

From Physics to Logic

This course aims to introduce you to the layers of abstraction of modern computer systems. We won't spend much time below the level of bits, bytes, words, and functional units, but I think you should at least be aware of the richness that exists below this level and what implications it has for the higher layers.

Simulating Classical Physics With Quantum Physics

Modern computer systems are based ultimately on quantum physics. However, they use the quantum effects of materials essentially to simulate classical physics. Most mathematical theories of computation also assume classical physics. How we built computers and how we think about computers have matched for over half of a century. To a reasonable approximation, all current computers can be thought of as Universal Turing Machines, we can think of UTMs in terms of a Newtonian game of billiards, and we can implement that game by exploiting the quantum physics of certain materials.

Starting in the early 1980s, models of computation based directly on quantum physics have been developed. These models may be more powerful than UTMs in a theoretical sense. In the mid 1990s, Peter Schor published an algorithm that factored large numbers into primes very quickly on a computational model known as a Quantum Turing Machine. That this task is believed to be difficult on classical computers under-girds almost all of cryptography, and thus Schor's result generated much interest. At this point, however, we don't know whether it is feasible to build a large scale quantum computer.

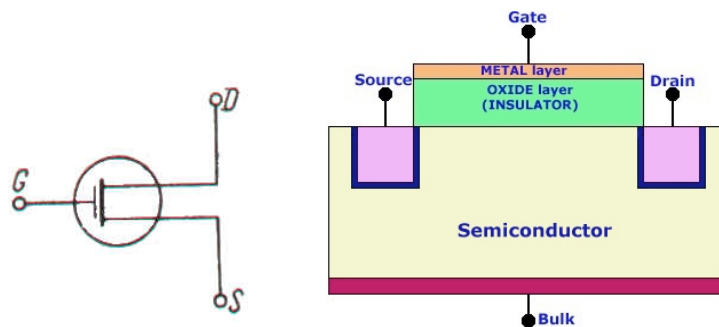
Semiconductors

The materials from which processors, memory, and other devices are constructed are known as semiconductors. While semiconductors are created using very mundane materials (silicon is the main ingredient of sand!), sophisticated chemistry is used to "grow" giant perfect silicon crystals and to subtly "dope" them with "impurities" in order to carefully control their properties. In addition to semiconductors, ultra-pure conductors, often created using copper, and insulators, typically created using silicon oxide (sand rust!) are the core materials of computers.

Transistors

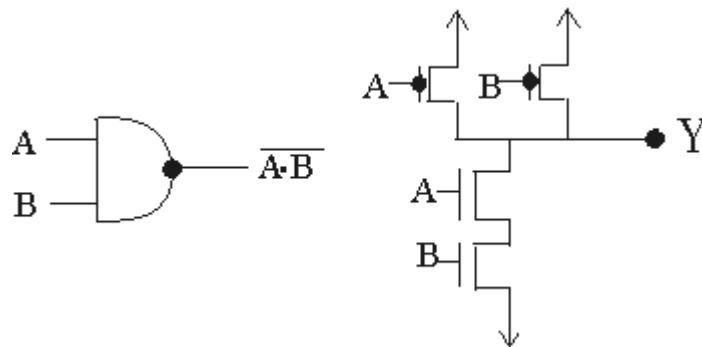
Semiconductors, insulators, and conductors are used to construct passive electronic components such as wires, resistors, capacitors, and inductors (rarely). However, the most important component built from semiconductors is an active one, the transistor. Transistors are three terminal devices, and there are many kinds. In the MOSFET, which is the typical kind of transistor in modern computer chips, the voltage on one terminal (the gate) controls how difficult it is to send electrical current between two other terminals (the source and drain).

In the nonlinear circuitry of computer chips (as opposed to the linear circuitry of, say, an audio amplifier), we usually use the gate to “turn on” and “turn off” the flow of current, like a simple valve. Below, you can see the symbol of a MOSFET and a cross-section of its typical construction. On a current processor chip, a MOSFET transistor has a size measured in the 10s of nanometers, about 1,000 to 10,000 times smaller than the width of a human hair.



Logic and Memory

Using transistors and capacitors, we can create the combinational logic and memory that are the basis of computers. For example, the following shows the symbol for a NAND (“Not And”) at the left and its implementation using four MOSFETs at right.

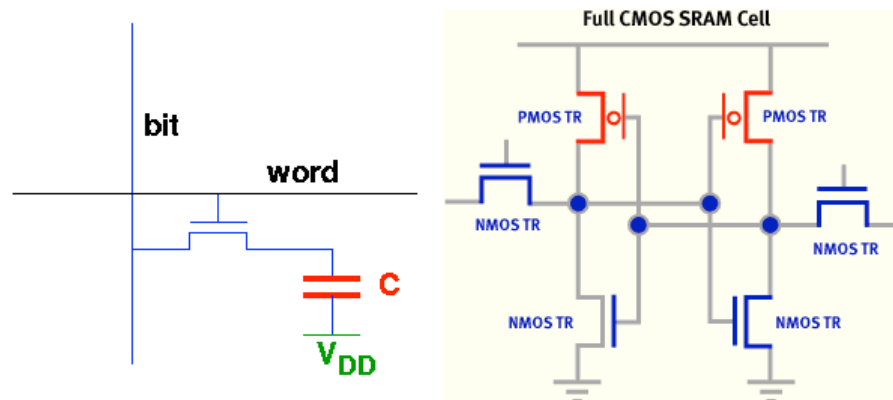


The solid circles (often shown as open circles as well) represent inversion. Basically, a circle at the gate of a MOSFET indicates that as gate voltage decreases, it gets easier to

send current from source to drain. Without a circle, decreasing voltage makes it gets harder.

The two inputs to the NAND represent truth values (true or false). The output of the NAND is true unless both inputs are true, in which the output is false. While this might seem like a strange logical operation, NAND is powerful because it is a *universal* logical operation – any other combinational logic operation (and, or, not, xor, etc) or logic circuit built out of them can be implemented by wiring NANDs together.

We can also use transistors (and other components) to build memory cells. Below, at left we see a DRAM cell (1 bit of memory), which is what main memory on most computers consists of. At right, we see an SRAM cell (again, 1 bit of memory), which is the typical cell in cache memories, processor register files, and other components. There is a third kind of memory cell, a latch, which is widely used in implementing processor pipelines.



An important thing to notice is that an SRAM cell uses no capacitors but uses six transistors instead of just one. SRAM cells are much faster than DRAM cells, but also much more expensive. That's why they tend to be put in caches. An SRAM cell will also hold its bit as long as there is power. In a DRAM cell, because the bit is held in a capacitor and capacitors always leak charge, the bit must be frequently refreshed by external circuitry. In your typical computer, all of the main memory is read and rewritten 10-20 times per second to keep the contents from being lost.

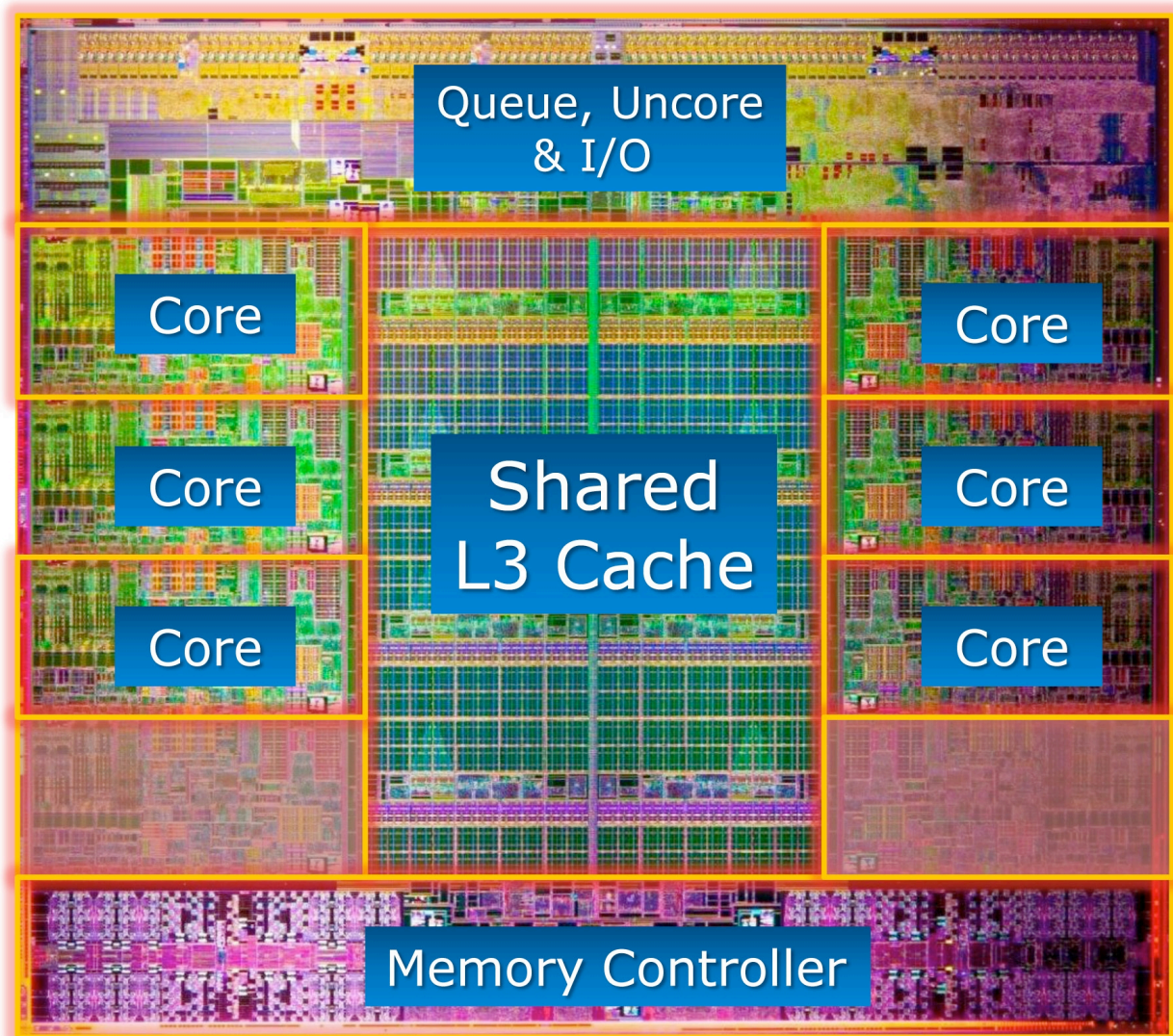
Photolithography and Chips

Transistors were first developed in the late 1940s. By the late 1950s, engineers had learned how to mass-produce individual transistors inexpensively. Around this time, Jack Kilby, an engineer at Texas Instruments, had a vision of putting building multiple transistors and passive devices on a single semiconductor substrate. He developed the first monolithic integrated circuit (IC or "chip"), inventing a technology that is still going strong today.

A chip is produced using a process called photolithography, which you can think of as a strange kind of photographic printing. If you've ever done black-and-white printing in a darkroom, you understand the gist of it already. Essentially, a negative called a mask is

produced that is the inverse of the layer of material to be put on the semiconductor substrate. The substrate is coated with the material that will be deposited and then covered in a material called photoresist. The mask is then placed over the photoresist and a bright light (ultraviolet light these days) is shown on it. Next, the photoresist is developed, and then the unexposed photoresist is etched away (like fixing a photographic print). Then the next layer can be created over the current one. Chips are built up out of multiple layers like this.

Below you'll see a bird's eye view of the Intel Core i7-3930K processor, which costs about \$500. This whole chip has an area of only 430 square millimeters (about double the size of the nail on your big toe) yet contains about 2.2 billion transistors.



Moore's Law, Its Limits, And So Many Transistors

In the late 1960s, Gordon Moore made a plot of the number of transistors on a chip as a function of the year the chip came out and found that it was exponential. Every 18 months, the number of transistors would double, and chips would become much faster because the transistors were smaller. He predicted that this would continue, and he was right. Moore's Law still holds true, except now the doubling period is just one year. Experts expect that it will continue for at least another 10 years. The current ways that we build processors will not be able to make use of this many transistors, and thus there is much current research on how to actually exploit having billions of transistors on a chip.

Multicore Processors, GPUs, and Parallelism

An important recent shift in processors has been the advent of the multicore era. In a multicore processor, the chip contains multiple copies of the processor. Even commonplace desktop microprocessors today have four cores (four copies of the processor). Our class's server machine has 40 cores, while a \$10,000 server can easily have 64 cores today. The Intel Phi within the class server has over 240. It is expected that since the number of transistors on chip will continue to scale exponentially for some time, the consequence is that the number of cores (copies of the processor) per chip will also grow exponentially for some time, probably doubling every year or two.

In addition to such completely general purpose processor cores, manufacturers of graphics hardware have also been making their hardware more general purpose, creating "GPUs" (Graphics Processing Units) that have a much larger number of cores, each of which is much "smaller" than a general purpose core. For example, the \$3,000 NVIDIA Tesla K20, which fits into a typical desktop or server computer, sports almost 2,500 cores. Our class's server machine also has one of these.

Having k cores, no matter what kind, does not mean that existing programs will get k times faster. A multicore chip is a parallel computer, meaning that it can do several things simultaneously at the hardware level. Parallel computers have been around for over 30 years, but the focus of research in parallel systems to this point has largely been to support scientific computations. With parallel computers on everyone's desktop and in everyone's cell phone, existing and new ideas from parallel systems will need to be employed. The challenge is how to feed the parallelism. Currently it is *the programmer* who needs to think about k things the chip should be doing simultaneously. This is unlikely the change anytime soon.