

# Digital Signal Processing at 1GHz in a Field-Programmable Object Array

Dirk R. Helgemo

**Abstract**—Autonomous MAC and ALU processors and register files (three types of Silicon Objects) are implemented with custom logic to achieve 1GHz fixed-point multiply and accumulate. Synchronous programmable interconnect and embedded storage reduces the need for difficult index calculation and the use of external memory for intermediate values. The flexibility of the objects and their interconnect allows the level of parallelism to be chosen freely based on performance requirements and resource constraints. Arraying hundreds of objects in parallel in a single chip enables incredible DSP performance from a flexible, in-circuit reprogrammable architecture.

For example, a 1024-point FFT with (16+16)-bit complex samples can be completed every 160 clock cycles (i.e., every 160 nanoseconds) using 64 butterflies (128 MAC, 128 ALU, and 64 RF objects) assisted by 128 ALU and 64 RF objects for inter-stage data routing.

**Index Terms**—Digital Signal Processing (DSP), Application-Specific Programmable Product (ASPP), Reconfigurable Architecture, Field-Programmable Object Array (FPOA).

## I. INTRODUCTION

MathStar is offering a massively parallel high-performance computation fabric. Individual processing units, called Silicon Objects, are programmed individually and act autonomously. Each object is less than 400x400 micrometers square, implemented in custom logic, allowing hundreds of high-speed objects to be tiled on a single chip. Silicon Objects and their interconnect are programmed to construct computation macro blocks – composing simple scalar operations (addition, multiplication, logic, storage) into complex functions (e.g., 1024-point FFT). Interconnect and instructions are configured after fabrication via PROM, resulting in a field-programmable object array (FPOA). All communication and processing is synchronized to a global clock (up to 1GHz), removing the design issue of analog timing closure altogether.

## II. SILICON OBJECT COMMUNICATION

Silicon Objects communicate via 21-bit buses composed of the following: sixteen bits of data, one bit indicating the validity of the data (e.g., for event-driven programming), and four bits of user-defined side-band control signals.

Communication proceeds synchronously and cooperatively. Buses are driven directly by registers (i.e., no intervening logic) for the most aggressive digital timing between objects. Values of interest are read by a cooperating receiving object; thus data is

pulled rather than pushed through the architecture. Objects synchronize to the same digital clock cycle (phase) via user programming of control signals and/or data patterns.

The communication topology is a hybrid: objects can read registers from adjacent neighbors, or from any distance via pipelined “party lines.” Neighbor registers in diagonal and Manhattan directions are observed with no latency (the same as local registers). Party lines can turn, pass, land, and/or launch at every object hop. The land/launch combination can be chosen to insert a pipeline delay and restore digital coherency, thereby enabling communication at any distance (at the expense of latency and party line landing registers). The communication topology thus facilitates the programming of high-speed computation kernels of arbitrary size and shape.

## III. SILICON OBJECT TYPES

Whereas the communication infrastructure across a given fabricated Silicon Object array is uniform, the silicon implementation of each element can be unique, yielding a heterogeneous array. The following are available element implementations, known as Silicon Object types.

### A. Multiply-Accumulate Object (MAC)

The MAC object type accepts two 16-bit signed integer inputs every clock cycle, multiplies them together, and adds or subtracts the product into the 32-bit accumulated result. The accumulator can be configured either to saturate, or to wrap into an 8-bit over-/underflow counter, tolerating up to 40-bit intermediate results. The entire accumulator is visible on object outputs and can be reset to zero or reloaded per control inputs, allowing either a new sequence to be started or a paused sequence to be resumed. The operation consumes fresh inputs and generates a result every clock cycle, with a processing latency of two clock cycles.

### B. Arithmetic-Logic Unit Object (ALU)

The ALU is the most general-purpose object type. It employs a 16-bit add, shift, and logic operator controlled by an 8-instruction state machine. Each instruction selects up to three 16-bit input words and a carry input bit, configures the operator (a.k.a. opcode), selects one or more result destination registers, and specifies conditional execution and branching options. This object type contains nine working registers (four for neighbors, five for party lines); two programmable constant registers; and two wired constants. Thus there are twenty-one possible inputs and nine possible outputs. In a single clock cycle, the current instruction is fetched and decoded, the operator is executed, the result is stored (subject to conditional execution), and the next instruction is selected per branching.

Unclassified manuscript sent May 30, 2003, to the HPEC 2003 Conference.

Dirk R. Helgemo is Chief Architect of MathStar, Inc., 5900 Green Oak Drive, Minneapolis, MN 55343; phone 952-746-2200; fax 952-746-2201; email Dirk.Helgemo@MathStar.com.

### C. Register File (RF)

The RF object type provides fast storage within the array. Up to two 20-bit values can be read and two 20-bit values written simultaneously every clock cycle, with an access latency of two clock cycles. Storage capacity is 64 20-bit words, also configurable as 32 40-bit words.

The read and write ports can each be configured for random or sequential access. Thus, an RF can be configured as a dual-port RAM, a FIFO, or random-write sequential-read. The last combination, also known as “sort” mode, allows values to be written in an arbitrary index order, but then retrieved as a sequence without the burden of address generation. That is, values can be written to arbitrary addresses in anticipation of the order in which they will be read out.

## IV. FAST FOURIER TRANSFORM (FFT) VIA OBJECTS

### A. Complex Multiplication

Two MAC objects can be efficiently ganged to multiply two complex numbers. Four products are generated, two of which are differenced, two of which are summed. Thus two MAC objects can generate a complex result every two clock cycles, with a latency of three clock cycles.

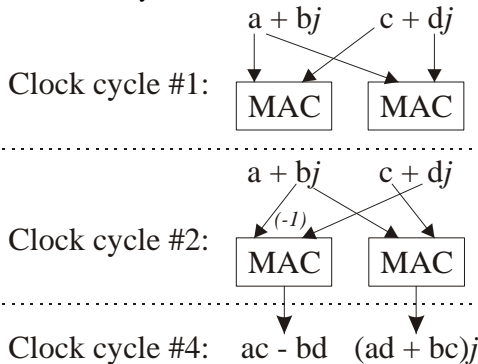


Figure 1: Complex Multiplication

### B. Radix-2 Butterfly

The butterfly kernel within the FFT algorithm accepts two complex numbers from a previous FFT stage, multiplies one of the inputs with a twiddle factor (a complex constant), and performs a complex sum and difference, yielding two complex numbers for the next FFT stage.

Consider the implementation of a decimation-in-time butterfly:  $out_1 = in_1 + W^k in_2$ ,  $out_2 = in_1 - W^k in_2$ , where  $W^k$  is the complex twiddle factor. Due to the predictability of the twiddle factors, they are precalculated and stored into an RF object configured for sequential read mode – and address generation is not required.

Thus the butterfly algorithm is as follows: 1. Fetch the precalculated complex  $W^k$  from the RF object. 2. Multiply  $W^k in_2$  via two MAC objects (as described above). 3. Use two ALU objects to both sum and difference the complex product against  $in_1$  to generate outputs  $out_1$  and  $out_2$ , respectively, over the next two clock cycles. Thus, butterfly outputs can be calculated every two clock cycles, with a latency of five clock cycles.

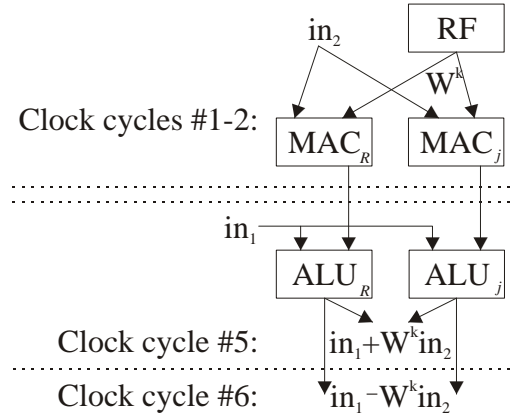


Figure 2: Decimation-In-Time Butterfly

### C. Fast Fourier Transform (FFT)

Each stage of butterflies chooses a different pairing of the previous stage’s results (as well as different twiddle factors) until all of the FFT inputs affect all of the FFT outputs (i.e.,  $2^n$  points require  $n$  stages). While a single butterfly can be leveraged to any size FFT, multiple parallel instantiations of the butterfly (in powers of two) increase the theoretical computational performance dramatically:

# Butterflies	MACs	ALUs	RFs	Rate	Latency
1	2	2	1	$1/(n2^{n+1})$	$3+n2^n$
2	4	4	2	$1/(n2^n)$	$3+n2^{n-1}$
$2^{n-2}$	$2^{n-1}$	$2^{n-1}$	$2^{n-2}$	$1/4n$	$3+4n$
$2^{n-1}$	$2^n$	$2^n$	$2^{n-1}$	$1/2n$	$3+2n$
$n2^{n-1}$	$n2^n$	$n2^n$	$n2^{n-1}$	$1/2$	$3+2n$

Figure 3: Butterfly Parallelism for  $2^n$ -point FFT ( $n$  stages) (Rate is results per clock cycle. Latency is clock cycles.)

Fortunately, practical performance does not substantially lag the theoretical ceiling. Butterflies are kept 100% utilized by providing two new complex inputs every two clock cycles. Either an RF object or two ALU objects can sustain this bandwidth indefinitely. The trick lies in efficient transitions between FFT stages.

Every butterfly result is used precisely twice as an input into the next stage. Therefore, the butterfly results (two complex result every two clock cycles) are routed via ALU objects (with stage-specific directions) toward the two butterflies for the next FFT stage. An RF object sorts the complex data values into the correct order for the next stage using a nearby ALU object to generate stage-specific write addresses into the RF object.

Performance is lost between stages only if the RF object cannot be loaded in time to start the next stage. In practice, index analysis of the data dependencies between stages allows the next stage to be started while the previous stage completes. (Ironically, fully parallelized butterflies cannot avoid stalling between FFT stages because none of them can start until the previous stage completes.)

## V. RESULTS

A 1024-point FFT with (16+16)-bit complex samples can be completed every 160 nanoseconds using 64 butterflies (128 MAC, 128 ALU, and 64 RF objects) assisted by 128 ALU and 64

RF objects for inter-stage data routing. An array of 25x25 objects provides the required number and arrangement (with over 100 objects remaining for control sequencing), yet fits within a 10x10 millimeter square of silicon. Note that the object commonality allows larger, smaller, and different mixes of object types, I/O, and on-chip RAM to be readily constructed according to specific application requirements.