

## A 1.5 GHz 90 nm Embedded Microprocessor Core

Franco Ricci, Lawrence T. Clark\*, Tim Beatty, Wing Yu, Alex Bashmakov, Shay Demmons,  
Eric Fox, Jay Miller, Manish Biyani, and Jon Haigh

Intel Corporation, Chandler, AZ 85226, USA

\* Department of Electrical Engineering, Arizona State University, Tempe, AZ 85287, USA

### Abstract

A 90 nm ARM™ V5TE compatible microprocessor core intended for high performance and low power embedded applications is described. The core includes an ECC protected 2<sup>nd</sup> level 512 MB cache, high-bandwidth single-cycle L1 cache line fill and evict, and cache coherency. Circuit design for high speed and low power are described, as well as their impact on the micro-architecture. Features to support low standby power modes and embedded test are also described.

**Keywords:** VLSI, microprocessor, power, embedded and cache

### Pipeline and Micro-architecture

The die plot and floor plan comprises figure 1(a). The core contains 39M transistors and is 2.5 x 7.5 mm with the L2 data array occupying approximately ½ the area as shown. Key blocks are labeled. The execution datapath (DDP) runs vertically under the L1 data cache, while the instruction pipeline (IDP) is alongside and includes a 128 entry branch target buffer. The execution pipeline is 7 stages, with loads and stores requiring 8 stages as in the previous generation design [1]. Substantial circuit and micro-architectural changes were required to reduce the gates per stage from an average of 27 and peak of 32 in a previous design to about 24 maximum. The target process technology is described in [2].

The pipelines are illustrated in Fig. 2, with the execution pipeline similar to our previous design [1]. The MAC requires from two to four stages depending on data and operation. Details of the integer pipeline are described in subsequent sections. The bus pipeline depends on the bus multiplier. The L2 pipeline requires 8 core clock cycles.

### Clock Distribution and Latch Design

Pulse-clocked latches minimize the delay through and the power consumption of the latch circuits. While over 40% sequential circuit and clock tree energy reduction have been demonstrated [1][3], increasing variation with process generations makes it increasingly difficult to use pulse-clocked latches. However, the power savings and lower delay sequential elements facilitates high performance block design. A high quality (low skew) clock distribution network is critically important.

The clock distribution network is divided into four regions. The pre-delivery region (CDC) consists of testing features and the bypass/ring oscillator muxes. The CDC contains 3 inversion stages. The delivery region (CDP) is embedded within the functional blocks and spans the entire die, comprising 9 inversion stages. The CDP contains shorting bars at multiple levels at the receiver inputs similar to the scheme described in [4]. The final stage of the delivery region, EGCLK, is gridded across five vertical trunks in the data flow direction (see Fig 1(b)). The third region comprises the Unit Clock Buffers (UCBs) as two inversion stages. The UCBs, which allow unit level clock gating, are located directly under the EGCLK trunk and are also embedded with

the functional blocks. The last region of the clock distribution network provides the final clock drivers, i.e., local clock buffers (LCBs).

Within the design, approximately 2.5k pulse LCBs drive 16k pulse latches. 4k non-pulse LCBs drive transparent latches and flip-flops. The former primarily drive dynamic circuits, but are also used in test and debug blocks. 780 unit level clock buffers (UCBs) comprise the preceding clock stage. 80% of the UCBs are gated to reduce unit level global clock activity. Simulations show final clock net activity factor over multiple tests is 0.12. Each LCB consists of two inversion stages (see Fig. 3). The Pulse LCB includes three stages of delay to generate the pulse width. The third delay stage is a NOR which allows enabling the pulse for reducing power consumption or functional gating. The enable signal hold time is approximately the width of the pulse and therefore does not require a latch to capture state.

The use of pulsed clocks requires special attention to asymmetries in the latch internal clocking and parasitics. The necessary pulse width is determined by the writing of a 0 into the state node, since writing a 1 through the P pass device P1 is helped by the N pass transistor N1 turning on early and the P pass transistor turning off late (see Fig. 4). The width requirement increases for the 0 case since the feedback transistor does not turn on until an additional inverter after the NMOS pass device is turned off. PMOS P1 does not aid the high to low transition as the N device aids the 0 to 1 transition since the state node is already a  $V_t$  below the gate voltage when the P device is enabled.

Scan is supported on the rising clock edge of all flip-flops and latches. All latches support a low standby power mode whereby all state is moved into high- $V_t$ , thick gate storage located within each latch (or flip-flop slave) [5]. Power can then be gated to the entire core while retaining state. The design is illustrated in Fig. 4 where the thick gate shadow latch also operates as the scan slave for LSSD scan clocking. Scan does not run at the full operating frequency.

### Execution Unit

The Execution pipeline stage contains the domino shifter and the static adder. The shifter critical path is three domino gates, each containing a multiplexor and following complex gate providing sign extension. Only right shifts are supported, providing rotate and sign extend capability. The shifter terminates in a multiplexing pulse-clocked latch for adder source select. It is fully bypassed and clock gated when not in use. The 32-b static adder occupies the second clock phase, and is completed after the next clock rising edge by the TLB CAM driver circuit. The adder critical path is complicated by a zero-detect controlling multiplexing of "Process ID" (PID) on the upper 7-b as well as recirculate and bypass features.

The adder is comprised of two sections. A modulo-3 Ling adder performs the lower 24-b addition. The upper 8 bits use

a carry-select scheme with a zero-detector on each output to identify whether the PID insertion is necessary. Carry select CARRY24 for the upper 8 bit adder and zero detect set up to the CAM driver clock edge so both possible values are carried through the recirculate and bypass multiplexors. The PID is inserted into the pulse-clocked differential TLB CAM driver latch as shown in Fig. 5.

### Multiply-Accumulator (MAC)

The pipelined MAC supports early termination and can multiply two 16-b operands every 2 clocks and 32-b operands every 4 clocks. A radix-2 Booth encoding algorithm is used. The Booth encoder works on 16 bits at a time producing 9 partial products. The feedback path from the previous pipeline stage generates two more partial products, which along with the accumulator operand, comprise the 12 to be added in one cycle. A fully static carry save adder (CSA) is used in the design. The CSA is implemented in three reduction stages, using static pass-gate circuits as shown in Fig. 6. The first stage uses three 4 to 2 reducers to reduce the 12 partial products down to six. The second stage employs 1-bit full adders as 3 to 2 reducers to produce four results. The final stage uses a single 4 to 2 reducer to produce a sum and a carry output. An adder sums the results in the next cycle.

The CSA is optimized to balance speed and power consumption. The CSA is 51 bits wide, presenting a physical design challenge in routing the partial products at various stages of the tree. To alleviate layout effort all reducer stages are of uniform width, enough to accommodate the 12 vertical partial product lines in one bitslice. The CSA is thus metal-limited in the vertical direction. An additional challenge was routing the Booth select lines to the 9 rows of the partial product mux. Due to this, the partial product multiplexor is metal-limited in the horizontal direction. Since the Partial Product/CSA speed paths are critical, the selects and partial product wires are widened to lower resistance. These constraints set the minimum height and width of the multiply-accumulator unit.

### L1 Caches

The virtually tagged, parity protected L1 caches are 4-way set associative with separate sense amplifier enables gated by the tag hit signals so that only 32 or 64 sense amps are enabled. A domino L1 data array output bus provides high speed. The 64-b bus flows directly from the cache, acting as the way multiplexor by combining data from different banks as well as data returned from the memory buffers. The latter are searched in a fully associative fashion for each access.

A separate 256-b bus supports fill and eviction data. This bus is time-multiplexed using a tri-state design. It provides L2 to L1 bandwidth of 40GB/s at 1.25 GHz and supports single-cycle fills and evictions. To allow the fast turnaround, this bus is driven only in one phase, with full keepers providing storage in the other phase. The physical organization of the cache was tailored to allow this high bandwidth, as 256 sense amplifiers or write drivers must be accessible from the bus as illustrated in Fig. 7. The words in each line are split over multiple banks as shown. For a fill/evict, each array (eight in all) has one word line enabled. So that two words may be read from a selected way during a load (64-b load double operation) the upper array way order

in each bank differs from the bottom way order to allow simultaneous access.

The tags use conventional NOR comparators with split match lines to limit the fan-in per side. The comparator is implemented as D2 domino, with the precharged sense amplifiers acting as the (footed) D1 stage. The exclusive-or function is implemented in a static OAI gate that allows a single pull-down transistor per bit, saving 4x diffusion loading per bit and limiting leakage currents on the comparator match nodes.

By checking if an instruction cache request is to the same line as the previous request, the access to the instruction tag array, including decode, read, sense and compare is avoided. By squashing accesses, all tag energy, save that from clock loading and leakage, can be avoided. Architectural level simulation shows that approximately 75% of instruction cache tag accesses are eliminated during normal operation. There is insufficient time in the pipeline to squash data cache tag accesses, but provision is made to avoid redundant TLB physical address lookups.

A separate tag array is used to allow snooping despite a virtually addressed cache. This allows multiprocessor operation and the increased micro-architectural parallelism afforded by a virtual data cache. The snoop tag is synchronized with the virtual tag, with each pointing to the other as data array when accessed as the cache tag as shown in Fig. 8.

### L2 Cache

The 8-way set associative L2 cache consists of a 512KB data array with corresponding tag and state arrays, and is re-sizeable to 256 kB with the removal of the upper 4 64kB data array blocks. Block level redundancy is supported. The replacement policy is not recently used. The L2 cache is a fully pipelined, non-blocking, and operates at  $\frac{1}{2}$  the core frequency. It is physically addressed using a 36-b address, providing 64 GB of addressable memory. All accesses to and from the L2 array occur at a full cache line width. The L2 cache always operates in write-back mode and is inclusive of the L1 cached lines. Hardware cache coherency is supported using the *MOESI* protocol. Line locking allows improved real-time response. Finally, the L2 supports a *push-cache* capability, where specially tagged bus transactions push data directly into the L2.

While the L2 pipeline operates at half the frequency of the core, it receives the same clock. This allows greater flexibility in generating internal signal edges. Clocks are only generated when valid operations are outstanding, as shown in Fig. 9. Array clocks are gated by a valid signal generated by the control block. Support for the  $\frac{1}{2}$  size option makes accesses to these halves mutually exclusive. Gating unit level clock buffers in the half not being accessed affords additional clock power savings. The L2 arrays are read with a dynamic sense amp activated by a following clock edge. This removes one of the basic races but does dissipate extra power due to greater than necessary bit line development when operating at low frequencies. ECC bits are read and written concurrently with data. The SEC/DED ECC scheme uses 10 bits to protect the 256-b lines. All data accesses are corrected as the syndrome generation and bit repair occupies the fourth L2 pipeline stage, as part of the data delivery to the L1 buffers. The L2 tags are parity protected.

### Design for Testability (DFT)

The core includes a combination of structural and functional test approaches. Scan insertion achieved over 98% latch coverage, assisting validation and debug as well as production test. A digital programmable clock stretch capability allows speed path debug using the scan chains.

Structural based functional testing (SBFT) allows core functional test independent of the SOC environment. To this end the core is isolated during test by a boundary scan-like ring. This not only allows the same core test content across products, but also allows for test parallelism within the product since SOC blocks can be running scan tests, while the core is running functional tests. Observability is provided by software controlled MISRs, which record all core outputs. All data flowing into the MISRs are fully qualified to avoid unknown state propagation. Small DFT “test agent” blocks test ensure coverage of the core interfaces.

The L2 SRAM arrays (data, tag and state) are tested using a programmable hardware BIST block providing maximal algorithmic coverage of the arrays early in the product life cycle and then a reduced test in the product life cycle to reduce test time. L1 caches are tested using the main execution pipe (decoder, register file, ALU) augmented by DFT registers and comparators. The test program is held in a 2kB SRAM attached to the L1 instruction cache bus. The fill/evict logic and bus provides high-bandwidth array read and write during test. This approach significantly reduces L1 cache test cost by using core datapath circuits.

### Measured Results

Silicon results at room temperature comprise Fig. 10 demonstrating 1.87 GHz operation at 1.5V. The minimum operating voltage is limited by a custom latch write failure in the core datapath. Static power is 58 mW @ 1.2V, room temperature and dynamic power is 770 mW @ 1.5 GHz and 1.3V as determined by Powermill simulations.

### Conclusions

The highest performance ARM compatible embedded microprocessor core to date has been described. The core features an L2 cache, cache coherency, and substantial testability features. Diverse circuit techniques, including pulse-clocked latches simulating master-slave flip-flops and combinations of static and dynamic circuits produce high performance at low power without resort to deep pipelines. Testability in an SOC environment is facilitated by combinations of BIST, signature analysis, and full scan capabilities. Reliability is enhanced by inline L2 error detection and correction and L1 parity protection. Greater than 1 GHz clock rates are achieved, while extensive clock gating and inclusion of low standby power modes improve applicability to hand-held battery powered SOC designs.

### References

- [1] L. Clark, et al., “An embedded 32-b microprocessor core for low-power and high-performance applications,” *IEEE JSSC*, 36, pp. 1599-1608, 2001.
- [2] K. Kuhn, et al., “A 90 nm communication technology featuring SiGe HBT transistors, RF CMOS, precision R-L-C RF elements and 1 $\mu\text{m}^2$  6-T SRAM cell,” *IEDM Tech. Dig.*, pp. 73-76, 2002.
- [3] J. Tschanz, et al., “Comparative delay and energy of single edge-triggered and dual edge-triggered pulsed flip-flops for high-performance microprocessors,” *Proc. ISLPED*, pp. 147-152, 2001.

- [4] N. Bindal, et al., “Scalable sub-10ps skew global clock distribution for a 90 nm multi-GHz IA microprocessor,” *ISSCC Dig.*, pp. 346-347, 2003.
- [5] L. Clark, F. Ricci, and M. Biyani, “Low standby power state storage for sub 130 nm technologies,” *to appear in IEEE JSSC*, Feb., 2005.

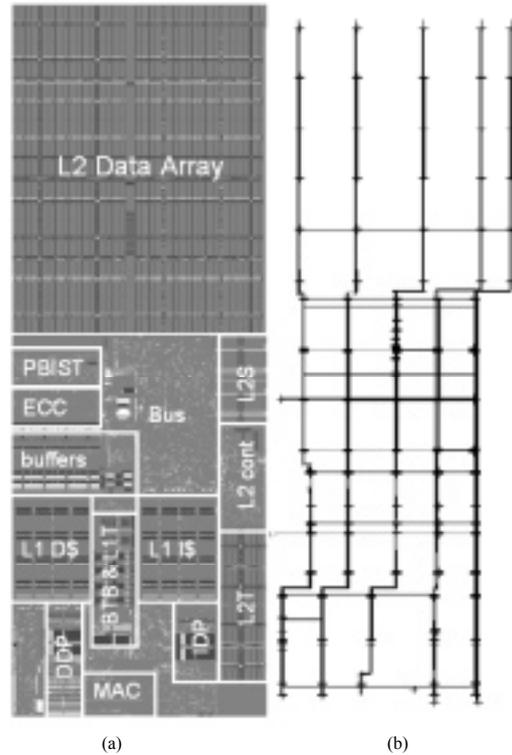


Figure 1. (a) Die micrograph and (b) clock tree structure.

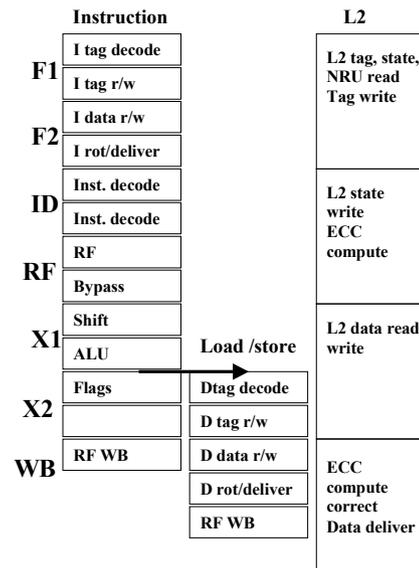


Figure 2. Instruction execution and L2 pipeline organization.

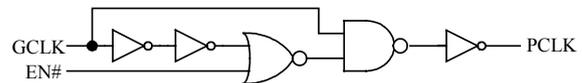


Figure 3. Pulse clock generating local clock buffer.

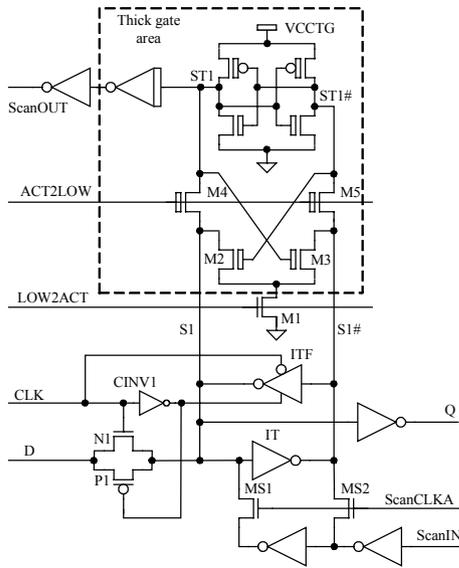


Figure 4. Pulse clocked latch with thick gate state retention.

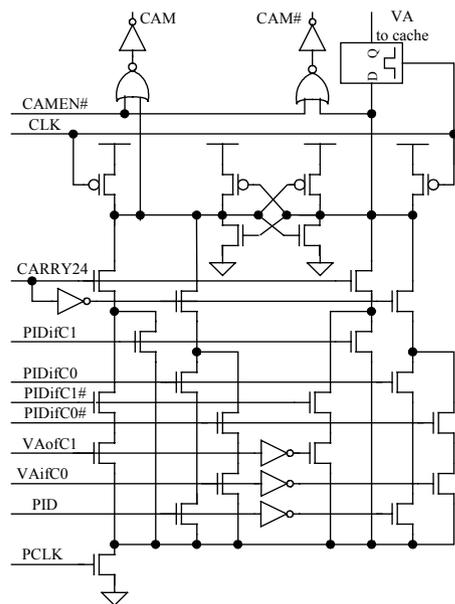


Figure 5. Pulse clocked latch with thick gate state retention.

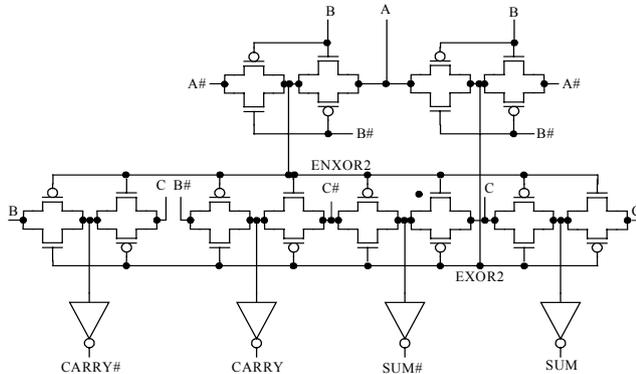


Figure 6. Static passgate implementation of the full adder used in the 3-to-2 reducer.

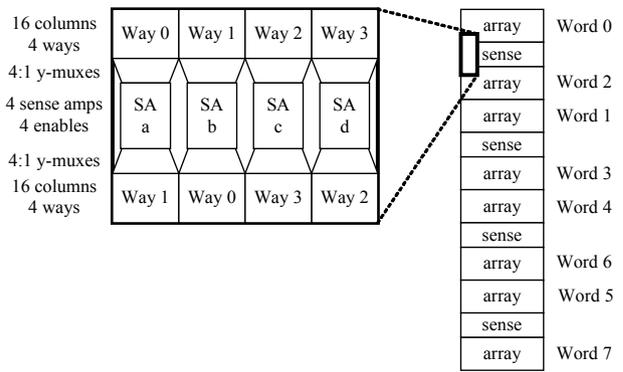


Figure 7. L1 sense amplifier organization to allow 256-b fill /evicts.

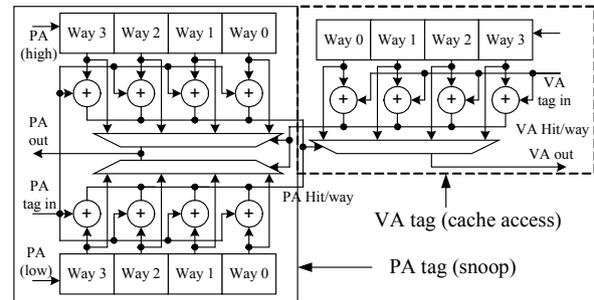


Figure 8. L1 Cache tag organization.

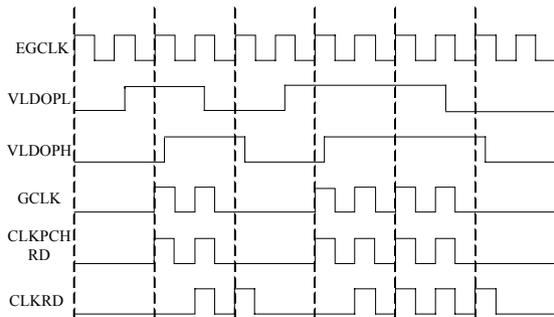


Figure 9. L2 clocking. Clocks are generated only when valid operations are ready, with a single and double read shown.

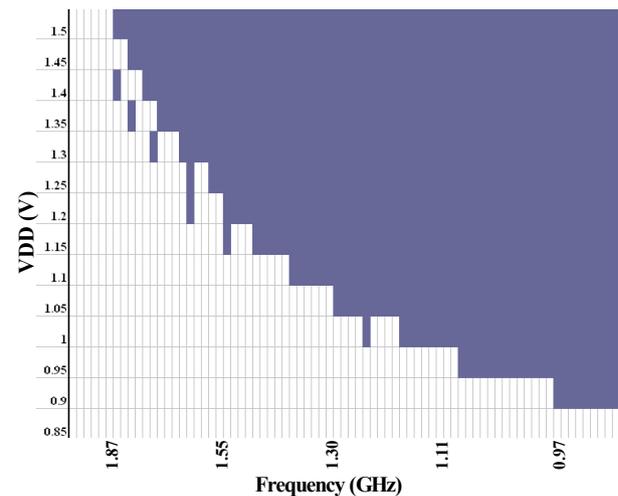


Figure 10. Shmoo plot of measured frequency vs. voltage.