

INTERNATIONAL  
TECHNOLOGY ROADMAP  
FOR  
SEMICONDUCTORS 2.0

2015 EDITION

SYSTEM INTEGRATION

## 2 ITRS 2.0 System Integration Chapter

THE ITRS IS DEvised AND INTENDED FOR TECHNOLOGY ASSESSMENT ONLY AND IS WITHOUT REGARD TO ANY COMMERCIAL CONSIDERATIONS PERTAINING TO INDIVIDUAL PRODUCTS OR EQUIPMENT.

# ITRS 2.0 SYSTEM INTEGRATION CHAPTER

## 1. MISSION

The mission of the System Integration (SI) chapter in ITRS2.0 is to establish a top-down, system-driven roadmapping framework for key market drivers of the semiconductor industry in the 2015-2030 period. The SI focus team is currently developing and constructing roadmaps of relevant system metrics for mobile, datacenter and Internet of Things (IoT) drivers. The mobile driver, embodied by the smartphone product, has redefined the existing ITRS SOC-CP (consumer portable system-on-chip) driver with richer feature requirements. As a fast-growing aspect of the datacenter, microservers have been separated out from the conventional server market segment [16]. IoT, as one of the fastest-growing market segments of electronic devices [15], imposes significantly different design considerations from conventional electronics designs due to low-power and ubiquitous deployment requirements. For these new drivers, the SI focus team seeks to describe new indicators (e.g., power management, bandwidth and integration) as functionalities expand, architectures evolve, and heterogeneous integration soars.

## 2. SCOPE AND TERMINOLOGY

**Changes in the semiconductor industry supply chain.** The 1980s and 1990s saw a semiconductor industry dominated by integrated device manufacturers (IDMs). During this period, the architecture of the main driver in the ITRS, the microprocessor unit (MPU), was not application-driven. Standard components in PC and server systems, e.g., memories and microprocessors, scaled their densities and operating frequencies continuously to meet aggressive performance and cost requirements. Applications had to be designed based on these components. However, in the past ten years, fabless design houses have changed the industry landscape. Design teams have been building customized system-on-chip (SOC) and system-in-package (SIP) products, rather than building standard components, to address specific application requirements. As applications evolve, they drive further requirements for heterogeneous integration, outside system connectivity, etc. A key goal of the SI focus team is to extract the technology requirements hidden behind the evolution of end products such as mobility, datacenter/microservers and IoT. In Table SYSINT1, both the near-term and long-term challenges to system integration are summarized.

Table SYSINT1: Summary of system integration challenges.

<i>Near term (within 5 years)</i>	<i>Sub-challenges</i>	<i>Relation between drivers</i>
Design productivity	System integration, AMS/MEMS co-design and design automation SIP and 3D (TSV-based) planning and implementation flows Heterogeneous integration (optical, mechanical, chemical, biomedical, etc.)	Mobile/IoT: Beneficial for system dimension scaling, performance improvement, and cost. Datacenter: Beneficial for system-wide bandwidth and power efficiency
Power management	Dynamic and static, system- and circuit-level power optimization	Mobile/ IoT: Beneficial for battery life Datacenter: Beneficial for cooling cost and energy fee
Manufacturability	Performance/power variability, device parameter variability, lithography limitations impact on design, mask cost, quality of (process) models	Mobile/datacenter/IoT: Beneficial for cost reduction and reliability improvement
Bandwidth / service latency	High performance memory / NVM interfaces, memory / processor stacking	Mobile: Beneficial for improving display capacity and developing more sophisticated services Datacenter: Beneficial for faster responses
Cooling	Temperature-constrained physical implementation, 3D integration/packaging	Mobile/datacenter: Avoiding heating issues

## 4 ITRS 2.0 System Integration Chapter

<i>Long term (&gt; 5 years)</i>	<i>Sub-challenges</i>	<i>Relation between drivers</i>
Design productivity	System-Level Design Automation (SDA) Executable Specification	Mobile/datacenter/IoT: Beneficial for faster design turnaround-time and less design effort
Power management	On-die power sensors, silicon photonics, novel transistors and memory	Mobile/ IoT: Beneficial for battery life Datacenter: Beneficial for cooling cost and energy fee
Manufacturability	Sequential 3D integration 3D transistors (LGAA, VGAA, CNT) Novel memory technologies	Mobile/datacenter/IoT: Beneficial for cost reduction and reliability improvement
Bandwidth / service latency	High radix networks, interfaces with novel memory devices	Mobile: Beneficial for improving display capacity and developing more sophisticated services Datacenter: Beneficial for faster responses
Cooling	Microfluidic cooling (single-phase / two-phase)	Mobile/datacenter: Avoiding heating issues

AMS—analogue/mixed signal  
NVM—non-volatile memory

MEMS—micro-electro-mechanical systems  
LGAA/VGAA—lateral/vertical gate-all-around

TSV—through silicon via  
CNT—carbon nanotube

**Motivations and distinctions between ITRS 2.0 system drivers and ITRS 1.0 “system drivers”.** Historically, in its 1998-2013 editions, the ITRS has used metrics such as transistor density, number of cores, power, etc., to roadmap technology evolution of integrated circuits (ICs). These metrics are essentially driven by the physical-dimension scaling as predicted by Moore’s Law. The current (2013 edition) ITRS System Drivers Chapter roadmaps key IC products that drive process and design technologies. However, new requirements from applications such as mobile devices, datacenters/microservers, etc. require a new, system-level roadmapping approach, as these applications imply roadmaps for system-level metrics (e.g., the number of sensors, memory bandwidth, etc.). The ITRS roadmapping process as previously seen in the System Drivers Chapter has not explicitly incorporated these system-level product requirements. Therefore, a crucial goal of “ITRS 2.0” is to connect emerging system product drivers, along with corresponding metrics, into the ITRS’s semiconductor roadmapping methodology.

### Terminology and definitions of drivers.

A mobile device is a computing device of which the form factor could be carried by users in daily life and has the capacity to connect to other devices, display information, and execute software.

A datacenter is a facility that centralizes an organization’s IT operations and equipment, and where it stores, manages, and disseminates its data. Data centers house a network’s most critical systems and are vital to the continuity of daily operations. Consequentially, the service latency, power management, scalability, dependability, and security of data centers and their information are top priorities for organizations [41] [65] [58].

An IoT device (a.k.a. smart object (SO) in other literature) is an autonomous, physical digital object augmented with sensing/actuating, processing, storing, and networking capabilities. It is able to sense/actuate, store, and interpret information created within itself and around the neighboring external world where it is situated, acts on its own, cooperates with other IoT devices, and exchanges information with other kinds of electronic devices (e.g., mobile device and datacenter) and human users [40].

**Driver roadmapping methodology used by system integration.** The roadmap process in ITRS2.0 is summarized in Figure SYSINT1. (i) Calibration data comes from sources such as published data from web searches, specification documents, datasheets and whitepapers from IC companies, teardown reports, and high-level comments from industry collaborators. (ii) Function categories are obtained by clustering the analysis of IC components. Based on the categorization, we create abstract block diagrams as system models. We also analyze the components and predict how metrics such as maximum operating frequency, die area, number of antennas, number of sensors, etc. evolve over the roadmap’s 15-year horizon. Finally, we produce a roadmap for system-level metrics based on the projected metrics and the abstract block diagrams.

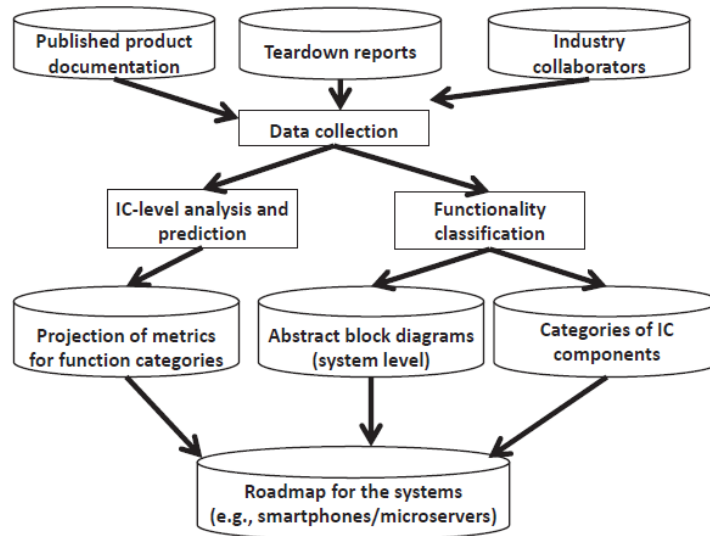


Figure SYSINT1: Flow of data collection, analysis, and metric projection in the ITRS 2.0 roadmapping methodology.

### 3. MOBILE DRIVER

In this section, we describe the mobile driver. We describe main features driven by future applications, along with projections of key system metrics. We then describe key technology challenges and potential solutions to which we map the projected system metrics.

#### 3.1 KEY METRICS TO OF MOBILE DRIVER

In recent years, mobile devices, notably smartphones, have shown significant expansion of computing capabilities. Since smartphone systems are built with multiple heterogeneous ICs (e.g., logic, memory, microelectromechanical systems (MEMS), and radio-frequency (RF)), we must understand tradeoffs at the system level. Beyond the current ITRS SOC-CP roadmap, ITRS 2.0 introduces a new mobile driver to comprehend and roadmap metrics at a higher, system level for mobility applications. Figure SYSINT2, based on the Qualcomm Snapdragon family of SOCs [1], illustrates the growth of features and degree of integration in recent application processors (APs). Each new technology generation (aka “node”), which enables reduced computation power (e.g., new instruction set architecture (ISA), new devices, new low-power techniques) or the introduction of new features (e.g., graphic processing unit (GPU) or 1080p video), brings an increased number of vertically-stacked bars in the plot. Figure SYSINT2 shows that the degree of integration after 2008 keeps increasing to meet the demands of (i) higher computation performance, (ii) faster wireless connections, and (iii) richer multimedia capabilities. The increasing number of heterogeneous components (RF, logic, memory and MEMS) complicates the system design and blocks form factor reductions, while increasing the smartphone design cost and power budget.

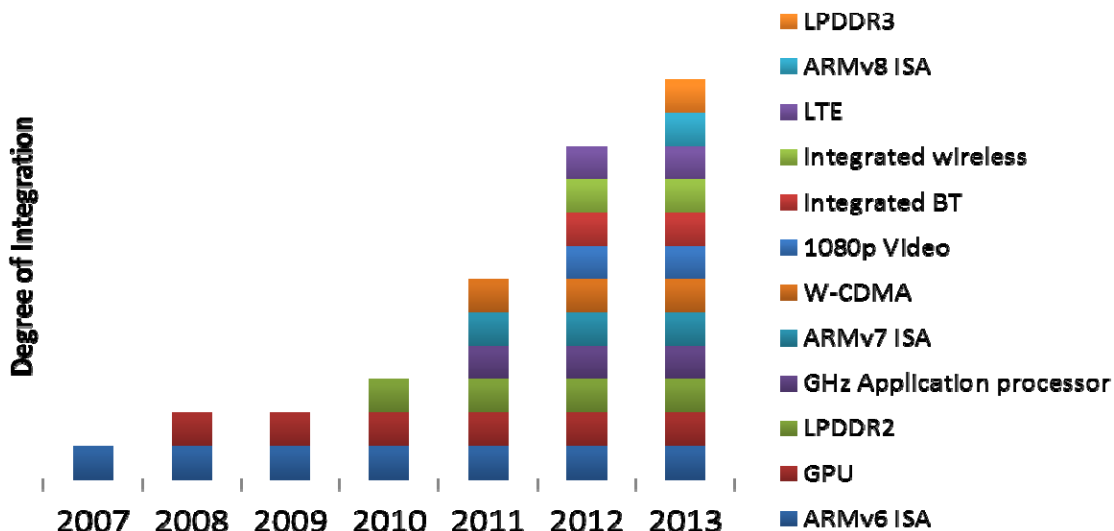


Figure SYSINT2: Increasing degree of integration in mobile application processors (Qualcomm Snapdragon™ family) [1].

Table SYSINT2 projects key system-level metrics of the mobile driver. Input Metrics in the table correspond to system metrics that are projected into the future. Output Metrics are implied by the trajectories of the Input Metrics. Baseline power growth for each IC component in the mobile driver is the same 7% per year that is specified for IC-level power in the 2013 ITRS roadmap of the SOC-CP product. A system (board-level) power projection (5% growth in power per year) is shown in Figure SYSINT3(a).<sup>1</sup> A 4.5W power management gap, relative to a system maximum power requirement of 4W, is projected to exist at the 15-year horizon. The power management gap for board-level power leads to a number of design challenges (heat and thermal/thermomechanical design, battery life, etc.). We expect that extremely aggressive low-power design techniques will need to be applied to IC components in the mobile driver to address power management challenges. Figure SYSINT3(b) shows a projection for another output metric in Table SYSINT2, namely, board area. An area gap of up to 46cm<sup>2</sup> (relative to a 60cm<sup>2</sup> limit)<sup>2</sup> is seen by the end of the roadmap.

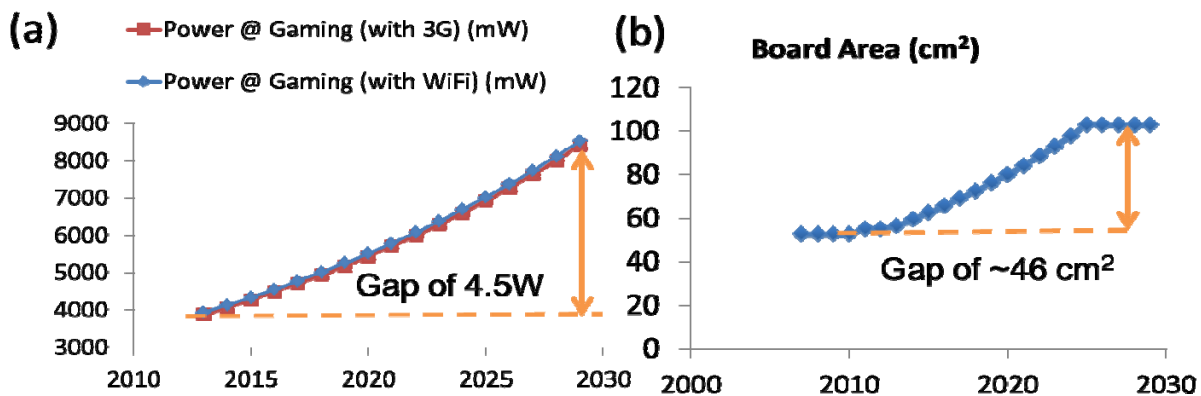


Figure SYSINT3: Implied requirements for mobile driver board area and system power.

Figure SYSINT4(a) shows the scaling of the number of pixels in the mobile driver displays. Display pixels of this driver are driven by high definition standards (e.g., 720p, 1080p, 4K, etc.). Increase in the display size as well as scaling of GPU cores increase the memory bandwidth requirement as shown in Figure SYSINT4(b). By 2029, ultra HD resolutions of 7680 × 4320 and large number of GPU cores could potentially increase memory BW requirements to 61.9GB/s. The rapid growth of bandwidth demands for system level interconnects and off-device interconnects is considered to be a challenge for the mobile driver design.

<sup>1</sup> The annual power growth rate of each component is assumed to be fixed. The SI focus team assumes the power of communication modules (RF, modem, WiFi, etc.) does not increase rapidly after new communication standard is introduced since new low power technologies should control the power under budget.

<sup>2</sup> Board area limit is calibrated with Apple iPhone 5.

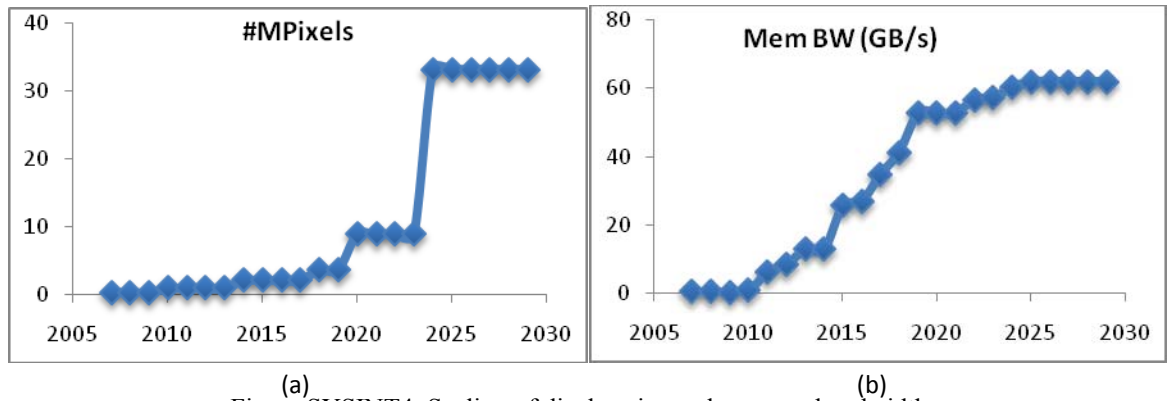


Figure SYSINT4: Scaling of display size and memory bandwidth.

The SI focus team has selected the metrics listed in Table SYSINT2 to develop a technology roadmap for the mobile driver. Through the projection in Table SYSINT2, the SI focus team is able to discover the major technology challenges that must be overcome to address the diverse feature requirement of this mobile driver.

Table SYSINT2: Summary of scaling trends of the mobile driver.

	Year	2007	2010	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029
Input Metrics	#AP cores	1	2	4	4	4	6	9	12	18	18	18	25	28	32	36	33	30	27	25
	#GPU cores	-	-	-	-	6	12	19	30	49	58	69	110	141	189	247	259	273	290	303
	Max freq. (GHz) <sup>3</sup>	0.6	1.5	2.5	2.6	2.7	2.8	2.9	3	3.2	3.3	3.4	3.6	3.7	3.8	4	4.2	4.3	4.5	4.7
	#MPixels <sup>4</sup>	0.307	0.922	0.922	2.1	2.1	2.1	2.1	3.7	3.7	8.8	8.8	8.8	8.8	33.2	33.2	33.2	33.2	33.2	33.2
	Mem BW (GB/s) <sup>5</sup>	-	-	12.8	12.8	25.6	26.9	34.8	41.3	52.6	52.6	52.6	56.7	57.3	60.2	61.9	61.9	61.9	61.9	61.9
	#Sensors	4	8	10	12	14	14	16	16	20	20	20	20	21	21	21	21	22	22	22
	#Antennas	6	8	10	11	11	11	13	13	13	13	14	14	15	15	15	15	15	15	15
	#ICs	8	12	9	7	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10	7-10
	Cellular data rate (MB/s) <sup>6</sup>	0.048	1.70	12.50	12.50	12.50	12.50	12.50	21.63	21.63	40.75	40.75	40.75	40.75	40.75	40.75	40.75	40.75	40.75	40.75
	WiFi data rate (Mb/s) <sup>7, 8</sup>	6.75	75	75	867	867	867	867	867	867	867	7000	7000	7000	28000	28000	28000	28000	28000	28000
Output Metrics	Power (w/ 3G) (W)	3.52	3.67	3.82	4.00	4.2	4.42	4.64	4.87	5.12	5.37	5.64	5.92	6.22	6.53	6.86	7.20	7.56	7.94	8.48
	Board area (cm <sup>2</sup> ) <sup>9</sup>	53	53	57	59	62	66	69	73	76	80	84	89	93	98	103	103	103	103	103

### 3.2 KEY CHALLENGES AND PROMISING SOLUTIONS

Several challenges exist in the development of the mobile driver, based on the projection of system metrics. In Table SYSINT3, the mapping between these challenges and potential solutions are summarized. In addition, the corresponding quantitative metrics are noted in Column 2 of the same table. Table SYSINT4 shows the timeline for the potential solutions. We note that timeline of potential solutions are not very clear beyond 2020.

<sup>3</sup> Max. frequency is defined as the highest operation frequency of components in the system.

<sup>4</sup> The #MPixel is defined as the pixel number (in million) of display. It is modeled based on the following speculative timeline on display formats [24] [25]: VGA in 2007, HD720 in 2010, HD1080 in 2014, WQHD in 2018, 4K in 2020, 8K in 2024.

<sup>5</sup> The memory bandwidth is defined as the bandwidth between APs and the main memory system. It will be driven by bandwidth-hungry applications, such as 3D display with double (120Hz) refresh rate, high-resolution imaging, GPU GFLOPS (follows from the projections from the SOC-CP model) and display, and multimedia features. The steps of memory bandwidth scaling are synchronous with #MPixel scaling due to the correlation.

<sup>6</sup> The cellular standard is modeled based on a speculative timeline on communication standards: 3G in 2007 [20], HSPA in 2010 [20], LTE in 2013 [21], LTE with 2x2 MIMO in 2018 [20], and 4x4 MIMO in 2020 [20]. The average growth rate is ~1.3x per year.

<sup>7</sup> The WiFi data rate is modeled based on the following speculative timeline on industrial standards: 802.11a/b/g in 2007, 802.11n in 2010, 802.11ac in 2014 [20], 802.11ad in 2021 [20], and WirelessHD 1.1 in 2024 [20]. The average growth rate is ~1.4x per year.

<sup>8</sup> The increasing bandwidth of cellular WiFi connections is a new challenge to power management since the transmission power is expected to increase. Multiple-input and multiple-output (MIMO) technology is expected to address this power challenge by improving the transmission power efficiency.

<sup>9</sup> Board area is defined as the total area of the PCB boards where major components are mounted.

## 8 ITRS 2.0 System Integration Chapter

(i) The form factor challenge. As the size of the mobile driver shrinks, especially in thickness, adding new functionalities within a compact form factor becomes very challenging. To address this challenge, the SI focus team has identified two roadblocks in technology development.

- a. The PCB footprint occupied by connectors and components should keep shrinking even though the memory bandwidth requirement and #ICs predicted in Table SYSINT2 increase.
- b. The degree of integration of heterogeneous components, such as logic, memories, non-volatile memories (NVMs), MEMs, RF/analog/mixed-signal (RF/AMS), should keep increasing to reduce the required footprint.

(ii) The system-level power management challenge. Since the predicted board power of the mobile driver will be beyond its 4W target limitation in 2018, system-level power management is an emerging challenge. The roadblocks to address this challenge are as follows.

- a. The increasing memory bandwidth requirement shown in Table SYSINT2 relies on faster signaling and wider system buses, which will increase the board-level power consumption.
- b. Increasing the number of sensors and other IC components require more PCB traces. Shrinking mobile driver form factors are expected to worsen this problem since routing traces will be more complicated, causing interference and higher power.

(iii) The system-wide bandwidth challenge. System-wide bandwidth refers to the bandwidth between application processors and memories or application processors and other peripherals. As the requirements of higher compute performance, #functionalities, and display bandwidth keep growing (as indicated by the scaling of #APs, #GPUs, #sensors, #pixels, and the communication bandwidth), delivering proportionate system-wide bandwidth will become challenging. Another aspect of this challenge will be the tradeoffs between power management and bandwidth.

(iv) Communication bandwidth scaling. This challenge refers to the gaps between required cellular data rate or WiFi data rate and achievable data rates. As the required communication standards supported by a single RF module keep increasing, improvement in transistor scaling should provide the technological capability for the mobile driver to integrate more bands and communication standards within a limited PCB footprint budget.

Table SYSINT3: Key challenges and potential solutions of the mobile driver.

Challenges	Metrics (Description)	Roadblocks	Potential solutions
Form Factor Challenge	#Sensors, #ICs, #Antennas, (#Components ↑)  Memory bandwidth (PCB routing complexity ↑, #connectors ↑)	Increasing PCB footprint occupied by connectors and components	1. Package-level integration 2. Through-silicon via (TSV)-based 3D integration 3. Sequential 3D integration <sup>10</sup>
		Integration of heterogeneous components	1. TSV-based 3D integration <sup>11</sup> 2. Unified logic/RF technology
		Die area explosion due to more functionalities	1. Technology scaling 2. Sequential 3D integration
System-Level Power Management	Max freq., #AP cores, #GPU cores, Memory bandwidth (Power consumption ↑)	High-speed off-processor memory buses	1. TSV-based 3D integration 2. Advanced DRAM (HBM, HMC) <sup>12</sup>
		Increasing #sensors and #IC components	1. TSV-based 3D integration 2. Sensor/MEMS/logic integration
System-wide Bandwidth Scaling	Memory bandwidth, #MPixel, Cellular data rate, WiFi data rate (Bandwidth requirement ↑)	High-speed off-processor memory buses	1. TSV-based 3D 2. Advanced DRAM (HBM, HMC)
		Increasing inter-component bandwidth requirement	1. TSV-based and sequential 3D integration 2. Integrated multi-standard comm. circuits
Communication Bandwidth Scaling	Cellular data rate, WiFi data rate (Bandwidth requirement ↑)	Increasing communication modes/bandwidth requirement (2015) for cellular phone and WiFi	1. Unified logic/RF technology 2. Integrated multi-standard communication circuits
Sensor Pixel Scaling	#MPixel (Pixel density ↑, Optical design complexity ↑)	Pixel dimension scaling limited by optical performance	1. Sensor/MEMS/logic integration (e.g., back-side illumination [12])

<sup>10</sup> Sequential 3D integration refers to 3D integration with fine-pitch TSV (that is close to the gate pitch) while TSV-based integration refers to 3D integration with coarse-pitch TSV's at function block level.

<sup>11</sup> 3D integration refers to the superset of sequential 3D, TSV-based 3D, memory/logic stacking, sensor/logic stacking, etc.

<sup>12</sup> HBM denotes high bandwidth memory; HMC denotes hybrid memory cube.



Table SYSINT4: Timeline for potential solutions for the mobile driver.

	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030
<b>Package-level integration</b>															
Stack thickness < 1mm (0.45mm in 2013) [87]															
<b>TSV-based 3D integration</b>															
Qualification: 3D contact > 100K/mm <sup>2</sup> [88]															
<b>Sequential 3D integration</b>															
Qualification: 3D contact > 5M/mm <sup>2</sup> [88]															
<b>Unified logic/RF technology</b>															
Qualification: Integration of multiband CMOS RF with APs [90] [91]															
<b>Technology scaling</b>															
16/14nm foundry node in 2016															
<b>Advanced DRAM (HBM, HMC)</b>															
Qualification: 100GB/sec [92] [93]															
<b>Integrated multi-standard comm. circuits</b>															
Qualification: Single-chip RF/TX/RX (> 40 bands) [94]															
<b>Sensor/MEMS/logic integration</b>															
Qualification: Multi-MEMS integrated with processing logic [12] [89]															

This legend indicates the time during research, development, and qualification pre-production should be taking place for the solution.

Qualification: the criteria to reach pre-production

Research Required

Development Underway

Qualification/Pre-production

Continuous Improvement



The solution timeline is explained as follows.

**Package-level integration** is a near-term solution for form factor scaling. The thickness of logic/memory package was 0.45mm in 2013 [87] and improvement will continue.

**TSV-based integration** is expected to research 100K/mm<sup>2</sup> 3D contact density by 2020. It will provide higher integration capacity for memory, logic, AMS, etc. [81]

**Sequential 3D** Beyond 2021, sequential 3D integration is expected to continue the Moore's Law scaling by enabling fine-grained interconnection (>5M/mm<sup>2</sup>) between stacked dies [84].

**Unified logic/RF circuits** [51]. This is a long-term trend to tightly integrate multi-standard RF circuits and application processors beyond 2023.

**Technology scaling.** 2016 onwards, technology scaling will continue as the main thrust for the semiconductor industry throughout the roadmap, along with the new device development predicted in the More Moore Chapter.

**Advanced DRAM** will provide more bandwidth than conventional LPDDRx interfaces. We expected to see its deployment in the mobile driver if bandwidth demand increases due to gaming and other visual applications. The bandwidth target is 100GB/sec by 2020 [92] [93].

**Sensor/MEMS/logic integration** will continue to be merged with each other through package-level integration, 3D, or single-chip solutions [89]. Since this involves several technologies to be developed (integration and devices), it is expected around 2020.

**Integrated multi-standard communication circuits.** Due to the increasing cellular bandwidth and multi-standard support, we expect the **integrated multi-standard communication circuits** in 2018 (this does not include application processors, but only modem logics).

## 4. DATACENTER AND MICROSERVER DRIVERS

In this section, we describe the main features, key metrics, key challenges, and potential solutions for the challenges of datacenter and microserver drivers.

### 4.1 KEY METRICS OF DATACENTER AND MICROSERVER DRIVERS

Recent studies of datacenters (e.g., by Doller et al. [2]) suggest that high-performance MPU (MPU-HP) and networking SOC (SOC-NW) products are the main components in datacenters. These products may be implemented either in a single chip or in a multichip module (MCM). Optimized datacenter architecture cannot be achieved with a single chip as its key building block; rather, a co-optimization of storage, interconnects and software is required. Since the raw data stored in datacenters is usually sparse, pre-processing that is typically executed in traditional server cores are precluded, due to energy budget. Besides integrating power-efficient cores to be within an energy budget, datacenters require high

## 10 ITRS 2.0 System Integration Chapter

bandwidth and accessibility for local memories (mostly non-volatile memories) to execute data-intensive operations. Datacenters are a driver for functionality scaling, lithography and device scaling, high-density integration and packaging, and advanced interconnect solutions. Due to datacenter-specific optimizations and system-level design requirements such as high rack density and cooling capacity, the metrics of servers in datacenters are different from those of server chips in existing products which are comprehended by ITRS.

Some new design challenges to the microserver driver are introduced by their deployments in datacenters. Big data computing requires a drastic reduction in communication latencies to meet an under-100ms requirement, meaning data must be increasingly localized. The collected data suggests that the microserver driver addresses the cost issue by limiting the number of cores per rack unit and the latency issue by localizing user-specific search data. The volume of information in datacenters is anticipated to grow at a very high rate (e.g., double every two years, or even faster). When users search for specific information, latencies can be on the order of tens of milliseconds because datacenters typically store information in a highly distributed manner. As datacenters grow in size, communication latencies increase along with power consumption. To limit power and temperature of datacenters, companies are forced to invest huge amounts of money to establish and maintain power plants adjacent to datacenters, and to construct datacenters in geographies with “natural refrigeration”. There is a limit to such investment in power plants and cooling. Cooling costs, which can reach over 35% of electricity costs, continue to rise in server farms and datacenters. This creates a need to reduce the number of cores and operating frequencies to limit this cost.

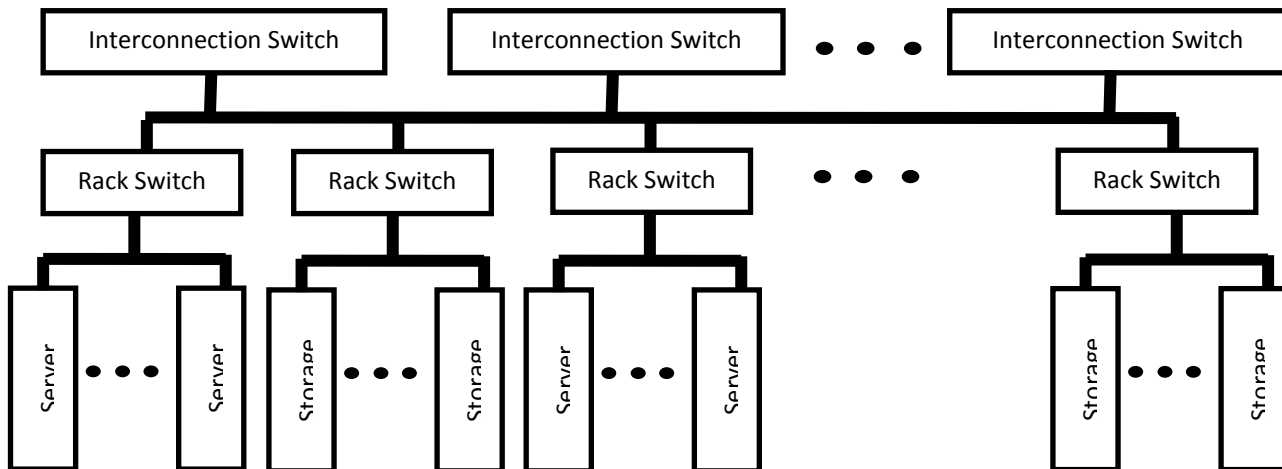


Figure SYSINT5: Datacenter block diagram.

Figure SYSINT5 shows a datacenter block diagram. As discussed above, latency and bandwidth dominate the datacenter metrics roadmap. Electrical networking components will eventually become optical, non-volatile memory will take over storage from mechanical disks, and the storage hierarchy and overall server topology are expected to flatten as the roadmap advances.

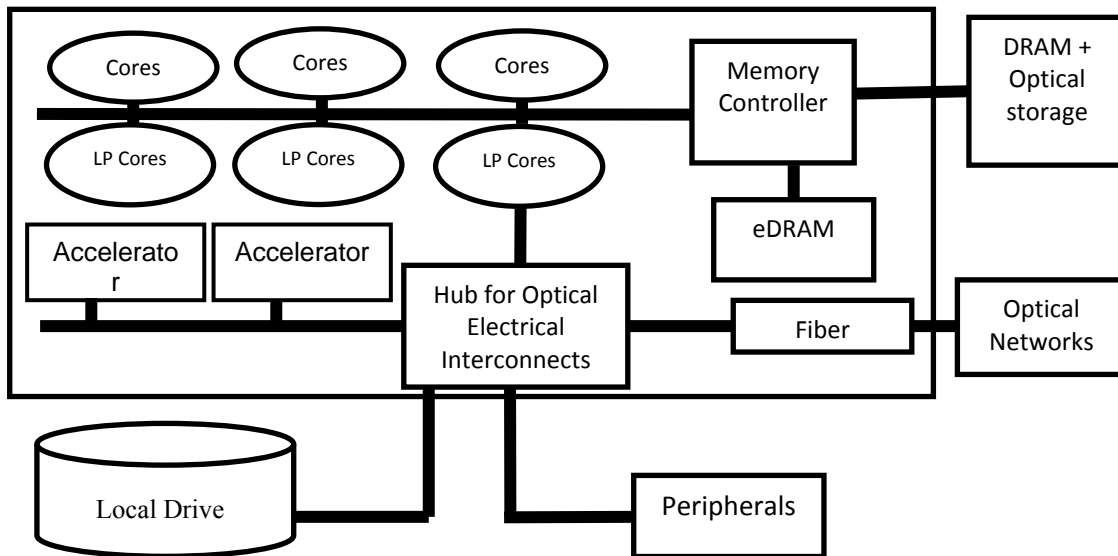


Figure SYSINT6: Microserver block diagram.

Computational power density challenges of microservers may imply SOC/ASIC-like trajectories with more light / low-power (LP) cores, peripherals, application-specific accelerators, complicated inter-IP communication fabrics, and higher integrated network bandwidth to address system throughput requirements under strict form-factor and power constraints. Figure SYSINT6 shows a block diagram of microservers in 2020 and beyond. The microserver organization integrates optical interconnects and low-power and heterogeneous cores.

To reduce operating costs, datacenters and microservers must maximize the number of cores in a rack unit subject to power and thermal constraints. Form factor, energy-efficiency, and networking throughput are important for these drivers owing to strict latency requirements. The microserver design is further challenged by form factor needs. As a consequence, demand for reduced form factor and system design efforts drive the integration of the MPU and the chipset. Compared to a 1U server (MPU-HP in ITRS), a microserver has a higher degree of integration as it includes on-chip Ethernet and peripheral hubs. Recent MPUs for microservers integrate application-specific accelerators to improve energy efficiency. Hence, high integration of functionalities is also a challenge for both datacenters and microservers. The SI focus team has selected the metrics listed in Tables SYSINT5 and SYSINT6 to develop technology roadmaps of datacenters and microservers, respectively.

Table SYSINT5: Summary of scaling trends of datacenters.

Year	2010	2014	2015	2017	2019	2021	2023	2025	2027	2029
# Cores (K)	64	300	<b>360</b>	1044	3008	4935	5825	<b>7578</b>	8967	10602
Storage (PB) <sup>13</sup> [68] [26]	20	100	<b>300</b>	1559	4676	14029	42088	<b>126264</b>	378792	1136377
Area (MSF) [66] <sup>14</sup>	0.2	0.5	<b>0.5</b>	0.9	1.6	2.2	2.2	<b>2.42</b>	2.42	2.42
Power [27] (MkWh) <sup>15</sup>	777.2	770.6	<b>779.8</b>	839	1004.7	1137.2	1226.1	<b>1380.7</b>	1635.6	2044.3
Switch [28] BW(Tb/s) <sup>16</sup>	100	631	<b>1000</b>	2512	6309	10000	15849	<b>25119</b>	39811	63096
GFLOPS/W <sup>17</sup>	0.4	1.7	<b>2.4</b>	4.9	10	17	24	<b>33.9</b>	47.9	67.8
IU/rack <sup>18</sup>	40	40	<b>40</b>	40	40	40	40	<b>40</b>	40	40
Cores/socket <sup>19</sup>	8	15	<b>18</b>	29	47	59	74	<b>93</b>	117	147
Power/IU (W) <sup>20</sup>	700	700	<b>700</b>	700	700	700	700	<b>700</b>	700	700

<sup>13</sup> The total data amount hosted in the datacenter.

<sup>14</sup> The area of one server building. A datacenter is composed with multiple server buildings.

<sup>15</sup> Total power consumption of a datacenter.

<sup>16</sup> The total switching capability of a datacenter.

<sup>17</sup> Giga flops per watt, which measures the energy efficiency of a datacenter.

<sup>18</sup> Server unit density of a rack.

<sup>19</sup> #Processor cores installed to a single socket on the server main board.

<sup>20</sup> Power consumed by a single server unit.

## 12 ITRS 2.0 System Integration Chapter

Memory / 1U (GB) <sup>21</sup>	8	24	<b>32</b>	45	64	76	91	<b>108</b>	129	154
Power/socket [29] <sup>22</sup>	180	168	<b>165</b>	159	153	149	145	<b>141</b>	137	133
NW BW/1U (Gb/s) [30] <sup>23</sup>	1	10	<b>40</b>	40	100	100	100	<b>400</b>	400	400
NW BW/rack (Gb/s)	40	400	<b>1600</b>	1600	4000	4000	4000	<b>16000</b>	16000	16000
Rack switch BW (Gb/s) [31] <sup>24</sup>	10	400	<b>1200</b>	1200	3000	3000	3000	<b>12000</b>	12000	12000
Power efficiency <sup>25</sup>	0.5	0.54	<b>0.55</b>	0.57	0.59	0.61	0.63	<b>0.65</b>	0.67	0.69
Communication power (MkWh) <sup>26</sup>	2.19	13.83	<b>21.91</b>	55.05	138.26	87.66	138.93	<b>220.19</b>	348.97	553.08
Storage power (MkWh) [32] <sup>27</sup>	0.001	0.00876	<b>0.0657</b>	0.259	0.449	0.778	1.347	<b>2.334</b>	4.043	7.004
Cooling power (MkWh) <sup>28</sup>	18.51	34.36	<b>38.99</b>	55.02	86.13	237.31	264.26	<b>307.03</b>	374.89	482.56
#Users (B) [33] <sup>29</sup>	0.44	1.32	<b>1.74</b>	3.01	5.21	9.03	15.63	<b>27.07</b>	46.89	81.21
Data Upload (GB/month/user) [34] [67] <sup>30</sup>	4	6	<b>14</b>	43	75	129	224	<b>389</b>	673	1166

Table SYSINT6: Summary of scaling trends of microservers.<sup>31</sup>

Year		2010	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029
Input Metrics	#MPU cores/rack unit	8	16	19	<b>23</b>	27	32	38	45	<b>51</b>	57	64	72	81	<b>91</b>	102	114	128	<b>144</b>
	Max freq. (GHz)	2.8	3.46	3.6	<b>3.74</b>	3.89	4.05	4.21	4.38	<b>4.56</b>	4.74	4.93	5.13	5.34	<b>5.55</b>	5.77	6	6.24	<b>6.49</b>
	DRAM cap. (GB)/rack unit	32	128	202	<b>319</b>	504	796	1258	1988	<b>3141</b>	4963	7842	12390	19576	<b>30930</b>	48869	77213	121997	<b>192755</b>
	DRAM BW (GB/s) <sup>32</sup>	10.6	51.2	64.8	<b>82</b>	103.7	131.2	166	210	<b>265.6</b>	336	425	537.6	680	<b>860.1</b>	1088	1376.2	1740.8	<b>2202</b>
	Off-MPU BW (GB/s) <sup>33</sup>	25.6	64	82	<b>105</b>	135	173	222	285	<b>337</b>	399	472	558	660	<b>781</b>	924	1093	1293	<b>1530</b>
	MPU freq. × #Cores (GHz)/rack unit	22	55	68	<b>86</b>	105	130	160	197	<b>233</b>	270	316	369	433	<b>505</b>	589	684	799	<b>935</b>

The projections in Tables SYSINT5 and SYSINT6 imply technology challenges to deployment of datacenter and microservers for big data and server-type applications.

### 4.2 KEY CHALLENGES AND PROMISING SOLUTIONS

A major trend of datacenters is consolidation of distributed smaller datacenters to centralized gigantic datacenters. The first observation is the recent migration of data hosts from local storage to cloud operated by major cloud computing service providers (e.g., Google and Amazon AWS) [52] [53]. Major technology requirements in datacenter applications, such as latency, dependability, and scalability, have been driving the consolidation. Due to these technology requirements, a qualified local datacenter might be too expensive to maintain, which is also accelerating the consolidation and the migration to cloud. As an example, we have observed significant effort toward datacenter consolidation in government

<sup>21</sup> Main memory installed to a server unit.

<sup>22</sup> Power consumed the cores on a single socket.

<sup>23</sup> Network bandwidth connected to a server unit.

<sup>24</sup> Switching capacity of within a single rack.

<sup>25</sup> Ratio of power consumed within a datacenter to the power delivered from grid to the datacenter.

<sup>26</sup> Power consumed by networking and switching.

<sup>27</sup> Power consumed by storage.

<sup>28</sup> Power consumed by cooling facility.

<sup>29</sup> Users served by a datacenter.

<sup>30</sup> Monthly data amount uploaded by user.

<sup>31</sup> Metrics are defined for a single server unit.

<sup>32</sup> Bandwidth between main memory and MPUs.

<sup>33</sup> Bandwidth from MPU to peripheral. E.g., the BW through PCI Express to network or storage.

segments.<sup>34</sup> We believe the trend of consolidation will continue as the Internet services and IoT services demand lower latency, dependable, and scalable datacenters.

However, datacenter consolidation is faced with some showstoppers. The analysis from [57] indicates that building diverse new services is slowing down the consolidation because the service providers may intensively rely on localized and customized datacenters to enable differentiated services from competitors. The technology challenges of hardware are also slowing down the consolidation. To continue the consolidation, the following challenges need to be addressed: service latency, space density, power, and integration. These challenges are also confirmed by several observations regarding the major cloud-service companies: **Rapid growth of server nodes in a datacenter.** We have noticed that the number of servers (in terms of rack) of Amazon AWS has grown 27% per year recently [59]. This is mapped to the space density and the integration challenges. **Better power efficiency.** The power consumption of Google data center has decreased by 10% per year recently [60]. **Service latency gap between cloud and local datacenters.** The performance comparison between local storage and Amazon AWS in [61] shows the latency overhead (1~2ms) of cloud datacenters could be a significant performance degradation when the expected service latency is of several millisecond level.

Table SYSINT7 lists the challenges of the datacenter and microserver drivers, organized into latency, power management and integration challenges. A mapping of the challenges to quantitative metrics (Table SYSINT5 and Table SYSINT6) is given in Column 2 of Table SYSINT7, and Table SYSINT8 shows the timeline of the potential solutions. We note that timeline of potential solutions are not very clear beyond 2020.

- (i) The service latency challenge rises in datacenter/microserver design because of the crucial requirement for service latency. The research in [3] proposes much more pessimistic metrics (from 50<sup>th</sup> percentile to 99<sup>th</sup> percentile latency) to ensure service quality could be guaranteed when “Big Data” is hosted. To address this application requirement, solutions are expected from a wide range of providers.
  - a. Since network performance will dominate service latency, high-radix photonics switching networks are expected to be introduced to address the internode bandwidth requirement.
  - b. To host Big Data, conventional memory architectures will be unable to address access time requirements for Big Data. Spindle-based hard drives will be replaced by storage-class memories (SCM).
  - c. To improve the intra-node communication performance (e.g., for MPU to memories in Table SYSINT6 or memories to NVMs), better solutions for heterogeneous integration are expected.
- (ii) To provide sufficient computing resources with MPU cores and application-specific accelerators,
  - a. Moore’s Law should continue the transistor scaling so that more functionalities could be hosted in the same die area while avoiding power increases that result in too much overhead to the cooling equipment.
  - b. Better memory integration (e.g., memory-over-logic) in each computing node is expected to ease the power management challenge by reducing the power impact.
  - c. Advanced power management techniques such as adaptive power management with on-die power sensors [6] are expected to be developed to address the power management issue.
- (iii) The electro-optical integration challenge. Since the power and performance requirements of datacenter are both crucial, highly-integrated photonic inter-node networks are expected by 2020 [3]. Since the electro-optical interfaces are distributed all over the datacenter, it is necessary to develop on-chip light sources and on-chip photonic modulators and detectors to reduce the power, space, and performance overhead due to off-chip converters for electro-optical interfaces.

Table SYSINT7: Key challenges and potential solutions of the datacenter and microserver drivers.

Challenges	Metrics	Roadblocks	Potential solutions
Service Latency Challenge	#MPU cores/rack unit, DRAM cap./ rack unit, Max Freq., DRAM bandwidth, Off-MPU bandwidth (Performance requirement ↑)	Low #hop connections	High-radix network [3]
		Low bit/J transmission	Silicon photonics [3] which can deliver higher switch BW and lower pJ/bit
		High performance memory architecture	SCM to replace hard drives [3]
		High storage bandwidth	Distributed storage nodes [3]
		Encrypted data processing	Distributed compression and encryption engines [3]
		NVM reliability	Novel memory devices [3] Control algorithm [4]
Node Density/Cooling/Power Management	#MPU cores/rack unit, DRAM cap./ rack unit, Max Freq., DRAM bandwidth, Off-MPU	Die areas increase due to more functionalities	Moore’s Law scaling

<sup>34</sup> The Department of Health and Human Services (HHS) of the United States proposed a 20% reduction plan for their datacenters from 2010 to 2015 [55].

## 14 ITRS 2.0 System Integration Chapter

Challenge	bandwidth (Power consumption ↑)	Low power processor architecture	64-bit ARM core [5]
		Lack of one-fits-all processor architecture	Modularized processor 3D stacks (TSV-based) [5]
		Power management for different application context	Integrated on-die power sensors [6]
Electro-Optical Integration Challenge	DRAM bandwidth, Off-MPU bandwidth (Bandwidth requirement ↑)	On-chip light source	Silicon compatible laser source [100]
		On-chip detector / modulator	Silicon compatible laser source [100]

Table SYSINT8: Timeline of potential solutions of the datacenter and microserver drivers.

	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030
<b>Silicon photonics</b>															
Qualification: (1) Monolithic integration of silicon photonics [96] [95] (2) > 200GB/s directly off MPU															
<b>Silicon compatible laser source</b>															
Qualification: Light source on silicon substrate [100]															
<b>High-radix networks</b>															
Qualification: Optical switches with > 1K ports [97]															
<b>Storage class memory (SCM)</b>															
Qualification: sub-us latency [98]															
<b>Distributed compression, encryption engines</b>															
Qualification: built-in 10Gbps-level IPsec/SSL accelerator															
<b>Novel memory devices</b>															
Qualification: achieving < sub-us latency [98] [99]															
<b>64-bit ARM processor core</b>															
Ready solution															
<b>Modularized processor 3D stacks (using TSV-based 3D)</b>															
Qualification: 3D contact > 100K/mm <sup>2</sup> [88]															
<b>On-die power sensors</b>															
Qualification: 25x power efficiency improvement [101]															

This legend indicates the time during research, development, and qualification pre-production should be taking place for the solution.

Qualification: the criteria to reach pre-production

Research Required

Development Underway

Qualification/Pre-production

Continuous Improvement



The solution timeline is explained as follows.

**Silicon photonics.** Optical device development and integration challenges of silicon photonics will be the main showstoppers for deployment of this technology. By 2020, we expect monolithic integration could provide more than 200GB/s directly between the cores/cache [95] [96]. The requirements of high-bandwidth and low-energy to process huge amounts of data will drive pre-production by 2020 [97] [104].

**Silicon compatible laser source.** Prior to deployment of silicon photonics, silicon compatible laser source needs to be ready by 2018.

**High-radix network.** Since the service latency is critical to the datacenter design, high-radix networks provide fewer hops among server nodes. We expect more than 1000 ports in a single optical switch by 2020.

**Storage-class memory.** Replacing mechanical hard drives with non-volatile storage is required for bandwidth and energy efficiency requirements (sub- $\mu$ s latency [98] [99]). The full deployment of SCM is expected in pre-production by 2017.

**Distributed compression, encryption engines.** The datacenter is expected to enhance data security and compress sparse data to save both bandwidth and storage. We have observed the existence of the encryption engines within processor cores. We expect higher throughput in the future.

**Novel memory devices.** New memory technologies provide promising solutions for service latency (sub- $\mu$ s latency), energy efficiency, and server node density scaling. Flash memory, as the current mainstream NVM storage media, will need to be replaced due to durability, integration, and performance issues. Novel memory devices are promising candidates for the datacenter driver. However, the deployment may be constrained by device research progress. We expect novel memory devices to be in pre-production by 2019.

**64-bit ARM cores.** The ARM architecture is popular in the mobile market for its energy efficiency. This solution is ready for the datacenter.

**Modularized processor 3D stacks [5].** From [65], we expect datacenter design would require diverse types of server nodes (e.g., computation and storage). Modularized 3D stacking (based on TSV-based 3D) provides a solution to reduce the NRE cost due to small processor module amount for each node type. The interposer-based integration in [5] is ready, but we require development and optimization of the system. Deployment is constrained by the components to be integrated (e.g., novel memory devices and silicon photonics) and is expected to be in pre-production by 2020.

**On-die power sensor.** This technology combines thermal sensors and DVFS information for real-time power monitoring. Although monitor circuits are ready, aggressive and pervasive self-management is expected in 2019 to achieve 25x power efficiency improvement [6].

## 5. IoT DRIVER

The IoT driver represents the extreme low power requirement in the semiconductor industry. The operations of IoT devices are constrained by battery sizes and long operation life time. Instead of roadmapping computation performance, we predict power efficiency instead. The key metrics are operation duty cycles, suspend current, Ion/core frequency, and life time.

### 5.1 KEY METRICS OF IOT DRIVER

The exemplar block diagram of the IoT driver is shown in Figure SYSINT7. The architecture of the IoT driver is mission-oriented. The hardware implementation could be a subset of the listed blocks, or constructed by programmable logic (e.g. FPGA). The integration approach varies from board level (e.g., batteries), package level (sensor/MCU package), 3D stacking (MCU and memory stacking), or chip level integration (baseband and MCU).

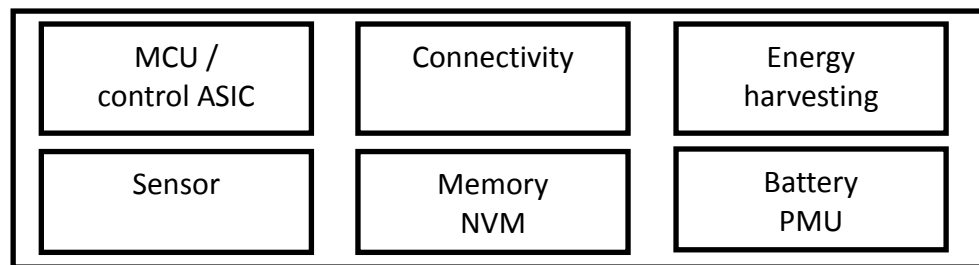


Figure SYSINT7: Exemplar block diagram of IoT device.

Based on the baseline configuration shown in Figure SYSINT7, the SI focus team has identified four categories of metrics to be tracked for IoT Driver: power, form factor, performance, and peripheral. The categories and metrics tracked by the SI focus team are listed in Table SYSINT9. **Power:** The trend of VDD follows latest ORTC voltage scaling. The SI focus team roadmaps the mass deployment of energy harvesting to start at 2019. Due to the constrained energy source, the SI focus team roadmaps the suspend current of the MCU to scale at 0.85x per year. The efficiency and power density of integrated DC-DC converters are roadmapped to scale at 1.01x and 1.08x per year, respectively. The connectivity power consumption, constrained by energy source and communication throughput, is also roadmapped to scale down per year. The peak current (mainly constrained by the storage and source of energy) scales at 0.62x per year and the transmission power per bit scales at 0.63x per year. **Form factor:** the system form factors of IoT devices are mainly constrained by the battery and the passive components. The SI focus team roadmaps the form factor will gradually scale down at 0.8x per year after 2018 due the mass deployment of energy harvesting and higher integration of components. **Performance:** the SI focus team expects MCU performance to scale conservatively relative to MPU or SOC because the strict power efficiency constraint and form factor constraint. The MPU clock scales at 1.085x per year before 2020 and 1.017x per year afterward. The power efficiency metric in terms of  $I_{ON}$  per MHz, scales at 0.85x per year until 2022 and 0.93x afterward. The number of MCU core is fixed at single cores, and the performance improvement (MCU DMIPS) increases linearly. **Peripheral:** the SI focus team roadmap expects the number of sensors to increase rapidly in the 2015 to 2017 timeframe. The increase will saturate afterward due to the limitation of integration effort. The power consumed by sensors will scale at 0.7x from 2015 to 2017, 0.85x from 2018 to 2021, and 0.92x afterward, due to the energy constraint. Conversely, the density of battery or other form of energy storage will increase at 1.07x per year to support the increase of connectivity throughput, MCU performance, and integrated sensors.

Table SYSINT9: Summary of scaling trends of IoT driver.

Categories	Year	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029
Power	Energy source (B = battery; H = energy harvesting)	B	B	B	B + H	B + H	B + H	B + H	B + H	B + H	B + H	B + H	B + H	B + H	B + H	B + H
	Lowest VDD (V) <sup>35</sup>	0.8	0.8	0.75	0.75	0.7	0.7	0.65	0.65	0.65	0.55	0.55	0.55	0.45	0.45	0.45
	Deep suspend current [35] (nA) <sup>36</sup>	100	85	72	61	52	44	38	32	27	23	20	17	14	12	10
	DC-DC efficiency (%) <sup>37</sup>	80%	81%	82%	85%	86%	87%	88%	88%	89%	90%	91%	92%	93%	94%	95%
	DC-DC power density (W/mm <sup>2</sup> ) <sup>38</sup>	10.00	10.80	11.66	12.60	13.60	14.69	15.87	17.14	18.51	19.99	21.59	23.32	25.18	27.20	29.37
	Peak Tx/Rx current (mA) <sup>39</sup>	50.00	31.05	19.28	11.97	7.44	4.62	2.87	1.78	1.11	0.69	0.43	0.26	0.16	0.10	0.06
	Tx/Rx power per bit ( $\mu$ W/bit) <sup>40</sup>	2.480	1.552	0.972	0.608	0.381	0.238	0.149	0.093	0.058	0.037	0.023	0.014	0.009	0.006	0.004
Form factor	Module footprint (mm <sup>2</sup> ) <sup>41</sup>	500	500	500	350	280	224	179	143	115	92	73	59	47	38	30
Performance	MCU #Cores	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	MCU I <sub>ON</sub> / Operation frequency [36] ( $\mu$ A/MHz) <sup>42</sup>	30.0	25.5	21.7	18.4	15.7	13.3	11.3	9.6	8.9	8.3	7.7	7.2	6.7	6.2	5.8
	Max MCU Frequency (MHz) [37]	200	217	235	255	277	301	306	311	316	322	327	333	338	344	350
	MCU Flash Size (KB) [38] <sup>43</sup>	1024	1024	1024	1024	2048	2048	4096	4096	4096	4096	8192	8192	8192	8192	8192
	MCU DMIPS <sup>44</sup>	200	220	242	266	293	322	354	390	429	472	519	571	628	690	759
Peripheral	#Sensors <sup>45</sup>	4	4	8	8	10	10	12	12	12	12	13	13	13	13	13
	Max Sensor Power [39] ( $\mu$ W) <sup>46</sup>	2850	1995	1397	1187	1009	858	729	671	617	568	522	480	442	407	374
	Battery Power Density (Watt-Hr/Liter)	561	600	642	687	735	787	842	901	964	1031	1104	1181	1263	1352	1447

### 5.2 Key challenges and promising solutions.

From the current roadmapped metrics, the SI focus team has identified the following technology challenges for the IoT driver. The key challenges and their technology solutions are summarized in Table SYSINT10, and the timeline for these solutions are summarized in Table SYSINT11.

(i) Transistor device design and scaling (“reverse” Moore) challenge: due to the extremely low power requirements for both communication and computing, the conventional Moore’s Law scaling of transistors is not seen in the IoT driver. The leading technology node of the IoT driver falls behind other drivers (e.g., conventional MPU and SOC) by more than two nodes. This reverse trend is constrained by (i) the IoT device’s sensitivity to energy loss during its suspend mode; (ii) design technology that has not provided enough matched leakage reduction for the increasing transistor number. This reverse trend is a showstopper for increasing the IoT performance to meet the application requirement. Promising solutions to the transistor device scaling challenge include:

<sup>35</sup> Lowest VDD consumed by the components in the system.

<sup>36</sup> Suspend current of MCU, of which the always-on blocks will dominate the system suspend current.

<sup>37</sup> Conversion efficiency of the integrated DC-DC converter at nominal voltage.

<sup>38</sup> The spatial efficiency of the converter, which is defined as the output power divided by the circuit area. For the DC-DC converter specification, we obtain the data from a major circuit conference [51]. However, the data points are still sparse so we will keep tracking the trend.

<sup>39</sup> Peak current consumed by the connectivity interface.

<sup>40</sup> Transmission power for each bit, which is the energy efficiency metric for communication.

<sup>41</sup> Physical footprint of the system.

<sup>42</sup> Current consumption normalized to the operation frequency, which is the energy efficiency metric for the computation.

<sup>43</sup> Flash (or over NVM) size to store programs, configuration, and data

<sup>44</sup> DMIPS benchmark as the performance metric of MCU.

<sup>45</sup> Number of sensors integrated to the system.

<sup>46</sup> Total power consumed by sensors.



(a) Development of emerging devices is the primary solution for the device scaling with leakage current under control. In the recent nodes, FinFET has shown manageable leakage current while providing improved performance. As the technology nodes advance, lateral and vertical gate-all-around (LGAA and VGAA) transistors are promising since they improve the effective gate length of transistors while the device footprint is shrinking. Carbon nanotube (CNT) is the long term solution to the device scaling since its novel physical structure provides better electrical characteristics.

(b) Device scaling could be relaxed by 3D integration [49]. Emerging memory devices, such as resistive RAM (RRAM) [50], provide better durability against wearing, better integration density, and ease of integration to present back end-of-line (BEOL) technology. By means of the device innovation, the scaling challenge could be addressed.

(ii) IP/sensor integration and scaling challenge: the IoT driver is cost sensitive due to the huge amount of deployment. In order to reduce the system cost and integration effort, we have observed the requirement to reduce the component number in the system (i.e., total cost of bill of material, BOM). However, integrating extra components, such as baseband processors for communication, is constrained by the reversed Moore's Law mentioned above. Power conversion circuits, such as integrated DC-DC converters significantly improve the power efficiency and relax the power system design effort. We have observed the tendency of integrating power conversion circuits into the IoT MCUs, which also has also raised interest in the circuit design community in the recent years [51]. Promising solutions to the IP/sensor integration challenge include:

(a) Heterogeneous 3D addresses different requirements for electrical characteristics among analog, logic, MEMS, and imaging pixels, which allows system designers to integrate different substrates in limited package footprints and simplify the system level design effort.

(b) Design technology, such as on-chip passive components (e.g. inductors) over BEOL and configurable/fine-grained regulation are promising technologies to realize high efficiency on-die regulation, which will improve the power efficiency constrained by harsh limits on energy sources (i.e., energy harvesting or compact battery modules).

(iii) Supply voltage scaling challenge: the supply voltage of IoT drivers has been roadmapped to aggressively scale due to the extremely low power requirement. However, scaling is severely constrained by the two design challenges mentioned above. Since the scaling of device dimensions and threshold voltages are blocked by a lack of effective low power design technology, the supply voltage of the IoT driver is unable to be pushed lower. Meanwhile, the integrated blocks in the IoT driver, such as baseband processors and analog blocks, may have different supply voltage requirements. The constraints from device scaling and the fine-grained power domain may complicate the design of on-chip power distribution systems, which calls for improvement by design technology. Promising solutions to the supply voltage scaling challenge include several emerging computing paradigms.

Near-threshold [45] allows circuits to operate with limited headroom above the threshold voltage by compound approaches of architecture, device electrical characteristics optimization, or control over body biasing. Asynchronous computing replaces clock signals with handshaking mechanism to relax the timing requirement. Stochastic computing replaces parallel binary number representation with serial bit numbers to relax the timing requirement. Approximate computing allows bounded errors to relax timing requirement due to the redundant and error-tolerant nature of the network among IoT devices.

Table SYSINT10: Key challenges and potential solutions of the IoT driver.

Challenges	Metrics	Roadblocks	Potential solutions
Transistor device design and scaling	MCU #Cores, MCU ION / Operation frequency ( $\mu\text{A}/\text{MHz}$ ), MCU Flash Size (KB), Deep suspend current (nA)	Leakage current management	Device: FinFET, LGAA, VGAA, and CNT
		Reliability issues due to logic transistor scaling	
		Threshold voltage scaling	
		Data corruption due to NVM device scaling	Device: Emerging memory devices (e.g., RRAM) Design technology: 3D stacking
IP/Sensor integration and scaling (More than Moore) challenge	#Sensors, Max Sensor Power ( $\mu\text{W}$ ), DC-DC efficiency (%), DC-DC power density ( $\text{W}/\text{mm}^2$ )	Exclusive technology (analog, MEMS, logic... etc.)	Heterogeneous 3D integration
		On-die voltage regulation, scaling of passive components [42], conversion efficiency between different input/output of regulators	Integrated passive components [43], configurable switching capacitor [44]

Supply voltage scaling challenge	Lowest VDD (V), Battery Power Density (Watt-Hr/Liter), Peak Tx/Rx current (mA), Tx/Rx power per bit ( $\mu$ W/bit)	Threshold voltage scaling and performance requirement	Near-threshold computing [45], asynchronous computing [46], stochastic computing [47], approximate computing [48]
----------------------------------	--	---	---

Table SYSINT11: Timeline of potential solutions of the IoT driver.

	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030
<b>New devices: CNT, FinFET, LGAA, VGAA</b>	[Continuous improvement]														
Continuous improvement															
<b>New memory RRAM</b>	[Qualification: 28nm, 4F <sup>2</sup> area, sub-us latency [79]]														
Qualification: 28nm, 4F <sup>2</sup> area, sub-us latency [79]															
<b>Heterogeneous 3D integration (using TSV-based 3D)</b>	[Qualification: 3D contact > 100K/mm <sup>2</sup> [88]]														
Qualification: 3D contact > 100K/mm <sup>2</sup> [88]															
<b>On-chip Passive Components</b>	[Qualification: magnetic material and deep-trench cap.]														
Qualification: magnetic material and deep-trench cap.															
<b>Configurable/fine-grained regulation</b>	[Qualification: > on-chip regulation and multiple power domains [102]]														
Qualification: > on-chip regulation and multiple power domains [102]															
<b>Near-threshold computing</b>	[Qualification: 1.23x dynamic power improvement [103]]														
Qualification: 1.23x dynamic power improvement [103]															
<b>Asynchronous computing</b>	[Qualification: Similar to NTC]														
Qualification: Similar to NTC															
<b>Stochastic computing</b>	[Qualification: Similar to NTC]														
Qualification: Similar to NTC															
<b>Approximate computing</b>	[Continuous improvement]														
Continuous improvement															

This legend indicates the time during research, development, and qualification pre-production should be taking place for the solution.

Qualification: the criteria to reach pre-production

- Research Required
- Development Underway
- Qualification/Pre-production
- Continuous Improvement



The timeline explanations are as follows.

**New devices.** The development of new devices is under continuous development. FinFET is the current low-power solution. According to ITRS 2.0 More Moore Chapter, LGAA is expected in 2019; VGAA is expected in 2021; and CNT is beyond 2030 as the successor of CMOS.

**New memory RRAM.** RRAM (resistive RAM) is expected for relaxation of integration to compact product sizes, low power, high performance, and small footprint (4F<sup>2</sup>). It is expected to be in pre-production by 2020 to achieve the low power requirement when energy harvesting is introduced [78] [79].

**Heterogeneous 3D integration.** The integration of IoT includes heterogeneous IC components (logic, MEMS, sensors, etc.) By 2020, 3D heterogeneous integration is expected in pre-production using TSV-based integration [81].

**Integrated passive components.** Integrated switching regulators (DC-DC) avoid high energy loss due to linear regulation. However, integrated passive components are available now [80] and are expected to provide better quality (e.g., Q-factor) and compact die area.

**The configurable/fine-grained regulators** are expected to be in pre-production by 2018 to meet the low power requirement after energy harvesting is introduced. The IoT roadmap in [82] also points out the low power integration is expected by 2020 to address the application requirement.

**Near-threshold computing** relies on co-optimization of devices, circuits, and architectures. It is expected to be in pre-production by 2022.

**Asynchronous computing and stochastic computing** are expected to be in pre-production by 2020 and 2022, respectively, due to design and verification tool development [46] [47].

**Approximate computing** exists today in the form of variable-precision hardware. More aggressive approximations are expected for error-tolerant applications, such as neuromorphic processors or machine learning applications [83].

## 6. LOOKING AHEAD

In the current System Integration Chapter, three new drivers are identified based on the study and data collection of market-dominant electronics products (mobile, datacenter, and IoT). At the frontier of industry transformation, there are other new applications emerging that will probably become new technology drivers of the semiconductor industry. We summarize these promising applications as follows.

**Automobile Drivers (self-driving cars and drones.)** There are adequate indicators that new drivers will heavily influence the semiconductor industry in the future. Companies like Apple and Google who are typically associated with the mobile industry are putting efforts into the development of autonomous vehicles. Additionally, many traditional automobile companies (Audi, BMW, Ford, and Nissan) have put efforts into the development of self-driving cars. One estimate predicts 10 million autonomous vehicles will be on the road by the year 2020.

Google's self-driving cars use deep learning and convolutional neural networks (CNNs) to process a complex set of video feeds and make intelligent driving decisions. Processing is done both locally with onboard computers, as well as remotely at server farms [69]. However, for autonomous vehicles to fully replace human operators, improvements in the reliability of CNNs must be made. Statistics released by Google indicate that 13 "potential contacts" with other vehicles would have been made in a span of 14 months, had it not been for human intervention [70]. Another problem that must be handled in autonomous vehicles is simultaneous localization and mapping (SLAM). Autonomous vehicles first map a trajectory ahead of time, and then post-process the data to obtain a very accurate map of the environment. In subsequent drives, the vehicles are able to localize themselves within the map very quickly [71].

NVIDIA's Drive PX2 unit hopes to bring powerful neural network capabilities into autonomous vehicles [72]. Motivators for these actions include improvements in reliability, which can greatly benefit safety-criticality necessary for fully autonomous driving. The cost and size of the multitude of sensors necessary for autonomous driving, the ability for neural networks to perform as well as humans, and the capability to make decisions locally are all important requirements for self-driving cars.

Unmanned drones are another form of automated driving that is emerging. Presently, the drone market is restricted by regulations that greatly limit where fully autonomous drone operation is allowed. However, major companies like Amazon are working to create regulations friendlier to autonomous drone operation. Because drones pose a less-imminent threat to humans, and are more limited by weight and power consumption than the systems going into self-driving cars, companies have developed autonomous drone systems from older mobile SOCs. Qualcomm has released a reference platform for unmanned aerial vehicles (UAVs) that is based on the Snapdragon 801 [73].

**Graphic/visual drivers (derived from GPUs.)** Graphics processing units (GPUs) are cards with thousands of simple cores capable of efficiently completing highly parallel tasks. Used in conjunction with CPUs, these hybrid compute systems can offer significant performance improvements over traditional CPU-only systems. GPUs have already been put to use in deep learning, large-scale simulations, and imaging. NVIDIA offers GPUs that meet different metrics, and are largely driven by the specific market they are geared towards. For example, consumer desktops (GeForce), server/workstation (Tesla), and mobile (Tegra) all have different performance capabilities and power requirements. GPUs are ubiquitous in many of the traditional drivers, and will be an important part of using heterogeneous system integration to meet power efficiency and performance standards in the future. [74][75]

**Bio-Chip Drivers.** There is a large push for continuous monitoring devices in health and fitness. Continuous monitoring can improve the speed at which clinical drug trials are conducted, and provide better patient care through enhanced data collection. Key factors in this field are size, power consumption, energy harvesting, and many-channel ADC's and controllers. Brain machine interfaces can also be used to restore motor function after injury. **Lab on-a-chip.** The combination of conventional electronics devices and biomedical components is one of the promising applications. The concepts of "lab-on-a-chip" or "Organs-on-Chips" are expected to be a powerful tool for drug development [62] [63]. **Precision medical devices.** Implanted devices [64] enable tailored medical treatment for patients. These applications will create more challenges of system integration, power and reliability between electronics and biomedical components.

**Fabric Drivers.** Advances in high performance computing (HPC) have allowed it to enter the market for multiple use cases. The major driver is the fabric for these HPC systems. A system-level approach to the problem is being used by Intel and Cavium [76]. These companies had originally focused on being compute providers, but the expenses of the total

system increasingly moved away from the compute portion alone. Providers of HPC fabrics can no longer design a single homogenous product that appeals to many consumers. Instead, they are focusing on a platform that works with all of the divergent applications for HPC, such as modeling and simulation, visualization, machine learning, and data analytics. Additionally, the new fabrics are designed to support scalability from small racks to massive supercomputing centers, as well as both local and cloud-based modes of operation [77].

## 7. SUMMARY

The System Integration Chapter has developed a new roadmapping methodology to identify three new drivers (mobile, datacenter/microserver and IoT) for the semiconductor industry and extract the technology requirements from the new drivers. From the study, the scope of drivers expands from a single chip to the whole system, and new technology challenges such as system bandwidth, power management, integration, are explored.

The System Integration Chapter will continuously track the potential drivers and their impact on the semiconductor industry. The System Integration Chapter invites inputs from users, electronics product manufacturers, and academy for future revisions.

## REFERENCES:

- [1] [http://en.wikipedia.org/wiki/Qualcomm\\_Snapdragon](http://en.wikipedia.org/wiki/Qualcomm_Snapdragon)
- [2] E. Doller et al., "DataCenter 2020: Near-memory Acceleration for Data-oriented Applications", *Proc. Symposium on VLSI Circuits*, 2014
- [3] [https://www.usenix.org/sites/default/files/conference/protected-files/fast14\\_asanovic.pdf](https://www.usenix.org/sites/default/files/conference/protected-files/fast14_asanovic.pdf)
- [4] [http://www.snia.org/sites/default/education/tutorials/2009/spring/solid/JonathanThatcher\\_NandFlash\\_SolidState\\_Storage\\_ReliabilityV1-0-nc.pdf](http://www.snia.org/sites/default/education/tutorials/2009/spring/solid/JonathanThatcher_NandFlash_SolidState_Storage_ReliabilityV1-0-nc.pdf)
- [5] Y. Durand et al., "EUSERVER: Energy Efficient Node for European Micro-servers", *Proc. Euromicro Conference on Digital System Design*, 2014.
- [6] [http://www.theregister.co.uk/Print/2014/06/20/amd\\_25x20\\_power\\_efficiency\\_pledge/](http://www.theregister.co.uk/Print/2014/06/20/amd_25x20_power_efficiency_pledge/)
- [7] <http://www.intel.com/content/www/us/en/data-center/intel-labs-silicon-photonics-demo.html>
- [8] <http://spectrum.ieee.org/semiconductors/design/the-silicon-solution>
- [9] [http://www.hotchips.org/wp-content/uploads/hc\\_archives/hc22/HC22.23.330-1-Alduino-Intel-SiP-Link.pdf](http://www.hotchips.org/wp-content/uploads/hc_archives/hc22/HC22.23.330-1-Alduino-Intel-SiP-Link.pdf)
- [10] <http://www.intel.com/content/www/us/en/research/intel-labs-silicon-photonics-mxc-connector.html>
- [11] <http://isl.stanford.edu/~abbas/presentations/ICCP09.pdf>
- [12] <http://www.aptna.com/news/FSI-BSI-WhitePaper.pdf>
- [13] <http://www.datacenterknowledge.com/archives/2014/06/24/next-gen-intel-phi-coprocessor-to-use-silicon-photonics-interconnect/>
- [14] Hamzaoglu et al., "A 1Gb 2Ghz Embedded DRAM using a 22 nm tri-gate CMOS logic Technology", *Proc. ISSCC*, 2014.
- [15] <http://www.businessinsider.com/growth-in-the-internet-of-things-2013-10>
- [16] <http://www.transparencymarketresearch.com/pressrelease/microserver-market.htm>
- [17] [http://en.wikipedia.org/wiki/IEEE\\_802.11ac](http://en.wikipedia.org/wiki/IEEE_802.11ac)
- [18] [http://en.wikipedia.org/wiki/Wireless\\_Gigabit\\_Alliance](http://en.wikipedia.org/wiki/Wireless_Gigabit_Alliance)
- [19] <http://en.wikipedia.org/wiki/WirelessHD>
- [20] [http://en.wikipedia.org/wiki/List\\_of\\_device\\_bit\\_rates](http://en.wikipedia.org/wiki/List_of_device_bit_rates)
- [21] [http://www.itu.int/ITU-D/arb/COE/2010/4G/Documents/Doc4-LTE%20Workshop\\_TUN\\_Session3\\_LTE%20Overview.pdf](http://www.itu.int/ITU-D/arb/COE/2010/4G/Documents/Doc4-LTE%20Workshop_TUN_Session3_LTE%20Overview.pdf)
- [22] <https://www.comp.nus.edu.sg/~tulika/DAC14.pdf>
- [23] [http://mobiledevices.kom.aau.dk/research/energy\\_measurements\\_on\\_mobile\\_phones/results/](http://mobiledevices.kom.aau.dk/research/energy_measurements_on_mobile_phones/results/)
- [24] [http://en.wikipedia.org/wiki/Display\\_resolution#mediaviewer/File:Vector\\_Video\\_Standards8.svg](http://en.wikipedia.org/wiki/Display_resolution#mediaviewer/File:Vector_Video_Standards8.svg)
- [25] [http://en.wikipedia.org/wiki/8K\\_resolution](http://en.wikipedia.org/wiki/8K_resolution)
- [26] Storage is calibrated to 900PB in 2015 from <http://www.extremetech.com/computing/129183-how-big-is-the-cloud>
- [27] Power is calibrated to 780MkWh in 2015 from <http://www.datacenterknowledge.com/the-facebook-data-center-faq-page-2/>
- [28] Interconnection switch bandwidth is calibrated to 1000Tb/s in 2015 from <http://www.extremetech.com/computing/129183-how-big-is-the-cloud>; <https://storageservers.wordpress.com/2013/07/17/facts-and-stats-of-worlds-largest-data-centers/>
- [29] Power/socket is calibrated to 150W in 2015 from [http://ark.intel.com/products/84683/Intel-Xeon-Processor-E7-8880-v3-45M-Cache-2\\_30-GHz](http://ark.intel.com/products/84683/Intel-Xeon-Processor-E7-8880-v3-45M-Cache-2_30-GHz)
- [30] Network BW/1U is calibrated to 40Gb/s in 2015 from <http://www.dell.com/us/business/p/poweredge-r730xd/pd> ; <http://googlecloudplatform.blogspot.com/2015/06/A-Look-Inside-Google-Data-Center-Networks.html>
- [31] Rack switch BW is calibrated to 1.2Tb/s in 2015 from <http://www.redbooks.ibm.com/redbooks/pdfs/sg247960.pdf> ; <http://www.theplatform.net/2015/04/20/enterprises-gear-up-for-10ge-racks-100ge-spines/>
- [32] Storage power is assumed to be 25W/PB from <https://storageservers.wordpress.com/2013/07/17/facts-and-stats-of-worlds-largest-data-centers/>
- [33] #Users is calibrated to 1.32B in 2014 from <http://www.datacenterknowledge.com/archives/2009/05/14/whos-got-the-most-web-servers/>
- [34] Data upload per month is calibrated to 500KB/user/month in 2010 from <http://www.extremetech.com/computing/129183-how-big-is-the-cloud>
- [35] Deep suspend current is calibrated to 100nA in 2015 from <http://ambiqmicro.com/low-power-microcontroller>
- [36] MCU Ion/freq is calibrated to 30µA/Mhz in 2015 from <http://ambiqmicro.com/low-power-microcontroller>
- [37] Max MCU frequency is calibrated to 200MHz in 2015 from <http://www.ti.com/lit/ml/spr067/spr067.pdf>
- [38] MCU flash is calibrated to 1MB in 2015 from [http://www.ti.com/lit/ml/microcontrollers\\_16-bit\\_32-bit/c2000\\_performance/control\\_automation/f28m3x/overview.page](http://www.ti.com/lit/ml/microcontrollers_16-bit_32-bit/c2000_performance/control_automation/f28m3x/overview.page)
- [39] Sensor power is calibrated to 5.28mW in 2015 from [http://www.bosch-sensortec.com/en/homepage/products/3/6\\_axis\\_sensors\\_2/inertial\\_measurement\\_unit\\_1/bmi160/bmi160\\_1](http://www.bosch-sensortec.com/en/homepage/products/3/6_axis_sensors_2/inertial_measurement_unit_1/bmi160/bmi160_1)
- [40] G. Fortino et al., "An Agent-Based Middleware for Cooperating Smart Objects", Springer, 2013.
- [41] <https://www.paloaltonetworks.com/resources/learning-center/what-is-a-data-center.html>
- [42] [http://www.hotchips.org/wp-content/uploads/hc\\_archives/hc23/HC23.17.1-tutorial1/HC23.17.130.Switched-Cap-DC-DC-Alon\\_UCB.pdf](http://www.hotchips.org/wp-content/uploads/hc_archives/hc23/HC23.17.1-tutorial1/HC23.17.130.Switched-Cap-DC-DC-Alon_UCB.pdf)

- [43] N. Strurchen et al., "A 2.5D Integrated Voltage Regulator Using Coupled-Magnetic-Core Inductors on Silicon Interposer", *JSSC*, 2013.
- [44] E. Alon, "Fully Integrated Switched-Capacitor DC-DC Conversion", *Tutorial at Hotchips*, 2011.
- [45] R. G. Dreslinski, "Centip3De: A 64-Core, 3D Stacked, Near-Threshold System", *Tutorial at Hotchips*, 2012.
- [46] <http://www1.cs.columbia.edu/~nowick/async-overview-extended-10-10.pdf>
- [47] A. Alaghi et al., "Optimizing Stochastic Circuits for Accuracy-Energy Tradeoffs", *Proc. ICCAD*, 2015.
- [48] J. Han et al., "Approximate Computing: An Emerging Paradigm For Energy-Efficient Design", *Proc. ETS*, 2013.
- [49] <http://www.micron.com/about/innovations/3d-nand>
- [50] <http://www.computerworld.com/article/2859266/a-terabyte-on-a-postage-stamp-rram-heads-into-commercialization.html>
- [51] <http://isscc.org/trends/>
- [52] <http://www.datacenterdynamics.com/colo-cloud-/number-of-data-centers-to-decrease-after-2017/91495.fullarticle>
- [53] <http://www.datacenterknowledge.com/archives/2014/11/11/idc-amount-of-worlds-data-centers-to-start-declining-in-2017/>
- [54] <http://www1.unece.org/stat/platform/display/msis/Big+Data>
- [55] [http://www.hhs.gov/sites/default/files/ocio/ea/documents/hhs\\_datacenter\\_consolidation\\_plan.pdf](http://www.hhs.gov/sites/default/files/ocio/ea/documents/hhs_datacenter_consolidation_plan.pdf)
- [57] <http://searchdatacenter.techtarget.com/opinion/Following-both-sides-of-the-decentralized-vs-centralized-IT-debate>
- [58] <http://www.federaltimes.com/story/government/it/data-center/2015/02/17/data-center-consolidation-goals-not-aggressive/23556995/>
- [59] <https://huanliu.wordpress.com/2014/02/26/amazon-ec2-grows-62-in-2-years/>
- [60] <https://www.google.com/about/datacenters/efficiency/internal/>
- [61] <http://cdn2.hubspot.net/hub/525875/file-3621451007-png/blog-files/oracle-perf-aws-90-random-read-results.png?t=1453388256719>
- [62] P. Neuzil et al, Revisiting Lab on a Chip Technology for Drug Discovery, *Nature Reviews Drug Discovery*, 2012.
- [63] E. Esch et al., Organs-on-Chips at the Frontiers of Drug Discovery, *Nature Reviews Drug Discovery*, 2015.
- [64] J. A. V. Arx et al., Implantable Medical Device with Antenna, US Patent 8,615,305, 2013.
- [65] [https://www.usenix.org/sites/default/files/conference/protected-files/fast14\\_asanovic.pdf](https://www.usenix.org/sites/default/files/conference/protected-files/fast14_asanovic.pdf)
- [66] Datacenter layout with 42U racks. <http://www.42u.com/42U-cabinets.htm>
- [67] UNECE Global Data. <http://www1.unece.org/stat/platform/display/msis/Big+Data>
- [68] Datacenter consolidation. [http://www.hhs.gov/sites/default/files/ocio/ea/documents/hhs\\_datacenter\\_consolidation\\_plan.pdf](http://www.hhs.gov/sites/default/files/ocio/ea/documents/hhs_datacenter_consolidation_plan.pdf)
- [69] <http://www.theatlantic.com/technology/archive/2014/05/all-the-world-a-track-the-trick-that-makes-googles-self-driving-cars-work/370871/>
- [70] <http://www.wsj.com/articles/alphabet-argues-for-quicker-regulatory-path-for-fully-autonomous-cars-1454019334>
- [71] <http://vision.in.tum.de/research/vslam/lstdslam>
- [72] <http://wccftech.com/nvidia-drive-px-dual-chip-tegra-x1-revealed/>
- [73] <http://linuxgizmos.com/qualcomm-unveils-linux-based-uav-reference-platform/>
- [74] <http://www.nvidia.com/object/gpu-applications.html>
- [75] <http://www.nvidia.com/object/what-is-gpu-computing.html>
- [76] <http://www.networkcomputing.com/networking/caviums-xpliant-paves-way-multi-purpose-ethernet-switches/1227500873>
- [77] <http://www.nextplatform.com/2015/11/16/intel-rounds-out-scalable-systems-with-omni-path/>
- [78] [http://www.flashmemorysummit.com/English/Collaterals/Proceedings/2014/20140807\\_304C\\_Zitlaw.pdf](http://www.flashmemorysummit.com/English/Collaterals/Proceedings/2014/20140807_304C_Zitlaw.pdf)
- [79] [http://www.semicontaiwan.org/en/sites/semicontaiwan.org/files/data15/docs/2\\_5.\\_advances\\_and\\_trends\\_in\\_rram\\_technology\\_semicon\\_taiwan\\_2015\\_fi](http://www.semicontaiwan.org/en/sites/semicontaiwan.org/files/data15/docs/2_5._advances_and_trends_in_rram_technology_semicon_taiwan_2015_fi)  
nal.pdf
- [80] [http://www.hotchips.org/wp-content/uploads/hc\\_archives/hc23/HC23.17.1-tutorial1/HC23.17.121,Inductance-Gardner-Intel-DG%20081711-correct.pdf](http://www.hotchips.org/wp-content/uploads/hc_archives/hc23/HC23.17.1-tutorial1/HC23.17.121,Inductance-Gardner-Intel-DG%20081711-correct.pdf)
- [81] [http://4.bp.blogspot.com/\\_Bq3z0Dk2y8M/TK9MPbrFLiI/AAAAAAAAANG/93MpJOWFj-U/s1600/yole.jpg](http://4.bp.blogspot.com/_Bq3z0Dk2y8M/TK9MPbrFLiI/AAAAAAAAANG/93MpJOWFj-U/s1600/yole.jpg)
- [82] [http://www.yole.fr/iso\\_upload/News\\_Illustration/Illustration\\_IoT\\_TechnologyRoadmap\\_YOLE%20DEVELOPPEMENT\\_June%202014.jpg](http://www.yole.fr/iso_upload/News_Illustration/Illustration_IoT_TechnologyRoadmap_YOLE%20DEVELOPPEMENT_June%202014.jpg)
- [83] H. Esmailzadeh, "Neural Acceleration for General-Purpose Approximate Programs", *ISCA*, 2012.
- [84] <http://electroiq.com/wp-content/uploads/2014/07/fig-11.png>
- [85] <http://img.deusm.com/eetimes/3p-0005-monolithic-3d-ic-01-1g.jpg>
- [86] <http://www.3dincites.com/2014/02/411-cea-letis-interposer-roadmap/>
- [87] <http://www.dailytech.com/Report+Samsung+Semiconductor+Bounces+Back+w+iPhone+SoC+and+Memory+Orders/article37193.htm>
- [88] M. Scannell, "3D Integration Activities at Leti", Talk Slides, 2012.
- [89] A. C. Fischer et al., "Integrating MEMS and ICs", *Microsystems & Nanoengineering*, 2015.
- [90] M. Ireland, "RF SOI Redefining Mobility in the Front End", SOI Consortium, 2014.
- [91] J. Moreira et al., "A Single-Chip HSPA Transceiver with Fully Integrated 3G CMOS Power Amplifiers", *ISSCC*, 2015.
- [92] H. Jun, "HBM (High Bandwidth Memory) for 2.5D", *Semicon*, 2015.
- [93] J. Hruska, "Beyond DDR4: The differences between Wide I/O, HBM, and Hybrid Memory Cube", 2015.
- [94] <http://gigaom.com/2013/02/21/qualcomm-new-radio-chip-gets-us-one-step-closer-to-a-global-4g-phone/>
- [95] J. L. Malinge, "A View on the Silicon Photonics Trends and Market Prospective", 2014.
- [96] <http://arstechnica.com/information-technology/2015/05/ibm-demos-first-fully-integrated-monolithic-silicon-photonics-chip/>
- [97] firebox
- [98] Storage Class Memory: Towards a Disruptively Low-Cost Solid-State Non-Volatile Memory, *IBM*, 2013.
- [99] <http://www.kitguru.net/components/memory/anton-shilov/micron-readies-second-gen-3d-xpoint-memory-working-on-all-new-memory-tech/>
- [100] Z. Zhou et al., "On-chip light sources for silicon photonics", *Light: Science & Applications*, 2015.
- [101] <https://www.amd.com/Documents/energy-efficiency-whitepaper.pdf>
- [102] W. Godycki et al., "Enabling Realistic Fine-Grain Voltage Scaling with Reconfigurable Power Distribution Networks", 2014.
- [103] G. Smith, "Updates of the ITRS Design Cost and Power Models", *ICCD*, 2014.
- [104] "Roadmap Report: Photonics for Disaggregated Data Centers", *OSA Industry Development Associates*, March 2015.