

3.1 A 22nm IA Multi-CPU and GPU System-on-Chip

Satish Damaraju¹, Varghese George¹, Sanjeev Jahagirdar¹,
Tanveer Khondker¹, Robert Milstrey¹, Sanjib Sarkar¹, Scott Siers¹,
Israel Stoloro², Arun Subbiah¹

¹Intel, Folsom, CA

²Intel, Haifa, Israel

This paper describes the 22nm Intel® processor codenamed Ivy Bridge that integrates up to four high-performance Intel Architecture (IA) cores, a power/performance optimized graphics/media processing unit (GPU), as well as memory, PCIe, and display controllers in the same die. The Ivy Bridge architecture is derived from the second-generation Intel Core™ processor [2], seen in Fig. 3.1.1. The processor has about 1.4 billion transistors in about 160mm² in its largest incarnation. It introduces several enhancements in power, performance, and features over its predecessor. The IA core adds a pipelined divider, a next page prefetcher, additions to the ISA for 16b floating point conversion, fast string moves, and fast access of the FS/GS base registers. The Graphics/Media block provides significantly improved performance along with DX11 API support. Further, display capabilities have been augmented with a third independent display pipeline. The on-die power management control unit (PCU) and its associated firmware have added several power and thermal optimizations to improve performance and yield within the existing platform power envelopes, as well as to improve idle power relative to its predecessor. Power gates are distributed throughout the cores, enabling the PCU to independently either reduce voltage to a state retention voltage, or turn off voltage to a given core, depending on the current core usage conditions. The processor also implements power gating for portions of the DDR I/O buffers, reducing CPU power consumption when memory is in “self-refresh” mode. To optimize core sleep time, Ivy Bridge incorporates smart interrupt routing logic, which sends interrupts to active cores. The PCU also analyzes the processor’s inherent voltage-frequency dependency to determine the optimum operating voltages across the entire dynamic range of operation (Fig. 3.1.2). The processor uses die temperature to estimate power and energy consumption in order to maximize performance.

The Ivy Bridge processor is the first microprocessor designed on Intel’s latest 22nm process technology. This technology is the first process to use 3D tri-gate transistors for improved channel control and low-voltage operation. In addition to delivering a ~2× improvement in transistor density, the process allows a much lower V_t to deliver similar speed at 200mV lower operating voltage, or up to 37% faster speed at low operating voltages, as shown in Fig. 3.1.3. The process also offers three transistor types with varying speed/leakage tradeoffs, with fast, nominal leakage devices, medium-speed, “quarter-leakage” devices, and slower, “tenth-leakage” devices. High-frequency areas of the processor use ~30% nominal and ~70% quarter-leakage devices, while low speed areas use ~25% quarter-leakage and ~75% tenth-leakage devices. While the new process provides improved transistors, it did pose some design challenges. First, the 3D “fin” defines the transistor width and inherently limits small changes in transistor width. Losing fine granularity adjustment of transistor width in turn has an impact on device sizing, particularly for small devices. To overcome this problem and improve scaling, small devices are deliberately over scaled (that is, converting any portion of a fin after scaling to a full fin). This yields much better timing with little power impact. When possible during design closure, tools exchange excessive timing margin for the lower power options of either downsizing the device or converting it to a lower leakage variant.

Design convergence duration was reduced relative to previous generations by utilizing a holistic design migration methodology. This migration accounted for schematic and layout process migration of the register file and custom design blocks, as well as for clean-up of timing, noise, and design rule violations. This resulted in about a 15% reduction in execution duration relative to the preceding project. Further, in order to ensure adequate timing margin for high-speed/high-voltage operation, an additional 7% was removed from the timing targets.

Pre-silicon power analysis was performed via design-methodology-driven modeling including modeling the new lower leakage devices, systematic algorithmic solutions (clock, sequential, random logic sizing), binning strategies for fixed

frequency domains, focused tests, and a detailed accounting of the power consumed on I/O and PLL power planes, in addition to the traditional digital estimations. Post-silicon results showed good correlation to this analysis, with measured C_{dyn} within 5% of pre-silicon estimates.

Minimum operating voltage (V_{ccmin}) is one of the key parameters impacting battery life. Several circuit techniques were developed to optimize V_{ccmin} for the L3 cache and register files (RF) with minimal area/power increase. Specifically, write/read-assist circuit features have been incorporated to minimize bitcell passgate-to-pullup contention during write, while maintaining cell stability [4]. The write-assist circuitry used in GPU RFs is illustrated in Fig. 4.1.4. Post-silicon knobs to control pulses on inputs C0-C4 have been added to determine the optimal CVCC node drop and duration for achieving the best V_{ccmin} results. Similar techniques are employed in other parts of the processor’s cache design. Since cache-specific V_{ccmin} failures are spread across the cache array randomly, a Dynamic Cache Shrink architecture feature has been introduced to further improve V_{ccmin} , whereby L3 cache size is dynamically reduced from 2MB per slice to 256KB using the PCU. In addition, randomly failing bits in the L3 cache are replaced by a per-die post-silicon redundancy flow. These techniques together resulted in >250mV V_{ccmin} benefit for the L3 cache.

In the Ivy Bridge processor, DDR3 performance has been extended from 1333MHz to 1600MHz with up to two DIMMs per channel, as the design supports training routines during boot for input reference voltage, on-die termination resistances, and buffer strength. The buffers also support the 1.35V DDR3L specification, which enables lower platform power consumption compared to DDR3. In addition, when the DDR I/O power is partially gated, the DDR I/O buffer power consumption is about 50mW. The processor’s PCIe receiver consists of a high-bandwidth CTLE with automatic gain control, integration-DFE, and a Mueller-Müller CDR circuit [3][1]. PCIe lane power consumption is about 10mW in L1 and 100mW in L0. The design is capable of tolerating wide process variations and contains on-die dfx features to test for jitter tolerances and voltage margins.

The basic clocking structure in the processor is derived from its predecessor’s implementation [2]. The clocking is supported by 14 PLLs and an embedded clock compensation scheme, which achieves the required skew across different power planes and frequency domains. The stringent jitter/bandwidth and fixed-frequency requirements of serial I/Os are supported by LC PLLs. For variable-frequency PLLs, in order to mitigate the high variability of P1270 analog device parameters, two key features were introduced: auto-VCO frequency folding to optimize operating conditions on non-tunable VCOs, and a tunable SBPLL design which eliminates VCO F_{max} excursion across target frequencies. Fig. 3.1.5 demonstrates both concepts, where, for example, target frequencies in the 1-2GHz range are achieved by running the VCO at 2x frequency (2-4GHz) and applying a post divider (“folded”) so that the VCO runs at twice the target frequency for optimal VCO performance. The optimal frequency-folding regions are determined by an on-die state machine that hunts for the closed-loop VCO F_{max} during PLL training. A simplified tunable VCO circuit architecture is shown in Fig. 3.1.6. In VCO oscillator instantiations, tunable discrete loads are added or removed from the circuit to optimize VCO target frequency. The red curve (label ‘100’) in Fig. 3.1.5 indicates the VCO tuning target, with the family of target frequencies associated with added/removed loads shown in black. Final VCO targeting is accomplished on a per-PLL basis by configuring the appropriate loads to bring the VCO as close as possible to the target VCO frequency.

References:

- [1] F. Spagna, et al., “A 78mW 11.8Gb/s serial link transceiver with adaptive RX equalization and baud-rate CDR in 32nm CMOS,” *ISSCC Dig. Tech. Papers*, pp. 366-367, 2010.
- [2] M. Yuffe, et al., “A Fully Integrated Multi-CPU, GPU and Memory,” *ISSCC Dig. Tech. Papers*, pp. 264-266, 2011.
- [3] K. Mueller, et al., “Timing Recovery in Digital Synchronous Data Receivers”. *IEEE Trans. on Communications*, vol. 24, no. 5, pp. 516-531, 1976.
- [4] M. Khellah, et al., “Process, Temperature, and Supply-Noise Tolerant 45 nm Dense Cache Arrays With Diffusion-Notch-Free (DNF) 6T SRAM Cells and Dynamic Multi-Vcc” *IEEE J. Solid-State Circuits*, vol. 44, no. 4, pp. 1199-1208, 2009.

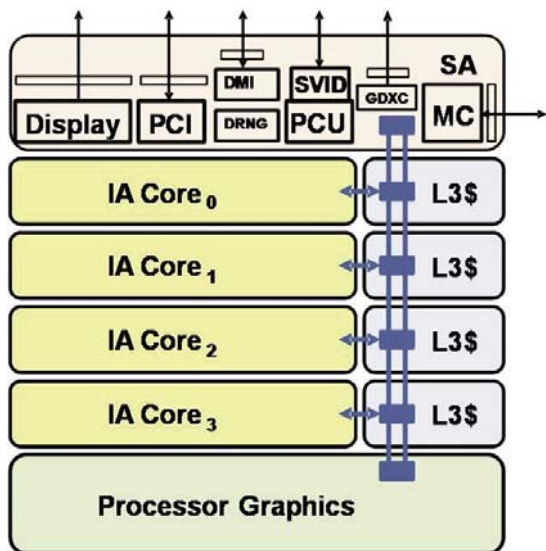


Figure 3.1.1: Ivy Bridge processor block diagram.

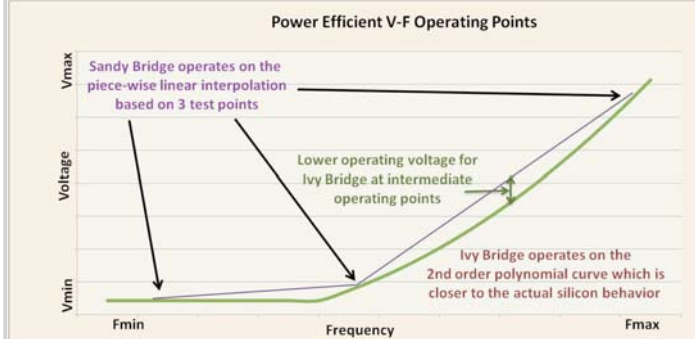


Figure 3.1.2: Second order model for transistor behavior.

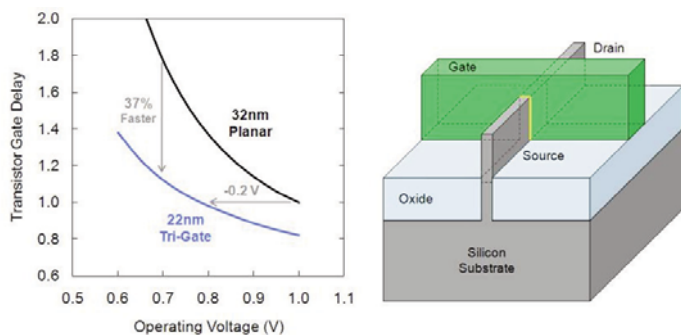


Figure 3.1.3: Tri-gate's low voltage performance as a function of voltage.

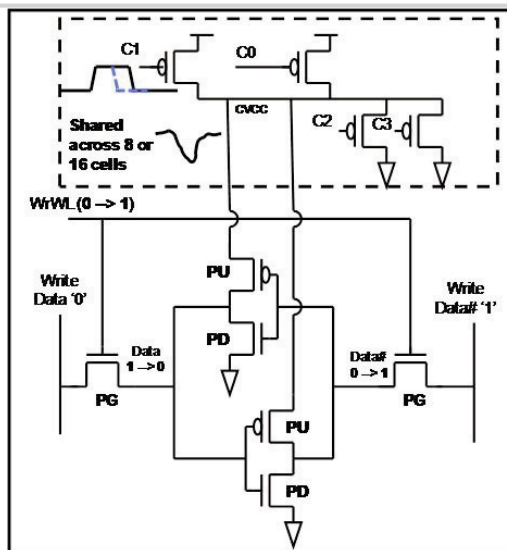


Figure 3.1.4: RF Write Assist Circuit to improve V_{ccmin} .

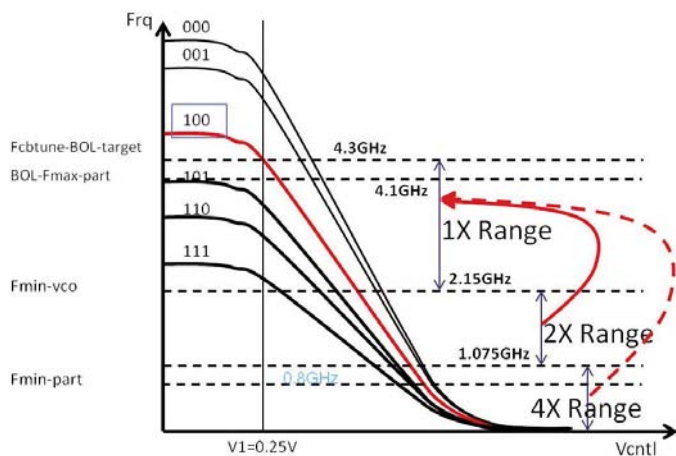


Figure 3.1.5: PLL frequency folding and tuning illustration.

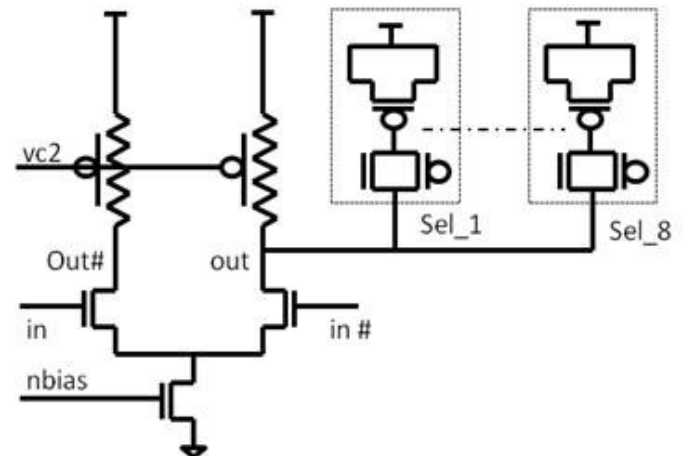


Figure 3.1.6: Simplified tunable VCO circuit.