

# THE P VERSUS NP PROBLEM

STEPHEN COOK

## 1. STATEMENT OF THE PROBLEM

The **P** versus **NP** problem is to determine whether every language accepted by some nondeterministic algorithm in polynomial time is also accepted by some (deterministic) algorithm in polynomial time. To define the problem precisely it is necessary to give a formal model of a computer. The standard computer model in computability theory is the Turing machine, introduced by Alan Turing in 1936 [37]. Although the model was introduced before physical computers were built, it nevertheless continues to be accepted as the proper computer model for the purpose of defining the notion of *computable function*.

Informally the class **P** is the class of decision problems solvable by some algorithm within a number of steps bounded by some fixed polynomial in the length of the input. Turing was not concerned with the efficiency of his machines, rather his concern was whether they can simulate arbitrary algorithms given sufficient time. It turns out, however, Turing machines can generally simulate more efficient computer models (for example, machines equipped with many tapes or an unbounded random access memory) by at most squaring or cubing the computation time. Thus **P** is a robust class and has equivalent definitions over a large class of computer models. Here we follow standard practice and define the class **P** in terms of Turing machines.

Formally the elements of the class **P** are languages. Let  $\Sigma$  be a finite alphabet (that is, a finite nonempty set) with at least two elements, and let  $\Sigma^*$  be the set of finite strings over  $\Sigma$ . Then a *language over  $\Sigma$*  is a subset  $L$  of  $\Sigma^*$ . Each Turing machine  $M$  has an associated *input alphabet*  $\Sigma$ . For each string  $w$  in  $\Sigma^*$  there is a computation associated with  $M$  with input  $w$ . (The notions of Turing machine and computation are defined formally in the appendix.) We say that  $M$  *accepts*  $w$  if this computation terminates in the accepting state. Note that  $M$  fails to accept  $w$  either if this computation ends in the rejecting state, or if the computation fails to terminate. The *language accepted by  $M$* , denoted  $L(M)$ , has associated alphabet  $\Sigma$  and is defined by

$$L(M) = \{w \in \Sigma^* \mid M \text{ accepts } w\}.$$

We denote by  $t_M(w)$  the number of steps in the computation of  $M$  on input  $w$  (see the appendix). If this computation never halts, then  $t_M(w) = \infty$ . For  $n \in \mathbb{N}$  we denote by  $T_M(n)$  the *worst case run time* of  $M$ ; that is,

$$T_M(n) = \max\{t_M(w) \mid w \in \Sigma^n\},$$

where  $\Sigma^n$  is the set of all strings over  $\Sigma$  of length  $n$ . We say that  $M$  *runs in polynomial time* if there exists  $k$  such that for all  $n$ ,  $T_M(n) \leq n^k + k$ . Now we

define the class **P** of languages by

$$\mathbf{P} = \{L \mid L = L(M) \text{ for some Turing machine } M \text{ that runs}$$

in polynomial time}\}.

The notation **NP** stands for “nondeterministic polynomial time”, since originally **NP** was defined in terms of nondeterministic machines (that is, machines that have more than one possible move from a given configuration). However, now it is customary to give an equivalent definition using the notion of a *checking relation*, which is simply a binary relation  $R \subseteq \Sigma^* \times \Sigma_1^*$  for some finite alphabets  $\Sigma$  and  $\Sigma_1$ . We associate with each such relation  $R$  a language  $L_R$  over  $\Sigma \cup \Sigma_1 \cup \{\#\}$  defined by

$$L_R = \{w\#y \mid R(w, y)\}$$

where the symbol  $\#$  is not in  $\Sigma$ . We say that  $R$  is *polynomial-time* iff  $L_R \in \mathbf{P}$ .

Now we define the class **NP** of languages by the condition that a language  $L$  over  $\Sigma$  is in **NP** iff there is  $k \in \mathbb{N}$  and a polynomial-time checking relation  $R$  such that for all  $w \in \Sigma^*$ ,

$$w \in L \iff \exists y (|y| \leq |w|^k \text{ and } R(w, y)),$$

where  $|w|$  and  $|y|$  denote the lengths of  $w$  and  $y$ , respectively.

**Problem Statement.** Does  $\mathbf{P} = \mathbf{NP}$ ?

It is easy to see that the answer is independent of the size of the alphabet  $\Sigma$  (we assume  $|\Sigma| \geq 2$ ), since strings over an alphabet of any fixed size can be efficiently coded by strings over a binary alphabet. (For  $|\Sigma| = 1$  the problem is still open, although it is possible that  $\mathbf{P} = \mathbf{NP}$  in this case but not in the general case.)

It is trivial to show that  $\mathbf{P} \subseteq \mathbf{NP}$ , since for each language  $L$  over  $\Sigma$ , if  $L \in \mathbf{P}$  then we can define the polynomial-time checking relation  $R \subseteq \Sigma^* \cup \Sigma^*$  by

$$R(w, y) \iff w \in L$$

for all  $w, y \in \Sigma^*$ .

Here are two simple examples, using decimal notation to code natural numbers: The set of perfect squares is in **P**, since Newton’s method can be used to efficiently approximate square roots. The set of composite numbers is in **NP**, where (denoting the decimal notation for a natural number  $c$  by  $\bar{c}$ ) the associated polynomial time checking relation  $R$  is given by

(1)

$$R(\bar{a}, \bar{b}) \iff 1 < b < a \text{ and } b|a.$$

(Recently it was shown that in fact the set of composite numbers is also in **P** [1], answering a longstanding open question.)

## 2. HISTORY AND IMPORTANCE

The importance of the **P** vs **NP** question stems from the successful theories of **NP**-completeness and complexity-based cryptography, as well as the potentially stunning practical consequences of a constructive proof of  $\mathbf{P} = \mathbf{NP}$ .

The theory of **NP**-completeness has its roots in computability theory, which originated in the work of Turing, Church, Gödel, and others in the 1930s. The computability precursors of the classes **P** and **NP** are the classes of decidable and c.e. (computably enumerable) languages, respectively. We say that a language  $L$  is c.e. (or *semi-decidable*) iff  $L = L(M)$  for some Turing machine  $M$ . We say that  $L$

is *decidable* iff  $L = L(M)$  for some Turing machine  $M$  that satisfies the condition that  $M$  halts on all input strings  $w$ . There is an equivalent definition of c.e. that brings out its analogy with **NP**, namely  $L$  is c.e. iff there is a computable “checking relation”  $R(x, y)$  such that  $L = \{x \mid \exists y R(x, y)\}$ .

Using the notation  $\langle M \rangle$  to denote a string describing a Turing machine  $M$ , we define the Halting Problem  $HP$  as follows:

$$HP = \{\langle M \rangle \mid M \text{ is a Turing machine that halts on input } \langle M \rangle\}.$$

Turing used a simple diagonal argument to show that  $HP$  is not decidable. On the other hand, it is not hard to show that  $HP$  is c.e.

Of central importance in computability theory is the notion of reducibility, which Turing defined roughly as follows: A language  $L_1$  is *Turing reducible* to a language  $L_2$  iff there is an oracle Turing machine  $M$  that accepts  $L_1$ , where  $M$  is allowed to make membership queries of the form  $x \in L_2$ , which are correctly answered by an “oracle” for  $L_2$ . Later, the more restricted notion of many-one reducibility ( $\leq_m$ ) was introduced and defined as follows.

**Definition 1.** Suppose that  $L_i$  is a language over  $\Sigma_i$ ,  $i = 1, 2$ . Then  $L_1 \leq_m L_2$  iff there is a (total) computable function  $f : \Sigma_1^* \rightarrow \Sigma_2^*$  such that  $x \in L_1 \iff f(x) \in L_2$ , for all  $x \in \Sigma_1^*$ .

It is easy to see that if  $L_1 \leq_m L_2$  and  $L_2$  is decidable, then  $L_1$  is decidable. This fact provides an important tool for showing undecidability; for example, if  $HP \leq_m L$ , then  $L$  is undecidable.

The notion of **NP**-complete is based on the following notion from computability theory:

**Definition 2.** A language  $L$  is *c.e.-complete* iff  $L$  is c.e., and  $L' \leq_m L$  for every c.e. language  $L'$ .

It is easy to show that  $HP$  is c.e.-complete. It turns out that most “natural” undecidable c.e. languages are, in fact, c.e.-complete. Since  $\leq_m$  is transitive, to show that a c.e. language  $L$  is c.e.-complete it suffices to show that  $HP \leq_m L$ .

The notion of polynomial-time computation was introduced in the 1960s by Cobham [8] and Edmonds [13] as part of the early development of computational complexity theory (although earlier von Neumann [38], in 1953, distinguished between polynomial-time and exponential-time algorithms). Edmonds called polynomial-time algorithms “good algorithms” and linked them to tractable algorithms.

It has now become standard in complexity theory to identify polynomial-time with feasible, and here we digress to discuss this point. It is of course not literally true that every polynomial-time algorithm can be feasibly executed on small inputs; for example, a computer program requiring  $n^{100}$  steps could never be executed on an input even as small as  $n = 10$ . Here is a more defensible statement (see [10]):

**Feasibility Thesis.** *A natural problem has a feasible algorithm iff it has a polynomial-time algorithm.*

Examples of natural problems that have both feasible and polynomial-time algorithms abound: Integer arithmetic, linear algebra, network flow, linear programming, many graph problems (connectivity, shortest path, minimum spanning tree, bipartite matching), etc. On the other hand, the deep results of Robertson and Seymour [29] provide a potential source of counterexamples to the thesis: They

prove that every minor-closed family of graphs can be recognized in polynomial time (in fact, in time  $O(n^3)$ ), but the algorithms supplied by their method have such huge constants that they are not feasible. However, each potential counterexample can be removed by finding a feasible algorithm for it. For example, a feasible recognition algorithm is known for the class of planar graphs, but none is currently known for the class of graphs embeddable in  $\mathbb{R}^3$  with no two cycles linked. (Both examples are minor-closed families.) Of course the words “natural” and “feasible” in the thesis above should be explained; generally we do not consider a class with a parameter as natural, such as the set of graphs embeddable on a surface of genus  $k$ ,  $k > 1$ .

We mention two concerns related to the “only if” direction of the thesis. The first comes from randomized algorithms. We discuss at the end of Section 3 the possibility that a source of random bits might be used to greatly reduce the recognition time required for some language. Note, however, that it is not clear whether a truly random source exists in nature. The second concern comes from quantum computers. This computer model incorporates the idea of superposition of states from quantum mechanics and allows a potential exponential speed-up of some computations over Turing machines. For example, Shor [32] has shown that some quantum computer algorithm is able to factor integers in polynomial time, but no polynomial-time integer-factoring algorithm is known for Turing machines. Physicists have so far been unable to build a quantum computer that can handle more than a half-dozen bits, so this threat to the feasibility thesis is hypothetical at present.

Returning to the historical treatment of complexity theory, in 1971 the present author [9] introduced a notion of **NP**-completeness as a polynomial-time analog of c.e.-completeness, except that the reduction used was a polynomial-time analog of Turing reducibility rather than of many-one reducibility. The main results in [9] are that several natural problems, including Satisfiability and 3-SAT (defined below) and subgraph isomorphism are **NP**-complete. A year later Karp [21] used these completeness results to show that 20 other natural problems are **NP**-complete, thus forcefully demonstrating the importance of the subject. Karp also introduced the now standard notation **P** and **NP** and redefined **NP**-completeness using the polynomial-time analog of many-one reducibility, a definition that has become standard. Meanwhile Levin [23], independently of Cook and Karp, defined the notion of “universal search problem”, similar to the **NP**-complete problem, and gave six examples, including Satisfiability.

The standard definitions concerning **NP**-completeness are close analogs of Definitions 1 and 2 above.

**Definition 3.** Suppose that  $L_i$  is a language over  $\Sigma_i$ ,  $i = 1, 2$ . Then  $L_1 \leq_p L_2$  ( $L_1$  is p-reducible to  $L_2$ ) iff there is a polynomial-time computable function  $f : \Sigma_1^* \rightarrow \Sigma_2^*$  such that  $x \in L_1 \iff f(x) \in L_2$ , for all  $x \in \Sigma_1^*$ .

**Definition 4.** A language  $L$  is **NP**-complete iff  $L$  is in **NP**, and  $L' \leq_p L$  for every language  $L'$  in **NP**.

The following proposition is easy to prove: Part (b) uses the transitivity of  $\leq_p$ , and part (c) follows from part (a).

**Proposition 1.** (a) If  $L_1 \leq_p L_2$  and  $L_2 \in \mathbf{P}$ , then  $L_1 \in \mathbf{P}$ .  
(b) If  $L_1$  is **NP**-complete,  $L_2 \in \mathbf{NP}$ , and  $L_1 \leq_p L_2$ , then  $L_2$  is **NP**-complete.

(c) If  $L$  is NP-complete and  $L \in P$ , then  $P = NP$ .

Notice that parts (a) and (b) have close analogs in computability theory. The analog of part (c) is simply that if  $L$  is c.e.-complete then  $L$  is undecidable. Part (b) is the basic method for showing new problems are NP-complete, and part (c) explains why it is probably a waste of time looking for a polynomial-time algorithm for an NP-complete problem.

In practice, a member of NP is expressed as a decision problem, and the corresponding language is understood to mean the set of strings coding YES instances to the decision problem using standard coding methods. Thus the problem Satisfiability is: Given a propositional formula  $F$ , determine whether  $F$  is satisfiable. To show that this is in NP, we define the polynomial-time checking relation  $R(x, y)$ , which holds iff  $x$  codes a propositional formula  $F$  and  $y$  codes a truth assignment to the variables of  $F$  that makes  $F$  true. This problem was shown in [9] to be NP-complete essentially by showing that, for each polynomial-time Turing machine  $M$  that recognizes a checking relation  $R(x, y)$  for an NP language  $L$ , there is a polynomial-time algorithm that takes as input a string  $x$  and produces a propositional formula  $F_x$  (describing the computation of  $M$  on input  $(x, y)$ , with variables representing the unknown string  $y$ ) such that  $F_x$  is satisfiable iff  $M$  accepts the input  $(x, y)$  for some  $y$  with  $|y| \leq |x|^{O(1)}$ .

An important special case of Satisfiability is 3-SAT, which was also shown to be NP-complete in [9]. Instances of 3-SAT are restricted to formulas in conjunctive normal form with three literals per clause. For example, the formula

$$(2) \quad (P \vee Q \vee R) \wedge (\bar{P} \vee Q \vee \bar{R}) \wedge (P \vee \bar{Q} \vee S) \wedge (\bar{P} \vee \bar{R} \vee \bar{S})$$

is a YES instance to 3-SAT since the truth assignment  $\tau$  satisfies the formula, where  $\tau(P) = \tau(Q) = True$  and  $\tau(R) = \tau(S) = False$ .

Many hundreds of NP-complete problems have been identified, including Sub-setSum (given a set of positive integers presented in decimal notation, and a target T, is there a subset summing to T?), many graph problems (given a graph  $G$ , does  $G$  have a Hamiltonian cycle? Does  $G$  have a clique consisting of half of the vertices? Can the vertices of  $G$  be colored with three colors with distinct colors for adjacent vertices?). These problems give rise to many scheduling and routing problems with industrial importance. The book [15] provides an excellent reference to the subject, with 300 NP-complete problems listed in the appendix.

Associated with each decision problem in NP there is a search problem, which is, given a string  $x$ , find a string  $y$  satisfying the checking relation  $R(x, y)$  for the problem (or determine that  $x$  is a NO instance to the problem). Such a  $y$  is said to be a *certificate* for  $x$ . In the case of an NP-complete problem it is easy to see that the search problem can be efficiently reduced to the corresponding decision problem. In fact, if  $P = NP$ , then the associated search problem for every NP problem has a polynomial-time algorithm. For example, an algorithm for the decision problem Satisfiability can be used to find a truth assignment  $\tau$  satisfying a given satisfiable formula  $F$  by, for each variable  $P$  in  $F$  in turn, setting  $P$  to True in  $F$  or False in  $F$ , whichever case keeps  $F$  satisfiable.

The set of complements of NP languages is denoted coNP. The complement of an NP-complete language is thought not to be in NP; otherwise  $NP = coNP$ . The set TAUT of tautologies (propositional formulas true under all assignments) is the standard example of a coNP-complete language. The conjecture  $NP \neq coNP$

is equivalent to the assertion that no formal proof system (suitably defined) for tautologies has short (polynomial-bounded) proofs for all tautologies [12]. This fact has motivated the development of a rich theory of propositional proof complexity [22], one of whose goals is to prove superpolynomial lower bounds on the lengths of proofs for standard propositional proof systems.

There are interesting examples of **NP** problems not known to be either in **P** or **NP**-complete. One example is the graph isomorphism problem: Given two undirected graphs, determine whether they are isomorphic.

Another example, until recently, was the set of composite numbers. As mentioned in the first section, this set is in **NP**, with checking relation (1), but it is now known also to be in **P** [1]. However, the search problem associated with the checking relation (1) is equivalent to the problem of integer factoring and is thought unlikely to be solvable in polynomial time. In fact, an efficient factoring algorithm would break the RSA public key encryption scheme [28] commonly used to allow (presumably) secure financial transactions over the Internet.

There is an **NP** decision problem with complexity equivalent to that of integer factoring, namely

$$L_{\text{fact}} = \{\langle a, b \rangle \mid \exists d (1 < d < a \text{ and } d|b)\}.$$

Given an integer  $b > 1$ , the smallest prime divisor of  $b$  can be found with about  $\log_2 b$  queries to  $L_{\text{fact}}$ , using binary search. It is easy to see that the complement of  $L_{\text{fact}}$  is also in **NP**: a certificate showing  $\langle a, b \rangle$  is not in  $L_{\text{fact}}$  could be the complete prime decomposition of  $b$ . Thus  $L_{\text{fact}}$  is a good example of a problem in **NP** that seems unlikely to be either in **P** or **NP**-complete.

Computational complexity theory plays an important role in modern cryptography [16]. The security of the Internet, including most financial transactions, depends on complexity-theoretic assumptions such as the difficulty of integer factoring or of breaking DES (the Data Encryption Standard). If **P** = **NP**, these assumptions are all false. Specifically, an algorithm solving 3-SAT in  $n^2$  steps could be used to factor 200-digit numbers in a few minutes.

Although a practical algorithm for solving an **NP**-complete problem (showing **P** = **NP**) would have devastating consequences for cryptography, it would also have stunning practical consequences of a more positive nature, and not just because of the efficient solutions to the many **NP**-hard problems important to industry. For example, it would transform mathematics by allowing a computer to find a formal proof of any theorem that has a proof of reasonable length, since formal proofs can easily be recognized in polynomial time. Such theorems may well include all of the CMI prize problems. Although the formal proofs may not be initially intelligible to humans, the problem of finding intelligible proofs would be reduced to that of finding a recognition algorithm for intelligible proofs. Similar remarks apply to diverse creative human endeavors, such as designing airplane wings, creating physical theories, or even composing music. The question in each case is to what extent an efficient algorithm for recognizing a good result can be found. This is a fundamental problem in artificial intelligence, and one whose solution itself would be aided by the **NP**-solver by allowing easy testing of recognition theories.

Even if **P** ≠ **NP** it may still happen that every **NP** problem is susceptible to a polynomial-time algorithm that works on “most” inputs. This could render cryptography impossible and bring about most of the benefits of a world in which **P** = **NP**. This also motivates Levin’s theory [24], [18] of average case completeness,

in which the  $\mathbf{P} = \mathbf{NP}$  question is replaced by the question of whether every  $\mathbf{NP}$  problem with any reasonable probability distribution on its inputs can be solved in polynomial time on average.

In [34] Smale lists the  $\mathbf{P}$  vs  $\mathbf{NP}$  question as problem 3 of mathematical problems for the next century. However, Smale is interested not only in the classical version of the question, but also in a version expressed in terms of the field of complex numbers. Here Turing machines must be replaced by a machine model that is capable of doing exact arithmetic and zero tests on arbitrary complex numbers. The  $\mathbf{P}$  vs  $\mathbf{NP}$  question is replaced by a question related to Hilbert's Nullstellensatz: Is there a polynomial-time algorithm that, given a set of  $k$  multivariate polynomials over  $\mathbb{C}$ , determines whether they have a common zero? See [4] for a development of complexity theory in this setting.

The books by Papadimitriou [25] and Sipser [33] provide good introductions to mainstream complexity theory.

### 3. THE CONJECTURE AND ATTEMPTS TO PROVE IT

Most complexity theorists believe that  $\mathbf{P} \neq \mathbf{NP}$ . Perhaps this can be partly explained by the potentially stunning consequences of  $\mathbf{P} = \mathbf{NP}$  mentioned above, but there are better reasons. We explain these by considering the two possibilities in turn:  $\mathbf{P} = \mathbf{NP}$  and  $\mathbf{P} \neq \mathbf{NP}$ .

Suppose first that  $\mathbf{P} = \mathbf{NP}$  and consider how one might prove it. The obvious way is to exhibit a polynomial-time algorithm for 3-SAT or one of the other thousand or so known  $\mathbf{NP}$ -complete problems, and, indeed, many false proofs have been presented in this form. There is a standard toolkit available [7] for devising polynomial-time algorithms, including the greedy method, dynamic programming, reduction to linear programming, etc. These are the subjects of a course on algorithms, typical in undergraduate computer science curriculums. Because of their importance in industry, a vast number of programmers and engineers have attempted to find efficient algorithms for  $\mathbf{NP}$ -complete problems over the past 30 years, without success. There is a similar strong motivation for breaking the cryptographic schemes that assume  $\mathbf{P} \neq \mathbf{NP}$  for their security.

Of course, it is possible that some nonconstructive argument, such as the Robertson–Seymour proofs mentioned earlier in conjunction with the Feasibility Thesis, could show that  $\mathbf{P} = \mathbf{NP}$  without yielding any feasible algorithm for the standard  $\mathbf{NP}$ -complete problems. At present, however, the best proven upper bound on an algorithm for solving 3-SAT is approximately  $1.5^n$ , where  $n$  is the number of variables in the input formula.

Suppose, on the other hand, that  $\mathbf{P} \neq \mathbf{NP}$ , and consider how one might prove it. There are two general methods that have been tried: diagonalization with reduction and Boolean circuit lower bounds.

The method of diagonalization with reduction has been used very successfully in computability theory to prove a host of problems undecidable, beginning with the Halting Problem. It has also been used successfully in complexity theory to prove super-exponential lower bounds for very hard decidable problems. For example, Presburger arithmetic, the first-order theory of integers under addition, is a decidable theory for which Fischer and Rabin [14] proved that any Turing machine deciding the theory must use at least  $2^{2^cn}$  steps in the worst case, for some  $c > 0$ . In general, lower bounds using diagonalization and reduction relativize; that is,

they continue to apply in a setting in which both the problem instance and the solving Turing machine can make membership queries to an arbitrary oracle set  $A$ . However, in [3] it was shown that there is an oracle set  $A$  relative to which  $\mathbf{P} = \mathbf{NP}$ , suggesting that diagonalization with reduction cannot be used to separate these two classes. (There are nonrelativizing results in complexity theory, as will be mentioned below.) It is interesting to note that relative to a *generic* oracle,  $\mathbf{P} \neq \mathbf{NP}$  [5, 11].

A *Boolean circuit* is a finite acyclic graph in which each non-input node, or *gate*, is labelled with a Boolean connective; typically from {AND, OR, NOT}. The input nodes are labeled with variables  $x_1, \dots, x_n$ , and for each assignment of 0 or 1 to each variable, the circuit computes a bit value at each gate, including the output gate, in the obvious way. It is not hard to see that if  $L$  is a language over  $\{0, 1\}$  that is in  $\mathbf{P}$ , then there is a polynomial-size family of Boolean circuits  $\langle B_n \rangle$  such that  $B_n$  has  $n$  inputs, and for each bit string  $w$  of length  $n$ , when  $w$  is applied to the  $n$  input nodes of  $B_n$ , then the output bit of  $B_n$  is 1 iff  $w \in L$ . In this case we say that  $\langle B_n \rangle$  *computes*  $L$ .

Thus to prove  $\mathbf{P} \neq \mathbf{NP}$  it suffices to prove a super-polynomial lower bound on the size of any family of Boolean circuits solving some specific  $\mathbf{NP}$ -complete problem, such as 3-SAT. Back in 1949 Shannon [31] proved that for almost all Boolean functions  $f : \{0, 1\}^n \rightarrow \{0, 1\}$ , any Boolean circuit computing  $f$  requires at least  $2^n/n$  gates. Unfortunately, his counting argument gives no clue as to how to prove lower bounds for problems in  $\mathbf{NP}$ . Exponential lower bounds for  $\mathbf{NP}$  problems have been proved for restricted circuit models, including monotone circuits [26], [2] and bounded depth circuits with unbounded fan-in gates [17], [35] (see [6]). However, all attempts to find even super-linear lower bounds for unrestricted Boolean circuits for “explicitly given” Boolean functions have met with total failure; the best such lower bound proved so far is about  $4n$ . Razborov and Rudich [27] explain this failure by pointing out that all methods used so far can be classified as “natural proofs”, and natural proofs for general circuit lower bounds are doomed to failure, assuming a certain complexity-theoretic conjecture asserting that strong pseudo-random number generators exist. Since such generators have been constructed assuming the hardness of integer factorization, we can infer the surprising result that a natural proof for a general circuit lower bound would give rise to a more efficient factoring algorithm than is currently known.

The failure of complexity theory to prove interesting lower bounds on a general model of computation is much more pervasive than the failure to prove  $\mathbf{P} \neq \mathbf{NP}$ . It is consistent with present knowledge that not only could Satisfiability have a polynomial-time algorithm, it could have a linear time algorithm on a multitape Turing machine. The same applies for all 21 problems mentioned in Karp’s original paper [21]. There are complexity class separations that we know exist but cannot prove. For example, consider the sequence of complexity class inclusions

$$\text{LOGSPACE} \subseteq \mathbf{P} \subseteq \mathbf{NP} \subseteq \mathbf{PSPACE}.$$

A simple diagonal argument shows that the first is a proper subset of the last, but we cannot prove any particular adjacent inclusion is proper.

As another example, let **LINEAR-SIZE** be the class of languages over  $\{0, 1\}$  that can be computed by a family  $\langle B_n \rangle$  of Boolean circuits of size  $O(n)$ . It is not known whether either  $\mathbf{P}$  or  $\mathbf{NP}$  is a subset of **LINEAR-SIZE**, although Kannan [20] proved that there are languages in the polynomial hierarchy (a generalization of

**NP**) not in **LINEAR-SIZE**. Since if  $\mathbf{P} = \mathbf{NP}$ , the polynomial hierarchy collapses to  $\mathbf{P}$ , we have

**Proposition 2.** *If  $\mathbf{P} \subseteq \mathbf{LINEAR-SIZE}$ , then  $\mathbf{P} \neq \mathbf{NP}$ .*

This proposition could be interpreted as a method of proving  $\mathbf{P} \neq \mathbf{NP}$ , but a more usual belief is that the hypothesis is false.

A fundamental question in complexity theory is whether a source of random bits can be used to speed up substantially the recognition of some languages, provided one is willing to accept a small probability of error. The class **BPP** consists of all languages  $L$  that can be recognized by a randomized polynomial-time algorithm, with at most an exponentially small error probability on every input. Of course  $\mathbf{P} \subseteq \mathbf{BPP}$ , but it is not known whether the inclusion is proper. The set of prime numbers is in **BPP** [36], but it is not known to be in  $\mathbf{P}$ . A reason for thinking that  $\mathbf{BPP} = \mathbf{P}$  is that randomized algorithms are often successfully executed using a deterministic pseudo-random number generator as a source of “random” bits.

There is an indirect but intriguing connection between the two questions  $\mathbf{P} = \mathbf{BPP}$  and  $\mathbf{P} = \mathbf{NP}$ . Let **E** be the class of languages recognizable in exponential time; that is the class of languages  $L$  such that  $L = L(M)$  for some Turing machine  $M$  with  $T_M(n) = O(2^{cn})$  for some  $c > 0$ . Let **A** be the assertion that some language in **E** requires exponential circuit complexity. That is,

**Assertion A.** *There is  $L \in \mathbf{E}$  and  $\epsilon > 0$  such that, for every circuit family  $\langle B_n \rangle$  computing  $L$  and for all sufficiently large  $n$ ,  $B_n$  has at least  $2^{\epsilon n}$  gates.*

**Proposition 3.** *If **A** then  $\mathbf{BPP} = \mathbf{P}$ . If not **A** then  $\mathbf{P} \neq \mathbf{NP}$ .*

The first implication is a lovely theorem by Impagliazzo and Wigderson [19] and the second is an intriguing observation by V. Kabanets that strengthens Proposition 2. In fact, Kabanets concludes  $\mathbf{P} \neq \mathbf{NP}$  from a weaker assumption than **not A**; namely that every language in **E** can be computed by a Boolean circuit family  $\langle B_n \rangle$  such that for at least one  $n$ ,  $B_n$  has fewer gates than the maximum needed to compute any Boolean function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$ . But there is no consensus on whether this hypothesis is true.

We should point out that Proposition 3 relativizes, and, in particular, relative to any **PSPACE**-complete oracle Assertion A holds and  $\mathbf{BPP} = \mathbf{P} = \mathbf{NP}$ . Thus a nonrelativizing construction will be needed if one is to prove  $\mathbf{P} \neq \mathbf{NP}$  by giving small circuits for languages in **E**. Nonrelativizing constructions have been used successfully before, for example in showing **IP** (Interactive Polynomial time) contains all of **PSPACE** [30]. In this and other such constructions a key technique is to represent Boolean functions by multivariate polynomials over finite fields.

#### APPENDIX: DEFINITION OF TURING MACHINE

A Turing machine  $M$  consists of a finite state control (i.e., a finite program) attached to a read/write head moving on an infinite tape. The tape is divided into squares, each capable of storing one symbol from a finite alphabet  $\Gamma$  that includes the blank symbol  $b$ . Each machine  $M$  has a specified input alphabet  $\Sigma$ , which is a subset of  $\Gamma$ , not including the blank symbol  $b$ . At each step in a computation,  $M$  is in some state  $q$  in a specified finite set  $Q$  of possible states. Initially, a finite input string over  $\Sigma$  is written on adjacent squares of the tape, all other squares are blank (contain  $b$ ), the head scans the left-most symbol of the input string, and  $M$  is in

the initial state  $q_0$ . At each step  $M$  is in some state  $q$  and the head is scanning a tape square containing some tape symbol  $s$ , and the action performed depends on the pair  $(q, s)$  and is specified by the machine's transition function (or program)  $\delta$ . The action consists of printing a symbol on the scanned square, moving the head left or right one square, and assuming a new state.

Formally, a Turing machine  $M$  is a tuple  $\langle \Sigma, \Gamma, Q, \delta \rangle$ , where  $\Sigma, \Gamma, Q$  are finite nonempty sets with  $\Sigma \subseteq \Gamma$  and  $b \in \Gamma - \Sigma$ . The state set  $Q$  contains three special states  $q_0, q_{\text{accept}}, q_{\text{reject}}$ . The *transition function*  $\delta$  satisfies

$$\delta : (Q - \{q_{\text{accept}}, q_{\text{reject}}\}) \times \Gamma \rightarrow Q \times \Gamma \times \{-1, 1\}.$$

If  $\delta(q, s) = (q', s', h)$ , the interpretation is that, if  $M$  is in state  $q$  scanning the symbol  $s$ , then  $q'$  is the new state,  $s'$  is the symbol printed, and the tape head moves left or right one square depending on whether  $h$  is  $-1$  or  $1$ .

We assume that the sets  $Q$  and  $\Gamma$  are disjoint.

A *configuration* of  $M$  is a string  $xqy$  with  $x, y \in \Gamma^*$ ,  $y$  not the empty string, and  $q \in Q$ .

The interpretation of the configuration  $xqy$  is that  $M$  is in state  $q$  with  $xy$  on its tape, with its head scanning the left-most symbol of  $y$ .

If  $C$  and  $C'$  are configurations, then  $C \xrightarrow{M} C'$  if  $C = xqsy$  and  $\delta(q, s) = (q', s', h)$  and one of the following holds:

$$C' = xs'q'y \text{ and } h = 1 \text{ and } y \text{ is nonempty.}$$

$$C' = xs'q'b \text{ and } h = 1 \text{ and } y \text{ is empty.}$$

$$C' = x'q'as'y \text{ and } h = -1 \text{ and } x = x'a \text{ for some } a \in \Gamma.$$

$$C' = q'bs'y \text{ and } h = -1 \text{ and } x \text{ is empty.}$$

A configuration  $xqy$  is *halting* if  $q \in \{q_{\text{accept}}, q_{\text{reject}}\}$ . Note that for each non-halting configuration  $C$  there is a unique configuration  $C'$  such that  $C \xrightarrow{M} C'$ .

The *computation* of  $M$  on input  $w \in \Sigma^*$  is the unique sequence  $C_0, C_1, \dots$  of configurations such that  $C_0 = q_0w$  (or  $C_0 = q_0b$  if  $w$  is empty) and  $C_i \xrightarrow{M} C_{i+1}$  for each  $i$  with  $C_{i+1}$  in the computation, and either the sequence is infinite or it ends in a halting configuration. If the computation is finite, then the number of steps is one less than the number of configurations; otherwise the number of steps is infinite. We say that  $M$  *accepts*  $w$  iff the computation is finite and the final configuration contains the state  $q_{\text{accept}}$ .

#### ACKNOWLEDGMENTS

My thanks to Avi Wigderson and Hugh Woodin for many helpful suggestions for improving an earlier version of this paper.

#### REFERENCES

- [1] M. Agrawal, N. Kayal, and N. Saxena, *Primes is in P*, Ann. Math. **160** (2004), 781–793.
- [2] N. Alon and R.B. Boppana, *The monotone circuit complexity of boolean functions*, Combinatorica **7** (1987), 1–22.
- [3] T. Baker, J. Gill, and R. Solovay, *Relativizations of the P =? NP question*, SICOMP: SIAM Journal on Computing, 1975.
- [4] L. Blum, F. Cucker, M. Shub, and S. Smale, *Complexity and Real Computation*, Springer-Verlag, New York, 1998.
- [5] M. Blum and R. Impagliazzo, *Generic oracles and oracle classes*, in Proceedings of the 28th Annual Symposium on Foundations of Computer Science, A.K. Chandra, ed., IEEE Computer Society Press, Los Angeles, 1987, 118–126.

- [6] R.B. Boppana and M. Sipser, *The complexity of finite functions*, Handbook of Theoretical Computer Science, Volume A: Algorithms and Complexity, J. Van Leeuwen, ed., Elsevier and The MIT Press, Cambridge, MA, 1990, 759–804.
- [7] T. Cormen, C. Leiserson, R. Rivest, and C. Stein, *Introduction to Algorithms*, 2nd edition, McGraw Hill, New York, 2001.
- [8] A. Cobham, *The intrinsic computational difficulty of functions*, in Proceedings of the 1964 International Congress for Logic, Methodology, and Philosophy of Science, Y. Bar-Hillel, ed., Elsevier/North-Holland, Amsterdam, 1964, 24–30.
- [9] S. Cook, *The complexity of theorem-proving procedures*, in Conference Record of Third Annual ACM Symposium on Theory of Computing, ACM, New York, 1971, 151–158.
- [10] S. Cook, *Computational complexity of higher type functions*, in Proceedings of the International Congress of Mathematicians, Kyoto, Japan, Springer-Verlag, Berlin, 1991, 55–69.
- [11] S. Cook, R. Impagliazzo, and T. Yamakami, *A tight relationship between generic oracles and type-2 complexity theory*, Information and Computation **137** (1997), 159–170.
- [12] S. Cook and R. Reckhow, *The relative efficiency of propositional proof systems*, J. Symbolic Logic **44** (1979), 36–50.
- [13] J. Edmonds, *Minimum partition of a matroid into independent subsets*, J. Res. Nat. Bur. Standards Sect. B **69** (1965), 67–72.
- [14] M.J. Fischer and M.O. Rabin, *Super-exponential complexity of Presburger arithmetic*, in Complexity of Computation **7**, AMS, Providence, RI, 1974, 27–41.
- [15] M.R. Garey and D.S. Johnson, *Computers and Intractability, a Guide to the Theory of NP-Completeness*, W.H. Freeman and Co., San Francisco, 1979.
- [16] O. Goldreich, *The Foundations of Cryptography — Volume 1*, Cambridge University Press, Cambridge, UK, 2000.
- [17] J. Hastad, *Almost optimal lower bounds for small depth circuits*, in Randomness and Computation, Advances in Computing Research **5**, JAI Press Inc., Greenwich, CT, 1989, 143–170.
- [18] R. Impagliazzo, *A personal view of average-case complexity*, in 10th IEEE Annual Conference on Structure in Complexity Theory, IEEE Computer Society Press, Washington, DC, 1995, 134–147.
- [19] R. Impagliazzo and A. Wigderson, *P = BPP if E requires exponential circuits: Derandomizing the XOR lemma*, in ACM Symposium on Theory of Computing (STOC), ACM, New York, 1997, 220–229.
- [20] R. Kannan, *Circuit-size lower bounds and non-reducibility to sparse sets*, Information and Control **55** (1982), 40–56..
- [21] R.M. Karp, *Reducibility among combinatorial problems*, in Complexity of Computer Computations, R. E. Miller and J. W. Thatcher, eds., Plenum Press, New York, 1972, 85–103.
- [22] J. Krajicek, *Bounded Arithmetic, Propositional Logic, and Complexity Theory*, Cambridge University Press, Cambridge, 1995.
- [23] L. Levin, *Universal search problems* (in Russian), Problemy Peredachi Informatsii **9** (1973), 265–266. English translation in B. A. Trakhtenbrot, *A survey of Russian approaches to Perebor (brute-force search) algorithms*, Annals of the History of Computing **6** (1984), 384–400.
- [24] L. Levin, *Average case complete problems*, SIAM J. Computing **15** (1986), 285–286.
- [25] C. Papadimitriou, *Computational Complexity*, Addison-Wesley, Reading, MA, 1994.
- [26] A.A. Razborov, *Lower bounds on the monotone complexity of some boolean functions*, Soviet Math. Dokl. **31** (1985), 354–357.
- [27] A.A. Razborov and S. Rudich, *Natural proofs*, Journal of Computer and System Sciences **55** (1997), 24–35.
- [28] R.L. Rivest, A. Shamir, and L. Adleman, *A method for obtaining digital signatures and public-key cryptosystems*, Comm. ACM **21** (1978), 120–126.
- [29] N. Robertson and P.D. Seymour, *Graph minors i-xiii*, Journal of Combinatorial Theory B, 1983–1995.
- [30] A. Shamir, IP = PSPACE, J.A.C.M. **39** (1992), 869–977.
- [31] C. Shannon, *The synthesis of two-terminal switching circuits*, Bell System Technical Journal **28** (1949), 59–98.
- [32] P. Shor, *Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer*, SIAM J. Computing **26** (1997), 1484–1509.
- [33] M. Sipser, *Introduction to the Theory of Computation*, PWS Publ., Boston, 1997.

- [34] S. Smale, *Mathematical problems for the next century*, Math. Intelligencer **20**, no. 2, 1998, 7–15.
- [35] R. Smolensky, *Algebraic methods in the theory of lower bounds for boolean circuit complexity*, in ACM Symposium on Theory of Computing (STOC) **19**, ACM, New York, 1987, 77–82.
- [36] R. Solovay and V. Strassen, *A fast Monte-Carlo test for primality*, SIAM Journal on Computing **6** (1977), 84–85.
- [37] A. Turing, *On computable numbers with an application to the entscheidungsproblem*, Proc. London Math. Soc. **42** (1936), 230–265.
- [38] J. von Neumann, *A certain zero-sum two-person game equivalent to the optimal assignment problem*, in Contributions to the Theory of Games II, H.W. Kahn and A.W. Tucker, eds. Princeton Univ. Press, Princeton, NJ, 1953, 5–12.