

Limits on Fundamental Limits to Computation

Igor L. Markov, the University of Michigan (currently at Google)

Abstract

An indispensable part of our lives, *computing* has also become essential to industries and governments. Steady improvements in computer hardware have been supported by periodic doubling of transistor densities in integrated circuits over the last fifty years. Such *Moore scaling* now requires increasingly heroic efforts, stimulating research in alternative hardware and stirring controversy. To help evaluate emerging technologies and enrich our understanding of integrated-circuit scaling, we review fundamental limits to computation: in manufacturing, energy, physical space, design and verification effort, and algorithms. To outline what is achievable *in principle* and *in practice*, we recall how some limits were circumvented, compare loose and tight limits. We also point out that engineering difficulties encountered by emerging technologies may indicate yet-unknown limits.

1 Introduction

Emerging technologies for computing promise to outperform conventional integrated circuits in computation bandwidth or speed, power consumption, manufacturing cost, or form factor [1,2]. However, razor-sharp focus on any one nascent technology and its benefits sometimes neglects serious limitations or discounts ongoing improvements in established approaches. To foster a richer context for evaluating emerging technologies, we review limiting factors and salient trends in computing that determine what is achievable *in principle* and *in practice*. Several fundamental limits remain substantially loose, possibly indicating viable opportunities for emerging technologies. To clarify this uncertainty, we study *limits on fundamental limits*.

Universal and general-purpose computers. Viewing *clocks and watches* as early computers, it is easy to see the importance of *long-running calculations that can be repeated with high accuracy by mass-produced devices*. The significance of *programmable* digital computers became clear at least 200 years ago, as illustrated by *Jacquard looms* in textile manufacturing. However, the existence of *universal computers* that can efficiently simulate (almost) all other computing devices — analog or digital — was only articulated in the 1930s by Church and Turing (Turing excluded quantum physics when considering universality) [3]. Efficiency was studied from a theoretical perspective at first, but strong demand in military applications in the 1940s lead Turing and von Neumann to develop detailed hardware architectures for universal computers — Turing’s design (Pilot ACE) was more efficient, but von Neumann’s was easier to program. The *stored-program architecture* made universal computers practical in the sense that a single computer design could be effective in many diverse applications. Such practical universality thrives (*i*) in economies of scale in computer hardware and (*ii*) among extensive software stacks. Not surprisingly, the most sophisticated and commercially successful computer designs and components, such as Intel and IBM CPUs, were based on the von Neumann’s paradigm. The numerous uses and large markets of general-purpose chips, as well as the exact reproducibility of their results, justify the enormous capital investment in the design, verification and manufacturing of leading-edge integrated circuits. Today general-purpose CPUs power cloud server-farms and displace specialized (but still universal) mainframe processors in many supercomputers. Emerging universal computers based on field-programmable gate-arrays (FPGAs) and general-purpose graphics processing units (GPGPUs) outperform CPUs in some cases, but their efficiencies remain complementary to those of CPUs. The success of deterministic general-purpose computing manifests in the convergence of diverse functionalities in portable inexpensive smartphones. After steady improvement, general-purpose computing displaced entire industries (newspapers, photography, etc) and launched new applications (video conferencing, GPS navigation, online shopping, networked entertainment, etc) [4]. Application-specific integrated circuits (ASICs) streamline input-output and networking, or

optimize functionalities previously performed by general-purpose hardware. They speed up biomolecular simulation 100-fold [5,6] and improve the efficiency of video decoding 500-fold [7], but require design effort with keen understanding of specific computations, impose high costs and financial risks, need markets where general-purpose computers lag behind, and often cannot adapt to new algorithms. Recent techniques for *customizable domain-specific computing* [8] offer better tradeoffs, while many applications favor the combination of *general-purpose hardware and domain-specific software*, including specialized programming languages [9,10] such as Erlang used in Whatsapp.

Limits as aids to evaluating emerging technologies. Without sufficient history, we cannot extrapolate *scaling laws* for emerging technologies, yet expectations run high. For example, new proposals for analog processors appear frequently (as illustrated by adiabatic quantum computers), but fail to address concerns about analog computing, such as its limitations on scale, reliability, and long-running error-free computation. General-purpose computers meet these requirements with digital integrated circuits (IC) and now command the electronics market. In comparison, quantum computers — both digital and analog — hold promise *only in niche applications* and do not offer *faster general-purpose computing* as they are no faster for *sorting* and other specific tasks [11–13]. In exaggerating the engineering impact of quantum computers, popular press missed this nuance. But in scientific research, building quantum computers may help simulating quantum-chemical phenomena and reveal new fundamental limits. Sections 4 and 6 discuss limits on emerging technologies.

Technology extrapolation versus fundamental limits. The scaling of commercial computing hardware regularly runs into formidable obstacles [1], but near-term technological advances often circumvent them. The International Technology Roadmap for Semiconductors (ITRS) [14] keeps track of such obstacles and possible solutions with a focus on frequently-revised consensus estimates. For example, consensus estimates initially predicted 10 GHz CPUs for the 45 nm technology node [15], versus the 3–4 GHz range seen in practice. In 2004, the unrelated Quantum Information Science and Technology Roadmap [16] forecast 50 physical qubits by 2012. Such optimism arose by assuming technological solutions long before they were developed and validated, and by overlooking important limits. The authors of [17,18] classify limits to device and interconnect as *fundamental*, *material*, *device*, *circuit*, and *system limits*. These categories define the rows of Table 1, and the columns reflect sections of this paper where specific limits are examined for tightness.

2 Engineering obstacles

Engineering obstacles limit specific technologies and choices. For example, a key bottleneck today is IC manufacture, which packs billions of transistors and wires in several cm^2 of silicon with astronomically low defect rates. Layers of material are deposited on silicon and patterned with lasers, fabricating all circuit components simultaneously. Precision optics and photochemical processes ensure accuracy.

Limits on manufacturing. No account of limits to computing is complete without the Abbe diffraction limit: light with wavelength λ , traversing a medium with refractive index η , and converging to a spot with angle θ (perhaps, focused by a lens) creates a spot with diameter $d = \lambda/NA$, where $NA = \eta \sin \theta$ is the numerical aperture. NA reaches 1.4 for modern optics, so it would seem that semiconductor manufacturing is limited to feature sizes $\lambda/2.8$, hence ArF lasers with 193nm wavelength should not support photolithographic manufacturing of transistors with 65nm features. Yet, they supported *sub-wavelength lithography* for the 45nm–22nm technology nodes using *asymmetric illumination* and *computational lithography* [19]. Here one starts with optical masks that look like the intended image, but when the image gets blurry, alter masks by gently shifting edges to improve the image, possibly giving up the semblance between the two. Clearly, some limits are formulated to be broken! Ten years ago, researchers demonstrated patterning of nanomaterials by live viruses [20]. Known virions exceed 20nm in diameter, whereas subwavelength lithography with 193nm-wavelength ArF laser recently extended to 14nm semiconductor manufacturing [14]. Hence, viruses and microorganisms are no longer at the forefront of semiconductor manufacturing. Extreme ultra-violet (X-ray) lasers have been energy-limited, but are improving. Their use requires changing refractive optics to reflective. Additional progress in *multiple patterning* and *directed self-assembly* promises to support photolithography beyond the 10nm technology node.

Limits	Engineering	Design and Validation	Energy, time	Space, time	Information, Complexity
Funda- mental	Abbe (diffraction) Amdahl Gustafson	Error-corr. & dense codes Fault- tolerance thresholds	Einstein $E=mc^2$ Heisenberg $\Delta E\Delta t$ Landauer $kT \ln 2$ Bremermann Adiabatic thrms	Speed of light Planck scale Bekenstein Fisher $T(n)^{1/(d+1)}$	Shannon Holevo NC, NP, #P Turing (decidability)
Mate- rial	Dielectric constant Carrier mobility Surface morphology Fabrication-related	Analytical & numerical modeling	Conductivity Permittivity Bandgap Heat flow	Propagation speed Atomic spacing No gravitational collapse	Information transfer between carriers
Device	Gate dielectric Channel charge ctrl Leakage, Latency Crosstalk, Aging	Compact modeling Parameter selection	CMOS, quantum Charge-centric Signal to noise Energy conversion	Entropy density Entropy flow Interfaces & contacts Size & delay variation	Universality
Circuit	Delay, Inductance Thermal-related Yield, Reliability, IO	Interconnect Test Validation	Dark, darker, dim and gray silicon Cooling efficiency Power density/supply	Interconnect 2D or 3D	Circuit complexity bounds
System +SW	Specification, Implementation Validation, Cost		Synchronization, Physical integration Parallelism, <i>Ab initio</i> limits (Lloyd)		The CAP theorem

Table 1: Some of the known limits to computation [5, 13–15, 17, 18, 22–24, 27, 32, 40, 42, 43, 47, 49–51, 54, 55, 58–61, 63, 64, 66, 75–77, 79, 88, 97].

Limits on individual interconnects. Despite the doubling of transistor density with Moore’s law [21], semiconductor integrated circuits (ICs) would not work without fast and dense interconnects. Metallic wires can be either fast or dense, but not both at the same time — smaller cross-section increases electrical resistance, while greater height or width increase parasitic capacitance with neighboring wires (wire delay grows with RC). In 1995, an Intel researcher pointed out that *on-chip interconnect scaling* is the real limiter of high-performance ICs [22]. The scaling of interconnect is also moderated by *electron scattering against rough edges of metallic wires* [18, 24], inevitable with atomic-scale wires. Hence, IC interconnect stacks have evolved [15, 23] from four equal-pitch layers in 2000 to 16 layers with pitches varying by 32 times, including a large amount of dense (thin) wiring and fast (thick) wires used for global on-chip communication (Figure 3). Aluminum and copper remain unrivaled for conventional interconnects and can be combined in short wires [24]; *carbon-nanotube* and *spintronic* interconnects are also evaluated in [24]. *Photonic waveguides* and *RF links* offer alternative IC interconnect [25, 26], but obey fundamental limits derived from Maxwell’s equations, such as the *maximum propagation speed of EM waves* [18]. I/O links are limited by the perimeter or surface area of a chip, whereas chip capacity grows with area or volume, respectively.

Limits on conventional transistors. Transistors are limited by their tiniest feature — *the width of the gate dielectric*, — which recently reached the size of several atoms (Figure 1), creating problems: (i) a few missing atoms could alter transistor performance, (ii) manufacturing variation makes all transistors slightly different (Figure 2), (iii) electric current tends to leak through thin narrow dielectrics [17]. Instead of a *thinner* dielectric, transistors can be redesigned with *wider* dielectric layers [27] that surround a *fin shape* (Figure 4). Such configurations improve the control of electric field, reduce current densities and leakage, and diminish process variations. Each transistor can use several fins, extending transistor scaling by several generations. Semiconductor manufacturers adopted FinFETs for upcoming technology nodes. One step further, in *tunneling transistors* [28] a gate wraps around the channel to control tunnelling rate.

Limits on design effort. In the 1980s, Mead and Conway formalized IC design using a regular grid, enabling automated layout through algorithms. But resulting optimization problems

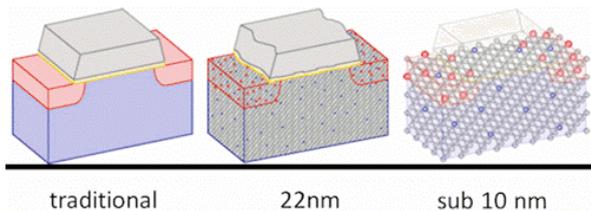


Figure 1: As MOSFET transistors shrink, gate dielectric (yellow) thickness approaches several atoms (0.5nm at the 22nm technology node). Atomic spacing limits device density to 1 device/nm even for radical devices.

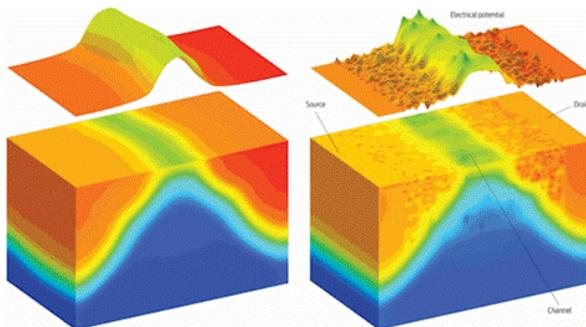


Figure 2: As MOSFET transistors shrink, the shape of electric field departs from basic rectilinear models, and level curves become disconnected. Atomic-level manufacturing variations, especially for dopant atoms, start affecting device parameters, making each transistor slightly different [97, 98]. IMAGE CREDIT: GOLD STANDARD SIMULATIONS.

remain hard, and heuristics are only good enough for practical use. Besides frequent algorithmic improvements, each technology generation alters circuit physics and requires new CAD software. The cost of design has doubled in a few years, becoming prohibitive for ICs with limited market penetration [14]. Emerging technologies, such as FinFETs and high-K dielectrics, circumvent known obstacles using forms of design optimization. Therefore, reasonably tight limits should account for potential future optimizations. Low-level technology enhancements, no matter how powerful, are often viewed as one-off improvements, in contrast to architectural redesigns that affect many processor generations. Between technology enhancements and architectural redesigns are global and local optimizations that alter “the texture” of IC design, such as *logic restructuring*, *gate sizing* and *device parameter selection*. Moore’s law promises higher transistor densities, but some transistors are designed to be 32 times larger than others. Large gates consume greater power to drive long interconnects at acceptable speed and satisfy performance constraints. Minimizing circuit area and power, subject to timing constraints (by configuring each logic gate to a certain size, threshold voltage, etc), is a hard but increasingly important optimization with a large parameter space. A recent convex optimization method [29] saved 30% power in Intel chips, and the impact of such improvements grows with circuit size. Many aspects of IC design are being improved, continually raising the bar for technologies that compete with CMOS.

Completing new IC designs, optimizing and verifying them requires great effort and continuing innovation, e.g., the lack of scalable design automation is a limiting factor for analog ICs [30, 31]. In 1999, bottom-up analysis of digital IC technologies [15, 32] outlined design scaling up to self-contained modules with 50K standard cells (each cell contains 1-3 logic gates), but further scaling was limited by global interconnect. In 2010, physical separation of modules became less critical, as large-scale placement optimizations assumed greater responsibility for IC layout and learned to blend nearby modules [33, 34]. In a general trend, powerful design automation [35] frees circuit engineers to focus on microarchitecture [34], but increasingly relies on algorithmic optimization. Until recently, this strategy suffered significant losses in performance [36] and power [37] compared to ideal designs, but has now become both successful and indispensable due to rapidly increasing complexity

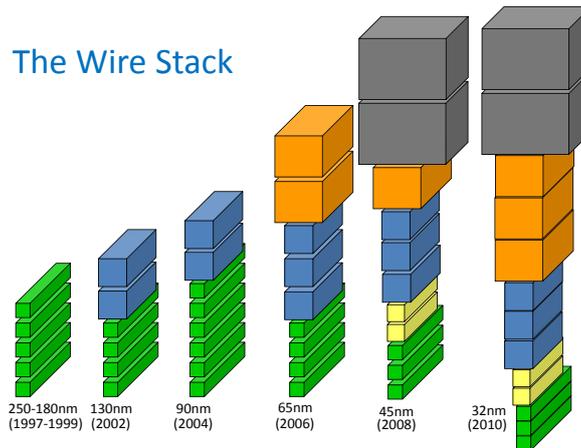


Figure 3: The evolution of metallic wire stacks from 1997 to 2010 by semiconductor technology nodes. IMAGE CREDIT: IBM RESEARCH (modified).

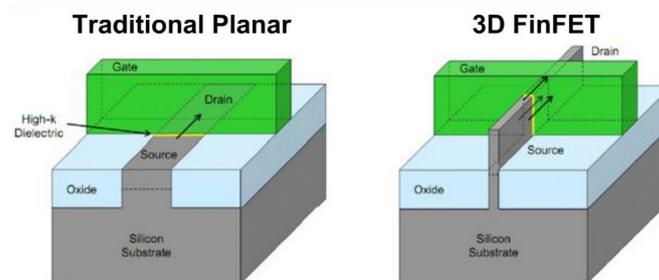


Figure 4: FinFET transistors possess much wider gate dielectric (surrounding the fin shape) than MOSFET transistors and can use multiple fins.

of digital and mixed-signal electronic systems. Hardware and software must now be co-designed and co-verified, with software efforts increasing at a faster rate. *Platform-based design* combines high-level design abstractions with effective reuse of components and functionalities in engineered systems [38]. Customizable domain-specific computing [8] and domain-specific programming languages [9, 10] offload specialization to software running on reusable hardware platforms.

3 Energy-time limits

In predicting the main obstacles to improving modern electronics, the International Technology Roadmap for Semiconductors (ITRS) highlights *the management of system power and energy* as the dominant Grand Challenge [14]. The faster the computation, the more energy it consumes, but actual power-performance tradeoffs depend on the physical scale. While the ITRS, by its charter, focuses on near-term projections and IC design techniques, fundamental limits reflect available energy resources, properties of the physical space, power-dissipation constraints, and energy waste.

3.1 Reversibility

A 1961 result by Landauer [39] shows that erasing one bit of information entails an energy loss $\geq kT \ln 2$ (*thermodynamic threshold*), where k is the Boltzmann constant and T is the temperature in Kelvin. This principle was validated empirically in 2012 [40] and seems to motivate *reversible computing* [41], where all input information is preserved, incurring additional costs. Formally speaking, zero-energy computation is prohibited by *the energy-time form of the Heisenberg uncertainty principle* ($\Delta t \Delta E \geq \hbar$): faster computation requires greater energy [42, 43]. However, recent work in applied su-

perconductivity [44] demonstrates “highly exotic” *physically-reversible* circuits operating at 4°K with energy dissipation below the thermodynamic threshold. They apparently fail to scale to large sizes, run into other limits, and remain no more practical than “mainstream” superconducting circuits and refrigerated low-power CMOS circuits. Technologies that implement *quantum circuits* [45] can *approximate* reversible Boolean computing, but currently do not scale to large sizes, are energy-inefficient at the system level, rely on fragile components, and require heavy fault-tolerance overhead [13]. Conventional ICs also do not help obtaining energy savings from reversible computing because they dissipate 30-60% of all energy in (reversible) wires and repeaters [23]. At room temperature, Landauer’s limit amounts to 2.85×10^{-21} J — a very small fraction of the total, given that modern ICs dissipate 0.1-100 Watts and contain $< 10^9$ logic gates. With the increasing dominance of interconnect (Section 4), more energy is spent on communication than on computation. Logically-reversible computing is important for reasons other than energy — in cryptography, quantum information processing, etc [46].

3.2 Power constraints and CPUs

The end of CPU frequency scaling. In 2004, Intel Corp. abruptly cancelled a 4GHz CPU project because high power density required awkward cooling technologies. Other CPU manufacturers kept clock-frequencies in the 1-6GHz range, but also resorted to multicore CPUs [47]. Since dynamic circuit power grows with clock frequency and supply voltage squared [48], energy can be saved by distributing work among slower, lower-voltage parallel CPU cores *if parallelization overhead is small*.

Dark, darker, dim, gray silicon. A companion trend to Moore’s law — the Dennard scaling theory [49] — shows how to keep power consumption of semiconductor ICs constant while increasing their density. But Dennard scaling broke down ten years ago [49]. Extrapolation of semiconductor scaling trends for CMOS — the dominant semiconductor technology for 20 years — shows that the power consumption of transistors available in modern ICs reduces more slowly than their size (which is subject to Moore’s law) [50, 51]. To ensure the performance envelope of transistors, chip power density must be limited, and a fraction of transistors must be kept dark at any given time. Modern CPUs have not been able to use all their circuits at once, but this asymptotic effect — termed the *utilization wall* [50] — will soon black out 99% of the chip, hence the term *dark silicon* and a reasoned reference to the apocalypse [50]. Saving power by slowing CPU cores down is termed *dim silicon*. Detailed studies of dark silicon [51] show similar results. To this end, executives from Microsoft and IBM have recently proclaimed an end to the era of multicore microprocessors [52]. Two related trends appeared earlier: (i) increasingly large IC regions remain transistor-free to aid routing and physical synthesis, to accommodate power-ground networks, etc [53, 54] — we call them *darker silicon*, (ii) increasingly many gates do not perform useful computation but reinforce long, weak interconnects [55] or slow down wires that are too short — call them *gray silicon*. Today, 50-80% of all gates in high-performance ICs are repeaters.

Limits for power supply and cooling. Data centers in the US consumed 2.2% of total U.S. electricity in 2011. As powerplants take time to build, we cannot sustain past trends of doubled power consumption per year. It is possible to improve the efficiency of transmission lines (using high-temperature superconductors [56]) and power conversion in datacenters, but the efficiency of on-chip power-networks may soon reach 80-90%. Modern IC power management includes *clock* and *power gating* [47], per-core voltage scaling [57], *charge recovery* [58] and, in recent processors, a CPU core dedicated to power scheduling. IC power consumption depends quadratically on supply voltage, which has decreased steadily for many years, but recently stabilized at 0.5-2V [48]. Supply voltage typically exceeds the *threshold voltage* of field-effect transistors by a safety margin that ensures circuit reliability, fast operation and low leakage. Threshold voltage depends on the thickness of gate dielectric, which reached a practical limit of several atoms (Section 2). Supply voltage is limited by around 200mV [17] — five times below current practice — and simple circuits reach this limit. With slower operation, *near-* and *sub-threshold circuits* may consume 100 times less energy [59]. Cooling technologies can improve too, but fundamental quantum limits bound the efficiency of heat removal [60–62].

3.3 Broader limits

The study in [63] explores a general *binary-logic switch model* with binary states represented by two *quantum wells separated by a potential barrier*. Representing information by electric charge requires energy for binary switching and thus limits the logic-switching density, if a significant fraction of the chip can switch simultaneously. To circumvent this limit, one can encode information in *spin-states, photon polarizations, super-conducting currents, or magnetic flux*, noting that these carriers have already been in commercial use. Spin-states are particularly attractive because they promise high-density nonvolatile storage [64] and scalable interconnects [24]. More powerful limits are based on the amount of material in Earth’s crust (where silicon is the second most common element after oxygen), on atomic spacing (Section 2), radii and energies, bandgaps, as well as the wavelength of the electron. We are currently using only a tiny fraction of Earth’s mass for computing, and yet various limits could be circumvented if new particles are discovered. Beyond atomic physics, some limits rely on basic constants: the speed of light, the gravitational constant, the quantum (Planck) scale, the Boltzmann constant, etc. Lloyd [43], as well as Kraus [65] extend well-known bounds by Bremermann and Bekenstein, and give Moore’s law 150 and 600 years, respectively. These results are too loose to obstruct the performance of practical computers. In contrast, current consensus estimates from the ITRS [14] give Moore’s law only 10-20 years, due to both *technological* and *economic* considerations [2].

4 Asymptotic space-time limits

Engineering limits for deployed technologies can often be circumvented, while first-principles limits on energy and power are very loose. Reasonably tight limits are rare.

Limits to parallelism. Suppose we wish to compare a parallel and sequential computer built from the same units, to argue that a new parallel algorithm is many times faster than the best sequential algorithm (the same reasoning applies to logic gates on an IC). Given N parallel units and an algorithm that runs K times faster on sufficiently large inputs, one can *simulate* the parallel system on the sequential system by dividing its time between N computational slices. Since this simulation is roughly N times slower, it runs K/N times faster than the original sequential algorithm. If this algorithm was best possible, we have $K \leq N$. The bound is reasonably tight in practice for small N and can be violated slightly since N CPUs include more CPU cache, but such violations do not justify parallel algorithms — one could instead buy/build one CPU with a larger cache. Such linear speedup is optimistically assumed for *the parallelizable component* in the 1988 Gustafson’s law that suggests scaling the number of processors with input size (as illustrated by instantaneous Web search queries over massive data sets) [5]. Also in 1988, Fisher [66] employed *asymptotic runtime estimates* instead of numerical limits and avoided the breakdown into parallel and sequential runtime components, assumed in Amdahl’s [67] and Gustafson’s laws [5]. Asymptotic estimates neglect leading constants and offer a powerful way to capture nonlinear phenomena occurring at large scale.

Fisher [66] assumes a sequential computation with $T(n)$ elementary steps for input of size n , and limits the performance of its parallel variants that can use an unbounded d -dimensional grid of finite-size computing units (electrical switches on a semiconductor chip, logic gates, CPU cores, etc) communicating at a finite speed, say, bounded by the speed of light. We highlight only one aspect of this four-page work — the parallel computation requires $\Omega(\sqrt[d+1]{T(n)})$ steps. This result undermines the N -fold speedup assumed in Gustafson’s law for N processors on appropriately sized input data [5]. A more realistic speedup from $\sim n^k$ to $\sim \log n$ can be achieved in an abstract model of computation for matrix multiplication and fast Fourier transforms. But not in physical space [66]. Surprising as it may seem, after reviewing many loose limits to computation, we have identified a reasonably tight limit (the impact of I/O — a major bottleneck today — is also covered in [66]). Indeed, many parallel computations today (excluding multimedia processing and Web search) are limited by several forms of communication and synchronization, including network and storage access. The billions of logic gates and memory elements in modern ICs are linked by up to 16 levels of wires (Figure 3), longer wires are segmented by repeaters. Most of the physical volume and circuit delay are attributed to

interconnect [23]. This is relatively new, as gate delays were dominant until 2000 [14], but wires get slower relative to gates at each new technology node. This uneven scaling has compounded in ways that would surprise Turing and von Neumann — a single clock cycle is now far too short for a signal to cross the entire chip, and even the distance covered in 200 ps (5 GHz) at light-speed is close to chip size. Yet, most electrical engineers and computer scientists continue to focus on gates.

Implications to 3D ICs and other emerging technologies. The promise of 3D integration for improving IC performance can be contrasted with technical obstructions to its industry adoption. To derive limits on possible improvement, we use the result from [66] sensitive to the dimension of the physical space: a sequential computation with $T(n)$ steps requires $\Omega(\sqrt[3]{T(n)})$ steps in 2D and $\Omega(\sqrt[4]{T(n)})$ in 3D. Letting $t = \sqrt[3]{T(n)}$, shows that 3D integration asymptotically reduces t to $t^{3/4}$ — a significant but not dramatic speedup. This speedup requires an unbounded number of 2D device layers, otherwise there is no asymptotic speedup [68]. For 3D ICs with 2-3 layers, the main benefits of 3D IC integration today are in improving manufacturing yield, improving I/O bandwidth, and combining 2D ICs that are optimized for random logic, dense memory, FPGA, analog, MEMS, etc. Ultra-high density CMOS logic ICs with *monolithic* 3D integration [69] suffer higher routing congestion than traditional 2D ICs. Emerging technologies promise to improve device parameters, but often remain limited by scale, faults, and interconnect, e.g., *quantum dots* enable Terahertz switching but hamper nonlocal communication [70]. CNT-FETs [71] leverage extraordinary carrier mobility in semiconducting carbon nanotubes to use interconnect more efficiently by improving drive strength, while reducing supply voltage. Emerging interconnects include *silicon photonics*, shown by Intel in 2013 [72] as a 100Gb/s replacement of copper cables connecting adjacent chips. It promises to reduce power consumption and form factor. Quantum physics alters the nature of communication with Einstein’s “spooky action at a distance” facilitated by entanglement [13]. However, the flows of information and entropy are subject to quantum limits [60, 61]. Several quantum algorithms run asymptotically faster than best conventional algorithms [13], but fault-tolerance overhead offsets their potential benefits in practice, and empirical evidence of quantum speedups has not been compelling so far [73, 74]. Several stages in the development of quantum information processing remain challenging [99], and the surprising difficulty of *scaling up* reliable quantum computation could stem from limits on *communication* and *entropy* [13, 60, 61]. In contrast, Lloyd [43] notes that *individual* quantum devices now approach energy limits for switching, whereas nonquantum devices remain orders of magnitude away. This suggests an obstacle to simulating quantum physics on conventional computers (abstract models aside). In terms of computational complexity though, quantum computers *cannot* attain significant advantage for many problem types [11–13]. Such lack of *consistent general-purpose speedup* limits the benefits of several emerging technologies in mature applications with diverse algorithmic steps, e.g., computer-aided design and Web search. Accelerating one step usually does not greatly speed up the entire application, as noted by Amdahl in 1967 [67]. Figuratively speaking, *the most successful computers are designed for the decathlon, rather than for sprint only*.

5 Complexity-theoretic limits

Section 4 enabled tighter limits by neglecting energy and using asymptotic rather than numeric bounds — a more abstract model focuses on the impact of scale, and recurring trends quickly overtake one-off device-specific effects. Next, we neglect spatial effects and focus on the nature of computation in an abstract model (used by software engineers) that represents computation by elementary steps with input-independent runtimes. Such limits survive many improvements in computer technologies, and are often stronger for specific problems. For example, the best-known algorithms for multiplying large numbers are only slightly slower than reading the input (an obvious speed limit), but only in the asymptotic sense — for numbers with <1000 bits, those algorithms lag behind simpler algorithms in actual performance. To focus on what matters, we now do not just track asymptotic worst-case complexity of best algorithms for a given problem, but merely distinguish *polynomial* asymptotic growth from *exponential*. Limits formulated in such crude terms (unsolvability in polynomial time *on any computer*) are powerful [75]: the hardness of number-factoring underpins Internet commerce, while

the $P \neq NP$ conjecture explains the lack of satisfactory, scalable solutions to important algorithmic problems, e.g., in optimization and verification of IC designs [76]. A similar conjecture $P \neq NC$ seeks to explain why many algorithmic problems that can be *solved* efficiently have not *parallelized* efficiently [77]. Most of these limits have not been proven. Some can be circumvented by using radically different physics, e.g., quantum computers solve number factoring in polynomial time (in theory). But quantum computation does not affect $P \neq NP$ [78]. The lack of proofs, despite heavy empirical evidence, requires faith and is an important limitation of many nonphysical limits to computing. This faith is not universally shared — Donald Knuth argues¹ that $P=NP$ would not contradict anything we know today. A rare *proven* result by Turing (also invulnerable to quantum physics) states that checking if a given program ever halts is *undecidable*: no algorithm solves this problem in all cases regardless of runtime. Yet, software developers solve this problem during peer code reviews, and computer science teachers — when grading exams in programming courses. *Worst-case analysis* is another limitation of nonphysical limits to computing, but suggests potential gains through approximation and specialization. For some NP-hard optimization problems, such as the *Euclidean Travelling Salesman Problem* (EucTSP), polynomial-time approximations exist, but in other cases, such as the *maximum clique problem*, accurate approximation is as hard as finding optimal solutions [79]. For some important problems and algorithms, such as the Simplex algorithm for *linear programming*, few inputs lead to exponential runtime, and minute perturbations reduce runtime to polynomial [80].

6 Conclusions

The death march of Moore’s law [1, 2] invites discussions of fundamental limits and alternatives to silicon semiconductors [71]. Near-term constraints invariably tie to *costs* and *capital*, but are explained away by new markets for electronics, increasing Earth population, and growing world economy [2]. Such economic pressures emphasize the value of *computational universality* and broad applicability of IC architectures to solve multiple tasks under conventional environmental conditions. In a likely scenario, only CPUs, GPUs, FPGAs and dense memory ICs will remain viable at the end of Moore’s law, while specialized circuits will be manufactured with less advanced technologies. Indeed, memory chips have lead Moore scaling by leveraging their simpler structure, modest interconnect, and more controllable manufacturing, but their scaling is slowing down [2]. The decelerated scaling of CMOS ICs still outperforms the scaling of the most viable emerging technologies. Empirical scaling laws describing the evolution of computing are well-known [81]. In addition to Moore’s law, Dennard scaling, as well as Amdahl’s and Gustafson’s laws reviewed earlier, Metcalfe’s law [82] states that the value of a computer network, such as the Internet or Facebook, scales as the number of user-to-user connections that can be formed. Grosch’s law [83] ties N -fold improvements in computer performance to N^2 -fold cost increases (in equivalent units). Applying it in reverse, we can estimate acceptable performance of cheaper computers. But such laws only capture *ongoing scaling* and will break down in the future.

The *roadmapping process* represented by the International Technology Roadmap for Semiconductors (ITRS) [14] relies on consensus estimates and works around engineering obstacles. It tracks improvements in materials and tools, collects best practices and outlines promising design strategies. As suggested in [17, 18], it can be enriched by analysis of limits. We additionally focus on how closely such limits can be approached. Aside from historical “wrong turns” recalled in Sections 2 and 3, we find interesting effects when examining the tightness of individual limits. While energy-time limits are most critical in computer design [14, 84], space-time limits appear tighter [66] and capture bottlenecks formed by interconnect and communication. They suggest optimizing gate locations and sizes, and placing gates in three dimensions. One can also adapt algorithms to spatial embeddings [85, 86] and seek space-time limits. But the gap between current technologies and energy-time limits hints at greater rewards. *Charge recovery* [58], *power management* [47], voltage scaling [57], and *near-threshold computing* [59] reduce energy waste. Optimizing algorithms and circuits simultaneously for energy and spatial embedding [87] gives biological systems an edge (from the 1D worm *C. elegans* with 302

¹See Question 17 in <http://www.informit.com/articles/article.aspx?p=2213858>

neurons to the 3D human brain with 86 billion neurons) [1]. Yet, using mass-energy to compute can be a veritable *nuclear option*. In a 1959 talk, which predated Moore’s law, Richard Feynman suggested that there was “plenty of room at the bottom,” forecasting the miniaturization of electronics. Today, with relatively little physical room left, *there is plenty of energy at the bottom*. If this energy is tapped for computing, how can resulting heat be removed? Recycling heat into mass or electricity seems ruled out by limits to energy conversion and the acceptable thermal envelope.

Technology-specific limits for modern computers tend to express tradeoffs, especially for systems with conflicting performance parameters and properties [88]. Little is known about limits on *design technologies*. Given that large-scale complex systems are often designed and implemented hierarchically [53] with multiple levels of abstraction, it would be valuable to capture losses incurred at abstraction boundaries and between levels of design hierarchies. It is common to estimate resources required for a subsystem and then implement the subsystem to satisfy resource budgets. Underestimation is avoided because it leads to failures, but overestimation results in overdesign. Inaccuracies in estimation and physical modeling also lead to losses during optimization, especially in the presence of uncertainty. Clarifying engineering limits gives hope to circumvent them.

Technology-agnostic limits look simple and have had significant impact in practice, for example Aaronson explains why NP-hardness is unlikely to be circumvented by through physics [78]. Limits to *parallel computation* became prominent after CPU speed levelled off ten years ago. They suggest using faster interconnect [18], local computation that reduces communication [89], time-division multiplexing of logic [90], architectural and algorithmic techniques [91], solving larger problem instances, and altering applications to embrace parallelism [5]. John Gustafson advocates a *natural selection*: the survival of applications fittest for parallelism. In another twist, the performance and power consumption of industry-scale distributed systems is often described by probability distributions, rather than single numbers [92, 93], making it harder to even formulate appropriate limits. We also cannot yet formulate fundamental limits related to the complexity of the software-development effort, the efficiency of CPU caches [94], and computational requirements of incremental functional verification, but we have noticed that many known limits are either loose or can be circumvented, leading to *secondary limits*. To wit, the $P \neq NP$ limit is worded in terms of worst-case rather than average-case performance, and has not been proven despite heavy evidence. Researchers have ruled out entire categories of proof techniques as insufficient to complete such a proof [76, 95]. While esoteric, such *tertiary limits* can be effective in practice — in August 2010, they helped researchers quickly invalidate Vinay Deolalikar’s highly-technical attempt at proving $P \neq NP$. On the other hand, the correctness of lengthy proofs for some key results could not be established with acceptable level of certainty by reviewers, prompting efforts in verifying mathematics by computation [96].

In summary, we have reviewed what is known about limits to computation, including existential challenges arising in the sciences, design and optimization challenges arising in engineering, as well as current state of the art. These categories are closely linked due to the rapid pace of technology development. When a specific limit is approached and obstructs progress, understanding its assumptions is a key to circumventing it. Some limits are hopelessly loose and can be ignored, while other limits remain conjectured based on empirical evidence and may be very difficult to establish rigorously. Such *limits on limits to computation* deserve further study.

Acknowledgments. This work was supported in part by the Semiconductor Research Corporation (SRC) Task 2264.001 (funded by Intel and IBM), US Airforce Research Laboratory Award FA8750-11-2-0043, and US National Science Foundation (NSF) Award 1162087.

References

- [1] Cavin, R. K., Lugli, P., and Zhirnov, V. V. Science and Engineering Beyond Moore's Law, *Proc. IEEE* 100:1720-1749 (2012).
- [2] Chien, A. A., and Karamcheti, V. Moore's Law: The First Ending and a New Beginning, *IEEE Computer* 46(13):48-53 (2013).
- [3] Herken, R., ed. The Universal Turing Machine: A Half-Century Survey, 2nd ed., *Springer* (2013).
- [4] Andreesen, M. Why Software Is Eating The World, *Wall Street Journal* August 11 (2011).
- [5] Padua, D. A., editor, Encyclopedia of Parallel Computing (*Springer* 2011).
- [6] Shaw, D. E., Anton: a Special-Purpose Machine that Achieves a Hundred-Fold Speedup in Biomolecular Simulations, In *Proc. Int'l Symp. on High Performance Distributed Computing (HPDC)* 129-130 (2013).
- [7] Hameed, R., Qadeer, W., Wachs, M., Azizi, O., Solomatnikov, A., Lee, B. C., Richardson, S., Kozyrakis, C., Horowitz, M. Understanding Sources of Inefficiency in General-Purpose Chips, *Commun. ACM* 54(10): 85-93 October (2011).
- [8] Cong, J., Reinman, G., Bui, A. T., Sarkar, V. Customizable Domain-Specific Computing, *IEEE Design & Test of Computers* 28(2):6-15 (2011).
- [9] Mernik, M., Heering, J., Sloane, A. M. When and How to Develop Domain-Specific Languages, *ACM Comput. Surv.* 37(4): 316-344 (2005).
- [10] Olukotun, K., Beyond Parallel Programming with Domain Specific Languages, In *Proc. Symposium on Principles & Practice of Parallel Programming (PPOPP)* 179-180 (2014).
- [11] Aaronson, S., Shi, Y., Quantum Lower Bounds for the Collision and the Element Distinctness Problems, *J. ACM* 51(4): 595-605 (2004).
- [12] Jain, R., Ji, Z., Upadhyay, S., Watrous, J., "QIP = PSPACE," *Commun. ACM* 53(12): 102-109 (2010).
- [13] Nielsen, M. A., and Chuang, I. L. Quantum Computation and Quantum Information, *Cambridge University Press* (2011).
- [14] International Technology Roadmap for Semiconductors (ITRS), <http://www.itrs.net/>
- [15] Sylvester, D., Keutzer, K. A Global Wiring Paradigm for Deep Submicron Design, *IEEE Trans. on CAD* 19(2):242-252 (2000).
- [16] A Quantum Information Science and Technology Roadmap, *Los Alamos Technical Report LA-UR-04-1778*, 2004 <http://qist.lanl.gov>
- [17] Meindl, J., Low Power Microelectronics: Retrospective and Prospect, *Proc. IEEE* 83(4):619-635, 1995.
- [18] Davis, J. A., Venkatesan, R., Kaloyeros, A., Beylansky, M., Souri, S. J., Banerjee, K., Saraswat, K. C., Rahman, A., Reif, R., Meindl, J. D. Interconnect Limits on Gigascale Integration (GSI) in the 21st Century, *Proc. IEEE* 89(3):305-324 (2001).
- [19] Ma, X., and Arce, G. R., Computational Lithography (*Wiley* 2011).
- [20] Mazzola, L. Commercializing Nanotechnology, *Nature Biotechnology* 21(10):1137-1143 (2003).
- [21] Moore, G. E. Cramming More Components onto Integrated Circuits, *Electronics* 38:1-4 (1965).

- [22] Bohr, M., Interconnect Scaling — The Real Limiter to High Performance ULSI, In *Proc. Int'l Elec. Device Meeting (IEDM)* 241-244 (1995).
- [23] Shelar, R., Patyra, M. Impact of Local Interconnects on Timing and Power in a High Performance Microprocessor, *IEEE Trans. CAD* 32(10): 1623-1627 (2013).
- [24] Naeemi, A., Ceyhan, A., Kumar, V., Pan, C., Iraei, R. M., Rakheja, S. BEOL Scaling Limits and Next Generation Technology Prospects. In *Proc. Design Automation Conf. (DAC)* 2014: 1-6.
- [25] Almeida, V. R., Barrios, C. A., Panepucci, R. R., Lipson, M. All-optical Control of Light on a Silicon Chip, *Nature* 431: 1081-1084 (2004).
- [26] Chang, M.-C. F., Roychowdhury, V. P., Zhang, L., Shin, H., Qian, Y., RF/wireless Interconnect for Inter- and Intra-chip Communications, *Proc. IEEE* 89(4):455-66 (2002).
- [27] Hisamoto, D., Lee, W.-C., Kedzierski, J., Takeuchi, H., Asano, K., Kuo, C., Anderson, E., King, T.-J., Bokor, J., Hu, C. FinFET-a Self-aligned Double-gate MOSFET Scalable to 20 nm, *IEEE Trans. on Electron Devices* 47(12):2320-2325 (2002).
- [28] Seabaugh, A. The Tunneling Transistor, *IEEE Spectrum*, September (2013).
- [29] Ozdal, M. M., Burns, S. M., and Hu, J. Algorithms for Gate Sizing and Device Parameter Selection for High-Performance Designs, *IEEE Trans. CAD* 31(10):1558-1571 (2012).
- [30] Rutenbar, R. A. Design Automation for Analog: the Next Generation of Tool Challenges, In *Proc. Int'l Conf. Computer-Aided Design of Integrated Circuits (ICCAD)* 458-460 (2006).
- [31] Rutenbar, R. A. Analog Layout Synthesis: What's Missing? In *Proc. Int'l Symp. Physical Design of Integrated Circuits (ISPD)* 43 (2010).
- [32] Ho, R., Mai, K., Kapadia, H., Horowitz, M. Interconnect Scaling Implications for CAD, In *Proc. Int'l Conf. Computer-Aided Design of Integrated Circuits (ICCAD)* 425-429 (1999).
- [33] Markov, I. L., Hu, J., Kim, M.-C. Progress and Challenges in VLSI Placement Research, In *Proc. Int'l Conf. Computer-Aided Design of Integrated Circuits (ICCAD)* 2012: 275-282.
- [34] Puri, R. Opportunities and Challenges for High-Performance CPU Designs and Design Automation, In *Proc. Int'l Symp. Physical Design of Integrated Circuits* 179 (2013).
- [35] Lavagno, L., Martin, G., Scheffer, L. Electronic Design Automation for Integrated Circuits Handbook, (*CRC Press* 2006).
- [36] Chinnery, D. G., and Keutzer, K. Closing the Gap Between ASIC and Custom - Tools and Techniques for High-Performance ASIC Design (*Springer* 2004).
- [37] Chinnery, D. G., and Keutzer, K. Closing the Power Gap between ASIC and Custom - Tools and Techniques for Low Power Design (*Springer* 2007).
- [38] Sangiovanni-Vincentelli, A. L., Carloni, L. P., De Bernardinis, F., Sgroi, M. Benefits and Challenges for Platform-Based Design, In *Proc. Design Automation Conf. (DAC)* 409-414 (2004).
- [39] Landauer, R. Irreversibility and Heat Generation in the Computing Process, *IBM J. of Research and Development* 5: 183-191 (1961).
- [40] Bérut, A., Arakelyan, A., Petrosyan, A., Ciliberto, S., Dillenschneider R., and Lutz, E. Experimental Verification of Landauer's Principle Linking Information and Thermodynamics, *Nature* 483: 187-189 (2012).
- [41] Bennett, C. H., and Landauer, R. The Fundamental Limits of Computation, *Scientific American* 253, July (1985).

- [42] Aharonov, Y., Bohm, D. Time in the Quantum Theory and the Uncertainty Relation for Time and Energy, *Physical Review* 122 (5): 1649-1658 (1961).
- [43] Lloyd, S. Ultimate Physical Limits on Computation, *Nature* 406:1046-1054 (2000).
- [44] Ren, J., Semenov, V. K. Progress With Physically and Logically Reversible Superconducting Digital Circuits, *IEEE Trans. on Applied Superconductivity* 21(3): 780-786 (2011).
- [45] Monroe, C., Raussendorf, R., Ruthven, A., Brown, K. R., Maunz, P., Duan, L.-M., and Kim, J. Large-scale Modular Quantum-Computer Architecture with Atomic Memory and Photonic Interconnects, *Phys. Rev. A* 89, 022317.
- [46] Saeedi, M., and Markov, I. L. Synthesis and Optimization of Reversible Circuits — a Survey, *ACM Comput.Surv.* 45(2):21 (2013).
- [47] Borkar, S. Thousand-Core Chips: A Technology Perspective, In *Proc. Design Automation Conf. (DAC)* 746-749 (2007).
- [48] Rabaey, J. M., Chandrakasan, A., Nikolic, B. Digital Integrated Circuits A Design Perspective, *Pearson Education, Inc* (2003).
- [49] Bohr, M. A 30 Year Retrospective on Dennard’s MOSFET Scaling Paper, *IEEE Solid-State Circuits Society Newsletter* 12(1): 11-13 (2007).
- [50] Taylor, M. B. Is Dark Silicon Useful? Harnessing the Four Horsemen of the Coming Dark Silicon Apocalypse, In *Proc. Design Automation Conf. (DAC)* 1131-1136 (2012).
- [51] Esmailzadeh, H., Blem, E. R., St.-Amant, R., Sankaralingam, K., Burger, D., Power Challenges May End the Multicore Era, *Comm. ACM* 56(2): 93-102 (2013).
- [52] Yeraswork, Z. 3D Stacks and Security Key for IBM in Server Market, *EE Times* 12/17/2013.
- [53] Caldwell, A. E., Kahng, A. B., and Markov, I. L. Hierarchical Whitespace Allocation in Top-down Placement, *IEEE Trans. on CAD*, vol. 22(11): 716-724, November (2003).
- [54] Adya, S. N., Markov, I. L., Villarrubia, P. G., On Whitespace and Stability in Physical Synthesis, *Integration: the VLSI Journal* 39(4): 340-362 (2006).
- [55] Saxena, P., Menezes, N., Cocchini, P., Kirkpatrick, D. Repeater Scaling and Its Impact on CAD, *IEEE Trans. CAD* 23(4): 451-463 (2004).
- [56] Oestergaard, J., Okholm, J., Lomholt, K., Toennesen, G. Energy Losses of Superconducting Power Transmission Cables in the Grid, *IEEE Transactions on Applied Superconductivity* 11: 2375 (2001).
- [57] Pinckney, N. R., Dreslinski, R. G., Sewell, K., Fick, D., Mudge, T. N., Sylvester, D., Blaauw, D. Limits of Parallelism and Boosting in Dim Silicon, *IEEE Micro* 33(5): 30-37 (2013).
- [58] Kim, S., Ziesler, C.H., Papaefthymiou, M.C. Charge-Recovery Computing on Silicon, *IEEE Trans. Comput.* 54(6):651-659, (2005).
- [59] Dreslinski, R. G., Wieckowski, M., Blaauw, D., Sylvester, D., Mudge, T. Near-Threshold Computing: Reclaiming Moore’s Law Through Energy Efficient Integrated Circuits, *Proc. IEEE* 98(2): 253-266 (2010).
- [60] Pendry, J. B. Quantum Limits to the Flow of Information and Entropy, *J. Phys. A: Math. Gen.* 16 2161 (1983).
- [61] Blencowe, M. P., Vitelli, V. Universal Quantum Limits on Single-channel Information, Entropy, and Heat flow, *Phys. Rev. A* 62 052104 (2000).

- [62] Whitney, R. S. Most Efficient Quantum Thermoelectric at Finite Power Output, *Phys. Rev. Lett.* 112: 130601 (2014).
- [63] Zhirnov, V. V., Cavin, R. K., Hutchby, J. A., Bourianoff, G. I. Limits to Binary Logic Switch Scaling — a Gedanken Model, *Proc. IEEE* 91(11):1934-1939 (2003).
- [64] Wolf, S. A., Awschalom, D. D., Buhrman, R. A., Daughton, J. M., von Molnár, S., Roukes, M. L., Chtchelkanova, A. Y., Treger, D. M., Spintronics: A Spin-Based Electronics Vision for the Future, *Science* 294:1488-1494 (2001).
- [65] Krauss, L. M., and Starkman, G. D. Universal Limits on Computation [arXiv:astro-ph/0404510](https://arxiv.org/abs/astro-ph/0404510).
- [66] Fisher, D. Your Favorite Parallel Algorithms Might Not Be as Fast as You Think, *IEEE Trans. on Computers* 37(2): 211-213 (1988).
- [67] Amdahl, G. M. Computer Architecture and Amdahl’s Law, *IEEE Computer* 46(12): 38-46 (2013).
- [68] Mak, W.-K., Chu, C. Rethinking the Wirelength Benefit of 3-D Integration. *IEEE Trans. VLSI Syst.* 20(12): 2346-2351 (2012).
- [69] Lee, Y.-J., Morrow, P., Lim, S. K., “Ultra High Density Logic Designs Using Transistor-Level Monolithic 3D Integration,” In *Proc. Int’l Conf. Computer-Aided Design of Integrated Circuits ICCAD 2012*: 539-546.
- [70] Sherwin, M. S., Imamoglu, A., Montroy, Th. Quantum Computation with Quantum Dots and Terahertz Cavity Quantum Electrodynamics, *Phys. Rev. A* 60: 3508 (1999).
- [71] Shulaker, M., Hills, G., Patil, N., Wei, H., Chen, H.-Y., Wong H.-S. P., and Mitra, S. Carbon Nanotube Computer, *Nature* 501: 526-530 (2013).
- [72] Simonite, T., Intel’s Laser Chips Could Make Data Centers Run Better, *MIT Technology Review* 9/04/2013.
- [73] Rønnow, T. F., Wang, Z., Job, J., Boixo, S., Isakov, S. V., Wecker, D., Martinis, J. M., Lidar, D. A., and Troyer, M., Defining and Detecting Quantum Speedup, [arXiv:1401.2910](https://arxiv.org/abs/1401.2910) (2014).
- [74] Shin, S. W., Smith, G., Smolin, J. A., Vazirani, U., How ‘Quantum’ is the D-Wave Machine? [arXiv:1401.7087](https://arxiv.org/abs/1401.7087) (2014).
- [75] Sipser, M. Introduction to the Theory of Computation, 3rd ed. (*Cengage Learning* 2012).
- [76] Fortnow, L. The Status of the P versus NP problem, *Commun. ACM* 52(9): 78-86 (2009).
- [77] Markov, I. L. Know Your Limits: A Review of ‘Limits to Parallel Computation: P-Completeness Theory’, *IEEE Design and Test* 30(1): 78-83 (2013).
- [78] Aaronson, S., Guest Column: NP-Complete Problems and Physical Reality, *SIGACT News* 36(1): 30-52 (2005).
- [79] Vazirani, V. Approximation Algorithms (*Springer* 2002).
- [80] Spielman, D., Teng, S.-H. Smoothed Analysis of Algorithms: Why the Simplex Algorithm Usually Takes Polynomial Time, In *Proc. Symp. Theory of Computing* 296305 (2001).
- [81] Getov, V. Computing Laws: Origins, Standing, and Impact, *IEEE Computer* 46 (12): 24-25 (2013).
- [82] Metcalfe, B. Metcalfe’s Law after 40 Years of Ethernet, *IEEE Computer* 46(12): 26-31 (2013).

- [83] Ryan, P. S., Falvey, S., Merchant, R. When the Cloud Goes Local: The Global Problem with Data Localization, *IEEE Computer* 46(12):54-59 (2013).
- [84] Wenisch, T. F., Buyuktosunoglu, A. Energy-Aware Computing, *IEEE Micro* 32(5): 6-8 (2012).
- [85] Bachrach, J., and Beal, J. Developing Spatial Computers, *Technical Report MIT-CSAIL-TR-2007-017* (2007).
- [86] Rosenbaum, D. Optimal Quantum Circuits for Nearest-Neighbor Architectures, arXiv:1205.0036 (2012).
- [87] Patil, D., Azizi, O., Horowitz, M., Ho, R., Ananthraman, R. Robust Energy-Efficient Adder Topologies, *Computer Arithmetic (ARITH)* 16-28 (2007).
- [88] Brewer, E. CAP Twelve Years Later: How the ‘Rules’ Have Changed, *IEEE Computer* 45(2): 23-29 (2012).
- [89] Demmel, J. Communication-Avoiding Algorithms for Linear Algebra and Beyond, In *Proc. Int’l Parallel and Distributed Processing Symp.* 585 (2013).
- [90] Halfill, T. R., Tabula’s Time machine, *Microprocessor Report*, 3/29/2010.
- [91] Dror, R. O., Grossman, J. P., Mackenzie, K. M., Towles, B., Chow, E., Salmon, J. K., Young, C., Bank, J. A., Batson, B., Shaw, D. E., Kuskin, J., Larson, R. H., Moraes, M. A., Shaw, D. E. Overcoming Communication Latency Barriers in Massively Parallel Scientific Computation, *IEEE Micro* 31(3): 8-19 (2011).
- [92] Dean, J., and Barroso, L. A. The Tail at Scale, *Commun. ACM* 56(2):74-80 (2013).
- [93] Barroso, L. A., Clidaras, J., Hölzle, U. The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines, 2nd ed.: Synthesis Lectures on Computer Architecture (*Morgan & Claypool Publishers* 2013).
- [94] Balasubramonian, R., Jouppi, N. P., Muralimanohar, N. Multi-Core Cache Hierarchies, Synthesis Lectures on Computer Architecture (*Morgan & Claypool Publishers* 2011).
- [95] S. Aaronson, A. Wigderson, “Algebrization: A New Barrier in Complexity Theory,” *ACM Trans. Complexity Theory* 1(1), 2009.
- [96] Avigad, J., and Harrison, J. “Formally Verified Mathematics,” *Comm. ACM* 57(4): 66-75 (2014).
- [97] Asenov, A. Random Dopant Induced Threshold Voltage Lowering and Fluctuations in Sub-0.1 m MOSFETs: a 3-D “Atomistic” Simulation Study, *IEEE Trans. on Electron Devices* 45(12):2505-2513 (1998).
- [98] Miranda, M. The Threat of Semiconductor Variability, *IEEE Spectrum*, June (2012).
- [99] Devoret, M. H. and Schoelkopf, R. J. Superconducting circuits for quantum information: an outlook. *Science* 339, 11691173 (2013).