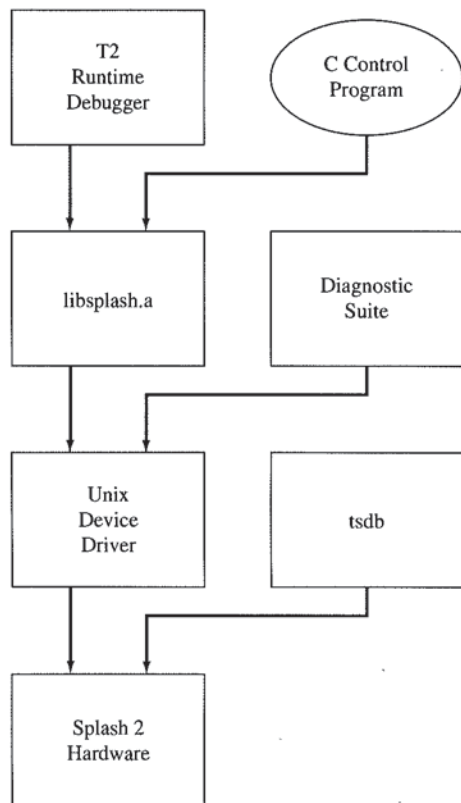


FIGURE 5.4 The Splash 2 Runtime Environment

management and several system calls (see Chapter 6). At the lowest level is a hardware debugger called "tsdb" (for Trivial SBus Debugger) that allows the designer to examine and set locations in the hardware based upon the physical address space of a single SBus slot.

The user interface of T2 is built upon the tool command language Tcl [9]. Tcl is an interpreted language with a C-like syntax which may be embedded into applications to provide an extensible user interface. The Tcl interface to T2 allows users to write simple programs to aid with debugging and experimentation on Splash 2. T2 also supports symbolic debugging by reading back the internal state of all FPGAs at the end of every clock cycle and associating the state of each flip-flop with the corresponding VHDL signal name. From T2 the user can step through the execution of the program, continuously displaying the contents of some or all of the registers in the design.

## 5.5 PROGRAMMER'S VIEW OF SPLASH 2

Every Splash 2 application may be divided into three main components: the portions that run on the Array Boards, the Interface Board, and the host computer. At the Splash 2 Array Board level, the programmable components consist of the Processing Elements, X1 through X16, the Control Element, X0, and the crossbar. At the Interface Board level, the Control Elements XL and XR are user-programmable, as are many of the control registers. The host interface must provide the input data streams, handle the output data streams, and control the operation of the Splash 2 system. In this section we discuss the process by which applications programs are developed, and then look at the programming model presented at each of the three levels. More details of the implementation may be found in Chapter 6 and in the Splash 2 Programmer's Manual [1].

### 5.5.1 Programming Process

Developing an application for Splash 2 is not unlike designing a program for a massively parallel computer. The programmer must choose an overall control paradigm, typically either data parallel (SIMD) or pipelined, and then plan the data flow among the Processing Elements, including the use of the crossbar and memories. On massively parallel computers the data layout among the processors is often critical to performance. In pipelined Splash 2 applications, it is the *control* layout that is critical. An algorithm must be partitioned carefully among the Processing Elements to maximize the efficiency of the inter-PE communication. Unfortunately, we know of no good automated tools that will find and exploit the structure of the design, so this partitioning must be performed manually by the programmer.

Once the basic control paradigm is chosen and the algorithm is partitioned, the communication and control protocols among the Processing Elements may be designed. Since Splash 2 is a globally synchronous system, these protocols are implemented as a set of finite state machines (VHDL *processes*) communicating through a set of *signals*. Input data from the Interface Board may be tagged as valid, or the clock may be controlled such that all input data seen by the Array Boards are valid.

A principal design goal in most applications is to maximize the utilization of one or more of the resources of Splash 2. For example, to maximize the memory bandwidth, many applications pipeline the accesses such that a new memory read or write operation occurs in every clock cycle. This is accomplished by registering the address and data within the IOBs of the Processing Elements, as is discussed in Chapter 6.

One major shortcoming of the programming methodology for FPGAs is the inability to determine the “size” (percent utilization) of a design without running through the entire compilation procedure. Splash 2 programmers have developed a crude “rule of thumb” to estimate the size of a design. The number of bits of register storage (including state machines) is summed, and if the number is within about 25 percent of the total number of flip-flops in the XC4010 (800 CLB flip-flops plus input and output flip-flops on the principal data ports), the design may be too large to fit. The 25 percent margin allows for inefficiencies in the placement and routing tools and the crude estimation heuristic. Many Splash 2 applications were in fact able to achieve or exceed a 98 percent CLB utilization rate.

### 5.5.2 Processing Element View

Programs for the individual Processing Elements of the Splash 2 Array Board are written in VHDL and must conform to the predefined Processing Element Entity declaration. The Processing Element Entity is essentially “boilerplate” code, common to all Splash 2 applications, and specifies the names and data types of the interface ports. The body of a PE program is a VHDL Architecture corresponding to the standard Entity. The interface ports include the data paths to the left- and right-hand neighbor PEs, the data path and control signals to the crossbar, the address and data path to the external memory, and a variety of control signals such as the global OR signals, the broadcast, interrupt, and handshake signals.

It is often important for timing considerations and CLB utilization to exploit the flip-flops in the input/output cells (IOBs) of the Processing Elements. To avoid long propagation delays between the logic core of one Processing Element to another, it is standard practice (although not required) to register data both entering and leaving the PE. Since the propagation delay on the major buses is significant, it is strongly recommended that input data from the SIMD Bus and output data to the RBus be registered. The timing of the external memory control requires that the address and control signals be registered in the Processing Element. This final constraint is enforced by the gate list postprocessor, `edif2xnf` or `xnfer`.

The set of configurations for the central crossbar is specified by an ASCII file that is interpreted by both the simulator and the runtime system. The configuration in use at any given time is selected by the Control Element (X0), but the output enable signals of the crossbar must still be set correctly by the individual Processing Elements. Another user-provided ASCII file may be used to specify initial contents of any of the external memories.

The Control Element (X0) is typically used to implement Array Board-level controller functions, such as SIMD instruction decode, and to store and broadcast common data tables. The Control Element has a different I/O interface than the Processing Element, and hence a unique Entity declaration. X0 has an input data port from the SIMD Bus and a bidirectional data port to the crossbar that is shared with

X16. Another three-bit output port is used to select the current crossbar configuration. The two global OR signals from each of the Processing Elements may be used to perform reduction and synchronization operations among all 16 PEs and the results combined with other Array Boards in the system via the systemwide global OR signals. Since the internal global OR signals are bidirectional, they may also be used to signal sequencing information from X0 to the individual PEs.

### 5.5.3 Interface Board View

The Development (or "Quick and Dirty") Board was built to assist in the debugging of the Splash 2 Array Boards while the final Interface Board was still being designed, and to provide an early application development environment. The Development Board maps every signal from the backplane side of the Interface Board to a host-accessible register, allowing the host to emulate in software the behavior of the Final Interface Board. The system clock is generated by host accesses to one of two special registers: the "Software Clock" register, which produces a fixed-width clock pulse, and the "SIMD Clock" register, which places the write data in the SIMD register and then generates a clock pulse. The functionality of the Development Board was retained in the Final Interface Board by incorporating a "bypass" mode that allowed applications and diagnostics written for the Development Board to run on the Final Interface Board by simply recompiling the code.

The Final Interface Board (IB) is responsible for controlling the data streams to and from Splash 2, the system clock, the RBus master and direction, and the FPGA configuration and state readback. The two FPGAs on the IB, XL and XR, are user-programmable, but the Splash 2 VHDL Library includes several standard designs that perform the most common control operations such as tagging input data with a "valid" indicator and only writing output data so tagged. More complicated designs can be implemented by modifying one of these programs. The input data source to the SIMD bus is selected from Channel B or C by XL, which may also perform preprocessing on the data, such as parallel-to-serial conversion. The output data is typically received by XR from the RBus and sent back to the host on Channel A, although XR may also send or receive data from XL. Both XR and XL have the ability to stop and restart the system clock, depending upon the state of the data channels. For example, if an input DMA channel is empty or an output channel is full, the clock may be stopped until the condition is cleared by the host (for instance, when another DMA operation occurs). The state machines in XL and XR that control the system clock may be clocked by a separate free-running clock signal. XR also controls the ownership of the RBus by setting the linear array size (RSize) and the direction (RDir) backplane signals.

### 5.5.4 Host View

A complete Splash 2 application includes a C program running on the host, which plays a pivotal role in the initialization and control of the hardware. This role includes downloading the configuration data to the FPGAs, establishing the input and output data streams, and controlling the system clock. The host program can also interact with the FPGA programs through a variety of both synchronous and asynchronous means.

The output of the compilation process is a set of configuration bitstream files, typically one for each FPGA. The host is responsible for merging the set of bitstream files for each Array Board into a single configuration stream called a "raw" file, which can then be downloaded directly to the Array Board. The host must also initialize the crossbar by reading and downloading the crossbar configuration file.

Symbolic debugging of running programs is supported by the state readback mechanism. To examine the internal state of a program, the host may stop the system clock and initiate a readback operation, which dumps the internal state of all of the FPGAs into a special buffer. The debugger may then extract the state of individual registers from the buffer, and associate the value with the VHDL symbol name.

The memory associated with the Processing Elements is mapped into the address space of the host program such that each PE memory appears as an array of integers. This allows the host program to read and write the memory using standard C data structures and pointer references. The kernel device driver is responsible for coordinating memory accesses with the FPGAs.

On systems with the Development Board, the SIMD and RBus data registers are mapped directly into the address space of the host program. To create an input stream to Splash 2 the host program simply writes data to the SIMD register. Likewise, an output stream is handled by reading from the RBus register. A set of library routines are available to facilitate these operations. The various asynchronous communications mechanisms such as the handshake registers may be accessed through C macros.

The Final Interface Board supports the use of standard Unix `read` and `write` system calls. An input data stream is created by writing the contents of an internal buffer or file to the device, while an output stream may be read from the device into a buffer or file. Higher-level library routines allow the concurrent handling of input and output streams. Another set of library routines permit the host to set the clock frequency, and start and stop the system clock.

## REFERENCES

- [1] J.M. Arnold and M.A. McGarry, "Splash 2 Programmer's Manual," Tech. Report SRC-TR-93-107, SRC, Bowie, Md., 1993.
- [2] R.A. Bergamaschi, "High-Level Synthesis in a Production Environment: Methodology and Algorithms," in J.P. Mermet, ed., *Fundamentals and Standards in Hardware Description Languages*, Kluwer Academic Publishers, Boston, 1993, pp. 195-230.
- [3] S.D. Brown et al., *Field-Programmable Gate Arrays*, Kluwer Academic Publishers, Boston, 1992.
- [4] M. Gokhale et al., "The Logic Description Generator," Tech. Report SRC-TR-90-011, SRC, Bowie, Md., 1990.
- [5] M. Gokhale and R. Minnich, "FPGA Programming in a Data Parallel C," *Proc. IEEE Workshop FPGAs for Custom Computing Machines*, CS Press, Los Alamitos, Calif., 1993, pp. 94-102.
- [6] *IEEE Standard VHDL Language Reference Manual*, Std 1076-1987, IEEE Press, New York, 1988.
- [7] *IEEE Standard VHDL Language Reference Manual*, Std 1076-1992, IEEE Press, New York, 1992.

- [8] P.J. Menchini, "An Introduction to VHDL," in J.P. Mermet, ed., *Fundamentals and Standards in Hardware Description Languages*, Kluwer Academic Publishers, Boston, 1993, pp. 359–384.
- [9] J.K. Ousterhout, *Tcl and the Tk Toolkit*, Addison-Wesley, Reading, Mass., 1994.
- [10] D.L. Perry, *VHDL*, McGraw-Hill, New York, 2nd ed., 1994.
- [11] *Reference Manual for the Ada Programming Language*, ANSI/MIL-STD-1815A-1983, U.S. Department of Defense, Washington, D.C., Feb. 1983.
- [12] Synopsys, Inc., *Design Compiler Reference Manual*, Synopsys, Inc., Mountain View, Calif., 1991.
- [13] Synopsys, Inc., *VHDL Compiler Reference Manual*, Synopsys, Inc., Mountain View, Calif., 1991.
- [14] Synopsys, Inc., *FPGA Compiler Reference Manual*, Synopsys, Inc., Mountain View, Calif., 1994.
- [15] D.E. Thomas and P.R. Moorby, *The Verilog Hardware Description Language*, Kluwer Academic Publishers, Boston, 1991.
- [16] S.M. Trimberger, ed., *Field Programmable Gate Array Technology*, Kluwer Academic Publishers, Boston, 1994.
- [17] Xilinx, Inc., *The XC4000 Data Book*, Xilinx, Inc., San Jose, Calif. 1994.

# CHAPTER 6

---

## Software Implementation

*Jeffrey M. Arnold*

### 6.1 INTRODUCTION

An important goal of the Splash 2 software effort was to provide a working programming environment as quickly as possible without sacrificing the ability to grow and evolve as the project progressed. We therefore chose to base the implementation on software standards and readily available tools as much as possible, allowing us to concentrate on the system integration and the development of applications. The standards we chose included the VHDL and C programming languages and the Unix operating system. This chapter shows how these standards and the tools that support them were assembled to produce a complete programming environment.

### 6.2 VHDL ENVIRONMENT

The Splash 2 VHDL programming environment consists of a library of useful VHDL constructs and a set of standard entity declarations for the various levels of the Splash 2 hierarchy. This section presents some details of that environment, and then discusses aspects of the VHDL programming style that was evolved by the Splash 2 programmers.

### 6.2.1 Splash 2 VHDL Library

The `Splash2 Library` contains a set of `Packages` that are used for the development of VHDL application code for Splash 2. The `TYPES` package contains definitions of the data types used for interchip communication and is essentially a superset of the IEEE 1164 Standard Logic data type package. All bidirectional interface ports are built upon a four-state subtype ('X', '0', '1', and 'Z') of the `StandardLogic` type, called `RBit3`.<sup>1</sup> Assignment of a value of 'Z' to a signal implies the synthesis of a tri-state driver. The 'X' state is used only in simulation, primarily to identify tri-state bus conflicts. All of the standard logical operators as well as signed and unsigned arithmetic operators are supported over vectors of `RBit3`.

The `SPLASH2` package contains a variety of constants, data types, and functions that are specific to either the Splash 2 architecture or the Splash 2 simulator. For example, constants are defined that specify the width and depth of the memories and the width of the linear data path. Subtypes are also defined to specify the Processing Element data ports.

The `COMPONENTS` package contains a set of components and procedures useful in writing applications. These include the "Pad" procedures, which can be used to interface between the tri-state (`RBit3` type) signals external to the Xilinx chips and the standard logic levels (`Bit` type) internal signals. There are four pad procedures, each overloaded to accept scalar and vector arguments of `Bit` and `RBit3` types. `Pad_Input` is used to receive inputs from off-chip; `Pad_Output` is used to drive off-chip; `Pad_InOut` is used to conditionally receive and drive off-chip signals; and `Pad_XBar` is used to receive and drive the crossbar data path, conditioned by the crossbar output enable signals.

The `HMACROS` package contains component declarations and simulation models for the set of *hard macros* [4] provided by Xilinx. Hard macros are logic modules that have been hand-optimized with fixed relative placement and routing for maximum efficiency. Until the release of the Xilinx XACT 5.0 tools in 1994, hard macros were the only mechanism for accessing special-purpose hardware such as the fast-carry-chain logic. The `HMACROS` package provides the means to structurally instantiate hard macros within an application.

### 6.2.2 Standard Entity Declarations

All Splash 2 applications programs must conform to the input and output behavior defined by the standard `ENTITY` declarations. There are four unique FPGA entities visible to the programmer: the Processing Element (X1 through X16) and the Control Element (X0) on the Array Board, and XL and XR on the Interface Board. The entity declaration for the Processing Element is shown in Figure 6.1. The generic parameters `BD_ID` and `PE_ID` are constant values, unique to each PE, provided by the software environment; they permit the application to customize each PE program to its physical position in the system. The port declarations represent the connections of the PE to its neighbors and the rest of the system. The ports of type `DataPath` represent the data paths to the left and right neighbors and the crossbar. For example,

<sup>1</sup>The original `RBit3` data type was a resolved three-state (0, 1, Z) logic developed before the release of IEEE 1164.



```

ENTITY Processing_Element IS
  GENERIC (BD_ID      : Integer;
           PE_ID      : Integer);
  PORT (XP_Left      : inout DataPath;
        XP_Right     : inout DataPath;
        XP_XBar      : inout DataPath;
        XP_XBar_EN_L : out   Bit_Vector(4 downto 0);
        XP_Mem_A     : inout MemAddr;
        XP_Mem_D     : inout MemData;
        XP_Mem_RD_L  : inout Bit;
        XP_Mem_WR_L  : inout Bit;
        XP_Int       : out   Bit;
        XP_Broadcast  : in   Bit;
        XP_Reset     : in   Bit;
        XP_HS0, XP_HS1 : inout RBit3;
        XP_GOR_Result : inout RBit3;
        XP_GOR_Valid  : inout RBit3;
        XP_LED       : out   Bit;
        XP_Clk       : in   Bit);
END Processing_Element;

```

FIGURE 6.1 Standard Processing Element Entity Declaration

XP\_Right of one PE is connected to XP\_Left of the next PE. The interface to the PE memory consists of an address bus (XP\_Mem\_A), a data bus (XP\_Mem\_D), and separate read and write control signals (XP\_Mem\_RD\_L and XP\_Mem\_WR\_L). The “\_L” appended to the name indicates the signal is active low. The user must set these signals to a ‘1’ when the memory is not in use. XP\_Int is the interrupt output signal. The interrupt signals from each of the FPGAs X0 through X16 are logically ANDed with the contents of the Array Board mask register, and then ORed to form the board-level interrupt. XP\_Broadcast is an input signal driven by X0 and is common to all 16 PEs. XP\_Reset is the systemwide reset signal, which may be set by the host. By default, XP\_Reset is automatically connected to the Global Set/Reset (GSR) signal of the Xilinx XC4010, but it is also available as a user input. The Array Board handshake registers appear to the PE as XP\_HS0 and XP\_HS1. Each PE is connected to a unique bit of the HS0 register, while HS1 is common to all PEs on the Array Board. XP\_GOR\_Result and XP\_GOR\_Valid are bidirectional signals between the Control Element (X0) and each of the Processing Elements. The signal names reflect their intended purposes (global AND/OR reduction and barrier synchronization), but the bidirectionality makes these ports useful for signaling state changes from X0 to individual PEs. The port XP\_LED is connected directly to a light-emitting diode (LED) on the front edge of the Array Board and is typically used for diagnostics. Finally, XP\_Clk is the global synchronous clock shared by every PE in the system.

The entity corresponding to the Control Element (X0) is similar to the Processing Element entity, although the port names are prefixed by X0\_ rather than XP\_ (see Figure 6.2). In place of the XP\_Left and XP\_Right data buses, the Control Element has 36-bit ports to the SIMD bus (X0\_SIMD) and to the crossbar (X0\_XB\_Data). X0\_GOR\_Result\_In and X0\_GOR\_Valid\_In are each 16-bit vectors of bidirectional

```

ENTITY Control_Element IS
  GENERIC (BD_ID      : Integer;
           PE_ID      : Integer);
  PORT (X0_SIMD      : inout DataPath;
        X0_XB_Data   : inout DataPath;
        X0_Mem_A     : inout MemAddr;
        X0_Mem_D     : inout MemData;
        X0_Mem_RD_L  : inout Bit;
        X0_Mem_WR_L  : inout Bit;
        X0_GOR_Result_In: inout RBit3_Vector(1 to 16);
        X0_GOR_Valid_In : inout RBit3_Vector(1 to 16);
        X0_GOR_Result  : out Bit;
        X0_GOR_Valid   : out Bit;
        X0_Clk         : in Bit;
        X0_XBar_Set    : out Bit_Vector(2 downto 0);
        X0_X16_Disable : out Bit;
        X0_XBar_Send   : out Bit;
        X0_Broadcast_In : in Bit;
        X0_Broadcast_Out: out Bit;
        X0_LED         : out Bit);
END Control_Element;

```

**FIGURE 6.2** Standard Control Element Entity Declaration

signals to each of the Processing Elements. `X0_GOR_Result` and `X0_GOR_Valid` are outputs connected to the wire-OR backplane signals. Access to the crossbar data is achieved by asserting `X0_X16_Disable`, which effectively isolates X16 from controlling the crossbar output enables, and then setting `X0_XBar_Send` high to transmit into the crossbar or low to receive.

The entity for the Interface Board part XL is shown in Figure 6.3. The principal data ports correspond to the data path shared by DMA channels B and C, the SIMD bus, and the data path to XR. There are separate input ports for the system clock and the free-running clock. A clock-enable output port is used to start and stop the system clock. Two separate channel control ports of 14 bits each convey the control and status

```

ENTITY XL IS
  PORT (XL_FIFO      : inout DataPath;
        XL_SIMD      : inout DataPath;
        XL_XR        : inout DataPath;
        XL_Free_Clk  : in Bit;
        XL_Splash_CLK : in Bit;
        XL_Enable_Clk : out Bit;
        XL_Chan_B    : inout ChanCtrl;
        XL_Chan_C    : inout ChanCtrl;
        XL_Ctrl      : inout Bit_Vector(4 downto 0);
        XL_GOR_Result : in Bit;
        XL_GOR_Valid  : in Bit;
        XL_BCast     : inout Bit);
END XL;

```

**FIGURE 6.3** Standard XL Entity Declaration

signals to and from the DMA channels. The global OR result and valid signals from the backplane are inputs to XL, as is a separate 5-bit handshake register, `XL_CTRL`.

The XR entity is very similar to XL. Its principal data ports correspond to data paths to DMA channel A, the RBus, and the XL-XR bus. It too has two clock inputs, a clock-enable output, and shares the same 5-bit handshake register. The RBus size and direction signals originate from ports on XR.

### 6.2.3 Programming Style

Over the course of the Splash 2 project a number of idiomatic VHDL constructs have evolved. Some of these constructs arose from requirements of the synthesis tools; others became a matter of programming style. In this section we examine a few of these idioms.

Signed and unsigned arithmetic is supported for `Integer`-derived types and for vectors of `Bit` and `RBit3` types. The default word size for both signed and unsigned `Integers` is 32 bits, which can lead to a tremendous waste of logic and routing resources when the data range is known to be small. Therefore, range constraints on integers are used to assist the synthesis tools in optimizing the width of operator logic. For example, the code in Figure 6.4 will synthesize to a 10-bit unsigned incrementer. Likewise, vector lengths may be used to control operator widths for arithmetic over bit vectors. Since the arithmetic operators are overloaded to work with either, the choice of whether to represent a value as an `Integer` or a `Bit_Vector` is one of programming style. It is often more convenient to specify ranges than vector lengths, but vectors allow easier expression of shifts, concatenation, and bitfield extraction.

```
SIGNAL i: Integer range (0 to 1023);
i <= i + 1;
```

FIGURE 6.4 Range Constrained Integer Assignment

All of the Processing Elements in Splash 2 receive a global synchronous clock, `XP_Clk`. Therefore, all `Processing_Element` architectures have one or more processes synchronized to this signal. As shown in Figure 6.5, a synchronous process has no *sensitivity list* to limit its execution, but rather contains a single `WAIT` statement conditioned to trigger execution of the process on the rising edge of the clock. Assignments to `SIGNAL` objects within the body of the synchronous process are used to imply registers, since assignment occurs only following the execution of the process body, effectively registering the result on the clock edge. Unregistered temporary values may be named within a process by assigning to `VARIABLE` objects.

The “Pad” procedures (`Pad_Input`, `Pad_Output`, `Pad_InOut`, and `Pad_XBar`) declared in the `COMPONENTS` package are typically used to connect logic within the

```
PROCESS BEGIN
  WAIT UNTIL XP_Clk'EVENT and XP_Clk = '1';
  -- Body of synchronous process
END PROCESS;
```

FIGURE 6.5 Synchronous Process

```

ARCHITECTURE Test OF Processing_Element IS
    SIGNAL Left : Bit_Vector(DATAPATH_WIDTH-1 downto 0);
    SIGNAL XBar_in : Bit_Vector(DATAPATH_WIDTH-1 downto 0);
    SIGNAL XBar_out : Bit_Vector(DATAPATH_WIDTH-1 downto 0);
    SIGNAL XBar_dir : Bit_Vector(4 downto 0);
BEGIN
    Pad_Input(XP_Left, Left);
    PROCESS BEGIN
        WAIT UNTIL XP_Clk'EVENT and XP_Clk = '1';
        Pad_XBar(XP_XBar, XBar_in, XBar_out, XBar_dir);
        Pad_Output(XP_XBar_EN_L, XBar_dir);
    END PROCESS;
END Test;

```

FIGURE 6.6 Example of Off-Chip Communication

Processing Element to the external Array Board environment. Figure 6.6 shows an example of the use of two of these procedures. The `Pad_Input` procedure receives data from the left-hand neighbor PE and makes it available on the internal `Bit_Vector` signal `Left`. Since this is a *concurrent* statement, there is no implicit registering of the data. In contrast, the call to `Pad_XBar` is a *sequential* statement within a synchronous process. Therefore, both the input to the PE (`XBar_in`) and the output to the crossbar (`XBar_out`) are registered. The 5-bit vector `XBar_dir` is used to control the direction of the five “bytes” of the crossbar and is also driven out to the crossbar to control the corresponding output enable pins.

Finite-state machines are implemented by embedding flow control constructs such as `IF` and `CASE` statements within synchronous processes. An enumerated type may be used to define the set of valid states and a `SIGNAL` object of this type declared to hold the current state. A `CASE` statement within a synchronous process is used to dispatch on the state variable. Within each `WHEN` clause, input conditions are tested, output signals are assigned, and the next state transition is computed.

Occasionally it is necessary to instantiate one or more components within a Splash 2 application program to create replicated structures or to gain access to specific FPGA features. VHDL provides several structural constructs, including component instantiation and conditional and iterative `Generate` statements. For example, access to the CLB RAM within the Xilinx PE may be accomplished by instantiating a special memory component. Generic parameters to this component are used to configure the width and depth of the memory as well as to specify any initial contents (such as for ROMs). When evaluated by the Splash 2 simulator, the model for this component uses these parameters to create an output file that can be read by the Xilinx-provided `MEMGEN` program. `MEMGEN` in turn creates a macro for the memory, which is incorporated into the FPGA load module by the place-and-route tools.

There are two standard modes of synchronizing the input data with the Splash 2 system. In the first mode, the XL chip on the Interface Board controls the system clock such that the Array Boards see a system clock pulse only when there is valid data on the SIMD bus. In the second mode, the system clock is allowed to run continuously while the presence of valid data is indicated by setting to 1 the most significant bit (bit 35) of the SIMD bus. More complex behavior can be achieved by modifying the XL program.

Access to the external memory is synchronous with the global clock. Memory read operations are performed by placing the address on the `XP_Mem_A` port in one clock cycle and reading the data from `XP_Mem_D` in the next cycle. The port `XP_Mem_Rd` is asserted with the address to indicate a read operation. The address may be changed every cycle to perform back-to-back reads. A write operation is performed by placing both the address and the data on the memory ports and simultaneously asserting `XP_Mem_WR`. The address and data may be changed every cycle to perform back-to-back writes.

### 6.3 SPLASH 2 SIMULATOR

The Splash 2 simulator itself is written in VHDL and consists of a hierarchical set of models of the various components of the system. An application program uses the Configuration statement to specify which architectures (models) to use at each of the levels of the hierarchy. The architectures specified at the leaves of the hierarchy may be any mix of user-provided VHDL code and predefined default models. The configuration statement also allows the user to customize individual models by setting the values of generic parameters. Among the parameters specified are the names of any files of test data. The configuration statement therefore specifies the construction of a complete model of the system. This model in turn is interpreted by a simulation engine, effectively executing the user's application.

This section begins by describing the structure of the simulator hierarchy. We then present the use of the configuration statement through a series of examples. Finally, we discuss some details of the system models.

#### 6.3.1 Structure

Figure 6.7 illustrates the structure of the Splash 2 simulator. The root of the hierarchy is the `System` model, which instantiates the `Interface` and the `S2Boards` models.

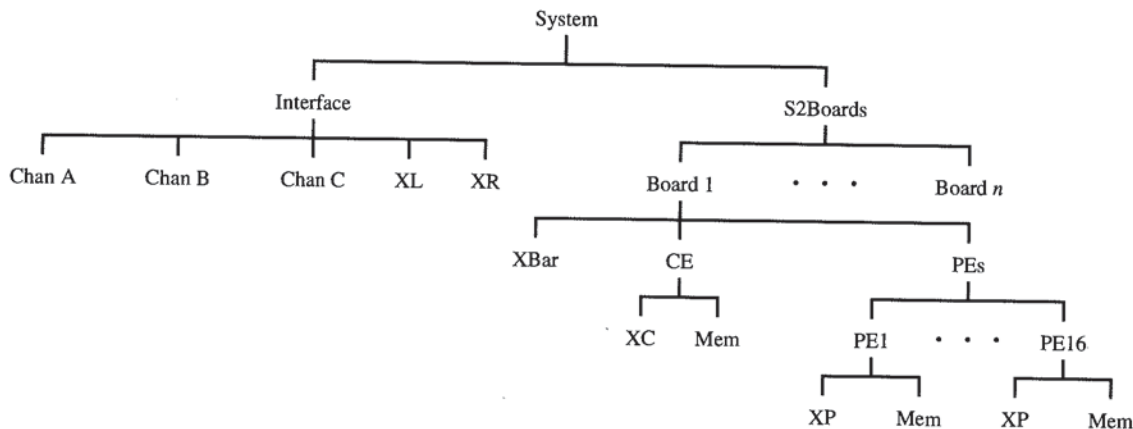


FIGURE 6.7 Structure of the Splash 2 Simulator

The `Interface` model is responsible for instantiating the DMA channels, XL and XR. Instances of the DMA channels are created only if there are values passed in the configuration statement for the corresponding input and/or output filenames. The `Interface` model is also responsible for generating the system clock, at a frequency determined by another generic parameter.

The number of Array Boards generated is controlled by the generic parameter `Number_Of_Boards`. The `S2Boards` model also passes the slot ID number to each Array Board component generated. The Array Board model in turn instantiates the crossbar (XBar) and the Control Element (CE) and contains a `Generate` statement that instantiates the Processing Elements (PEs). Each PE consists of a reference to the Processing Element component (XP) and the Memory component (Mem). The Processing Element reference passes the slot number and the PE number to the user's code through the generic parameters `BD_ID` and `PE_ID`, respectively.

### 6.3.2 Configuring the Simulator

The Splash 2 simulator assembles an application program according to the directions presented in the `VHDL Configuration` statement. This statement, typically stored in a file called `config.vhd`, identifies the architecture to use for each entity in the design and allows the user to specify values of generic parameters. Figure 6.8 shows the top, or outermost, level of a typical `config.vhd` file. Most of the file is common to all applications and may be copied from the Splash 2 library. In this example, `Top` is just a label to identify the configuration of the `Splash_System` entity. Within the `Structure` architecture there are two component instantiations: one for the `Interface_Board` entity and one for the `Splash2_Boards` entity. Any configuration information needed for the components would normally appear within the `for` clauses.

Figure 6.9 shows an example configuration for the `Interface_Board` component. The `Generic Map` construct is used to pass generic parameter values to the simulator. In this example, input for the application is taken from the file `test.dat` and output is written to the file `output.dat` in hexadecimal format. The clock model runs at a simulated frequency of 20 MHz. The user-programmable FPGAs XL and XR are both loaded with their `Valid` architectures from the `interface` library.

The rest of the Splash 2 simulator is configured in a similar manner. The complete `config.vhd` file contains places for specifying the number of Array Boards,

```
configuration TOP of Splash_System is
  for Structure
    for IFACE: Interface_Board
      -- Interface board configuration
    end for;
    for Splash: Splash2_Boards
      -- Configuration of array boards
    end for;
  end for;
end TOP;
```

FIGURE 6.8 Top Level of Simulator Configuration File

```

for IFACE: Interface_Board
  use entity interface.Interface_Board(Structure)
  Generic Map (input_file1 => "test.dat",
               output_file1 => "output.dat",
               File_Type    => Hex,
               Clock_Freq   => 20);
  for Structure
    for all: XL
      use entity interface.XL(Valid);
    end for;
    for all: XR
      use entity interface.XR(Valid);
    end for;
  end for;
end for;

```

**FIGURE 6.9** Interface Board Model Configuration

the crossbar configuration file, the Control Element and Processing Element architectures, and any initial memory tables.

### 6.3.3 Input and Output

Input to and output from the Splash 2 system are handled by the Interface Board model, which contains generic filename parameters for each of the three I/O channels. If the `config.vhd` file specifies a filename for a channel, the corresponding model opens the file for reading or writing. Input files are assumed to be ASCII containing one hexadecimal string per line, each line representing a new data value to be read from the channel. Output files are written in the same format, one value per line of the output file.

### 6.3.4 Crossbar and Memory Models

The crossbar model in the simulator is passed a generic parameter that contains the name of a file that contains up to eight settings. The first line of each setting consists of the keyword `configuration` followed by an integer from 0 to 7 (inclusive). The following lines of the setting are of the form:

output-port-number input-port-specifier

where “output-port-number” ranges from 1 to 16 and corresponds to the Processing Element number. The “input-port-specifier” is either a single integer (0 to 16) or five integers (0 to 16). If the input-port-specifier is a single integer in the range (1 to 16), it specifies a single source port for all 36 bits of the output. A value of 0 is used to indicate that port “output-port-number” is used as an input to the crossbar. If “input-port-specifier” consists of five integers, each integer specifies the source for one byte of the output, from most to least significant. If an output port number is missing from the configuration, it is assumed to be set to 0 (input to the crossbar). Note that simply setting “input-port-specifier” to 0 is not sufficient to disable the crossbar port; the corresponding `XP_XBar_EN_L` signals must be set to 1.

```

-- Odd PEs drive, even PEs receive
configuration 0
1 0
2 1
3 0
4 3
5 0
6 5
7 0
8 7
10 9
12 11
14 13
16 15

configuration 1
1 0
2 9 7 5 3 1

```

**FIGURE 6.10** Sample Crossbar Configuration File

Figure 6.10 shows an example crossbar configuration file with two settings defined. In the first setting, the odd-numbered PEs are driving into the crossbar while the even-numbered PEs are receiving. The crossbar ports corresponding to PEs X9, X11, X13, and X15 are implicitly set to 0 (disabled). In the second setting, PE X2 is receiving one byte each from X9, X7, X5, X3, and X1. Crossbar ports for X3 through X16 are implicitly disabled.

Users may choose from several predefined memory models available in the Splash 2 simulator, or they may add their own. The predefined memory architectures include the following:

- **None:** No memory is modeled. Access to the memory generates an error.
- **Zero:** Read-accesses from any location return a constant zero. Write-accesses are ignored.
- **Static:** Memory is modeled as a statically allocated, fixed-size array. The size of the array is determined by a generic parameter passed by the configuration file. This model is useful for lookup tables and sequentially accessed data.
- **Dynamic:** Storage for the memory model is allocated dynamically as needed. The first write to an address allocates a new storage cell, which subsequent reads will fetch. A read from an unwritten address generates an error. This model is useful for programs in which the access pattern is data-dependent (random).

The memory initialization, or “load,” file is an optional ASCII file that may be used to specify initial contents of the PE memories. The format of this file is simple. For each set of contiguous blocks of data, the base address of the block is given, followed by one or more data values. The base address is specified by the keyword `address` followed by an unsigned decimal integer. Subsequent (decimal) integers are interpreted as the 16-bit values to load into consecutive locations. A single load file may contain any number of blocks. Alternatively, the entire memory can be initialized to zero by including the keyword `clear` in the load file.



### 6.3.5 Hardware Constraints

One important role of the simulator is to verify that application programs satisfy certain hardware constraints that are difficult or impossible to verify through static analysis. The hardware constraints that can be checked in simulation include memory sequencing, reading uninitialized memory locations, and tri-state bus conflicts.

Since the data bus to the memory is bidirectional, the PE must ensure that it does not begin a write operation before the completion of a read. This constraint is verified within the simulator by checking that the memory read signal is deasserted at least one cycle before the memory write signal is asserted.

The major data buses on the Array Board are bidirectional, relying on the use of tri-state I/O drivers of the FPGAs to prevent conflict. The simulator models these tri-state pins with a four-state logic. The resolution function for the four-state logic detects any attempt to drive a signal to two different logic values simultaneously. When a conflict is found, a warning message is printed and the signal is set to the 'X', or unknown, state.

## 6.4 COMPILATION

A mix of logic synthesis and standard compilation techniques are used to compile VHDL programs into FPGA configurations. The logic synthesis tools from Synopsys Inc. [3] are used to map the VHDL code into a gate list. A custom peephole optimizer is then applied to the gate list to perform a variety of Xilinx-specific and Splash 2-specific optimizations. The resulting gate list is then mapped into the CLB structures and placed and routed using the Xilinx [4] tool package. The Xilinx tools are also used to extract the detailed timing information from the placed and routed design. This information may be used directly to manually optimize the design, or it may be used to construct a new structural VHDL model for each chip, which may be resimulated by the Splash 2 simulator to provide detailed timing analysis.

### 6.4.1 Logic Synthesis

In 1991, FPGA technology was still quite new and confined mainly to board-level "glue logic" applications. Consequently, very few commercial CAD tool vendors were targeting FPGAs for logic synthesis. After evaluating the few tools on the market, we chose to base our compiler on the Synopsys Design Compiler. This choice required the development of a custom technology library that allowed the Design Compiler to produce a technology-independent gate list. It was also necessary to write a net list conversion program to translate that generic gate list into a technology-dependent form suitable for the physical mapping tools.

The net list translator, called `edif2xnf`, parsed and flattened the hierarchical EDIF net list produced by the Design Compiler, creating another file in Xilinx Netlist Format (XNF). During the translation a number of Xilinx- and Splash 2-specific optimizations were also performed on the design, including:

- Flip-flops on the periphery of the logic were migrated to the I/O frame (IOBs) of the chip wherever possible.

- The board-level program reset and Xilinx Inhibit signals were connected to the internal GSR and GTS signals.
- Hard macro references were identified and marked as such.
- Port names were given specific pin assignments and pad slew-rate options.
- The clock signal was identified and a global buffer inserted in the clock net.
- A pattern-matching algorithm was also used to find opportunities to simplify the logic by exploiting the clock-enable feature of the Configurable Logic Block (CLB) flip-flops.

The output of the optimizer, in Xilinx Netlist Format (XNF), was fed into the Xilinx-provided placement and routing tools and static timing analyzer.

In late 1993 Synopsys released their “FPGA Compiler” product that incorporated much of the functionality of `edif2xnf` directly into the synthesizer. There were some minor problems, however, in the early releases of the FPGA Compiler, which necessitated our writing another program, `xnfer`, which was able to manipulate the XNF file. `xnfer` inserts the “drop in” logic that connects the internal GSR and GTS signals and moves peripheral flip-flops to the IOBs.

A Unix shell script (`vhdl2xnf`) presents the user with a simplified interface to the numerous controls of the FPGA Compiler and `xnfer`. `Vhdl2xnf` processes a number of command line options and constructs an execution script for the FPGA Compiler. This script, in turn, specifies any elaboration parameters (including the `BD_ID` and `PE_ID`), includes port mapping tables, and handles error conditions. The output of `vhdl2xnf` is an XNF file ready for physical mapping.

#### 6.4.2 Physical Mapping

The physical mapping of the design from XNF to a loadable bitstream is handled by the Xilinx-provided tools. The placement and routing tool, `PPR`, reads the XNF net list and produces an “LCA” file, which contains all of the configuration information in an ASCII format. The program `makebits` converts the LCA file into a bitstream format, called a “BIT” file. `Makebits` also produces an “LL” file that contains a table-mapping CLB and IOB flip-flops to positions in the readback bitstream. Another Unix shell script (`xnf2bit`) provides a convenient interface to the physical mapping tools and the symbol table creation.

#### 6.4.3 Debugging Support

To support the symbolic capabilities of the runtime debugger, a table is created associating the names of symbols with the location of the corresponding register bit in the readback stream. The information needed to build this table is extracted from the compiled design in two steps. First, the Xilinx tool `lca2xnf` is used to create an XNF file annotated with the location of each CLB flip-flop and its associated signal name.<sup>2</sup> The location information is then looked up in the LL file produced by `makebits` to produce an offset into the readback bitstream. A table of symbol names and offsets is then built for use by the debugger.

<sup>2</sup>This step is necessary to resolve ambiguities created by `PPR` through the use of feed-through CLBs.

## 6.5 RUNTIME SYSTEM

There are two host software interfaces to the Splash 2 system: a C language library that can be linked into an application-specific driving program, and an interactive symbolic debugger. Both interfaces are built upon the same underlying runtime system, and both provide the same basic functionality. The runtime system allows the user to open the Unix device, to map the Splash 2 memory into the host address space, to configure the FPGA devices and crossbar, to establish DMA data streams, and to control the system clock. The clock may be single-stepped, multiply-stepped, or allowed to run free. The user may also read and write various control registers, including the “handshake” registers.

### 6.5.1 T2: A Symbolic Debugger

To assist in the development of applications on the hardware, an interactive symbolic debugger, T2, was developed. The user interface to T2 is an interpreter executing the Tcl command language [1]. Tcl is a C-like language that provides a variety of control-flow mechanisms and allows the user to extend the command set by writing custom procedures. A set of built-in procedures provides access to the Splash 2 hardware resources and runtime software.

The built-in commands of T2 may be divided into three categories: hardware setup; program execution; and analysis. The setup commands include routines for hardware and software initialization, and configuration of the FPGAs and crossbar. To configure the system, the user specifies a map, which associates Processing Elements (individually or in groups) with bitstream files. The bitstream files for all of the PEs on an Array Board are then merged into a single image called a “raw” file. A given raw file may then be loaded to one or more Array Boards via the `ConfigArray` command. A raw file may be saved and reloaded on subsequent runs, obviating the need to associate and merge the bitstreams again. The crossbar is initialized from a crossbar configuration file by the `ConfigXBar` command.

Execution of an application program is controlled by the `Step` family of commands. There are a variety of these commands that allow the user to specify input and/or output files, file formats, and the interpretation of the “tag”: the most significant four bits of the 36-bit data word. For input files the tag may be set to a constant value or it may be taken from the input file. On output the tag is typically used to indicate valid data, so a mask may be provided to control which data are to be written to the output file. All of the `Step` commands allow the user to specify the number of clock cycles to execute.

The heart of the symbolic debugger is the Xilinx FPGA state readback mechanism. To trace a set of program variables, or symbols, the user issues the `AddReadBack` command after each `Step` command. `AddReadBack` adds the current state information to an internal history buffer. Another set of commands allows the user to inquire about the state of a particular symbol at a particular time. Symbols may be looked up individually, or an alias may be defined to aggregate multiple symbols.

The `wave` program allows users to view graphically the time-varying behavior of program symbols. The T2 command `Trace` adds a symbol (or alias) to a trace list. At the end of each clock cycle a readback is performed and the value of every

traced signal is written to a file. The `wave` program reads the trace file and paints a waveform display similar to a logic simulator or a hardware logic analyzer.

Finally, T2 provides a set of lower-level routines for reading and writing individual hardware registers. These routines are available for applications that require a level of control not provided by the higher-level interface.

### 6.5.2 Runtime Library

The Splash 2 runtime library, `libsplash`, which forms the foundation of T2, is also available to be linked into a user-written C program. For every built-in T2 command there is a corresponding entry point into `libsplash` that provides the same functionality. The routine `OpenAndInit` performs the basic hardware and software initialization, including opening the device, allocating and initializing the `Splash` device structure, and mapping the various pages of the physical address space into the user's address space. The `Splash` structure is also initialized with an array of pointers to each of the PE memories in the system. Separate minor devices corresponding to the DMA channels are also opened. The hardware initialization includes loading a passive, or idle, program into all of the FPGAs, setting the clock frequency, and resetting and disabling the DMA channels.

Application programs may be loaded and executed through `libsplash` in the same manner as from T2. Library routines exist to manipulate bitstream files and to create, save, and load raw files. The entire family of `Step` commands is also available as library routines.

In addition to the T2 commands, however, `libsplash` also provides a set of input and output routines based upon the standard Unix system calls `write` and `read`. The `Write` routine is a user-level interface to the `write` system call, which uses the DMA facility to transfer data from the user's address space to the Splash 2 Interface Board and XL. The `Read` routine similarly sets up a DMA transfer from the Interface Board to a buffer in the host memory.

Both `Write` and `Read` are *blocking* operations. That is, once called, these routines do not return control to the user program until the requested operation is completed. To implement two concurrent data streams, one input to Splash 2 and one output from Splash 2, the `WriteRead` routine is provided. `WriteRead` uses the first DMA controller to transfer data from one memory buffer to the Interface Board while simultaneously using the second DMA controller to move data from the Interface Board back to a different host memory buffer. This concurrency is accomplished by spawning a separate Unix process to perform the `Read` while the parent process proceeds with the `Write` operation. Once the `Write` has completed, the parent waits for the `Read` process to complete before returning control to the user program.

Both the output and input memory buffers used by the `WriteRead` routine are in the address space of the user program, but the data received from Splash 2 is in the address space of the `Read` process. Since Unix does not support shared memory very well, it is necessary to copy the received data from the `Read` process back to the address space of the parent process. This copying is accomplished by passing the data from the child back to the parent in a memory-mapped temporary file. The parent process opens a temporary file prior to spawning the child process. The child then maps the file into its address space using the `mmap` system call, and passes it as

the target buffer for the `Read` operation. Upon completion, the parent also memory-maps the temporary file and copies the data to the user's buffer. The temporary file is truncated to zero length prior to closing, to avoid any writes to the disk.

The buffer copying performed by the `WriteRead` routine can be avoided by forcing the user to manage the memory-mapped temporary files. The `WriteReadFD` routine allows the user to pass the file descriptors of memory-mapped files in place of the memory buffers. The output data is taken directly from the "write" temporary file, while the input data is written directly to the "read" temporary file.

The user program also has direct access to the various device registers in the Splash 2 system. The register-access commands of T2 are available to the C programmer, but for efficiency reasons a separate set of C macros is also provided. These macros typically accept a symbolic register identifier and a value and perform any necessary data alignment prior to reading or writing a register.

### 6.5.3 Device Driver

The interface between `libsplash` and the Splash 2 hardware is handled by the *device driver* [2]. A device driver is a body of code written for a particular physical device which executes within the protected domain of the operating system itself. The Splash 2 driver provides entry points for the various operating system calls such as `open`, `close`, `mmap`, `read`, and `write`. The `open` call reserves a device for use by the user process and typically performs a variety of hardware and software initializations, while `close` frees the device for use by other processes. The `mmap` system call is the mechanism by which the operating system makes available to the user some portion of the physical address space of the device. Input to and output from a device are done with the `read` and `write` calls. The driver is also responsible for handling system interrupts caused by the device.

The physical address space of the Splash 2 device is composed of several distinct segments corresponding to the registers and memory on the Interface and Array Boards. The register space of each board is further divided into two pieces: user mode space and kernel mode space. User mode space contains those registers which may be mapped directly into the address space of the user program, while kernel mode space contains registers reserved for use strictly within the device driver. As a rule of thumb, access to registers that may adversely affect the operation of the system, such as DMA and interrupt registers, is limited to the "trusted" device driver. The remaining registers and memories may be mapped into user space.

Since the Development Board does not support DMA-controlled input and output or interrupts, the device driver for systems with the Development Board relies entirely upon the `mmap` call. All the registers on the Development Board are mapped into the user's address space, and a set of user-level library routines is provided to support input and output. The `read` and `write` system calls are not supported.

The bank register is managed by the device driver and is transparent to the user program. Whenever a user reference crosses a 24-bit segment boundary, a memory fault is incurred that transfers control to the Splash 2 driver. The driver then unloads the mapping for the previous segment of memory, maps in the new segment, and updates the bank register. No further intervention by the driver is required until the next time a reference falls outside of the current segment.

Since the Processing Element memory is not truly dual-ported, special care must be taken to avoid simultaneous access from the host and from the FPGA. The device driver coordinates access to the memories through the operating system's page fault mechanism. Whenever the Splash 2 system clock is enabled, *all* of the Processing Element memories are unmapped from the user's address space. Therefore, any access to a PE memory from the host causes a fault, transferring control back to the driver. The driver then stops the system clock and unloads the mapping for the clock registers to prevent the user from inadvertently restarting the clock. After the clock has been stopped, the "Xilinx Disable" signal is asserted to passivate the FPGAs, a software timeout interrupt is scheduled for about 10 msec in the future, and the referenced memory segment is mapped in before control is returned to the user program. When either the timeout interrupt or a subsequent user reference to the clock registers occurs, the procedure is reversed by unloading the memory segment, deasserting the Xilinx Disable signal, and restarting the clock.

The device driver for systems containing the Interface Board supports DMA transfers through the `read` and `write` system calls. When a user-level input or output request is made, the driver must perform a variety of software bookkeeping operations before and after the actual data transfer. First, the user's data buffer is mapped into the kernel address space. Next, each page of the buffer is locked into physical memory to prevent the operating system from paging it to disk during the transfer. Then, if the buffer does not begin on a 16-word boundary, the transfer is aligned by manually copying data to or from the DMA channel. Once the buffer has been mapped, locked, and aligned, the DMA transfer is begun. When the transfer is complete, the Interface Board signals the driver by generating a hardware interrupt. The interrupt handler returns control to the driver, which reverses the process, copying any data remaining after the last 16-word block and then unlocking and unmapping the buffer.

The SBus hardware has a peak data bandwidth of nearly 60 MB/sec. Unfortunately, due to the software overhead associated with the DMA transfer, principally mapping the buffer into the kernel space and locking the pages in memory, the best transfer rate a user-level program can expect to achieve is about 23 MB/sec, or about 40 percent of the peak.<sup>3</sup> For small transfers, the software overhead of DMA can completely dominate the time to completion. Therefore, for requests of less than 1024 bytes, the `Read` and `Write` library routines handle the transfer entirely in user mode using slave read and write accesses to the FIFOs.

## 6.6 DIAGNOSTICS

The suite of diagnostic software for Splash 2 evolved from the need to test and debug the hardware. The diagnostics were originally written to support low-level hardware debugging and system software design, but as the project progressed they took on new roles in the postmanufacture testing of new boards and the routine health

<sup>3</sup>These values were empirically determined in our laboratory. A hardware logic analyzer was connected to the LED register and to the SBus grant signal on the Interface Board. A version of the device driver was written that marked events by writing to specific bits of the LED register. The logic analyzer then recorded the time spent in the various phases of the I/O transfer.

checkups of running systems. The principal components of the test suite are the `tsdb` debugger and the `robocop` diagnostic.

Support for the lowest level of hardware debugging is provided by the “trivial SBus debugger,” or `tsdb`. This tool is not specific to Splash 2, but rather operates on the physical address space of a given SBus slot. A simple command interpreter allows the user to examine and set locations by specifying an offset within the SBus slot space. Other commands include read and write loops to allow triggering of test equipment. The user can define a set of symbolic names to use in place of numeric values. These symbols can then be used in any command that expects a numeric value such as a physical address.

The main diagnostic program is called `robocop`. `Robocop` consists of a set of VHDL Processing Element programs and C host routines. The design philosophy of `robocop` is to test the functionality of the system in ever-increasing distance from the host, starting with the SBus interface and proceeding through the Interface Board eventually to the Array Boards. On the Interface Board, `robocop` begins by testing the various status and control registers, then the program and readback memory, the programmable clock, XL and XR, and finally each of the DMA controllers. Once the Interface Board passes all of the tests, `robocop` proceeds to test the Array Boards, starting with the PE memories, the FPGAs, and the Crossbar and data paths. Finally the data path between Array Boards is tested.

A simple menu-driven interface allows users to select tests to perform on individual components or run on the entire system at one of several levels of detail. Any errors discovered are logged on the host by both system name and by individual board serial number.

`Robocop` may also be configured to run in the background, automatically starting up whenever the Splash 2 system is not in use. This background mode is completely transparent to the user. If the diagnostics are running when a user attempts to start an application, the `OpenAndInit` routine in `libsplash` will send a signal to the `robocop` process causing it to gracefully shut down. Once `robocop` has exited, the `OpenAndInit` call returns control to the application in the normal manner. In addition to the normal error logging, when running in background mode, errors are also reported by sending electronic mail to a list of system maintainers.

## REFERENCES

- [1] J.K. Ousterhout, *Tcl and the Tk Toolkit*, Addison-Wesley, Reading, Mass., 1994.
- [2] J. Stigliani, *Writing SBus Device Drivers*, Sun Microsystems, Inc., Mountain View, Calif., 1990.
- [3] Synopsys, Inc., *FPGA Compiler Reference Manual*, Synopsys, Inc., Mountain View, Calif., 1994.
- [4] Xilinx, Inc., *The XC4000 Data Book*, Xilinx, Inc., San Jose, Calif., 1994.

# CHAPTER 7

---

## A Data Parallel Programming Model

*Maya Gokhale<sup>1</sup>*

### EDITORS' INTRODUCTION

The following chapter describes an alternative, data parallel, programming model suitable for some of the applications for Splash 2 or for CCMs in general. In the “standard” approach adopted for Splash 2, programmers must design the processor architecture, at least in concept, at the level of a block diagram of comparator boxes, adders, and such, for processing the data, and to a lesser extent for sequencing and control. Given the linear flow of data in many applications, many of the algorithms are most easily viewed as a series of processing boxes connected by lines representing the flow of data (see, for example, Figure 11.12 in Chapter 11). Although the programming of such an algorithm in VHDL in the Splash 2 programming environment is a relatively straightforward programming process, the fact remains that the hard work of determining the data processing steps needed, laying those steps out with the data flow, and partitioning the entire computation into chip-sized pieces has already been done before any programming ever takes place.

It is this process, the design of a processor architecture suitable for a given application, that the dbC approach suggests could be done automatically by the compiler. The underlying idea, at least part of which is certainly not new, is that the programmer is able to write code in a variant of C that supports both bit-oriented data types and massive parallelism of SIMD computation. The compiler then translates the dbC code into assembly language-level instructions and produces as output

<sup>1</sup>A version of this appeared as Gokhale and Schott [6] and is used with permission.



the VHDL code necessary to create an instance of a processor architecture capable of executing the specific assembly-level instructions needed for the particular computation at hand. It is in this last step that the Splash 2 dbC work differs from more ordinary SIMD machine computations. In the CM-2, for example, a specific target micro-architecture existed in hardware, and the Paris assembly-level instructions used that micro-architecture. Unlike some other approaches to computing on FPGAs, in which a specific micro-architecture is synthesized and instructions for that micro-architecture are used, with dbC in Splash 2 the architecture is defined by the instructions to be executed, and only as much architecture as is needed is eventually synthesized.

The key issue in the use of CCMs has always been the ability to produce working programs by programmers (rather than hardware engineers) with an expenditure of effort and time consistent with other programming tasks. A dbC approach, if it were to be successful, would go a long way toward resolving that key issue. There were three reasons, however, that precluded the consideration of dbC as the standard method by which Splash 2 was to be programmed.

- In its present form, dbC is suitable only for SIMD applications. Many of the Splash 2 applications simply are not suitable for a SIMD implementation (at least not a SIMD implementation on Splash 2), and it was necessary to have a programming environment that would accommodate those applications.
- The performance of dbC programs on Splash 2 and their use of the still relatively precious FPGA resources is not yet good enough that one could have demonstrated “success” on some important applications.
- Most important, and not unrelated to the previous point, dbC is a research project in its own right and, even now, has not come to closure. What was necessary for the Splash 2 demonstration project was a programming environment in which applications could be developed in a time frame consonant with the rest of the project. Although dbC had existed for other (standard) machines prior to the start of the Splash 2 project, dbC for Splash 2 did not exist and could not be predicted to exist in time for applications development. Further, the degree of risk concomitant with any real research project made the adoption of dbC as “the” programming mode for Splash 2 impractical.

In short, then, we offer this chapter as a suggestion of what the (near) future may hold for the programming of CCMs. Applications have been programmed using dbC, the research continues, and we would expect that future CCMs might rely on dbC or a similar language in a manner from which we were prevented by the sequence of events.

## 7.1 INTRODUCTION

The standard methodology for programming Splash 1 and Splash 2 was through hardware description languages. Splash 1 was programmed using the Logic Description Generator (LDG) described in Gokhale et al. [5], a textual HDL that facilitated the description of systolic, hierarchical designs. LDG was developed in-house to meet

the need for a low-level tool that nonetheless permitted the user to concisely describe a large amount of logic. On Splash 2, we were able to design at a much higher level. Behavioral VHDL approaches the expressive power of a parallel programming language. However, orchestrating a large number of concurrent event-driven loops is complicated, time-consuming, and error-prone. Splash 2 is a complex collection of devices. Although the simulator went a long way in helping to verify the correctness of a Splash 2 design, the programmer had full responsibility for creating the design in the first place, which required working out the timing of a multiplicity of interlocking events across the 17 FPGA chips, memories, FIFOs, crossbar, and host. VHDL programs were required for each distinct FPGA chip design. The crossbar program was in a separate ASCII file. A control program on the host was required to send data and control signals to the array and to read back results.

We knew that raising the conceptual level from hardware design to parallel programming would make Splash 2 (and custom computers in general) accessible to a much wider range of programmers. It would be ideal to write a single parallel program, with some portions executed on the Splash 2 Array Board and others on the host, with communication and coordination between the two (as well as among FPGA computing elements) managed automatically.

It was at this point that a related SRC project was synergistic to the problem of programming Splash 2. Another group at SRC had designed and built the TERASYS SIMD array, composed of custom Processor-in-Memory chips [3]. A new language, data-parallel bit C (dbC), was developed to program TERASYS [12]. Two features of the language and its implementation made it especially appealing for Splash 2.

First, in dbC, bits are first-class parallel objects. Variables of arbitrary bit length can be created, and operations over arbitrary bit length data objects are supported. On Splash 2, this allows us to create and operate with small (1-, 2-, and 4-bit) objects, saving valuable resource on the FPGA. For example, the genetic sequence comparison application uses 4-bit counters to record edit distances. It would not be possible to describe this structure accurately in conventional C, since bitfields are promoted to "int" for computation. In dbC, not only is the storage minimized, but the computation is over the actual bit length rather than a standard container size such as 32 bits.

The second enabling aspect of dbC was in its implementation: the dbC "compiler" is actually a translator from the parallel ANSI C superset to ANSI C. The parallel constructs are invoked as function (or macro) calls. For TERASYS, the parallel operations are implemented by a microcode library. For Splash 2, we synthesize logic on a program-by-program basis to support exactly those parallel operations that are required for a given program. The parallel operations are executed on the Array Board, with serial data manipulated on the host. Clock events, FIFOs, the crossbar, and FPGAs disappear from the programming model. A single dbC program controls both Splash 2 and the host. As an added advantage, a dbC simulator had been written for the TERASYS project and could be used by Splash programmers to debug their data parallel programs on a workstation prior to synthesis.

Thus, concurrent with application development in VHDL, we embarked on a research project to build a dbC-to-Splash 2 compiler. We realized at the outset that dbC, which follows the SIMD programming model, would not be suitable for all Splash 2 applications. Many applications are inherently MIMD: different processors perform different tasks. Some applications have real-time constraints, which might

not be met through high-level synthesis. Nevertheless, the key issue in the use of CCMs has always been the ability to produce working programs by programmers rather than by hardware designers, with an expenditure of effort and time consistent with other programming tasks. A dbC approach, even with the restricted application domain, would go a long way toward resolving that key issue.

In this chapter, we describe the dbC language and compiler, which translates programs written in a data parallel superset of ANSI C into high-level VHDL for the Splash 2 array of FPGA chips. The next section contains a brief introduction to dbC. Next, we describe the dbC/Splash 2 compiler and illustrate the compilation process with a simple example. Section 4 details how data parallel communication and global reduction operators are mapped onto Splash 2. Optimizations are described in Section 5. We evaluate our system in Section 6 by showing performance on a genetic database search problem coded in dbC. Finally, we summarize and sketch future directions.

## 7.2 DATA-PARALLEL BIT C

dbC is an ANSI C superset similar to MPL [10] and C\* [13]. The programming model is that of a SIMD processor array in which a host processor controls instruction sequencing of many Processing Elements (PEs) (see Figure 7.1). The PEs receive instructions from the host. A PE can be active, in which case it executes the current instruction, or inactive, in which case it ignores the instruction. The active state is controlled by a mask, the *context bit*. Each PE can communicate with its nearest-neighbor in the user-defined virtual topology (a linear topology is illustrated). PEs can also communicate in arbitrary any-to-any patterns through an interconnection network. Global combining operations (also called *reduction operations*) such as global OR, SUM, MAX can be performed over the entire PE array, with the result of the operation being returned to the host.

### 7.2.1 dbC Overview

The dbC programmer specifies the number of PEs by initializing two predefined variables, `DBC_net` and `DBC_net_shape`. `DBC_net` must be initialized to the number of dimensions in the PE array, and `DBC_net_shape` is a vector of rank `DBC_net`, each element of which gives the size of the corresponding PE dimension.

In dbC the programmer designates data which is to reside on the processor array with the attribute *poly*. Figure 7.1 shows a one-dimensional 12-processor array

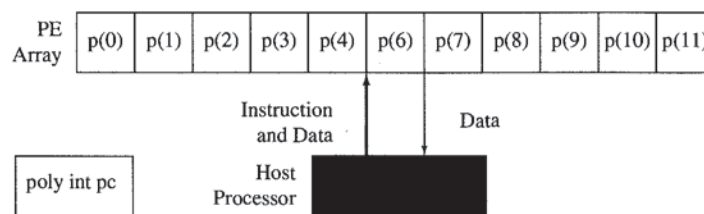


FIGURE 7.1 A SIMD Array

with a `poly int (P)` variable. All of the normal C operators can be used on poly data. In addition, special built-in operators that perform global combining operations are also defined. Processor activity is controlled by parallel control constructs such as “if,” “where,” or “while.” Interprocessor communication is initiated by the programmer with calls to intrinsic functions such as `DBC_net_send` for nearest-neighbor communication, `DBC_send` for arbitrary communication, and `DBC_read_from_proc` (`DBC_write_to_proc`) for PE-host (host-PE) communication.

As noted above, dbC was originally designed for the TERASYS SIMD array. It also runs on the Connection Machine-2. Both of these SIMD arrays have one-bit processors that perform arithmetic bit serially.

### 7.2.2 dbC Example

We show in Figure 7.2 a dbC program to compute the cross-correlation of two bitstreams. The program compares two bitstreams and accumulates a count of the number of times an individual bit in the bitstreams had the same value. The bitstreams are compared with a delay of zero bits, then with successively larger delays, usually one bit longer for each delay. For each of the delays a counter records the number of matches (see Figure 7.3). The delay is sometimes called a “lag.”

In a typical implementation, there are individual cells that perform the correlation between two streams of data for one value of the delay, that is, there is a cell for delay 0, delay 1, and so on. Each cell includes a comparator and a counter. The comparator compares the data in the bitstreams; if they are the same, the counter

```
#include <interproc.hd>
typedef poly unsigned Boolean:1;
Boolean a;
#define N 128
#define NPROC 64
unsigned DBC_net = 1;
poly unsigned int R:16 = 0;
unsigned DBC_net_shape [1] = {NPROC};
int right[1] = {1};
void main()
{
  all {
    int b;
    a = 0;
    for (b=0; b < N; b++) {
      DBC_write_to_proc(&a, 1, 0);
      R += (a ^ (Boolean) b);
      DBC_net_send(&a, a, right);
      printf("%d \n", DBC_read_from_proc(R, (b*NPROC)));
    }
  }
}
```

FIGURE 7.2 dbC Cross-Correlation Program