# Custom VLSI ASIC for Automotive Applications with Recurrent Networks

**R. Tawel and N. Aranki**
Jet Propulsion Laboratory
Pasadena, CA **911098099**


**G. V. Puskorius, K. A. Marko, L. A. Feldkamp,**
**J. V. James, G. Jesion, and T. M. Feldkamp**
Ford Research Laboratory
Dearborn, MI 48121-2053

## Abstract

Demands on the performance of vehicle control and diagnostic systems are steadily increasing as a consequence of stiff global competition **and** government **mandates.** Neural **networks provide a** means of creating control and diagnostic strategies that will help in meeting these demands efficiently and robustly. This paper describes a VLSI design that permits such networks to be executed in real time as well as the application, misfire detection, that served as a focus for the collaborative effort.

## 1 Introduction

The control system of a modern automobile involves several interacting subsystems, almost any one of which provides interesting theoretical and engineering challenges. Further, increasingly stringent emissions regulations require that any malfunctioning component or system with the potential to undermine the emissions control system be detected and identified. Neural networks have the potential for major impact in this work. Benefits may be anticipated in design time or performance of a control (Puskorius and Feldkamp, 1996) or diagnostic strategy (Marko et al., 1996). We have shown that both of these applications are suited to the use of recurrent multi-layer perceptrons, the architecture of which may be regarded as the joint generalization of a feedforward multi-layer perceptron and a o-lc-layer fully time-lagged recurrent network. A potential barrier to the use of such networks, however, arises from the considerable burden from other functions already carried by existing powertrain processors. This prompted an effort to develop a VLSI design that would

facilitate the implementation of recurrent networks in high-volume products.

## 2 Neuroprocessor Chip

The design constraints for this project called for the development of an inexpensive, fully autonomous, and commercially viable electronic chip. This single chip implement ation was required to be (1) extremel y compact in size (mass market potential) (2) flexible (several neural based applications would share the hardware and sequentially execute on it), and (3) accurate (no miscalls due to limited hardware resolution). Observing that combustion events occur, even at maximum engine speed, on a millisecond time scale, a novel and extremely compact and powerful layer-multiplexed bit-serial neuromorphic architecture was developed and exploited for the VLSI CMOS implementation.

### 2.1 Architect ure

The required computations can be summarized as a series of parallel multiply and accumulate (MAC) operations interspersed by an occasional nonlinear operation. We exploited five basic stategies to achieve our desired goals: (1) parallel intra-layer topology; (2) single-instruction-multipk-data (SIMD) architecture; (3) bit-serial fixed-point computational techniques; (4) inter-layer multiplexing of neuron resources; and (5) nonlinearities handled by look-up-tables.

The resulting architecture is shown schematically in Figure 1 and consists of (1) a global controller; (2) a pool of 16 bit-serial neurons; (3) a ROM based bipolar sigmoid activation look-up-table; (4) neuron state registers; and (5) a synaptic weight RAM. In this design, both inputs to the network as well as neuron outputs
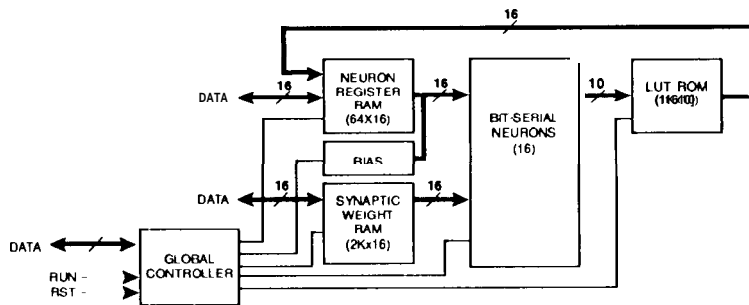


Figure 1: Schematic representation of forward propagation module.

arc stored in the neuron state RAM. When triggered by the global controller, each of the 16 neurons performs the multiply and accumulate (MAC) operation. They receive in a bit serial fashion as input the synaptic weights (from the synaptic weight RAM) and activations from either (a) input nodes or (b) outputs from other neurons and output the accumulated sum of partial products onto a tri-stated bus which is commonly shared by all 16 neurons. Because of the computat ional nature of neural networks where information is sequentially computed a layer at a time only enough neurons arc physically implemented in silicon as exist on the layer with the largest number of neurons for all applications of interest. As such, a candidate pool of 16 silicon neurons was chosen. This iutra-layer pool of neurons is organized in a SIMD configuration. Single-instruction (S1) means that all active neurons in the pool execute the same instruction at the same time. Multiple-data

(MD) means that each active neuron acts on its own slice of data, independently of all other processors. Thus the chip performs fully parallel computations under the supervision of the global controller.

A significant reduction in silicon real-estate was achieved by performing inter-layer multiplexing of the 16 neuron pool, i.e., the hardware used in calculating the activations of neurons in one layer is reused for the calculation of neurons in another layer, since neurocomputations are performed a layer at a time. We also used bit-serial algorithms extensively for arithmetic operations because their canonical nature and minimal interconnection requirements make them particularly suitable for efficient VLSI implementation.

## 2.2 Controller

At the heart of the neuroprocessor architecture is the global controller. The controller contains the logic to enable the neurochip to execute its task. This task is to load an architecture from RAM, and once triggered, to generate all necessary control signals in addition to orchestrate data movement on-chip and off-chip. When there are no computations being pm-formed, the global controller remains in the idle state, signaling its availability by having the active low BUSY flag set high. When a LOAD command is issued, the controller reads from RAM a neural network topology and goes into an idle state. When the RUN command is subsequently issued, the global controller is in charge of providing control signals to the RAM, ROM, and the 16 on-chip neurons, in order to proceed with the desired computation. Input activations are read out of the 64×16 Neurou State RAM, synaptic weights are read out of the 2K x 16 Synaptic Weight RAM, and both are propagated to the bank of 16 neurons. In this way, the global controller keeps track of both intra-layer operations as well as inter-layer operations. Upon completion of a forward pass through the network architecture, the global controller asserts the BUSY flag and returns to the idle state.
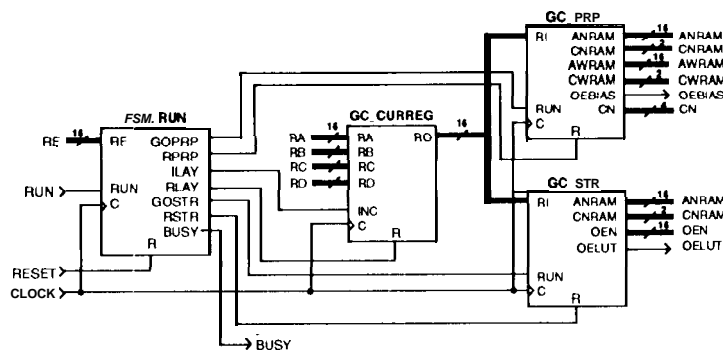


Figure 2: Run-time forward propagation controller.

## 2.3 Neurons

Fixed-point bit-serial algorithm for operations such as addition and mult iplication are uniquely suitable for efficient VLSI implementations because of their highly compact representations. For example, the size of an $n \times n$ bit multiplier scales quadratically $(O(n^2))$ for a bit-parallel implementation and linearly $(O(n))$ for a bit-serial one. ~lit-serial techniques were therefore used. A schematic representation

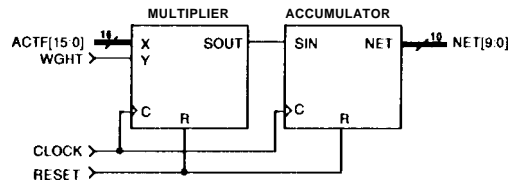of a bit-serial neuron is shown in Figure 3.



Figure 3: Ilit-serial neuron.

Precision constraints for the misfire problem called for the use of a $16 \times 16$ bit fixed-point multiplier. In operation, the multiplier accepts as input either an input stimulus to the neural network or an activation output from a neuron on a previous layer. It multiplies this quantity by the corresponding synaptic weight. The input stimulus (or activation output) is presented to the multiplier in a bit-parallel fashion, while 'the synaptic weights are-presented in a bit-serial fashion The serial output of the multiplier feeds directly into an accumulator.
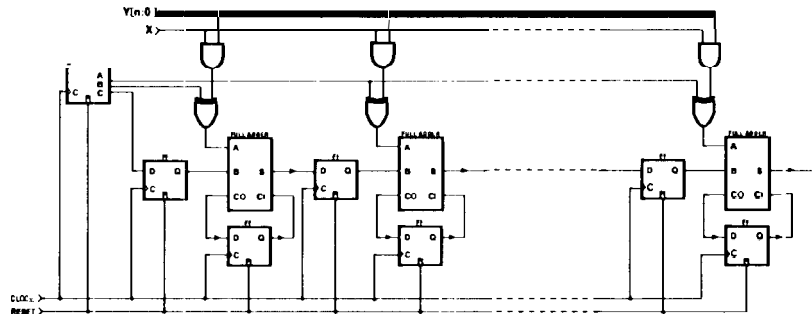


Figure 4: Hit-serial multiplier of length n

The multiplier shown in Figure 4 is a modified and improved version of a previously reported serial multiplier. Any size multiplier can be formed by cascading the basic multiplier cell. The bit-wise multiplication of the multiplier and multiplicand is performed by the AND gates. At each clock cycle, the bank of ANI) gates compute the partial product terms of the multiplier Y[15:0] and the serial multiplicand X (t). Two's complement multiplication is achieved by using XOR gates on the outputs of the AND gates. By controlling one of the inputs of the XOR gate, the finite state machine FSM can form the two's complement of selected terms based on its control ffow. In general, for an $n$ x n multiplier (resulting in a $2n$ bit product), the multiplier can be formed by using $2n$ basic cells and will perform the multiplication in $2n + 2$ clock cycles. Successive operations can be pipelined and the latency of the LSB of the product is $n + 2$ clock cycles.

The accumulator, shown in Figure 5, is also a bit-serial design. It is extremely compact as it consists of a single bit-serial adder linked to a chain of data registers. The length of the accumulator chain is governed by the multiplication length. The multiplier takes $2n + 2$ clock cycles to perform a complete $n$ x $n$ multiplication. At each clock cycle, the accumulator sums the bit from the input data stream with both the current contents of the data register on the circular chain as well as any carry bits that might have been generated from the addition in the previous clock
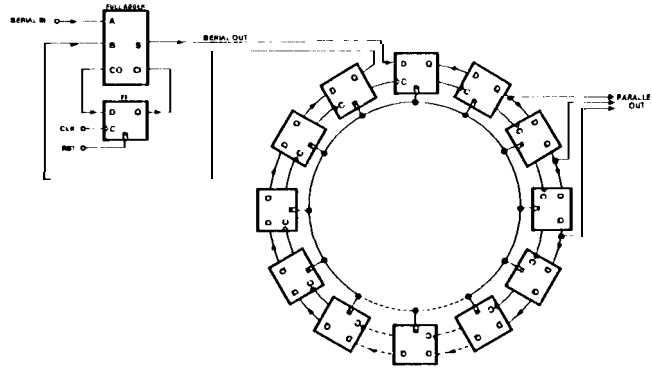
Figure 5: Ilit-serial accumulator of length n.

cycle. This value is subsequently stored onto the chain on the next clock cycle. This creates a circulating chain of data bits in the accumulator with period $2n+2$.

The neuroprocessor design was implemented using H P's 0.5 $\mu$m CMOS design rules. The first generation chip measured $8mm^2$ in size. *The* current design operates at a conservative 20 MHz clock speed. A neural application) can be loaded into the **hardware in** under 1 $\mu$s. Because of the S1 MD architecture, it takes 1.6 $\mu$s to simultaneousl y perform 16 multiply and accumulate operations. This translates into an effective computational throughput of 0.1 $\mu$s per MAC operation. The next generation processor will operate at 50 MHz.

## 3 The Misfire Diagnostic Problem

Because engine misfire can cause a significant, increase in tailpipe emissions and can damage the catalytic converter, it is a required diagnostic. Misfire detection must be performed between engine cylinder firings, which can occur at rates as high as 30,000 events per minute, so that approximately one billion events must be classified over the life of each vehicle. While there are many ways of detecting engine misfire, all currently practical methods rely on observing engine crankshaft rotational dynamics with a position sensor located at one end of the shaft,. Briefly stated, one looks for a crankshaft acceleration deficit following a cylinder firing and attempts to determine whether such a deficit is attributable to a lack of power provided on the most recent firing stroke. The method is complicated by several factors: 1) the crankshaft dynamics arc influenced by unregulated inputs from the driver and disturbances introduced through the driveshaft from road irregularities; 2) the dynamics arc obscured by measurement noise; 3) the crankshaft is not infinitely stiff and exhibits complex dynamics which mask the signature of the misfire event and which arc influenced by the event itself. In effect, we are observing the torsional oscillations of a nonlinear oscillator with driving forces applied at several locations along its main axis. While it is straightforward to write down dynamical equations that approximate the crankshaft rotational dynamics as a function of the combustion pressures applied to the piston faces, it is difficult, to solve those equations and even more difficult to solve the inverse inference problem associated with misfire diagnostics. Nonetheless, it was the expect ation of a discoverable dynamic relationship between the observed accelerations and the driving forces in this system, coupled with the absence of a satisfactory alternative approach, that motivated our exploration of recurrent networks as a solution to the problem.