

# Saving Energy and Increasing Density in Information Processing Using Photonics

David A. B. Miller

*Ginzton Laboratory, Spilker Building, 348 Via Pueblo Mall, Stanford CA 94305  
dabm@stanford.edu*

**Abstract:** We argue energy and interconnect density in information processing can be improved by orders of magnitude using parallel free-space optical channels inside and between racks, enabled by integrated waveguide photonics, and run synchronously without time-multiplexing. © 2020 D. A. B. Miller

Energy dissipation for information processing and communications consumes something of the order of 10% of all electricity (see [1] and references therein). If we do not reduce the energy per bit in such systems, then obviously we cannot readily continue scaling our use of information exponentially. The problem is obvious inside large systems like data centers, where power dissipation is a clear issue and limit. At the same time such systems are increasingly constrained by the density of connections – we have difficulty getting large enough amounts of information over wires and even fibers in our connections between or even within equipment racks.

The solution to such problems has been to use more optics, building on the technology developed first for long-distance communications, and evolving it for shorter distance use at lower power and cost. Optics avoids communications bandwidth and density limits of electrical wired connections over medium and long distances. Increasing bandwidth on optical fibers seems to offer a continuing path to higher densities. But there are problems with that approach; such increases in bandwidth can lead to increasing energy per bit, and even then we may still not have enough density of connections.

If we look microscopically down to the gate level in electronics, we can note that the energy to switch a CMOS gate is comparable to the energy required to charge a wire that could connect one gate to a nearby gate. Since most connections must go farther than this, we can conclude immediately that most energy is dissipated in interconnections, not in logic, even inside chips. Now, it would still be hard to argue for optical interconnects at short distances on chip, where wires are very inexpensive and the absolute interconnect energies are not yet very high (in the scale of femtojoules (fJ) to 100's of femtojoules per bit inside chips), but the farther we go, the more this electrical interconnect energy increases. The scale of this energy is set by the underlying capacitance of wires, which is quite generally of the order of picofarads (pF) per centimeter. By the time we reach the edge of the chip charging lines to go the next chip or further, we necessarily have ~ picojoule/bit (pJ/bit) or larger interconnect energies, just from charging such wiring. Changing to optical connections instead avoids such wire charging. So, if we can make the optical devices with low enough operating energy, less than the energy for charging such lengths of wiring, then we could reduce energy per bit in communications – for example, for all interconnects on and off chips and at longer distances.

Indeed, we already exploit this lower energy over distance in optical fiber communications and interconnects. Could we not simply continue to do this at every shorter distances, thereby progressively eliminating a large source of all power dissipation in information processing? The answer is yes, but there is a catch; in the way we do things now, and in evolutionary paths from our prior approaches, we reach multiple different limits at the level of picojoules per bit. This is frustrating because now we have viable approaches for optoelectronic devices themselves at the level of fJ's/bit or even lower [1].

The reasons for these pJ/bit energies are several. First, we do not integrate our optoelectronic devices well enough to be able to take advantage of these 10fJ/bit or lower device energies. We need integration of devices (especially photodetectors) within microns of the relevant electronics so that capacitances are in the fF range. This is possible, though it would require some technology development of integration technologies and the incorporation of materials that would allow very low device energies (e.g., Ge quantum wells for fJ/bit electroabsorption modulators).

A second reason for pJ/bit energies has nothing to do with the optoelectronic devices themselves; rather it is because of the way we design systems and the particular ways we try to use optics. These design approaches lead to the need for various interface circuits each with energies in scale of pJ's/bit. To get to the scale of 10 fJ/bit, we need to eliminate such circuits.

Integration of small (micron scale) photodetectors right beside transistors could allow us to eliminate much or all of the pJ/bit energies of receiver circuits; in the most extreme “receiverless” operating mode, the input received optical energy could then even be enough to swing the low total input capacitance by a logic voltage level.

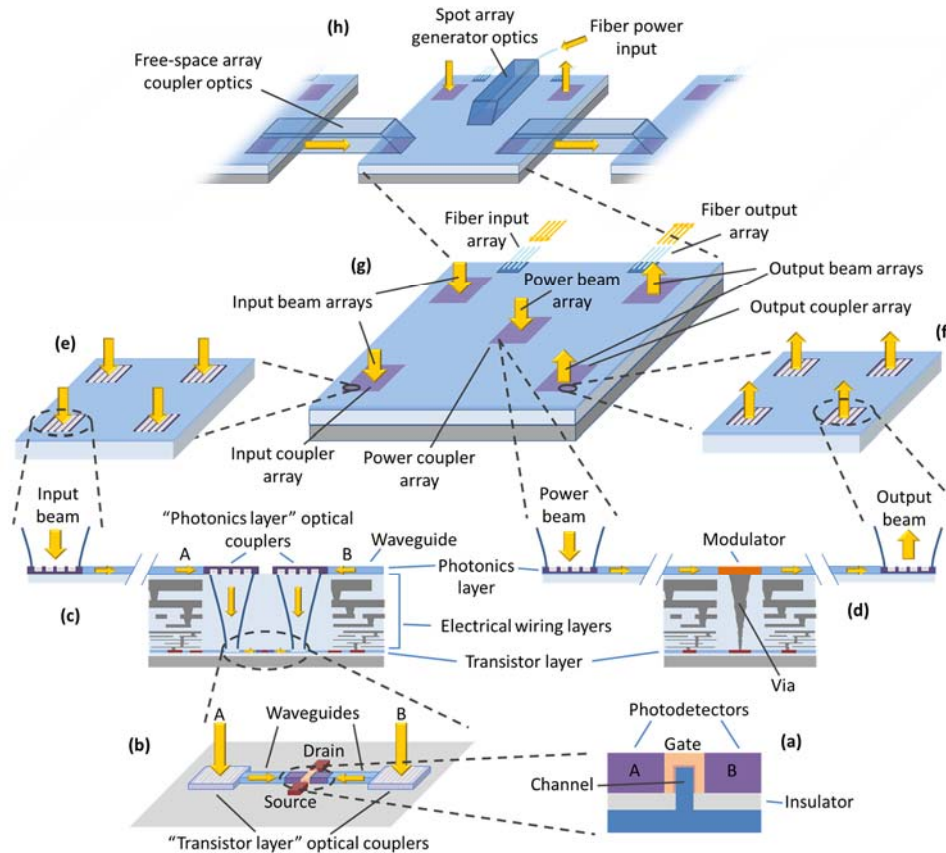


Fig. 1. Sketch of an optical platform for dense, low-energy interconnects, shown at multiple different length scales, from the transistors up to free-space arrays off a larger chip [1] (copyright IEEE, used with permission). (a) A pair of photodetectors is integrated beside the gate of the corresponding transistor input (here shown in the form of a FinFET structure). (b) Optical beams A and B are connected through “transistor-layer” couplers and short (e.g.,  $\sim 1\mu\text{m}$ ) waveguides to the photodetectors. (c) A photonics layer (e.g., as in silicon photonics) on top of the electrical wiring layers of the chip, with couplers between free space and waveguides in the photonics layer, and also from the photonics layer to detectors below. (d) Electrical “via” connections through the electrical wiring layer connect from output transistors to modulators in waveguides in the photonics layer. (e) and (f) show portions of input and output couplers and beam arrays. (g) shows a larger picture of the photonics layer on top of the entire chip, with various 2 dimensional coupler arrays: input and array coupler and beam arrays; a power array coupler and beam array; and linear arrays of fiber inputs and outputs. (h) shows spot array generator optics fed by input power from some central optical power source through a fiber, and how multiple chips might be connected latterly and vertically using free-space connections.

The need for the other pJ/bit circuits arises from our system approach. First, to get more information over fibers, we time-multiplex the signals. Such time multiplexing requires considerable energy because bits must be moved in registers, and at high speed. We might think electronic processing is “free” because of low transistor energies, but in fact, we should consider some amount of energy spent for every time we touch a bit. We cannot be precise about

that energy, but a general scale of it, given logic and interconnect, would be 1 – 100 fJ. If we think about the number of times we touch a bit in such time multiplexing processes, and we include some multiplying factor for running circuits at very high rates, we can see why pJ/bit energies arise in such serializer/deserializer (SERDES) circuits.

A related problem is that we run large systems with parts that are asynchronous with one another, so we need clock and data recovery circuits (CDR), with ~ pJ/bit energies, to line up the timing of communicated information. This has historically been unavoidable, in part because the time delay of signal propagation on electrical wires is not reliably predictable because of the temperature coefficient of the resistivity of copper.

Given these difficulties, we might conclude that even with our best use of optics, we are therefore stuck with pJ/bit energies for communications in and out of chips, in racks and between racks. The 10fJ/bit or below level that integrated optoelectronics itself might offer would seem to be unattainable given these circuit energies.

There is, however, a way out of this difficulty that could indeed let us progress to 10fJ/bit system energies. It involves using two other features of optics that we have so far not exploited. These features require no new physics, nor even any substantially new technology, but they do require that we change the way we physically build systems. The key feature of optics that has orders of magnitude of headroom available is using free-space parallelism. Optical systems can readily support thousands to millions of parallel channels (as evidenced by the optics inside every cell phone). To exploit this, we would move to free-space parallel arrays of light beams between chips and even racks. The key equation that shows the maximum number of free-space channels is [1,2]

$$N_{max} \approx \frac{A_T A_R}{\lambda^2 L^2}$$

assuming transmitter and receiver array areas  $A_T$  and  $A_R$  and separation  $L$ , and optical wavelength  $\lambda$ ; this supports 1000's to 10,000's of channels with millimeter areas and centimeter separations, or millions of channels with 10's of cm areas and meter separations. The optics to do this does exist and has been tested in laboratory situations. For short connections (chip to chip) these could be rigid optics. For longer connections (e.g. meters), array alignment could be servo-ed. The use of such massively parallel connections would mean time-multiplexing is not necessary – we could run the whole system at the efficient GHz clock rates of chips themselves. Fig. 1 shows an example physical system concept. The major required technologies are a photonics layer hybridized onto the chip, allowing dense, low-capacitance electrical connections to the electronics, and that supports efficient optoelectronic devices, waveguides, and free-space parallel coupler arrays and optics with 1000's to 10,000's of connections.

The other feature of optics to exploit is that, in optics, time delays are extremely predictable. Hence, we could contemplate running entire large systems entirely synchronous, at least within integer numbers of clock cycles of delay. Clock and data recovery is then unnecessary.

This proposal is, of course, radical. But (1) it is feasible with finite development of demonstrated technologies and (2) there does not seem to be any viable alternative solution that offers such substantially reduced energy and increased communications density inside information processing systems. Though such an approach is radical overall, many of the key technologies required – (a) denser, lower capacitance optoelectronic integration, (b) more efficient optical coupling, and (c) efficient, low-energy optical output devices, such as electroabsorption modulators in high performance materials – are ones that we would want in more conventional systems anyway, so it may be possible to approach such a radical change in smaller steps.

These arguments are discussed in detail in Ref. [1], and briefly in a broader context in [3]. The physical processes and system concepts are presented and used to propose a “straw man” system involving appropriate integrated optoelectronics and a combination of guided wave and free-space optical systems (Fig. 1). Specific arguments are made, with explicit numbers, as to how we might be able to reduce from 10's of pJ/bit to 10's of fJ/bit for all off-chip and even rack-to-rack connections. The talk will present and discuss these various arguments.

[1] D. A. B. Miller, “Attojoule Optoelectronics for Low-Energy Information Processing and Communications: A Tutorial Review,” *IEEE/OSA J. Lightwave Technology* **35** (3), 343-393 (2017) doi: 10.1109/JLT.2017.2647779

[2] D. A. B. Miller, “Waves, modes, communications, and optics: a tutorial,” *Adv. Opt. Photon.* **11**, 679-825 (2019) doi: 10.1364/AOP.11.000679

[3] J. M. Kahn and D. A. B. Miller, “Communications expands its space,” *Nature Photonics* **11**, 5 – 8 (2017) doi:10.1038/nphoton.2016.256