# The Physics
# of Computing

## Marilyn Wolf

# The Physics of Computing

**Marilyn Wolf**

**Notices**
Knowledge and best practice in this field are constantly changing. As new research and experience broaden our understanding, changes in research methods, professional practices, or medical treatment may become necessary.

Practitioners and researchers must always rely on their own experience and knowledge in evaluating and using any information, methods, compounds, or experiments described herein. In using such information or methods they should be mindful of their own safety and the safety of others, including parties for whom they have a professional responsibility.

To the fullest extent of the law, neither the Publisher nor the authors, contributors, or editors, assume any liability for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions, or ideas contained in the material herein.

For information on all Morgan Kaufmann publications
visit our website at https://www.elsevier.com/

Working together
to grow libraries in
developing countries

www.elsevier.com • www.bookaid.org

# Contents

# Preface

This book is an attempt to take a step back and provide a longer, more foundational view of computer engineering. While a lot of software design proceeded from the theoretical to the experimental—think about sorting, for example—computer system design has long taken an example-driven approach. Studying particular computer systems was useful in the early days when we were not quite sure what major concepts undergirded the design of these systems. But today we are much more capable of making generalized statements about many topics in computer engineering. Those foundational topics should form the core of our approach to computer engineering.

A foundational approach to digital systems is particularly important because it is so hard today to touch and see logic in operation. When I was a student, most systems were boards made of SSI and MSI logic. We had no choice but to deal with circuits and waveforms. Today, all that is hidden inside a chip; even boards have smaller wires and are harder to analyze. Many students today think that 1s and 0s float through chips—they have no conception of voltages and currents in computers.

This book does not ask students to design *anything*. Instead, we walk students through the principles that define the computer system design space. I like to think of the concepts in this course as knobs and meters on a control panel. Changing the settings of knobs change the values on the meters; the twist of one knob may change several meters. Consider, for example, the cascading effects of reducing MOS dielectric thickness: transistor transconductance changes, which in turn affects gate delay for the better and leakage current for the worse.

Engineers do not design systems in a vacuum. Instead, they design to a set of goals or requirements. The traditional goal for computer architects is performance, or more precisely throughput. But computers are actually designed with several other goals in mind, with energy/power and reliability paramount. Performance, energy, and reliability all have their roots in physical phenomenon. And they are inextricably linked by physics: improving one metric incurs costs in the other metrics. As in the rest of life, there is no free lunch in computer design.

Computer engineering is a relatively young field. Introductory courses have concentrated on building and taxonomy—build something, see the different ways to build something. As the field matures, it is time to start thinking about new ways to present the basics. The biologist E.O. Wilson said, "A field is initially defined by the questions it asks and eventually defined by the answers it provides." After 70 years of electronic computer engineering, it is time to start thinking of the field in terms of answers.

I have been thinking about computer design for a long time. Performance was for many years the paramount metric for computer design, although the ASIC world was always interested in area as a proxy for cost. I started to think more seriously about other fundamental limits, particularly power, for the class on pervasive information systems taught by Perry Cook and myself. A new course at Georgia Tech gave me the impetus to take this train of thought to its logical conclusion.

This course takes a broad view of both computer engineering and physics. Some of the most basic phenomena in computer architecture and even software design—the memory wall, the power wall, the race to dark—are due to fundamental physics. But the physics required to understand these major computational issues is much more than classical Shockley semiconductor theory. We need to understand thermodynamics, electrostatics, and good deal of circuit theory.

As I thought about this material, I came to regard Boltzmann's constant is a key concept. $k$ pops up everywhere: the diode equation, Arrhenius's equation, temperature, and the list goes on. Boltzmann's constant links temperature and energy, so it should not be a surprise that it is so tightly woven into the topics of this book.

Some of the material in this book is unique to modern CMOS—leakage mechanisms. But some of the topics can be applied to a wide range of circuit and device technologies. Delay through logic networks, metastability, and the fundamentals of reliability are examples of foundational concepts that computer engineers should understand even if CMOS is swept off the map by some other technology.

Parts of this book will probably strike some readers as hopelessly simplified and compressed. I hope that those same readers will find other parts of the book dense and hard to follow. The only way to understand how computers work is to understand the relationships between an entire collection of topics that do not always seem at first glance to be related. We know those knobs on the design process are linked because a change to one parameter that we think will help our design often results in side effects that negate much of the benefit we sought. I worked very hard to find the simplest possible description for a lot of concepts. Hopefully they are presented in a way that gives a basic understanding. The reader who is interested in a deeper understanding of one of these topics can pursue each one in more depth. But the singular goal of this book is to provide a unified description of the fundamental physical principles of computing machines.

Both computer engineers and electrical engineers are potential audiences for this book. These two groups come to the material with very different backgrounds. Computer engineers often have limited experience in circuit design. While they may have some knowledge of ideas like Kirchhoff's laws, they usually are not extremely comfortable with circuit analysis. Electrical engineers may not always have a detailed understanding of computer architecture. One of my challenges in writing this book was to provide enough background for each audience without overloading them.

I have mentioned a number of historical discoveries and inventions for several reasons. First, going through the sequence of ideas that led to today's design practices helps both to remind us that there is often more than one way to do things and to highlight the true advantages of the methods that were settled upon by history. Second, semiconductor physics and computer engineering have produced some of the most important inventions of the 20th century. These inventions will be used for centuries to come. We should not let ourselves become so accustomed to these ideas that we lose all capacity to marvel at them and at their inventors.

Richard Feynman's *Lectures on Computing* was an early inspiration for this class, but much of that book is devoted to quantum computing. That book does not consider many topics in traditional computing. For example, metastability is a fundamental physical concept in computer system design that Feynman does not touch. We salute him for his early recognition of the physical nature of computing. We also thank him for the *Lectures on Physics*, which provided a clear, concise way to think about many of the basic physical phenomena that underlie computing.

Several tips of the hat are in order. My friend and colleague Saibal Mukhopodhyay proposed the *Physical Foundations of Computer Engineering* course and provided many concrete suggestions, particularly on reliability and leakage. My thanks go out to him for his patience and insight. Dave Coelho graciously provided information on his power distribution system. Kees Vissers invented the welder current comparison. Alec Ishii advised me on clock distribution. Kevin Cao suggested how to make best use of the Predictive Technology Models. Bruce Jacob gave insight into DRAM. Srini Devadas's suggestions led to the advanced topics appendix. Tom Conte provided the core memory and Pentium Pro as well as many discussions on the present state and future of computing. Thanks to the reviewers for their helpful comments. And many thanks to my editor Nate McFadden for guiding this book through the long development and production process. Any faults in this book can be traced entirely and solely to me.

**Marilyn Wolf**
Atlanta, GA

# Electronic Computers

## 1.1 Introduction

We take computers for granted. But what is a computer? The computers we use every day can be summarized as **EBT**: *electronic, binary, Turing machines*. This combination of technologies has allowed us to conquer problems as diverse as weather forecasting, self-driving vehicles, and shopping.

But the path that led us to today's computers is long, and the EBT approach to computing is surprisingly recent. Some of the history helps us understand why we build computers the way we do. Machines that could compute—calculate with numbers and perform other multistep tasks—were originally conceived of and built with the only devices available, namely mechanical devices. These mechanical computers had severe limitations that were overcome only with the creation of electronic devices.

We will start with a brief introduction to the history of computing machines. We will then discuss the theory of computing that was developed based on this early experience and which underpins our modern ideas of computers. We then discuss the metrics we use to evaluate computer designs and the trade-offs between these goals that physics imposes on us. We will wrap up with a survey of the remainder of this book.

## 1.2 The long road to computers

Mechanical computation was inspired in large part by the need to control machines. People could not react quickly and reliably enough to maintain the proper operation of the equipment that appeared during the Industrial Revolution. Early mechanical computing devices were analog—they operated on continuous values. But machines that manipulated discrete values arrived surprisingly early. Eventually, theory was developed to describe these types of machines. One of those theories, the *Turing machine*, is the blueprint for the way in which we design computers today.
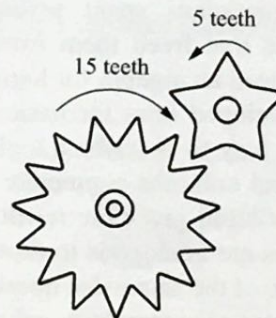
Weft

Warp

**FIGURE 1.2**

Operation of a Jacquard loom.

descend and grab its weft thread; no hole meant that the weft was not lifted at that point. Not only could the card represent complex two-dimensional patterns, but the entire fabric's pattern could be changed simply by replacing the card.

**Babbage machines**　　Charles Babbage started work in 1822 on a **difference machine** to calculate numbers, specifically to compute polynomial functions. His work was motivated by the need to tables of functions that were more accurate than those compiled by hand calculation, an error-prone process. At the heart of this machine's operation is the use of gears to count; this principle was also used in the mechanical calculators that were widely used well into the second half of the 20th century. The basic principle behind mechanical counting is simple and elegant, as shown in Fig. 1.3. We use a shaft whose position represents the value of a digit in the counting sum. We can put a radial mark on the end of each shaft, then a dial around the shaft with a position for each digit; the mark points the value of the shaft's digit. To build a two-digit counter, we connect the two shafts with gears. In this figure, the small gear has 5 teeth and the large gear has 15 teeth, giving a gearing ratio of 3:1. Three full rotations of the small gear produce one complete rotation of the large gear. The large gear remembers the number of teeth turned against it by the small gear up to a limit of 15. If we want to count in base 10, we gear the shafts at a 10:1 ratio: 10 rotations of the ones-digit



5 teeth

15 teeth

**FIGURE 1.3**

Using gears to count.