

INTERNATIONAL TECHNOLOGY ROADMAP FOR SEMICONDUCTORS

2001 EDITION

SYSTEM DRIVERS

TABLE OF CONTENTS

Scope	1
High-Volume Custom – Microprocessor (MPU).....	1
Analog/Mixed-Signal (AMS)	1
System-On-Chip (SOC)	2
Market Drivers.....	3
MPU System Driver	4
MPU Evolution	6
MPU Challenges	7
Mixed-Signal System Driver	8
Low-Noise Amplifier (LNA).....	8
Voltage-Controlled Oscillator (VCO)	9
Power Amplifier (PA).....	10
Analog-to-Digital Converter (ADC)	10
Mixed-Signal Evolution.....	11
Mixed-Signal Challenges.....	13
SOC System Driver	14
SOC-Multi-Technology.....	14
SOC-High Performance	15
SOC-Low Cost, Low Power.....	16
SOC Trends.....	17
SOC Challenges	21

LIST OF FIGURES

Figure 9 Recent ADC Performance Needs for Important Product Classes	12
Figure 10 First Integration of Technologies on SOC with Standard CMOS Process	15
Figure 11 Total Chip Power Trend for PDA Application.....	18
Figure 12 Power Gap Effect on Chip Composition	20
Figure 13 New and Reused Logic Content versus Memory Content with Constant Die Size and Insufficient (42% Per Node) Design Productivity Growth	22
Figure 14a Evolution of Maximum Logic Content with Different Rates of Design Productivity Improvement.....	22
Figure 14b 100% Productivity Improvement per Node Will Preserve Designer Freedom at the End of the ITRS Forecast Period,	22

LIST OF TABLES

Table 8 Major Product Market Segments and Impact On System Drivers.....	3
Table 9 Projected Mixed-signal Figures of Merit for Four Circuit Types	11
Table 10 System Functional Requirements for the PDA SOC-LP Driver.....	17
Table 11 Low Operating Power (LOP) and Low Standby Power (LSTP) Device and Process Attributes	19
Table 12 Power Management Gap	19

SYSTEM DRIVERS

SCOPE

Future semiconductor manufacturing and design technology capability is developed in response to economic drivers within the worldwide semiconductor industry. The ITRS must understand how technology requirements arise for product classes whose business and retooling cycles drive the semiconductor sector. Previous ITRS editions have focused on microprocessor (MPU), dynamic random-access memory (DRAM), and application-specific integrated circuit (ASIC) product classes, with some mention of system-on-chip (SOC) and analog/mixed-signal circuits. The unstated assumption was that technological advances need only be straight ahead and “linear”, and would be deployed in all semiconductor products. For this reason, specifics of the product classes (e.g., MPU or ASIC) were not required. Today, introduction of new technology solutions is increasingly application-driven, with products for different markets making use of different combinations of technologies at different times. General-purpose digital microprocessors for personal computers are being joined as drivers by mixed-signal systems for wireless communication and embedded applications. Wall-plugged servers are being replaced by battery-powered mobile devices. In-house, single-source chip designs are being supplanted by system-on-chip and system-in-package designs that incorporate building blocks from multiple sources.

The purpose of the 2001 ITRS System Drivers Chapter is to update and more clearly define the set of system drivers that has been used in previous ITRS editions. Together with the Overall Roadmap Technology Characteristics, the System Drivers Chapter seeks to provide a framework and motivation for specific technology requirements that is consistent across the respective ITRS technology areas and the 15-year span of the ITRS. The main contribution of the Chapter consists of quantitative, internally self-consistent models of the system drivers that support extrapolation into future technologies and adapt more smoothly to future technology developments. We focus on three system drivers: high-volume custom – microprocessor (MPU), analog/mixed-signal (AMS), and system-on-chip (SOC). A fourth system driver, high-volume custom – memory (DRAM), is not discussed due to its well-understood commodity nature.

HIGH-VOLUME CUSTOM – MICROPROCESSOR (MPU)

In high-volume custom designs, performance and manufacturing cost issues outweigh design or other non-recurring engineering (NRE) cost issues, primarily because of the large profits that these chips can potentially produce. These large profits result from very large sales volumes. Large volumes alone are neither necessary nor sufficient to warrant the custom design style, special process engineering and equipment, etc. often associated with such parts; the key is that the expected return on the combined NRE and manufacturing investment must be positive. Within the high-volume custom arena, the three dominant classes today are MPUs, memory¹ and reprogrammable (e.g., field-programmable gate array (FPGA)); the first of these is defined in the first section of this chapter, while the latter two are discussed only as “implementation fabrics” available to the SOC system driver. MPUs use the most aggressive design styles and manufacturing technologies to achieve their goals. It is for these high-volume parts that changes to the manufacturing flow are made, new design styles and supporting tools are created (the large revenue streams can pay for new tool creation), and subtle circuits issues are uncovered (not all risks taken by designers work out). Thus, while MPUs (and high-volume custom designs in general) are extremely labor-intensive, they create new technology and automation methods (in both design and fabrication) that are leveraged by the entire industry.

ANALOG/MIXED-SIGNAL (AMS)

AMS chips at least partially deal with input signals whose precise values matter. This broad class includes RF, analog, analog to digital and digital to analog converter, and, more recently, a large number of mixed-signal chips

¹ *Memory is a special class of high-volume custom design because of the very high replication rate of the basic memory cells and supporting circuits. Since these cells are repeated millions of times on a chip, and millions of chips are sold, the amount of custom design for these parts is extraordinary. This has led to separate fabrication lines for DRAM devices, with some of the most careful circuit engineering needed to ensure correct operation.*

2 System Drivers

where at least part of the chip design needs to measure signals with high precision. These chips have very different design and process technology demands than digital circuits. While technology scaling is always desirable for digital circuits due to reduced power, area and delay, it is not necessarily helpful for analog circuits since dealing with precision requirements or signals from a fixed voltage range is more difficult with scaled voltage supplies. Thus, scaling of analog circuits into new technologies is a difficult challenge. In general, AMS circuits (e.g., RF and embedded passives) and process technologies (e.g., silicon-germanium) present severe challenges to cost-effective CMOS integration.

The need for precision also affects tool requirements for analog design. Digital circuit design creates a set of rules that allow logic gates to function correctly: as long as these rules are followed, precise calculation of exact signal values is not needed. Analog designers, on the other hand, must be concerned with a number of “second-order effects” to obtain the required precision. Relevant issues include coupling (capacitance, inductance and substrate) and asymmetries (local variation of supplies, as well as implantation, alignment, etching, and other fabrication effects). Analysis tools for these issues are mostly in place but require expert users; synthesis tools are at best preliminary. Manufacturing test for AMS circuits is essentially unsolved.

SYSTEM-ON-CHIP (SOC)

SOC is a yet-evolving *product class and design style* that integrates pieces of technology from other system driver classes (e.g., MPU, memory, AMS, and reprogrammable) into a wide range of high-complexity, high-value semiconductor products. Manufacturing and design technologies for SOC are typically developed originally for high-volume custom drivers. The SOC driver class most closely resembles, and is evolved most directly from, the ASIC category since reduced design costs and higher levels of system integration are its principal goals.² The primary difference from ASIC is that in SOC design, the goal is to maximize *reuse* of existing blocks or “cores” – i.e., minimize the amount of the chip that is newly or directly created. Reused blocks in SOC include analog and high-volume custom cores, as well as blocks of *software* technology. A key challenge is to incent, create and maintain reusable blocks or cores so that they are available to SOC designers.³ The utility of SOC also depends on validation for reuse-based SOC designs being easier than for equivalent “from-scratch” designs.

SOC represents a confluence of previous product classes in several ways. As noted above, SOCs integrate building blocks from the other system driver classes, and are subsuming the ASIC category. The quality gap between full-custom and ASIC/SOC is diminishing: (i) the 2001 ITRS models overall ASIC and MPU logic densities as being equal; and (ii) “custom quality on an ASIC schedule” is increasingly achieved by on-the-fly (“liquid”) or tuning-based standard-cell methodologies. Finally, MPUs are evolving into SOCs: (i) MPUs are increasingly designed as cores to be included in SOCs, and (ii) MPUs are themselves designed as SOCs to improve reuse and design productivity (as discussed below, the 2001 ITRS MPU model has multiple processing cores and resembles an SOC in organization⁴). The most basic SOC challenge is presented by implementation productivity and manufacturing cost, which require greater reuse as well as platform-based design, silicon implementation regularity, or other novel circuit and system architecture paradigms. Another challenge is the heterogeneous integration of components from multiple implementation *fabrics* (e.g., RF, reprogrammable, MEMS, optoelectronic, and software).

² Most digital designs today are considered to be ASICs. ASIC connotes both a business model (with particular “handoff” from design team to ASIC foundry) and a design methodology (where the chip designer works predominantly at the functional level, coding the design at Verilog/VHDL or higher level description languages and invoking automatic logic synthesis and place-and-route with a standard-cell methodology). For economic reasons, custom functions are rarely created; reducing design cost and design risk is paramount. ASIC design is characterized by relatively conservative design methods and design goals (cf. differences in clock frequency and layout density between MPU and ASIC in previous ITRS editions) but aggressive use of technology, since moving to a scaled technology is a cheap way of achieving a better (smaller, lower power, and faster) part with little design risk (cf. convergence of MPU and ASIC process geometries in previous ITRS editions). Since the latter half of the 1990s, ASICs have been converging with SOCs in terms of content, process technology, and design methodology.

³ For example, reusable cores might require characterization of specific noise or power attributes (“field of use”, or “assumed design context”) that are not normally specified. Creation of an IC design artifact for reuse by others is substantially more difficult (by factors estimated at between 2X and 5X) than creation for one-time use.

⁴ The corresponding ASIC and structured-custom MPU design methodologies are also converging to a common “hierarchical ASIC/SOC” methodology. This is accelerated by customer-owned tooling business models on the ASIC side, and by tool limitations faced by both methodologies.

MARKET DRIVERS

Table 8 contrasts semiconductor product markets according to such factors as manufacturing volume, die size, integration heterogeneity, system complexity, and time-to-market. The influence on each class of system driver (with high-volume custom treated as a whole) is noted.⁵

Table 8 Major Product Market Segments and Impact On System Drivers

<i>MARKET DRIVERS</i>	<i>ASIC/SOC</i>	<i>ANALOG/MS</i>	<i>HIGH-VOLUME CUSTOM</i>
<i>I. Portable and Wireless</i>			
1. Size/weight ratio: peak in 2002 2. Battery life: peak in 2002 3. Function: 2× / 2 years 4. Time-to-market: ASAP 5. Time-in-market: decreasing	Low power paramount Need SOC integration (DSP, MPU, I/O cores, etc.)	Migrating on-chip for voice processing, RF A/D sampling, etc.	Specialized cores to optimize processing per microwatt.
<i>II. Broadband</i>			
1. Bandwidth: 2× / 9 months 2. Function: 20%/yr increase 3. Deployment/OperCost: flat 4. Reliability: asymptotic 99.999% 5. Time-in-market: long 6. Power: W/m ³ of system	Large gate counts. High reliability. Primarily SOC.	Migrating on-chip for signal recovery, RF A/D sampling, etc.	MPU cores and some specialized functions.
<i>III. Internet Switching</i>			
1. Bandwidth: 4× / 3-4 yrs. 2. Reliability 3. Time-to-market: ASAP 4. Power: W/m ³ of system	Large gate counts. High reliability. Primarily SOC, with more reprogrammability to accommodate custom functions.	Minimal on-chip analog. Migrating on-chip for I/O circuitry. MEMS for optical switching.	MPU cores, FPGA cores and some specialized functions.

⁵ Note that the driver classes are most clearly distinguished according to cost, time-to-market, and production volume. System cost is equal to Manufacturing cost + Design cost. Manufacturing cost breaks down further into non-recurring engineering (NRE) cost (masks, tools, etc.) and silicon cost (raw wafers + processing + test). The total system cost is correlated with function, #I/Os, package cost, power and speed. Hence, distinctions made in the 1999 ITRS between SOC-C (“cost-driven”) and SOC-P (“performance-driven”) simply reflect a cost continuum. Different regions of the (Manufacturing Volume, Time To Market, System Complexity) space are best served by ASIC, FPGA or HVC implementation fabrics, and by SOC or system-in-package (SIP) integration. This partitioning is continually evolving.

4 System Drivers

Table 8 Major Product Market Segments and Impact On System Drivers (continued)

MARKET DRIVERS	ASIC/SOC	ANALOG/MS	HIGH-VOLUME CUSTOM
<i>IV. Mass Storage</i>			
1. Density: 60% increase / yr 2. Speed: 2× by 2005 3. Form factor: shift toward 2.5"	High-speed front-end for storage systems. Primarily ASSP. Shift toward large FPGA and COT, away from ASIC costs and design flows	Increased requirement for higher precision position measurement, "inertia knowledgeable" actuator / power controllers integrated on-chip. MEMS on R/W head for sensing.	Demand for high-speed hardware for, e.g., "lookahead" in DB search, MPU instruction prefetch, data compression, S/N monitoring, failure prediction.
<i>V. Consumer</i>			
1. Cost: strong downward pressure 2. Time-to-market: <12 mos 3. Function: high novelty 4. Form factor 5. Durability / safety 6. Conservation / ecology	High-end products only. Reprogramability possible. Mainly ASSP; more SOC for high-end digital with cores for 3D graphics, parallel proc, RTOS kernel, MPU-MMU-DSP, voice synthesis and recognition, etc.	Increased integration for voice, visual, tactile, physical measurement (e.g., sensor networks). CCD or CMOS sensing for cameras.	For "long-life" mature products only. Decrease in long design cycles, and in use of high-cost non-prepackaged functions and design flows.
<i>VI. Computer</i>			
1. Speed: 2× / 2 yrs 2. Memory density: 2× / 2 yrs 3. Power: flat to decreasing, driven by cost and W/m ³ 4. Form factor: shrinking size 5. Reliability	Large gate counts. High speed. Drives demand for digital functionality. Primarily SOC integration of custom off-the-shelf MPU and I/O cores.	Minimal on-chip analog. Simple A/D and D/A. Video i/f for automated camera monitoring, video conferencing. Integrated high-speed A/D, D/A for monitoring, instrumentation, range-speed-position resolution.	MPU cores and some specialized functions. Increased industry partnerships on common designs to reduce development costs (requires data sharing and reuse across multiple design systems).
<i>VII. Automotive</i>			
1. Functionality 2. Ruggedness (external environment, noise) 3. Reliability and safety 4. Cost	Mainly entertainment systems. Mainly ASSP, but increasing SOC for high end using standard hardware platforms with RTOS kernel, embedded software.	Cost-driven on-chip ADC for sensor signals. Signal processing shifting to DSP for voice, visual. Physical measurement ("communicating sensors" for proximity, motion, positioning). MEMS for sensors.	

MPU SYSTEM DRIVER

The 2001 ITRS microprocessor (MPU) driver reflects general-purpose instruction-set architectures (ISAs) that are found standalone in desktop and server systems, and embedded as cores in SOC applications. MPUs are part of the high-volume custom segment that drives the semiconductor industry with respect to integration density and design complexity, power-speed performance envelope, large-team design process efficiency, test and verification, power management, and packaged system cost. The MPU system driver is subject to market forces that have historically led to (i) emergence of standard architecture platforms and multiple generations of derivatives, (ii) strong price sensitivities in the marketplace, and (iii) extremely high production volumes and manufacturing cost awareness. Key elements of the MPU driver model are as follows (studies in this Chapter can be run in the [GTX tool](#); MPU content is provided [in the following study](#).

(1) *Two types of MPU*—Historically, there have been two types of MPU: *cost-performance* (CP), reflecting “desktop”, and *high-performance* (HP), reflecting “server”. The CP versus HP taxonomy is retained in the 2001 ITRS—with each type parameterized with respect to volume production in a given node—largely to permit continuity with previous ITRS MPU models. Future MPU models will likely require a merged desktop-server category (this distinction is already blurred today) and a *mobile* category (essentially a low-power, high-performance SOC).

(2) *Constant die area*—Die areas are constant (140mm² for CP, 310mm² for HP) over the course of the roadmap, and are broken down into logic, memory, and integration overhead. Integration overhead reflects the presence of whitespace for interblock channels, floorplan packing losses, and potentially growing tradeoff of layout density for design turnaround time. The core message, in contrast to previous ITRS models, is that power and cost are strong limiters of die size. To first order, additional logic content would not be efficiently usable due to package power limits, and additional memory content (e.g., larger caches, more levels of memory hierarchy integrated on-chip) would not be cost-effective beyond a certain point.⁶

(3) *Multi-core organization*—MPU logic content reflects multiple processing units on-chip starting from the 130nm node. This integrates several factors: (i) organization of recent and planned commercial MPU products (both server and desktop); (ii) increasing need to reuse verification and logic design, as well as standard ISAs; (iii) ISA “augmentations” in successive generations (e.g., x86, MMX and EPIC) with likely continuations for encryption, graphics and multimedia, etc.; (iv) the need to enable flexible management of power at the architecture, OS and application levels via SOC-like integration of less efficient, general-purpose processor cores with more efficient, special-purpose “helper engines”⁷; (v) the limited size of processor cores (the estimate of a *constant* 20-25 million transistors per core⁸ is a conservative upper bound with respect to recent trends); and (vi) the convergence of SOC and MPU design methodologies due to design productivity needs. The number of cores on chip is projected to double with each successive technology node.

(4) *Memory content*—The MPU memory content is initially 512KBytes (512 × 1024 × 9 bits) of SRAM for CP and 2MBytes for HP in the 180nm node. Memory content, like logic content, is projected to double with each successive technology node, not with respect to absolute time intervals (e.g., every 18 months).^{9,10}

(5) *Layout density*—Due to their high levels of system complexity and production volume, MPUs are the driver for improved layout density.¹¹ Thus, MPU driver sets the layout densities, and hence the transistor counts and chip sizes, stated in the Overall Roadmap Technology Characteristics. The logic and SRAM layout densities in the 2001 ITRS ORTCs are analogous to the DRAM “A-factor”, and have been calibrated to recent MPU products. Logic layout densities reflect average standard-cell gate layouts of approximately 320F², where F is the minimum feature

⁶ *Multi-core organization (see Footnote 16) and associated power efficiencies may permit slight growth in die size, but the message is still that die areas are flattening out.*

⁷ *A “helper engine” is a form of “processing core” for graphics, encryption, signal processing, etc. The trend is toward architectures that contain more special-purpose, and less general-purpose, logic.*

⁸ *The CP core has 20 million transistors, and the HP core has 25 million transistors. The difference allows for more aggressive microarchitectural enhancements (trace caching, various prediction mechanisms, etc.) and other performance support.*

⁹ *The doubling of logic and memory content with each technology node, rather than with each 18- or 24-month time interval, is due to essentially constant layout densities for logic and SRAM, as well as conformance with other parts of the ITRS.*

Specifically, the ITRS remains planar CMOS-centric, with little or no acknowledgment of dual-gate FET, FinFET, etc. yet incorporated into the roadmap except as “research devices”. Adoption of such novel device architectures would allow improvements of layout densities beyond what is afforded by scaling alone.

¹⁰ *Deviation from the given model will likely occur around the 90nm node with adoption of denser embedded memories (eDRAM). Adoption of eDRAM, and integrated on-chip L3 cache, will respectively increase the on-chip memory density and memory transistor count by factors of approximately 3 from the given values. While this will significantly boost transistor counts, it is not projected to significantly affect the chip size or total chip power roadmap. Adoption of eDRAM will also depend strongly on compatibility with logic processes (notably the limited process window that arises from scaling of oxide thickness), the size and partitioning of memory within the individual product architecture, and density-performance-cost sensitivities.*

¹¹ *ASIC/SOC and MPU system driver products have access to similar processes, as forecast since the 1999 ITRS. This reflects emergence of pure-play foundry models, and means that fabric layout densities (SRAM, logic) are the same for SOC and MPU. However, MPUs drive high density and high performance, while SOCs drive high integration, low cost, and low power.*

6 System Drivers

size of the technology node.¹² As noted above, the logic layout density may improve significantly with the advent of novel devices. SRAM layout densities reflect use of a 6-transistor bitcell (the fitted expression for area per bitcell in units of $F^2 = 223.19F(\text{um}) + 97.74$) in MPUs, with 60% area overhead for peripheral circuitry.

(6) *Maximum on-chip (global) clock frequency*—MPUs also drive maximum on-chip clock frequencies in the Overall Roadmap Technology Characteristics; these in turn drive various aspects of the Interconnect, PIDS, FEP and Test roadmaps. The MPU maximum on-chip clock frequency has historically increased by a factor of 2 per generation. Of this, approximately 1.4× has been from device scaling (which runs into t_{ox} and other limits); the other 1.4× has been from reduction in number of logic stages in a pipeline stage (e.g., equivalent of 32 fanout-of-4 inverter (FO4 INV) delays¹³ at 180nm, and 26 FO4 INV delays at 130nm). There are several reasons why this historical trend will not continue: (i) well-formed clock pulses cannot be generated with period below 6-8 FO4 INV delays; (ii) there is increased overhead (diminishing returns) in pipelining (2–3 FO4 INV delays per flip-flop, 1–1.5 FO4 INV delays per pulse-mode latch); and (iii) around 14–16 FO4 INV delays is a practical lower limit for clock period given the latencies of L1 cache access, 64-bit integer addition, etc. The 2001 ITRS MPU model continues the historic rate of advance for maximum on-chip global clock frequencies, but flattens the clock period at 16 FO4 INV delays during the 90nm node ([a plot of historical MPU clock period data is provided](#)). The message is that from the 90nm node on, clock frequencies will advance only with device performance in the absence of novel circuit and architectural approaches.¹⁴

MPU EVOLUTION

An emerging “centralized processing” context integrates (i) centralized computing servers that provide high-performance computing via traditional MPUs (this driver), and (ii) *interface remedial processors* that provide power-efficient basic computing via, e.g., SOC integration of RF, analog/mixed-signal, and digital functions within a wireless handheld multimedia platform (see the low-power SOC PDA model, below). Key contexts for the future evolution of the traditional MPU are with respect to design productivity, power management, multi-core organization, I/O bandwidth, and circuit and process technology.

Design productivity—The complexity and cost of design and verification of MPU products has rapidly increased to the point where thousands of engineer-years (and a design team of hundreds) are devoted to a single design, yet processors reach market with hundreds of bugs.

Power management—Power dissipation limits of cost-effective packaging (estimated to reach 50 W/cm² for forced-air cooling by the end of the ITRS) cannot continue to support high supply voltages (historically scaling at 0.85× per generation instead of 0.7× ideal scaling) and frequencies (historically scaling by 2× per generation instead of 1.4× ideal scaling).¹⁵ Past clock frequency trends in the MPU system driver have been interpreted as future CMOS device performance (switching speed) requirements that lead to large off-currents and extremely thin gate oxides, as specified in the PIDS Chapter. Given such devices, MPUs that simply continue existing circuit and architecture

¹² A 2-input NAND gate is assumed to lay out in an 8x4 standard cell, where the dimensions are in units of contacted local metal pitch ($MP = 3.16 \times F$). In other words, the average gate occupies $32 \times (3.16)^2 = 320F^2$. For both semi-custom (ASIC/SOC) and full-custom (MPU) design methodologies, an overhead of 100% is assumed.

¹³ A FO4 INV delay is defined to be the delay of an inverter driving a load equal to 4 times its own input capacitance (with no local interconnect). This is equivalent to roughly 14 times the CV/I device delay metric that is used in the PIDS Chapter to track device performance. An explanation of the FO4 INV delay model used in the 2001 ITRS is provided in [supplemental material](#).

¹⁴ Unlike previous ITRS clock frequency models (e.g., Fisher/Nesbitt 1999), the 2001 model does not have any local or global interconnect component in its prototypical “critical path”. This is because local interconnect delays are negligible, and scale with device performance. Furthermore, buffered global interconnect does not contribute to the minimum clock period since long global interconnects are pipelined (cf. Intel Pentium-4 and Compaq Alpha 21264) – i.e., the clock frequency is determined primarily by the time needed to complete local computation loops, not by the time needed for global communication. Pipelining of global interconnects will become standard as the number of clock cycles required to signal cross-chip continues to increase beyond 1. “Marketing” emphases for MPUs necessarily shift from “frequency” to “throughput” or “utility”.

¹⁵ To maintain reasonable packaging cost, package pin counts and bump pitches for flip-chip are required to advance at a slower rate than integration densities (cf. the Assembly and Packaging Chapter). This increases pressure on design technology to manage larger wakeup and operational currents and larger supply voltage IR drops; power management problems are also passed to the architecture, OS and application levels of the system design.

techniques would exceed package power limits by factors of over 20× by the end of the ITRS; alternatively, MPU logic content and/or logic activity would need to decrease to match package constraints. Portable and low-power embedded contexts have more stringent power limits, and will encounter such obstacles earlier. Last, power efficiencies (e.g., GOps/mW) are up to four orders of magnitude greater for direct-mapped hardware than for general-purpose MPUs; this gap is increasing. As a result, traditional processing cores will face competition from application-specific or reconfigurable processing engines for space on future SOC-like MPUs.

Multi-core organization—In an MPU with multiple cores per die, the cores can be (i) smaller and faster to counter global interconnect scaling, and (ii) optimized for reuse across multiple applications and configurations. Multi-core architectures allow power savings as well as the use of redundancy to improve manufacturing yield.¹⁶ Organization of the MPU model also permits increasing amounts of the memory hierarchy on chip (consistent with processor-in-memory, or large on-chip eDRAM L3 starting in the 90nm generation). Higher memory content can, if only in a relatively trivial way, afford better “control” of leakage and total chip power. Evolutionary microarchitecture changes (superpipelining, superscalar, predictive methods) appear to be running out of steam. (“*Pollack’s Rule*” observes that in a given process technology, a new microarchitecture occupies 2–3× the area of the old (previous-generation) microarchitecture, while providing only 1.4–1.6× the performance.) Thus, more multithreading support will emerge for parallel processing, as well as more complex “hardwired” functions and/or specialized engines for networking, graphics, security, etc. Flexibility-efficiency tradeoff points shift away from general-purpose processing.

Input/output bandwidth—In MPU systems, I/O pins are mainly used to connect to memory, both high-level cache memory and main system memory. Increased processor performance has been pushing I/O bandwidth requirements. The highest-bandwidth port has traditionally been used for L2 or L3 cache, but recent designs are starting to integrate the memory controller on the processor die to reduce memory latency. These direct memory interfaces require more I/O bandwidth than the cache interface. In addition to the memory interface, many designs are replacing the system bus with high-speed point-to-point interfaces. These interfaces require much faster I/O design, exceeding Gbit/s rates. While serial links have achieved these rates for a while, integrating a large number of these I/O on a single chip is still challenging for design (each circuit must be very low power), test (need to have a tester that can run this fast) and packaging (packages must act as balanced transmission lines, including the connection to the chip and the board).

Circuit and process technology—Parametric yield (\$/wafer after bin-sorting) is severely threatened by the growing process variability implicit in feature size and device architecture roadmaps (Lithography and PIDS), including thinner and less reliable gate oxides, subwavelength optical lithography requiring aggressive reticle enhancement, and increased vulnerability to atomic-scale process variability (e.g., implant). This will require more intervention at the circuit and architecture design levels. Circuit design use of dynamic circuits, while attractive for performance in lower-frequency or clock-gated regimes, may be limited by noise margin and power dissipation concerns; less pass gate logic will be used due to body effect. Error-correction for single-event upset (SEU) in logic will increase, as will the use of redundancy and reconfigurability to compensate for yield loss. The need for power management will require a combination of techniques from several component technologies: (i) application-, OS- and architecture-level optimizations including parallelism and adaptive voltage and frequency scaling, (ii) process innovations including increased use of SOI, and (iii) circuit design techniques including the *simultaneous* use of multi- V_{th} , multi- V_{dd} , minimum-energy sizing under throughput constraints, and multi-domain clock gating and scheduling.

MPU CHALLENGES

The MPU driver strongly affects design and test technologies (distributed/collaborative design process, verification, at-speed test, tool capacity, power management), as well as device (off-current), lithography/FEP/interconnect (variability) and packaging (power dissipation and current delivery). The most daunting challenges are:

¹⁶ *Replication enables power savings through lowering of frequency and V_{dd} while maintaining throughput (e.g., two cores running at half the frequency and half the supply voltage will save a factor of 4 in CV^2f dynamic capacitive power, versus the “equivalent” single core). (Possibly, this could allow future increases in die size.) More generally, overheads of time-multiplexing of resources can be avoided, and the architecture and design focus can shift to better use of area than memory. Redundancy-based yield improvement occurs if, e.g., a die with $k-1$ instead of k functional cores is still useful.*

8 System Drivers

- *design and verification productivity* (e.g., total design cost, number of bug escapes) (Design),
- *power management and delivery* (e.g., GOps per mW) (Design, PIDS, Assembly and Packaging), and
- *parametric yield at volume production* (Lithography, PIDS, FEP, Design).

MIXED-SIGNAL SYSTEM DRIVER

In formulating an analog and mixed-signal (AMS) roadmap, simplification is necessary because there are many different circuits and architectures, and because the roadmap may be used by individuals not directly expert in mixed-signal design issues. We restrict our discussion to four basic analog circuits:

- Low-noise amplifier (LNA),
- Voltage-controlled oscillator (VCO),
- Power amplifier (PA), and
- Analog to digital converter (ADC).

The design and process technology used to build these four circuits will also determine the performance of many other mixed-signal circuits. Thus, the performance of these specific circuits, as described by figures of merit (FoMs), is a good basis for a mixed-signal roadmap.

The following discussion develops these FoMs in detail. By convention, all parameters (e.g., gain G) are given as absolute values instead of on a decibel scale. We also avoid preferences for specific solutions to given design problems. Indeed, we have sought to be as open as possible to different types of solutions since unexpected solutions have often helped to overcome barriers. (Competition, e.g., between alternative solutions, is a good driving force for all types of advances related to technology roadmapping.) Furthermore, we observe that a given type of circuit can have different requirements for different purposes; certain performance indicators might be contradictory for different applications.¹⁷ To avoid such situations, we adjust the figures of merit to a mainstream product. The economic regime of a mainstream product is usually highly competitive: it has a high production volume, and hence a high level of R&D investment by which its technology requirements can drive mixed-signal technology as a whole. The obvious mainstream product in this context is the *mobile phone*. Last, we evaluate the dependence of the FoMs on device parameters, so that circuit design requirements can lead to specific device and process technology specifications. Extrapolations are proposed that lead on the one hand to a significant advance of analog circuit performance and on the other hand to realistic and feasible technology advances. These parameters are given in the mixed-signal technology requirements table of the [PIDS](#) chapter.

LOW-NOISE AMPLIFIER (LNA)

Digital processing systems require interfaces to the analog world. Prominent examples for these interfaces are transmission media in wired or wireless communication. The LNA amplifies the input signal to a level which makes further signal processing insensitive to noise. The key performance issue for an LNA is to deliver the undistorted but amplified signal to downstream signal processing units without adding further noise.

LNA applications (GSM, CDMA, W-LAN, GPS, Bluetooth, etc.) operate in many frequency bands. The operating frequency and, in some cases, the operating bandwidth of the LNA will impact the maximum achievable performance; nonlinearity must also be considered to meet the specifications of many applications. These parameters must be included in the FoM. On the other hand, different systems are often not directly comparable and thus have diverging requirements. For example, very wide bandwidth is needed for high-performance wired applications, but this increases power consumption. Low power consumption is an important design attribute for low-bandwidth wireless applications. For wide-bandwidth systems, bandwidth may be more important than linearity

¹⁷ *Certain cases of application are omitted for the sake of simplicity, and arguments are given for the cases selected. In many cases, we have limited our considerations to CMOS since it is the prime technological driving force and in most cases the most important technology. Alternative solutions (especially other device families) and their relevance will be discussed for some cases, as well as at the end of this section.*

to describe the performance of an LNA. However, to avoid contradictory design constraints we focus on the *wireless* communication context.

The linearity of a low noise amplifier can be described by the output referenced third order intercept point ($OIP3 = G \times IIP3$ where G is the gain and $IIP3$ is the input referenced third order intercept point). A parameter determining the minimum signal that is correctly amplified by a LNA is directly given by the noise figure of the amplifier, NF . However, for consideration of the contribution of the amplifier to the total noise the value of $(NF-1)$ is a better measure, since it allows the ratio between the noise of the amplifier $N_{amplifier}$ and the noise already present at the input N_{input} to be directly evaluated. These two performance figures can be combined with the total power consumption P . The resulting figure of merit captures the dynamic range of an amplifier versus the necessary DC power. For roadmap purposes it is preferable to have available a performance measure that is independent of frequency and thus independent of the specific application. This can be achieved by assuming that the LNA is formed by a single amplification stage, so that the FoM scales linearly with operating frequency f . With these approximations and assumptions, a figure of merit (FoM_{LNA}) for LNAs is defined:

$$FoM_{LNA} = \frac{G \cdot IIP3 \cdot f}{(NF - 1) \cdot P} \quad (1)$$

Making further simplifying assumptions, and neglecting “design intelligence”, the evolution of the FoM with technology scaling can be extrapolated [1]¹⁸. Future trends of relevant device parameters for LNA design, including maximum oscillation frequency f_{max} , quality of inductors, inner gain of the MOSFETs ($g_m/g_{ds} |_{L_{min}}$), and RF supply voltages are shown in the mixed-signal technology requirements table of the [PIDS](#) chapter. The evolution of the FoM from recent best-in-class published accounts for CMOS LNAs shows a clear trend towards better performance for smaller device dimensions. This is in good agreement with the increase in the quality of the devices needed for LNA design. Extrapolating these data into the future, an estimate of future progress in LNA design is obtained as shown in Table 9.

VOLTAGE-CONTROLLED OSCILLATOR (VCO)

Another key component of RF signal processing systems is the VCO. The VCO is the key part of a phase-locked loop (PLL), which synchronizes communication between an integrated circuit and the outside world in high-bandwidth and/or high-frequency applications. The key design objectives for VCOs are to minimize the timing jitter of the generated waveform (or, equivalently, the phase noise) and to minimize the power consumption. From these parameters a figure of merit (FoM_{VCO}) is defined:

$$FoM_{VCO} = \left(\frac{f_0}{\Delta f} \right)^2 \frac{1}{L\{\Delta f\} \cdot P} \quad (2)$$

Here, f_0 is the oscillation frequency, $L\{\Delta f\}$ is the phase noise power spectral density measured at a frequency offset Δf from f_0 , and P is the total power consumption.

There is no clear correlation between the operating frequency and the figure of merit. However, a good value of the figure of merit is usually more difficult to achieve at higher frequencies. Therefore the figure of merit is not completely independent of the operating frequency. The definition also neglects the tuning range of the VCO since the necessary tuning range strongly depends on the application. Typically, however, the VCO’s phase noise or power consumption worsens if a larger tuning range is required.

By restricting to fully integrated CMOS tuned VCOs with on-chip load (LC-tank) and making further simplifications, the FoM can be linked to technology development [1]. The phase noise is mainly determined by the thermal noise and the quality factor of the LC-tank. Thermal noise versus power consumption is approximately constant over the technology nodes. Finally, the evolution of the figure of merit versus technology node mainly depends on the quality of the available inductors [1]. The evolution of FoMs from recent best in class published accounts for VCOs shows a

¹⁸ Reference [1] is the recent paper, R. Brederlow, S. Donnay, J. Sauerer, M. Vertregt, P. Wambacq, and W. Weber ‘A mixed signal design roadmap for the International Technology Roadmap for Semiconductors (ITRS)’, *IEEE Design and Test*, December 2001.

clear trend of increasing performance for decreasing CMOS minimum feature size. The FoMs are in good agreement with the data of the best available devices needed for VCO design in these technologies. Based on prediction of the relevant device parameters for future technology nodes (refer to the mixed-signal technology requirements table of the [PIDS](#) chapter), an extrapolation of the VCO FoM for future technology nodes is given in Table 9.

POWER AMPLIFIER (PA)

Power amplifiers are key components in the transmission path of wired or wireless communication systems. They deliver the transmission power required for transmitting information off-chip with high linearity to minimize adjacent channel power. For battery-operated applications in particular, minimum DC power at a given output power is required.

To establish a performance figure of merit, key parameters including output power P_{out} , power gain G , carrier frequency f , linearity (in terms of $IIP3$), and power-added-efficiency (PAE) must be taken into account. Unfortunately, linearity strongly depends on the operating class of the amplifiers, which makes it difficult to compare amplifiers of different classes. To remain independent of the design approach and the specifications of different applications we omit this parameter in our figure of merit. To compensate for the 20 dB/decade roll-off¹⁹ of the PA's RF-gain, a factor of f^2 is included into the figure of merit. This results in:

$$FoM_{PA} = P_{out} \cdot G \cdot PAE \cdot f^2 \quad (3)$$

Finally, restricting to the simplest PA architecture (class A operation) and making further simplifications enables correlation between the FoM and device parameters [1]. The key device parameters are seen to be the quality factor of the available inductors and f_{max} . Values for these parameters are mapped in the mixed-signal technology requirements table of the [PIDS](#) chapter. FoMs of best-in-class CMOS PAs have increased by approximately a factor of two per technology node in recent years, strongly correlated with progress in active and passive device parameters. From required device parameters for future technology nodes (see the mixed-signal technology requirements table of the [PIDS](#) chapter), we can deduce requirements for future PA FoM values, as shown in Table 9.

ANALOG-TO-DIGITAL CONVERTER (ADC)

Digital processing systems have interfaces to the analog world: audio and video interfaces, interfaces to magnetic and optical storage media, and interfaces to wired or wireless transmission media. The analog world meets digital processing at the analog-to-digital converter (ADC), where continuous-time and continuous-amplitude analog signals are converted to discrete-time (sampled) and discrete-amplitude (quantized). The ADC is therefore a useful vehicle for identifying advantages and limitations of future technologies with respect to system integration. It is also the most prominent and widely used mixed-signal circuit in today's integrated mixed-signal circuit design.

The main specification parameters of an ADC are related to sampling and quantization. The resolution of the converter, i.e., the number of quantization levels, is 2^n where n is the "number of bits" of the converter. This parameter also defines the maximum signal to noise level $SNR = n \cdot 6.02 + 1.76$ [dB]. The sampling rate of the converter, i.e., the number of n -wide samples quantized per unit time, is related to the bandwidth that needs to be converted and to the power consumption required for reaching these performance points. The Shannon/Nyquist criterion states that a signal can be reconstructed whenever the sample rate exceeds twice the converted bandwidth: $f_{sample} > 2 \times BW$.

To yield insight into the potential of future technology nodes, the ADC FoM should combine dynamic range, sample rate f_{sample} and power consumption P . However, these nominal parameters do not give accurate insight into the effective performance of the converter; a better basis is the effective performance extracted from measured data. Dynamic range is extracted from low frequency signal-to-noise-and-distortion ($SINAD_0$) measurement minus quantization error (both values in dB). From $SINAD_0$ an "effective number of bits" can be derived as

¹⁹ Most CMOS PAs are currently operated in this regime; using DC-gain for applications far below f_t would result in a slightly increased slope.

$ENOB_0 = (SINAD_0 - 1.76) / 6.02$. Then, the sample rate may be replaced by twice the effective resolution bandwidth ($2 \times ERBW$) if it has a lower value, to establish a link with the Nyquist criterion:

$$FoM_{ADC} = \frac{(2^{ENOB_0}) \times \min(\{f_{sample}\}, \{2 \times ERBW\})}{P} \quad (4)$$

For ADCs, the relationship between FoM and technology parameters is strongly dependent on the particular converter architecture and circuits used. The complexity and diversity of ADC designs makes it nearly impossible to come up with a direct relationship, as was possible for the basic RF circuits. Nevertheless, some general considerations regarding the parameters in the FoM are proposed in [1], and in some cases it is possible to determine performance requirements of the design from the performance requirements of a critical subcircuit. The device parameters relevant for the different ADC designs are stated in the mixed-signal technology requirements table of the [PIDS](#) chapter. The trend in recent years shows that the ADC FoM improves by approximately a factor of 2 every three years. Taking increasing design intelligence into account, these past improvements are in good agreement with improvements in analog device parameters. Current best-in-class is approximately 800G [conversion-step/Joule] for stand-alone CMOS/BiCMOS, and approximately 400G [conversion-step/Joule] for embedded CMOS. Expected future values for the ADC FoM are shown in Table 9. Major advances in design are needed to maintain performance increases for ADCs in the face of decreased voltage signal swings and supplies. In the long run, fundamental physical limitations (thermal noise) may block further improvement of the ADC FoM.

Table 9 Projected Mixed-signal Figures of Merit for Four Circuit Types

YEAR OF PRODUCTION	2001	2004	2007	2010	2013	2016	DRIVER
MPU $\frac{1}{2}$ PITCH	130	90	65	65	45	22	
F_oM_{LNA} [GHz]	10	15	25	30–40	40–50	50–70	PIDS*
F_oM_{VCO} [1/J] 10^{22}	5	6	7	8–9	10–11	12–14	PIDS
F_oM_{PA} [W•GHz ²] 10^4	6	12	24	40–50	80–90	100–130	PIDS
F_oM_{ADC} [1/J] 10^{12}	0.4	0.8	1-1.2	1.6–2.5	2.5–5	4–10	PIDS

*refer to the [Process Integration](#) chapter, table for Mixed-signal Technology Requirements

MIXED-SIGNAL EVOLUTION

Evolution of the mixed-signal driver, including its scope of application, is completely determined by the interplay between cost and performance. The above figures of merit measure mixed-signal *performance*. However, *cost of production* is also a critical issue for practical deployment of AMS circuits. Together, cost and performance determine the sufficiency of given technology trends relative to existing applications, as well as the potential of given technologies to enable and address entirely new applications.

Cost estimation. Unlike high-volume digital products where cost is mostly determined by chip area, in mixed-signal designs the area parameter is only one of several cost factors. The area of analog circuits in an SOC is typically in the range of 5-30%; economic forces to reduce mixed-signal area are therefore not as strong as for logic or memory. Related considerations include:

- analog area can sometimes be reduced by shifting the partitioning of a system between analog and digital parts (e.g. auto-calibration of A-to-D converters);
- process complexity is increased by introducing high-performance analog devices, so that solutions can have less area but greater total cost;
- technology choices can impact design cost by introducing greater risk of multiple hardware passes (tapeout iterations);
- manufacturing cost can also be impacted via parametric yield sensitivities; and
- a system-in-package solution with multiple die (e.g., large, low-cost digital and small, high-performance analog) can be cheaper than a single SOC solution.

Such considerations make cost estimation very difficult for mixed-signal designs. We may attempt to quantify mixed-signal cost by first restricting our attention to high-performance applications, since these also drive technology demands. Next, we note that analog features are embodied as high-performance passives or analog transistors, and that area can be taken as a proxy for cost.²⁰ Since scaling of transistors is driven by the need to improve density of the digital parts of a system, analog transistors can simply follow, thus rendering it unnecessary to specifically address their layout density. At the same time, total area in most current AMS designs is determined by embedded passives; their area consumption dominates the cost of the mixed-signal part of a system. Therefore, the mixed-signal technology requirements table of the *PIDS* chapter sets a roadmap of layout density for on-chip passive devices that is needed to improve the cost/performance ratio of high-performance mixed-signal designs.

Estimation of technology sufficiency—Figure 9 shows ADC requirements for recent applications in terms of a power/performance relationship. Under conditions of constant performance (resolution \times bandwidth), a constant power consumption is represented by a straight line with slope -1 . Increasing performance, achievable with better technology or circuit design, is equivalent to a shift of the power consumption lines toward the upper right. The data show a very slowly moving technological “barrier-line” for ADCs for a given ADC sampling rate and power consumption of 1W (Figure 9). Most of today’s ADC technologies (silicon, SiGe, and III-V compound semiconductor technologies and their hybrids) lie below the 1W barrier-line, and near-term solutions for moving the barrier-line more rapidly are unknown.

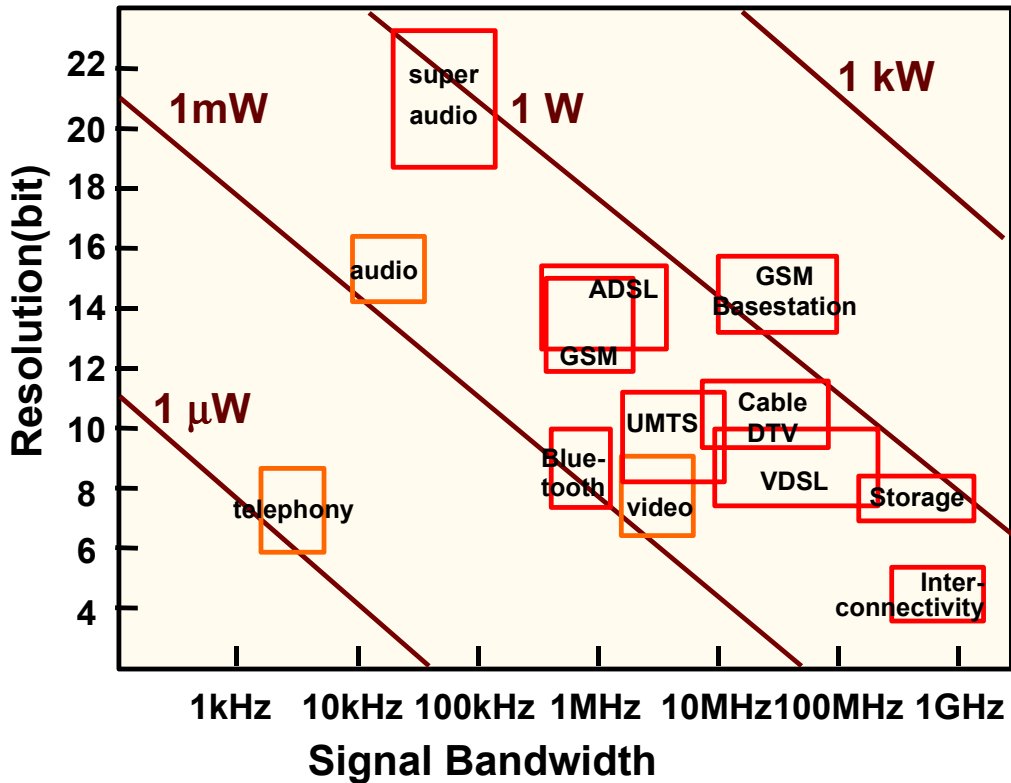


Figure 9 Recent ADC Performance Needs for Important Product Classes

²⁰ In analog designs, power consumption is often proportional to area, and since power is included in all four figures of merit, we have already implicitly considered area and cost criteria. Nonetheless, area requirements should be stated explicitly in a roadmap.

While the rate of improvement in ADC performance has been adequate for handset applications, this is clearly not the case for applications such as digital linearization of GSM base-stations, and handheld/mobile high-data rate digital video applications. For example, a multi-carrier GSM base-station with a typical setup of 32 carriers requires over 80dB of dynamic range. Implementing digital linearization in such a base-station with a 25 MHz transmitter band requires ADCs that have sampling rates of 300 MHz and 14 bits of resolution. According to Table 9 and assuming progress at recent rates, it will be perhaps until after 2010 before ADCs with such performance are manufactured in volume. While system designers would like to have such ADCs now, silicon and SiGe technologies have the necessary bit resolution (large numbers of devices per unit area) but not the speed; on the other hand, III-V compound semiconductor technologies have the speed but not the bit resolution. This motivates consideration of solutions that potentially increase the rate of ADC improvement at reasonable costs – e.g., use of compound semiconductors for their speed (perhaps combinations of HBTs, HEMTs, and resonant tunneling diodes), and hybrids of both CMOS and compound semiconductor technologies. The challenge for compound semiconductors is to increase the number of devices per unit area and to be co-integrated with CMOS processing.

Enabling new applications—For a given product, the usual strategy to increase unit shipments is cost reduction while increasing product performance. However, this is not the only driver for the semiconductor business, especially for products that include mixed-signal parts. Rather, improving technology and design performance enables *new* applications (comparable to the realization of the mobile handset in recent years), thus pushing the semiconductor industry into new markets. Analysis of mixed-signal designs as in Figure 9 can also be used to estimate design needs and design feasibility for future applications and new markets. We see that increasing performance is equivalent to the ability to develop new products which need higher performance or lower power consumption than available in today's technologies. Alternatively, when specifications of a new product are known, one can estimate the technology needed to fulfill these specifications, and/or the timeframe in which the semiconductor industry will be able to build that product with acceptable cost and performance. In this way, the FoM concept can be used to evaluate the feasibility and the market of potential new mixed-signal products. The ability to build high performance mixed-signal circuitry at low cost will continuously drive the semiconductor industry into such new products and markets.

MIXED-SIGNAL CHALLENGES

For most of today's mixed-signal designs, and particularly in classical analog design, the processed signal is represented by a voltage difference, so that the maximum signal is determined by the supply voltage. Decreasing supplies, a consequence of constant-field scaling, means decreasing the maximum achievable signal level. This has a strong impact on mixed-signal product development for SOC solutions. Typical development time for new mixed-signal parts is much longer than for digital and memory parts; sheer lack of design resources thus becomes another key challenge. An ideal design process would reuse existing mixed-signal designs and adjust parameters to meet interface specifications between a given SOC and the outside world, but such reuse depends on a second type of MOSFET that does not scale its maximum operating voltage. This has led to the specification in the PIDS Chapter of a mixed-signal CMOS transistor that uses a higher analog supply voltage and stays unchanged across multiple digital technology generations. Even with such a device, voltage reduction and development time of analog circuit blocks are major obstacles to low-cost and efficient scaling of mixed-signal functions. In summary, the most daunting mixed-signal challenges are:

- *decreasing supply voltage*, with needs including current-mode circuits, charge pumps for voltage enhancement, and thorough optimization of voltage levels in standard-cell circuits (PIDS, Design),
- *increasing relative parametric variations*, with needs including active mismatch compensation, and tradeoffs of speed versus resolution (PIDS, FEP, Lithography, Design),
- *increasing numbers of analog transistors per chip*, with needs including faster processing speed and improved convergence of mixed-signal simulation tools (Modeling and Simulation, Design),
- *increasing processing speed (clock frequencies)*, with needs including more accurate modeling of devices and interconnects, as well as test capability and package- and system-level integration (Test, Assembly and Packaging, Modeling and Simulation),

14 System Drivers

- *increasing leakage and crosstalk* arising from SOC integration, with needs including more accurate crosstalk and delay modeling, fully differential design for RF circuits, as well as technology measures outlined in the PIDS Chapter (PIDS, Modeling and Simulation, Design), and
- *shortage of design skills and productivity* arising from lack of training and poor automation, with needs including education and basic design tools research (Design).

SOC SYSTEM DRIVER

The system-on-chip (SOC) driver class is characterized by heavy reuse of intellectual property (IP) to improve design productivity, and by *system* integration that potentially encompasses heterogeneous technologies. SOCs exist to provide low cost and high integration. *Cost* considerations drive the deployment of low-power process and low-cost packaging solutions, along with fast-turnaround time design methodologies. The latter, in turn, require new standards and methodologies for IP description, IP test (including built-in self-test and self-repair), block interface synthesis, etc. *Integration* considerations drive the demand for heterogeneous technologies (flash, DRAM, MEMS, ferroelectric RAM (FRAM, MRAM), chemical sensors, etc.) in which particular system components (memory, sensors, etc.) are implemented, as well as the need for chip-package co-optimization. Thus, SOC is the driver for convergence of multiple technologies not only in the same system package, but potentially in the same manufacturing process. We discuss the nature and evolution of SOCs with respect to three variants driven respectively by multi-technology integration (MT), high performance (HP), and low power and low cost (LP). This partition is by no means disjoint, but rather reflects separate driving concerns (e.g., low-power design *is* high-performance design, but also reduces package and system cost).

SOC-MULTI-TECHNOLOGY

The need to build heterogeneous systems on a single chip is driven by such considerations as cost, form-factor, and reliability. Thus, process technologists seek to meld CMOS with MEMS, optoelectronics, and so on. Process complexity is a major factor in the cost of SOC-MT applications, since more technologies assembled on a single chip requires more complex processing. The total cost of processing is difficult to predict for future new materials and combinations of processing steps. However, at present cost considerations limit the number of technologies on a given SOC: processes are increasingly modular (e.g., enabling a flash add-on to a standard low-power logic process), but the modules are not generally “stackable”. Figure 10 shows how first integrations of each technology within standard CMOS processes – not necessarily together with other technologies, and not necessarily in volume production – might evolve. CMOS integration of the latter technologies (chemical sensing, electro-optical, electro-biological) is much less certain, since this depends not only on basic technical advances but also on SOC-MT being more cost-effective than multi-die system-in-package alternatives. Today, a number of technologies (flash, DRAM, GaAs) are more cost-effectively flipped onto or integrated side-by-side with silicon in the same module. Physical scale in system applications (e.g., ear-mouth = speaker-microphone separation, or distances within a car) also affect the need for single-die integration, particularly of sensors.

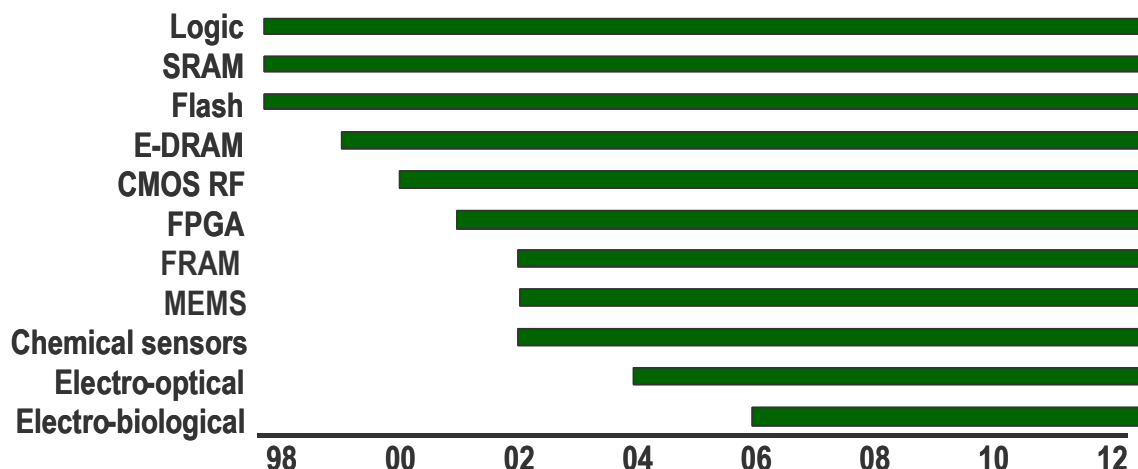


Figure 10 First Integration of Technologies on SOC with Standard CMOS Process

SOC-HIGH PERFORMANCE

Examples of SOC-HP include network processors and high-end gaming applications. Since it reflects MPU-SOC convergence, SOC-HP follows a similar trend as MPU and is not separately modeled here. However, one aspect of SOC-HP merits discussion, namely, that instances in the high-speed networking domain drive requirements for off-chip I/O signaling (which in turn create significant technology challenges to Test, Assembly and Packaging, and Design). Historically, chip I/O speed (per-pin bandwidth) has been scaling much more slowly than internal clock frequency. This is partly due to compatibility with existing slow I/O standards, but the primary limitation has been that unterminated CMOS signals on printed circuit boards are difficult to run at significantly greater than 100MHz due to slow settling times. During the past decade, high-speed links in technology initially developed for long-haul communication networks have found increasing use in other applications. The high-speed I/O eliminates the slow board settling problems by using point-to-point connections and treating the wire as a transmission line. Today the fastest of these serial links can run at 10Gbit/s.

A high-speed link has four main parts: a transmitter to convert bits to an electrical signal that is injected into the board-level wire, the wire itself, a receiver that converts the signal at the end of the wire back to bits, and a timing recovery circuit that compensates for the delay of the wire and samples the signal on the wire at the right place to get the correct data. Such links are intrinsically mixed-signal designs since receivers, transmitters, and timing recovery all require analog blocks (e.g., the VCO discussed as part of the Mixed-Signal driver is a key component of a timing recovery circuit). Broadly speaking, high-speed links are used in optical systems, chip-to-chip connections, and backplane connections. We now discuss each of these applications in slightly more detail.

Optical links generally push link performance the hardest; since there are generally a small number of optical signals, these links can tolerate relatively complex and power hungry interface circuits. Today, optical links run at 10Gbit/s per pin, and are expected to continue to scale up in frequency as projected in the Test Chapter (high-speed serial links discussion). Initially, electronics for these links were created in non-CMOS technologies, since CMOS was thought incapable of meeting the high-speed requirements. However, over the past five years, many researchers have developed circuits that can run at 10Gbit/s. While some papers have demonstrated links that run as fast as 1 FO4 delay per bit, most links run at 2–4 FO4 delays per bit, which yields 10Gbit/s in the 180nm node. Continuing to scale link speed with technology should be possible from the circuits standpoint, but will become difficult due to parasitics and packaging. Signals at this speed are highly sensitive to any discontinuities in their signal path. Even if controlled impedance packaging is used, vias in the package or board can cause impedance changes that will degrade the signal. The 1–2 pF parasitic capacitance from the ESD device will also significantly degrade the signal. Thus, continued performance scaling will require significant work in ESD, package and board design.

Chip-to-chip interconnections communicate information between two chips located on the same board, usually close to each other. The main metric driving the design of these links is not Gbit/s since it is generally possible to use a number of links in parallel to connect these chips. For example, if going twice as fast requires 10× the area and 10×

the power, it is better to use two links in parallel. Thus, these links are optimized for performance and cost, not just performance. In general, the highest chip-to-chip link speeds are 2–4 times slower than the highest optical link speeds. Bit times for these links vary dramatically, e.g., point-to-point links are available today with bit times ranging from about 2.5ns (400Mbit/s) to .4ns (2.5Gbit/s). This wide range of performance reflects dependencies on the number of IO required (higher IO counts have slower speeds), the degree of risk the designer is willing to take, and sometimes an existing I/O standard. Design of robust high-speed I/O is still a mixed-signal problem that cannot be automated or checked with current tools. Thus, many design teams are still conservative when choosing I/O rates. As technology scales and design tools become more robust, bit times should approach 4-8 FO4 delays, but this will require additional circuitry to compensate for package and other parasitic effects.

The last major application for high-speed links is in networking, where two chips on different boards must communicate. The signal path is still point-to-point, but travels from one chip through its package to the local board, through a connector to another board, through another connector to the destination board, and then through that board and receiver package to the receiver chip. For high bandwidth each chip generally has a large number of links, so that performance per unit cost is critical. The principal difference from chip-to-chip links is that the “wire” between the two chips has worse electrical properties. Wire issues are just becoming visible for current systems that run at 2.5Gbit/s; they will be a serious concern as speeds increase through 5Gbit/s and to 10Gbit/s, which is targeted for the 130nm node.

SOC-LOW COST, LOW POWER

Examples of SOC-LP include portable and wireless applications such as PDAs or digital camera chips. Table 10 sets requirements for various attributes of a low-power, consumer-driven, handheld wireless device (“PDA”) with multimedia processing capabilities, based in part on the model created by the Japan Semiconductor Technology Roadmap Working Group 1 and originally introduced in the 2000 ITRS update (Design Chapter). Key aspects of the model are as follows.²¹

- The system design consists of embedded blocks of CPU, DSP and other processing engines, and SRAM and embedded DRAM circuits. Processor core logic increases by 4× per node, and memory content increases by 2–4× per node.²²
- Die size increases on average by 20% per node through 2016 to accommodate increased functionality; this matches historical trends for the application domain.
- Layout densities for memory and logic fabrics are the same as for the MPU driver, with eDRAM density assumed to be 3× SRAM density.
- Maximum on-chip clock frequency is approximately 5–10% of the MPU clock frequency at each node.
- Peak power dissipation is limited to 0.1 W at 100°C, and standby power to 2.1 mW, due to battery life.²³

²¹ Other aspects of the model, which are not essential to the following analyses, address external communication speed (increasing by 6× per node in the near term, starting from 384 Kbps in 2001) and addressable system memory (increasing by 10× per node, starting from 0.1Gb in 2001).

²² The PDA contains approximately 20 million transistors in 2001. The model assumes that increasing parallel computation will be required in each generation of the device, to support video, audio and voice recognition functionality. This is reflected in CPU and DSP content (e.g., number of cores), which increases four-fold (4×) per technology node to match the processing demands of the corresponding applications. (By comparison, MPU logic content is projected to double with each node.) Overhead area (I/O buffer cells, pad ring, whitespace due to block packing, analog blocks, etc.) is fixed at 28% of the die. The 20M transistor count is broken down as follows. A typical CPU/DSP core (e.g., ARM) today is approximately 30-40K gates, or 125K transistors. We assume four such cores on chip in 2001, i.e., 500K CPU/DSP core transistors. In 2001, the “peripheral” logic transistor count is 11.5M transistors, and this count grows at 2X/node thereafter. SRAM transistor count is 8M in 2001, and grows at 2×/node thereafter. The composition of SRAM versus DRAM depends on the ratio of memory to logic. We assume that embedded DRAM (eDRAM) is cost effective when at least 30% of the chip area is memory. Its use is not invoked until the 30% trigger point, and begins at 16Mb in 2004. Once triggered, the eDRAM content quadruples every technology node. (While the SOC-LP PDA is a “single-chip design”, we do not imply any judgement as to whether multi-die or single-die implementation will be more cost-effective.)

²³ At 120Wh/kg in 2001, a 140g battery allows 0.1W operation for 7 days, 24 hours per day.

SOC TRENDS

Since SOC is aimed at low-cost and rapid system implementation, it is highly instructive to consider the implications of *power management* and *design productivity* on the achievable space of SOC designs. The following discussion develops trend analyses for the SOC-LP driver with respect to these issues.

Power—Two approaches can be used to derive the power dissipation for the SOC-LP model. The first approach is to accept the system specifications (0.1 W peak power, and 2.1 mW standby power) in a “top-down” fashion. The second approach is to derive the power requirements “bottom-up” from the implied logic and memory content, as well as process and circuit parameters. Logic power consumption is estimated based on $\alpha CV_{dd}^2 f + I_{off} V_{dd}$ model for dynamic plus static power, using area-based calculations similar to those in the MPU power analysis. The memory power consumption model also uses $\alpha CV_{dd}^2 f + I_{off} V_{dd}$ with a different factor for α .²⁴ For these calculations, we refer to the low-power device roadmap described in the PIDS Chapter. Table 11 lists key attributes for the low standby power (LSTP) and low operating power (LOP) devices, and contrasts these with the high-performance device model used for MPU power and frequency analyses. It is almost certain that future low-power SOCs will integrate multiple (LOP, LSTP, HP) technologies simultaneously within the same core, to afford greater control of dynamic power, standby power, and performance.

Figure 11 shows the “bottom-up” *lower bound* for total chip power at an operating temperature of 100°C, assuming that all logic is implemented with LOP or LSTP devices and operates as described in the previous footnote. We say that this is a lower bound since in practice some logic would need to be implemented with faster, higher-current devices. The figure indicates that SOC-LP power levels will substantially exceed the low-power requirements of the PDA application, and further provides a breakdown of power contributions for each case. As expected, LOP power is primarily due to standby power dissipation while LSTP power is primarily due to dynamic power dissipation²⁵. Total chip power using only LOP devices reaches 2.45W in 2016, mostly due to a sharp rise in static power after 2010. Total chip power using only LSTP devices reaches 1.5W in 2016; almost all of this is dynamic power.

Table 10 System Functional Requirements for the PDA SOC-LP Driver

YEAR OF PRODUCTION	2001	2004	2007	2010	2013	2016
Process Technology (nm)	130	90	65	45	32	22
Supply Voltage (V)	1.2	1	0.8	0.6	0.5	0.4
Clock Frequency (MHz)	150	300	450	600	900	1200
Application (maximum required performance)	Still Image Processing Web Browser Electric Mailer Scheduler	Real Time Video Codec (MPEG4/CIF)		Real Time Interpretation		
Application (other)		TV Telephone (1:1)		TV Telephone (>3:1)		
		Voice Recognition (Input)		Voice Recognition (Operation)		
		Authentication (Crypto Engine)				
Processing Performance (GOPS)	0.3	2	15	103	720	5042
Required Average Power (W)	0.1	0.1	0.1	0.1	0.1	0.1
Required Standby Power (mW)	2.1	2.1	2.1	2.1	2.1	2.1
Battery Capacity (Wh/Kg)	120	200		400		

²⁴ I_{off} denotes the NMOSFET drain current at room temperature, and is the sum of the NMOS sub-threshold, gate, and junction leakage current components, as described in the [PIDS](#) chapter. Details of active capacitance density calculations, dependences on temperature and threshold, etc. may be found in the PIDS Chapter documentation and in the following [supplemental file](#). The activity of logic blocks is fixed at 10%. The activity of memory blocks is estimated to be 0.4% based on the following analysis of large memory designs. We first assume that a memory cell contributes 2 gate capacitances of minimum size transistors for switching purposes, accounting for source/drain capacitances, contact capacitances and wiring capacitance along the bit lines. A write access requires power in the row/column decoders, word line and M bit lines, sense amplifiers and output buffers. We consider memory to be addressed with 2N bits and assume that memory power is due primarily to the column capacitances, and that $M \times 2^N$ bits are accessed simultaneously out of $2^N \times 2^N$ possible bits. Then $\alpha = M/2^N$ which is the ratio of accessed bit to total bits in the memory. For example, for a 16Mbit memory, $M=16$ and $N=12$; hence $\alpha=0.4\%$.

²⁵ At 25°C, dynamic power dissipation dominates the total power in both the LOP and LSTP cases.

18 System Drivers

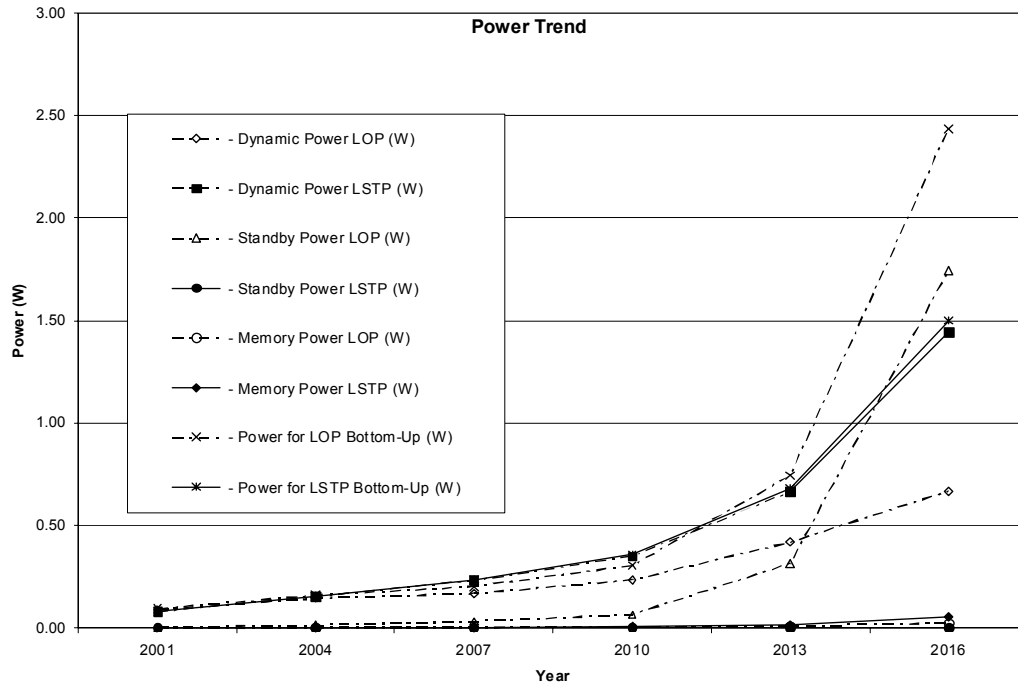


Figure 11 Total Chip Power Trend for PDA Application

Table 11 Low Operating Power (LOP) and Low Standby Power (LSTP) Device and Process Attributes

		99	00	01	02	03	04	05	06	07	10	13	16
Parameter	Type												
Tox (nm)	MPU	3.00	2.30	2.20	2.20	2.00	1.80	1.70	1.70	1.30	1.10	1.00	0.90
	LOP	3.20	3.00	2.2	2.0	1.8	1.6	1.4	1.3	1.2	1.0	0.9	0.8
	LSTP	3.20	3.00	2.6	2.4	2.2	2.0	1.8	1.6	1.4	1.1	1.0	0.9
Vdd	MPU	1.5	1.3	1.2	1.1	1.0	1.0	0.9	0.9	0.7	0.6	0.5	0.4
	LOP	1.3	1.2	1.2	1.2	1.1	1.1	1.0	1.0	0.9	0.8	0.7	0.6
	LSTP	1.3	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.2	1.1	1.0	0.9
Vth (V)	MPU	0.21	0.19	0.19	0.15	0.13	0.12	0.09	0.06	0.05	0.021	0.003	0.003
	LOP	0.34	0.34	0.34	0.35	0.36	0.32	0.33	0.34	0.29	0.29	0.25	0.22
	LSTP	0.51	0.51	0.51	0.52	0.53	0.53	0.54	0.55	0.52	0.49	0.45	0.45
Ion (uA/um)	MPU	1041	1022	926	959	967	954	924	960	1091	1250	1492	1507
	LOP	636	591	600	600	600	600	600	600	700	700	800	900
	LSTP	300	300	300	300	400	400	400	400	400	500	500	600
CV/I (ps)	MPU	2.00	1.64	1.63	1.34	1.16	0.99	0.86	0.79	0.66	0.39	0.23	0.16
	LOP	3.50	2.87	2.55	2.45	2.02	1.84	1.58	1.41	1.14	0.85	0.56	0.35
	LSTP	4.21	3.46	4.61	4.41	2.96	2.68	2.51	2.32	1.81	1.43	0.91	0.57
Ig (uA/um)	MPU	2e-5	1e-2	2e-2	2e-2	1e-1	2e-1	3e-1	3e-1	3e-4	1e-5	4e-9	2e-16
Ioff (uA/um)	MPU	0.00	0.01	0.01	0.03	0.07	0.10	0.30	0.70	1.00	3	7	10
	LOP	1e-4	1e-4	1e-4	1e-4	1e-4	3e-4	3e-4	3e-4	7e-4	1e-3	3e-3	1e-2
	LSTP	1e-6	1e-6	1e-6	1e-6	1e-6	1e-6	1e-6	1e-6	1e-6	1e-6	3e-6	7e-6
Gate L (nm)	MPU	100	70	65	53	45	37	32	30	25	18	13	9
	L(*P)	110	100	90	80	65	53	45	37	32	22	16	11
Gate cap (fF/um)	MPU	1.39	1.29	1.26	1.07	1.02	0.95	0.87	0.85	0.90	0.81	0.69	0.59
	LOP	1.43	1.39	1.28	1.23	1.10	1.00	0.95	0.85	0.89	0.75	0.63	0.53
	LSTP	1.43	1.39	1.15	1.10	0.99	0.89	0.84	0.77	0.82	0.71	0.61	0.51

Table 12 Power Management Gap

	2001	2004	2007	2010	2013	2016
Total LOP Dynamic Power Gap (X)	-0.06	0.59	1.03	2.04	6.43	23.34
Total LSTP Dynamic Power Gap (X)	-0.19	0.55	1.35	2.57	5.81	14.00
Total LOP Standby Power Gap (X)	0.85	5.25	14.55	30.18	148.76	828.71
Total LSTP Standby Power Gap (X)	-0.98	-0.98	-0.97	-0.88	-0.55	0.24

Table 12 shows the implied *power management gap*, i.e., the factor improvement in power management that must be achieved jointly at the levels of application, operating system, architecture, and IC design.²⁶ Required power reduction factors exceed 20× for dynamic power, and 800× for standby power. Here, the Total Power Gap is defined as (Total Power – 0.1W)/0.1W (the PDA total power requirement). Similarly, the Total Standby Power Gap is defined as (Total Standby Power – 2.1mW)/2.1mW (the PDA total standby power requirement). Negative values indicate the lack of any power management gap (i.e., existing techniques suffice).

²⁶ For HP MPUs implemented using high-performance devices, the ITRS model implies a nearly 30X power management gap by the end of the roadmap with respect to package power limits (this can be seen from the following [GTX study](#)). An alternative portrayal of the power management gap is that the maximum chip area that can contain logic decreases substantially if the chip is to remain within power constraints, and if we simply extrapolate current designs (without any improvement in application-level, OS-level, VLSI architecture, and IC design technology for power management).

20 System Drivers

Figure 12 projects logic/memory composition of SOC-LP designs, assuming that chip power is constrained according a power budget of 0.1W and that chip size is constrained to 100mm². Memory content outstrips logic content faster with LSTP devices since their operating power is much higher than that of LOP devices. Both models indicate that chips will be asymptotically dominated by memory by 2016 without substantial improvements in power management capability. Recall that the PDA chip size is projected to grow at approximately 20% per node even though power remains flat at 0.1W. This would lead to even more extreme memory-logic imbalances in the long-term years.

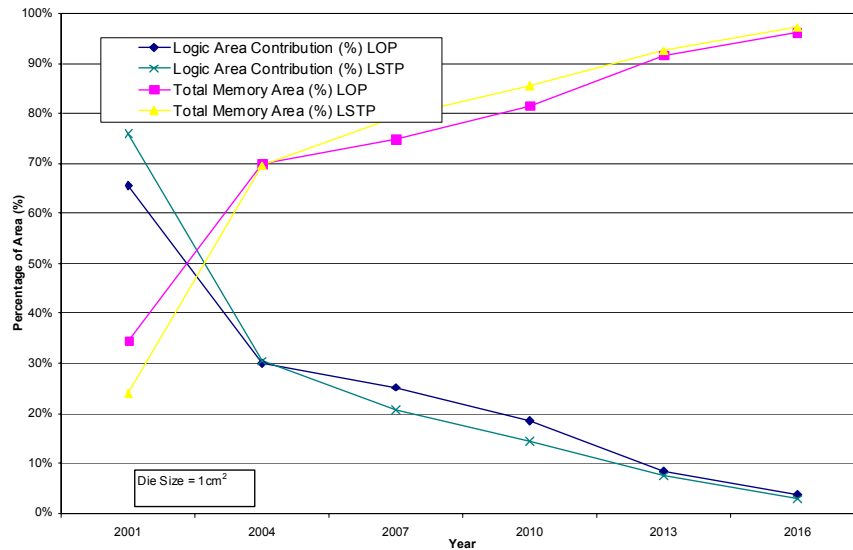


Figure 12 Power Gap Effect on Chip Composition

Productivity—The concept of (normalized) gates per designer-day may be used to measure the productivity of an integrated circuit designer. This figure of merit addresses the required improvement in designer productivity which must eventually be achieved by more reusable IP, better SOC design tools, etc. Normalized gates per designer-day can also be measured on a project by project basis to assess overall productivity of a design process. The importance of design productivity improvement can be seen from the effect of different rates of productivity improvements on the feasible logic and memory composition of the SOC. We assume that development of new logic blocks is expensive, but that reusable blocks also have overheads (learning curve, verification, integration, etc.) associated with their use. Memory design, through the use of compilers, is assumed to require minimal resources. These assumptions lead to the following model, an early version of which was discussed in the 2000 ITRS Design update, provided in the [supplemental material link](#) to the GTX study.

- Designer productivity for new logic is 360Kgates/year per designer in 1999 (aggregate 1000 gates per designer-day) and grows at a specified percentage per node thereafter.
- Designer productivity for reused logic is 720Kgates/year per designer in 1999 and grows at a specified percentage per node thereafter.
- Memory design has negligible cost due to the use of memory compilers.
- Available designer resources for the SOC are fixed at 10 person-years.
- Logic and memory densities are assumed to follow the MPU roadmap.

Figure 13 plots the maximum amount of chip area that can be occupied by logic (y-axis), given that a prescribed amount of chip area is already occupied by a given amount of memory (x-axis). Feasible solutions for the amounts

of reused logic and new logic – always maximizing new logic²⁷ – are plotted on the y-axis. We assume a constant die size of 1 cm², along with 42% productivity growth per node. The figure shows plots for four different nodes in the ITRS; valid solutions require 10 or fewer person-years of effort according to the y-axis scale on the right-hand side of each plot.²⁸ We see that in the year 2001 the designer can essentially dictate the memory/logic and new/reused ratios, since all solutions are within the 10 person-year design resource. On the other hand, in the year 2016 there is only a small feasible region, where roughly 95% of the chip must be memory and the remaining 5% can be constructed from new logic or reusable logic blocks. Figure 13 thus shows that without adequate productivity growth, using a constant designer resource with a constant die size will eventually force the SOC to contain only memory and (a small amount) of reusable blocks.

The design productivity requirement also follows from the growth in SOC logic content, or simply from the scaling of available transistor counts. If available transistor count doubles at each node, then keeping design content roughly the same requires at least 100% design productivity improvement per node. Figure 14a shows how the evolution of logic-memory balance changes with different rates of design productivity improvement; Figure 14b shows that a 100% improvement per node preserves the level of designer freedom enjoyed in the year 2001. Achieving such levels of productivity improvement is a key challenge for SOC. As noted in the Design Chapter, higher levels of reusability can be achieved for given applications using a *platform-based design* approach (whereby derivative designs are rapidly implemented from a single platform that has a fixed portion and a variable portion that permits proprietary or differentiated designs). Use of reprogrammable fabrics in such a platform can increase productivity further.

SOC CHALLENGES

SOC presents Design, Test, PIDS and other areas with a number of technology challenges, such as the development of reusable analog IP. The most daunting SOC challenges are:

- *design productivity improvement of > 100% per node*, with needs including platform-based design²⁹ and integration of programmable logic fabrics (Design),³⁰
- *management of power* especially for low-power, wireless, multimedia applications (Design, PIDS),
- *system-level integration of heterogeneous technologies* including MEMS and optoelectronics (PIDS, FEP, Design), and
- *development of SOC test methodology*, with needs including test reusability and analog/digital BIST.

²⁷ Chip area consists of memory plus new logic plus reused logic. A vertical line in the plot defines a possible combination of the three components. Only those solutions that require 10 person-years or less (right hand scale on y-axis) are considered to be valid.

²⁸ For example, in the year 2001 graph, if a chip has 20% memory area, then up to 18% new logic (and 62% reused logic) may be designed with a design resource of 10 person-years. However, if a chip has 80% memory in year 2001, then we can design all the way up to 20% new logic (and 0% reused logic) with only a 5 person-year design resource.

²⁹ Platform-based design is focused on a specific application domain. The platform embodies the hardware architecture, embedded software architecture, design methodologies for IP authoring and integration, design guidelines and modeling standards, IP characterization and support, and hardware/software verification and prototyping. Derivative designs may be rapidly implemented from a single platform that has a fixed portion and a variable portion that permits proprietary or differentiated designs. (See: H. Chang et al., *Surviving the SOC Revolution: A Guide to Platform-based Design*, Boston, Kluwer Academic, 1999.)

³⁰ A programmable logic core is a flexible logic fabric that can be customized to implement any digital logic function after fabrication. The structure of a programmable logic fabric may be similar to an FPGA capability within specific blocks of the SOC. They allow reprogrammability, adaptability and reconfigurability which greatly improves chip productivity. Applications include blocks that implement standards and protocols that continue to evolve, changing design specifications, and customization of logic for different, but related, applications and customers.

22 System Drivers

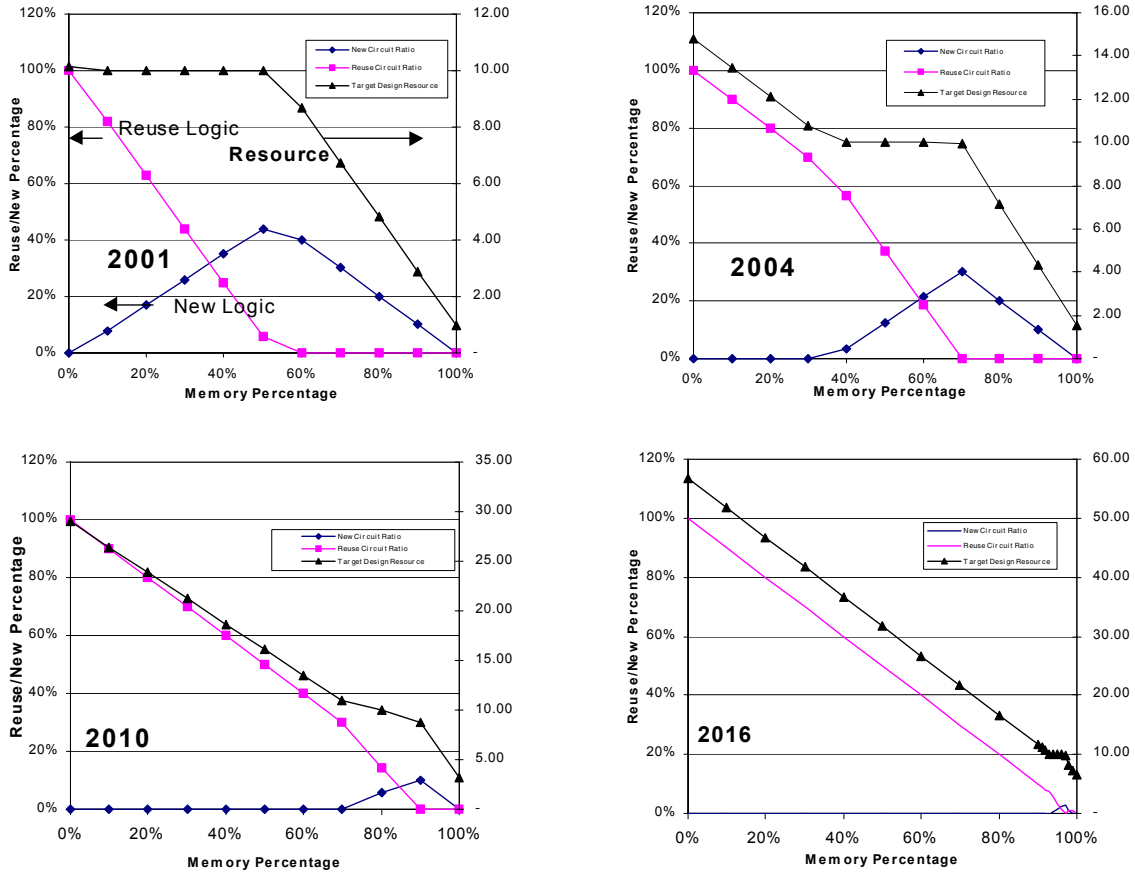


Figure 13 New and Reused Logic Content versus Memory Content with Constant Die Size and Insufficient (42% Per Node) Design Productivity Growth

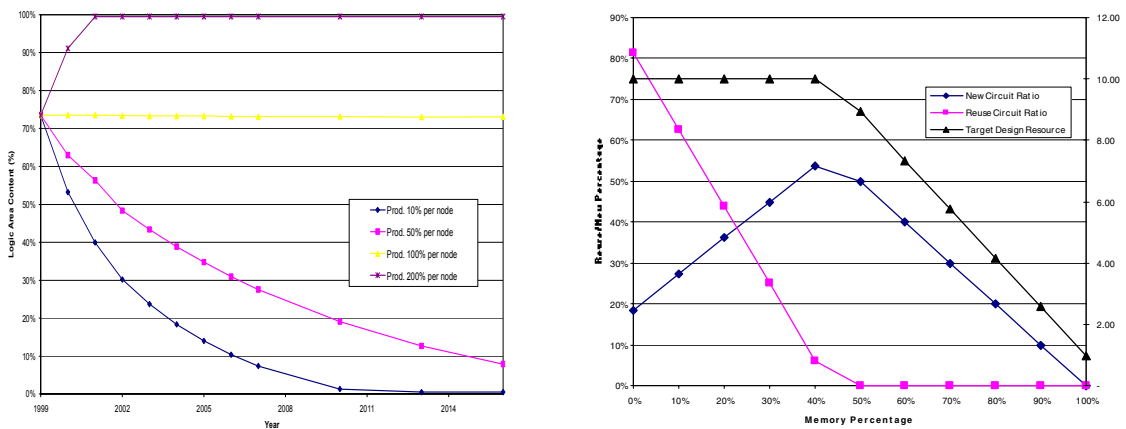


Figure 14a Evolution of Maximum Logic Content with Different Rates of Design Productivity Improvement

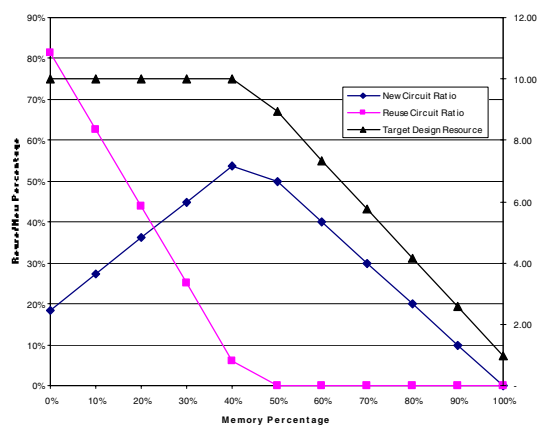


Figure 14b 100% Productivity Improvement per Node Will Preserve Designer Freedom at the End of the ITRS Forecast Period