# Computational Complexity:

# A Conceptual Perspective

Oded Goldreich

Department of Computer Science and Applied Mathematics
Weizmann Institute of Science, Rehovot, Israel.

October 12, 2006

# Chapter 3

# Variations on P and NP

In this chapter we consider variations on the complexity classes P and NP. We refer specifically to the non-uniform version of P, and to the Polynomial-time Hierarchy (which extends NP). These variations are motivated by relatively technical considerations; still, the resulting classes are referred to quite frequently in the literature.

**Summary:** Non-uniform polynomial-time (P/poly) captures efficient computations that are carried out by devices that can each only handle inputs of a specific length. The basic formalism ignore the complexity of constructing such devices (i.e., a uniformity condition). A finer formalism that allows to quantify the amount of non-uniformity refers to so called "machines that take advice."

The Polynomial-time Hierarchy (PH) generalizes NP by considering statements expressed by quantified Boolean formulae with a fixed number of alternations of existential and universal quantifiers. It is widely believed that each quantifier alternation adds expressive power to the class of such formulae.

The two different classes are related by showing that if NP is contained in P/poly then the Polynomial-time Hierarchy collapses to its second level. This result is commonly interpreted as supporting the common belief that non-uniformity is irrelevant to the P-vs-NP Question; that is, although P/poly extends beyond the class P, is is believed that P/poly does not contain NP.

Except for the latter result, which is presented in Section 3.2.3, the treatments of P/poly (in Section 3.1) and of the Polynomial-time Hierarchy (in Section 3.2) are independently of one another.

# 3.1   Non-uniform polynomial-time (P/poly)

In this section we consider two formulations of the notion of non-uniform polynomial-time, based on the two models of non-uniform computing devices that were presented in Section 1.2.4. That is, we specialize the treatment of non-uniform computing devices, provided in Section 1.2.4, to the case of polynomially bounded complexities. It turns out that both (polynomially bounded) formulations allow for solving the same class of computational problems, which is a strict superset of the class of problems solvable by polynomial-time algorithms.

The two models of non-uniform computing devices are Boolean circuits and "machines that take advice" (cf. §1.2.4.1 and §1.2.4.2, respectively). We will focus on the restriction of both models to the case of polynomial complexities, considering (non-uniform) polynomial-size circuits and polynomial-time algorithms that take (non-uniform) advice of polynomially bounded length.

The main motivation for considering non-uniform polynomial-size circuits is that their computational limitations imply analogous limitations on polynomial-time algorithms. The hope is that, as is often the case in mathematics and Science, disposing of an auxiliary condition (i.e., uniformity) that seems secondary[1] and is not well-understood may turn out fruitful. In particular, the (non-uniform) circuit model facilitates a low-level analysis of the evolution of a computation, and allow for the application of combinatorial techniques. The benefit of this approach has been demonstrated in the study of restricted classes of circuits (see Sections B.2.2 and B.2.3).

Polynomial-time algorithms that take polynomially bounded advice are useful in modeling auxiliary information available to possible efficient strategies that are of interest to us. Indeed, the typical cases are the modeling of adversaries in the context of cryptography and the modeling of arbitrary randomized algorithms in the context of derandomization. Furthermore, the model of polynomial-time algorithms that take advice allows for a quantitative study of the amount of non-uniformity, ranging from zero to polynomial.

## 3.1.1   Boolean Circuits

We refer the reader to §1.2.4.1 for a definition of (families of) Boolean circuits and the functions computed by them. For concreteness and simplicity, we assume throughout this section that all circuits has bounded fan-in. We highlight the following result stated in §1.2.4.1:

**Theorem 3.1** (circuit evaluation): *There exists a polynomial-time algorithm that, given a circuit $C : \{0,1\}^n \to \{0,1\}^m$ and an $n$-bit long string $x$, returns $C(x)$.*

Recall that the algorithm works by performing the "value-determination" process that underlies the definition of the computation of the circuit on a given input.

---

[1]The common belief is that the issue of non-uniformity is irrelevant to the P-vs-NP Question; that is, that resolving the latter question by proving that $\mathcal{P} \neq \mathcal{NP}$ is not easier than proving that NP does not have polynomial-size circuits. For further discussion see Appendix B.2 and Section 3.2.3.

This process assigns values to each of the circuit vertices based on the values of its children (or the values of the corresponding bit of the input, in the case of an input-terminal vertex).

**Circuit size as a complexity measure.** We recall the definitions of circuit complexity presented in to §1.2.4.1: The size of a circuit is defined as the number of edges, and the length of its description is almost linear in the latter; that is, a circuit of size $s$ is commonly described by the list of its edges and the labels of its vertices, which means that its description length is $O(s \log s)$. We are interested in families of circuits that solve computational problems, and thus we say that the circuit family $(C_n)_{n \in \mathbb{N}}$ computes the function $f : \{0,1\}^* \to \{0,1\}^*$ if for every $x \in \{0,1\}^*$ it holds that $C_{|x|}(x) = f(x)$. The size complexity of this family is the function $s : \mathbb{N} \to \mathbb{N}$ such that $s(n)$ is the size of $C_n$. The circuit complexity of a function $f$, denoted $s_f$, is the size complexity of the smallest family of circuits that computes $f$. An equivalent alternative follows.

**Definition 3.2** (circuit complexity): *The circuit complexity of $f : \{0,1\}^* \to \{0,1\}^*$ is the function $s_f : \mathbb{N} \to \mathbb{N}$ such that $s_f(n)$ is the size of the smallest circuit that computes the restriction of $f$ to $n$-bit strings.*

We stress that non-uniformity is implicit in this definition, because no conditions are made regarding the relation between the various circuits used to compute the function on different input lengths.

An interesting feature of Definition 3.2 is that, unlike in the case of uniform model of computation, the circuit complexity is the actual complexity of the function rather than an upper-bound on its complexity (cf. §1.2.3.4 and Section 4.2.1). This is a consequence of the fact that the circuit model has no "free parameters" (e.g., the finite algorithm in use).[2]

We will be interested in the class of problems that are solvable by families of polynomial-size circuits. That is, a problem is solvable by polynomial-size circuits if it can be solved by a function $f$ that has polynomial circuit complexity (i.e., there exists a polynomial $p$ such that $s_f(n) \leq p(n)$, for every $n \in \mathbb{N}$).

**A detour: uniform families.** A family of *polynomial-size* circuits $(C_n)_n$ is called uniform if given $n$ one can construct the circuit $C_n$ in poly($n$)-time. More generally:

**Definition 3.3** (uniformity): *A family of circuits $(C_n)_n$ is called uniform if there exists an algorithm $A$ that on input $n$ outputs $C_n$ within a number of steps that is polynomial in the size of $C_n$.*

We note that stronger notions of uniformity have been considered. For example, one may require the existence of a polynomial-time algorithm that on input $n$ and $v$, returns the label of vertex $v$ as well as the list of its children (or an indication that $v$ is not a vertex in $C_n$). For further discussion see Section 5.2.3.

---

[2]**Advanced comment:** Note that such "free parameters" underly linear speedup results such as Exercise 4.7, which in turn prevent the specification of the exact complexities of functions.

**Proposition 3.4** *If a problem is solvable by a uniform family of polynomial-size circuits then it is solvable by a polynomial-time algorithm.*

As was hinted in §1.2.4.1, the converse holds as well. The latter fact follows easily from the proof of Theorem 2.20 (see also the proof of Theorem 3.6).

**Proof:** On input $x$, the algorithm operates in two stages. In the first stage, it invokes the algorithm guaranteed by the uniformity condition, on input $n \overset{\text{def}}{=} |x|$, and obtains the circuit $C_n$. Next, it invokes the circuit evaluation algorithm (asserted in Theorem 3.1) on input $C_n$ and $x$, and obtains $C_n(x)$. Since the size and the description length of $C_n$ are polynomial in $n$, it follows that each stage of our algorithm runs in polynomial time (i.e., polynomial in $n = |x|$). Thus, the algorithm emulates the computation of $C_{|x|}(x)$, and does so in time polynomial in the length of its own input (i.e., $x$).  ∎

## 3.1.2   Machines that take advice

General (non-uniform) families of polynomial-size circuits and uniform families of polynomial-size circuits are two extremes with respect to the "amounts of non-uniformity" in the computing device. Intuitively, in the former, non-uniformity is only bounded by the size of the device, whereas in the latter the amounts of non-uniformity is zero. Here we consider a model that allows to decouple the size of the computing device from the amount of non-uniformity, which may indeed range from zero to the device's size. Specifically, we consider algorithms that "take a non-uniform advice" that depends only on the input length. The amount of non-uniformity will be defined to equal the length of the corresponding advice (as a function of the input length). Thus, we specialize Definition 1.12 to the case of polynomial-time algorithms.

**Definition 3.5** (non-uniform polynomial-time and $\mathcal{P}/\text{poly}$): *We say that a function $f$ is* computed in polynomial-time with advice of length $\ell : \mathbb{N} \to \mathbb{N}$ *if these exists a polynomial-time algorithm $A$ and an infinite* advice sequence $(a_n)_{n \in \mathbb{N}}$ *such that*

  1. *For every $x \in \{0, 1\}^*$, it holds that $A(a_{|x|}, x) = f(x)$.*
  2. *For every $n \in \mathbb{N}$, it holds that $|a_n| = \ell(n)$.*

*We say that a computational problem can be solved in polynomial-time with advice of length $\ell$ if a function solving this problem can be computed within these resources. We denote by $\mathcal{P}/\ell$ the class of decision problems that can be solved in polynomial-time with advice of length $\ell$, and by $\mathcal{P}/\text{poly}$ the union of $\mathcal{P}/p$ taken over all polynomials $p$.*

Clearly, $\mathcal{P}/0 = \mathcal{P}$. But allowing some (non-empty) advice increases the power of the class (see Theorem 3.7). and allowing advice of length comparable to the time complexity yields a formulation equivalent to circuit complexity (see Theorem 3.6). We highlight the greater flexibility available by the formalism of machines that take advice, which allows for separate specification of time complexity and advice length. (Indeed, this comes at the expense of a more cumbersome formulation, when we wish to focus on the case that both measures are equal.)

**Relation to families of polynomial-size circuits.**   As hinted before, the class of problems solvable by polynomial-time algorithms with polynomially bounded advice equals the class of problems solvable by families of polynomial-size circuits. For concreteness, we state this fact for decision problems.

**Theorem 3.6** *A decision problem is in $\mathcal{P}/\text{poly}$ if and only if it can be solved by a family of polynomial-size circuits.*

More generally, for any function $t$, the following proof establishes that equivalence of the power of machines having time complexity $t$ and taking advice of length $t$ versus families of circuits of size polynomially related to $t$.

**Proof Sketch:** Suppose that a problem can be solved by a polynomial-time algorithm $A$ using the polynomially bounded advice sequence $(a_n)_{n \in \mathbb{N}}$. We obtain a family of polynomial-size circuits that solves the same problem by adapting the proof of Theorem 2.20. Specifically, we observe that the computation of $A(a_{|x|}, x)$ can be emulated by a circuit of poly($|x|$)-size, *which incorporates $a_{|x|}$ and is given $x$ as input.* That is, we construct a circuit $C_n$ such that $C_n(x) = A(a_n, x)$ holds for every $x \in \{0,1\}^n$ (analogously to the way $C_x$ was constructed in the proof of Theorem 2.20, where it holds that $C_x(y) = M_R(x, y)$ for every $y$ of adequate length).

On the other hand, given a family of polynomial-size circuits, we obtain a polynomial-time algorithm for emulating this family *using advice that provide the description of the relevant circuits.* Specifically, we use the evaluation algorithm asserted in Theorem 3.1, while using the circuit's description as advice. That is, we use the fact that a circuit of size $s$ can be described by a string of length $O(s \log s)$, where the log factor is due to the fact that a graph with $v$ vertices and $e$ edges can be described by a string of length $2e \log_2 v$.   $\square$

**Another perspective.**   A set $S$ is called sparse if there exists a polynomial $p$ such that for every $n$ it holds that $|S \cap \{0,1\}^n| \leq p(n)$. We note that $\mathcal{P}/\text{poly}$ equals the class of sets that are Cook-reducible to a sparse set (see Exercise 3.2). Thus, SAT is Cook-reducible to a sparse set if and only if $\mathcal{NP} \subset \mathcal{P}/\text{poly}$. In contrast, SAT is Karp-reducible to a sparse set if and only if $\mathcal{NP} = \mathcal{P}$ (see Exercise 3.12).

**The power of $\mathcal{P}/\text{poly}$.**   In continuation to Theorem 1.13 (which focuses on advice and ignores the time complexity of the machine that takes this advice), we prove the following (stronger) result.

**Theorem 3.7** (the power of advice, revisited): *The class $\mathcal{P}/1 \subseteq \mathcal{P}/\text{poly}$ contains $\mathcal{P}$ as well as some undecidable problems.*

Actually, $\mathcal{P}/1 \subset \mathcal{P}/\text{poly}$. Furthermore, by using a counting argument, one can show that for any two polynomially bounded functions $\ell_1, \ell_2 : \mathbb{N} \to \mathbb{N}$ such that $\ell_2 - \ell_1 > 0$ is unbounded, it holds that $\mathcal{P}/\ell_1$ is strictly contained in $\mathcal{P}/\ell_2$; see Exercise 3.3.

**Proof:** Clearly, $\mathcal{P} = \mathcal{P}/0 \subseteq \mathcal{P}/1 \subseteq \mathcal{P}/\text{poly}$. To prove that $\mathcal{P}/1$ contains some undecidable problems, we review the proof of Theorem 1.13. The latter proof established the existence of uncomputable Boolean function that only depend on their input length. That is, there exists an undecidable set $S \subset \{0,1\}^*$ such that for every pair of equal length strings $(x,y)$ it holds that $x \in S$ if and only if $y \in S$. In other words, for every $x \in \{0,1\}^*$ it holds that $x \in S$ if and only if $1^{|x|} \in S$. But such a set is easily decidable in polynomial-time by a machine that takes one bit of advice; that is, consider the algorithm $A$ and the advice sequence $(a_n)_{n \in \mathbb{N}}$ such that $a_n = 1$ if and only if $1^n \in S$ and $A(a,x) = a$ (for $a \in \{0,1\}$ and $x \in \{0,1\}^*$). Note that indeed $A(a_{|x|}, x) = 1$ if and only if $x \in S$.   ■

## 3.2   The Polynomial-time Hierarchy (PH)

We start with an informal motivating discussion, which will be made formal in Section 3.2.1.

Sets in $\mathcal{NP}$ can be viewed as sets of valid assertions that can be expressed as quantified Boolean formulae using only existential quantifiers. That is, a set $S$ is in $\mathcal{NP}$ if there is a Karp-reduction of $S$ to the problem of deciding whether or not an existentially quantified Boolean formula is valid (i.e., an instance $x$ is mapped by this reduction to a formula of the form $\exists y_1 \cdots \exists y_{m(x)} \phi_x(y_1, ..., y_{m(x)})$).

The conjectured intractability of $\mathcal{NP}$ seems due to the long sequence of existential quantifiers. Of course, if somebody else (i.e., a "prover") were to provide us with an adequate assignment (to the $y_i$'s) whenever such an assignment exists then we would be in good shape. That is, we can efficiently verify proofs of validity of existentially quantified Boolean formulae.

But what if we want to verify the validity of a universally quantified Boolean formulae (i.e., formulae of the form $\forall y_1 \cdots \forall y_m \phi(y_1, ..., y_m)$). Here we seem to need the help of a totally different entity: we need a "refuter" that is guaranteed to provide us with a refutation whenever such exist, and we need to believe that if we were not presented with such a refutation then it is the case that no refutation exists (and hence the universally quantified formulae is valid). Indeed, this new setting (of a "refutation system") is fundamentally different from the setting of a proof system: In a proof system we are only convinced by proofs (to assertions) that we have verified by ourselves, whereas in the "refutation system" we trust the "refuter" to provide evidence against false assertions.[3] Furthermore, there seems to be no way of converting one setting (e.g., the proof system) into another (resp., the refutation system).

Taking an additional step, we may consider a more complicated system in which we use two agents: a "supporter" that tries to provide evidence in favor of an assertion and an "objector" that tries to refute it. These two agents conduct a debate (or an argument) in our presence, exchanging messages with the goal of

---

[3]More formally, in proof systems the soundness condition relies only on the actions of the verifier, whereas completeness also relies on the prover using an adequate strategy. In contrast, in "refutation system" the soundness condition relies on the proper actions of the refuter, whereas completeness does not depend on the refuter's actions.

making us (the referee) rule their way. The assertions that can be proven in this system take the form of general quantified formulae with alternating sequences of quantifiers, where the number of alternations equals the number of rounds of interaction in the said system. We stress that the exact length of each sequence of quantifiers of the same type does not matter, what matters is the number of alternations, denoted $k$.

The aforementioned system of alternations can be viewed as a two-party game, and we may ask ourselves which of the two parties has a $k$-move winning strategy. In general, we may consider any (0-1 zero-sum) two-party game, in which the game's position can be efficiently updated (by any given move) and efficiently evaluated. For such a fixed game, given an initial position, we may ask whether the first party has a ($k$-move) winning strategy. It seems that answering this type of question for some fixed $k$ does not necessarily allow answering it for $k + 1$. We now turn to formalize the foregoing discussion.

## 3.2.1 Alternation of quantifiers

In the following definition, the aforementioned propositional formula $\phi_x$ is replaced by the input $x$ itself. (Correspondingly, the combination of the Karp-reduction and a formula evaluation algorithm are replaced by the verification algorithm $V$ (see Exercise 3.7).) This is done in order to make the comparison to the definition of $\mathcal{NP}$ more transparent (as well as to fit the standard presentations). We also replace a sequence of Boolean quantifiers of the same type by a single corresponding quantifier that quantifies over all strings of the corresponding length.

**Definition 3.8** (the class $\Sigma_k$): *For a natural number $k$, a decision problem $S \subseteq \{0,1\}^*$ is in $\Sigma_k$ if there exists a polynomial $p$ and a polynomial time algorithm $V$ such that $x \in S$ if and only if*

$$\exists y_1 \in \{0,1\}^{p(|x|)} \forall y_2 \in \{0,1\}^{p(|x|)} \exists y_3 \in \{0,1\}^{p(|x|)} \cdots Q_k y_k \in \{0,1\}^{p(|x|)}$$
$$\text{s.t. } V(x, y_1, ..., y_k) = 1$$

*where $Q_k$ is an existential quantifier if $k$ is odd and is a universal quantifier otherwise.*

Note that $\Sigma_1 = \mathcal{NP}$ and $\Sigma_0 = \mathcal{P}$. The Polynomial-time Hierarchy, denoted $\mathcal{PH}$, is the union of all the aforementioned classes (i.e., $\mathcal{PH} = \cup_k \Sigma_k$), and $\Sigma_k$ is often referred to as the $k^{\text{th}}$ level of $\mathcal{PH}$. The levels of the Polynomial-time Hierarchy can also be defined inductively, by defining $\Sigma_{k+1}$ based on $\Pi_k \overset{\text{def}}{=} \text{co}\Sigma_k$, where $\text{co}\Sigma_k \overset{\text{def}}{=} \{\{0,1\}^* \setminus S : S \in \Sigma_k\}$ (cf. Eq. (2.4)).

**Proposition 3.9** *For every $k \geq 0$, a set $S$ is in $\Sigma_{k+1}$ if and only if there exists a polynomial $p$ and a set $S' \in \Pi_k$ such that $S = \{x : \exists y \in \{0,1\}^{p(|x|)} \text{ s.t. } (x,y) \in S'\}$.*

**Proof:** Suppose that $S$ is in $\Sigma_{k+1}$ and let $p$ and $V$ be as in Definition 3.8. Then define $S'$ as the set of pairs $(x, y)$ such that $|y| = p(|x|)$ and

$$\forall z_1 \in \{0,1\}^{p(|x|)} \exists z_2 \in \{0,1\}^{p(|x|)} \cdots Q_k z_k \in \{0,1\}^{p(|x|)} \text{ s.t. } V(x, y, z_1, ..., z_k) = 1 .$$

Note that $x \in S$ if and only if there exists $y \in \{0,1\}^{p(|x|)}$ such that $(x,y) \in S'$, and that $S' \in \Pi_k$ (see Exercise 3.6).

On the other hand, suppose that for some polynomial $p$ and a set $S' \in \Pi_k$ it holds that $S = \{x : \exists y \in \{0,1\}^{p(|x|)} \text{ s.t. } (x,y) \in S'\}$. Then, for some $p'$ and $V'$, it holds that $(x,y) \in S'$ if and only if $|y| = p(|x|)$ and

$$\forall z_1 \in \{0,1\}^{p'(|x|)} \exists z_2 \in \{0,1\}^{p'(|x|)} \cdots Q_k z_k \in \{0,1\}^{p'(|x|)} \text{ s.t. } V'(x,y,z_1,...,z_k) \neq 1$$

(see Exercise 3.6 again). By suitable encoding (of $y$ and the $z_i$'s as strings of length $\max(p(|x|), p'(|x|))$) and a trivial modification of $V'$, we conclude that $S \in \Sigma_{k+1}$. ∎

**Determining the winner in $k$-move games.**   Definition 3.8 can be interpreted as capturing the complexity of determining the winner in certain *efficient two-party game*. Specifically, we refer to two-party games that satisfy the following three conditions:

1. The parties alternate in taking moves that effect the game's (global) position, where each move has a description length that is bounded by a polynomial in the length of the current position.

2. The current position can be updated in polynomial-time based on the previous position and the current party's move.[4]

3. The winner in each position can be determined in polynomial-time.

A set $S \in \Sigma_k$ can be viewed as the set of initial positions (in a suitable game) for which the first party has a $k$-move winning strategy. Specifically, $x \in S$ if starting at the initial position $x$, there exists move $y_1$ for the first party, such that for every response move $y_2$ of the second party, there exists move $y_3$ for the first party, etc, such that after $k$ moves the parties reach a position in which the first party wins, where the final position as well as which party wins in it are determined by the predicate $V$ (in Definition 3.8). That is, $V(x,y_1,...,y_k) = 1$ if the position that is reached when starting from position $x$ and taking the move sequence $y_1,...,y_k$ is a winning position for the first party.

**The collapsing effect of some equalities.**   Extending the intuition that underlies the $\mathcal{NP} \neq \text{co}\mathcal{NP}$ conjecture, it is commonly conjectured that $\Sigma_k \neq \Pi_k$ for every $k \in \mathbb{N}$. The failure of this conjecture causes the collapse of the Polynomial-time Hierarchy to the corresponding level.

---

[4]Note that, since we consider a constant number of moves, the length of all possible final positions is bounded by a polynomial in the length of the initial position, and thus all items have an equivalent form in which one refers to the complexity as a function of the length of the initial position. The latter form allows for a smooth generalization to games with a polynomial number of moves (as in Section 5.4), where it is essential to state all complexities in terms of the length of the initial position.

**Proposition 3.10** *For every $k \geq 1$, if $\Sigma_k = \Pi_k$ then $\Sigma_{k+1} = \Sigma_k$, which in turn implies $\mathcal{PH} = \Sigma_k$.*

The converse also holds (i.e., $\mathcal{PH} = \Sigma_k$ implies $\Sigma_{k+1} = \Sigma_k$ and $\Sigma_k = \Pi_k$). Needless to say, Proposition 3.10 does not seem to hold for $k = 0$.

**Proof:** Assuming that $\Sigma_k = \Pi_k$, we first show that $\Sigma_{k+1} = \Sigma_k$. For any set $S$ in $\Sigma_{k+1}$, by Proposition 3.9, there exists a polynomial $p$ and a set $S' \in \Pi_k$ such that $S = \{x : \exists y \in \{0,1\}^{p(|x|)} \text{ s.t. } (x, y) \in S'\}$. Using the hypothesis, we infer that $S' \in \Sigma_k$, and so (using Proposition 3.9 and $k \geq 1$) there exists a polynomial $p'$ and a set $S'' \in \Pi_{k-1}$ such that $S' = \{x' : \exists y' \in \{0,1\}^{p'(|x'|)} \text{ s.t. } (x', y') \in S''\}$. It follows that

$$S = \{x : \exists y \in \{0,1\}^{p(|x|)} \exists z \in \{0,1\}^{p'(|(x,y)|)} \text{ s.t. } ((x, y), z) \in S''\}.$$

By collapsing the two adjacent existential quantifiers (and using Proposition 3.9 yet again), we conclude that $S \in \Sigma_k$. This proves the first part of the proposition.

Turning to the second part, we note that $\Sigma_{k+1} = \Sigma_k$ (or, equivalently, $\Pi_{k+1} = \Pi_k$) implies $\Sigma_{k+2} = \Sigma_{k+1}$ (again by using Proposition 3.9), and similarly $\Sigma_{j+2} = \Sigma_{j+1}$ for any $j \geq k$. Thus, $\Sigma_{k+1} = \Sigma_k$ implies $\mathcal{PH} = \Sigma_k$. ∎

**Decision problems that are Cook-reductions to NP.** The Polynomial-time Hierarchy contains all decision problems that are Cook-reductions to $\mathcal{NP}$ (see Exercise 3.4). As shown next, the latter class contains many natural problems. Recall that in Section 2.2.2 we defined two types of optimization problems and showed that under some natural conditions these two types are computationally equivalent (under Cook reductions). Specifically, one type of problems referred to finding solutions that have a value *exceeding some given threshold*, whereas the second type called for finding *optimal solutions*. In Section 2.3 we presented several problems of the first type, and proved that they are NP-complete. We note that corresponding versions of the second type are believed not to be in NP. For example, we discussed the problem of deciding whether or not a given graph $G$ has a clique of a given size $K$, and showed that it is NP-complete. In contract, the problem of deciding whether or not $K$ is the maximum clique size of the graph $G$ is not known (and quite unlikely) to be in $\mathcal{NP}$, although it is Cook-reducible to $\mathcal{NP}$. Thus, the class of decision problems that are Cook-reducible to $\mathcal{NP}$ contains many natural problems that are unlikely to be in $\mathcal{NP}$. The Polynomial-time Hierarchy contains all these problems.

**Complete problems and a relation to AC0.** We note that quantified Boolean formulae with a bounded number of quantifier alternation provide complete problems for the various levels of the Polynomial-time Hierarchy (see Exercise 3.7). We also note the correspondence between these formulae and (highly uniform) constant-depth circuits of unbounded fan-in that get as input the truth-table of the underlying (quantifier-free) formula (see Exercise 3.8).

### 3.2.2  Non-deterministic oracle machines

The Polynomial-time Hierarchy is commonly defined in terms of non-deterministic polynomial-time (oracle) machines that are given oracle access to a set in the lower level of the same hierarchy. Such machines are defined by combining the definitions of non-deterministic (polynomial-time) machines (cf. Definition 2.7) and oracle machines (cf. Definition 1.11). Specifically, for an oracle $f : \{0,1\}^* \to \{0,1\}^*$, a non-deterministic oracle machine $M$, and a string $x$, one considers the question of whether or not there exists an accepting (non-deterministic) computation of $M$ on input $x$ and access to the oracle $f$. The class of sets that can be accepted by non-deterministic polynomial-time (oracle) machines with access to $f$ is denoted $\mathcal{NP}^f$. (We note that this notation makes sense because we can associate the class $\mathcal{NP}$ with a collection of machines that lends itself to be extended to oracle machines.) For any class of decision problems $\mathcal{C}$, we denote by $\mathcal{NP}^{\mathcal{C}}$ the union of $\mathcal{NP}^f$ taken over all decision problems $f$ in $\mathcal{C}$. The following result provides an alternative definition of the Polynomial-time Hierarchy.

**Proposition 3.11** *For every $k \geq 1$, it holds that $\Sigma_{k+1} = \mathcal{NP}^{\Sigma_k}$.*

**Proof:** The first direction (i.e., $\Sigma_{k+1} \subseteq \mathcal{NP}^{\Sigma_k}$) is almost straightforward: For any $S \in \Sigma_{k+1}$, let $S' \in \Pi_k$ and $p$ be as in Proposition 3.9; that is, $S = \{x : \exists y \in \{0,1\}^{p(|x|)} \text{ s.t. } (x,y) \in S'\}$. Consider the non-deterministic oracle machine that, on input $x$, non-deterministically generates $y \in \{0,1\}^{p(|x|)}$ and accepts if and only if (the oracle indicates that) $(x,y) \in S'$. This machine demonstrates that $S \in \mathcal{NP}^{\Pi_k} = \mathcal{NP}^{\Sigma_k}$, where the equality holds by letting the oracle machine flip each (binary) answer that is provided by the oracle.[5]

For the opposite direction (i.e., $\mathcal{NP}^{\Sigma_k} \subseteq \Sigma_{k+1}$), let $M$ be a non-deterministic polynomial-time oracle machine that accepts $S$ when given oracle access to $S' \in \Sigma_k$. Note that (unlike the machine constructed in the foregoing argument) machine $M$ may issue several queries to $S'$, and these queries may be determined based on previous oracle answers. To simplify the argument, we assume, without loss of generality, that at the very beginning of its execution machine $M$ guesses (non-deterministic) all oracle answers and accepts only if the actual answers match its guesses. Thus, $M$'s queries to the oracle are determined by its input, denoted $x$, and its non-deterministic choices, denoted $y$. We denote by $q^{(i)}(x,y)$ the $i^{\text{th}}$ query made by $M$ (on input $x$ and non-deterministic choices $y$), and by $a^{(i)}(x,y)$ the corresponding (a priori) guessed answer (which is a bit in $y$). Thus, $M$ accepts $x$ if and only if there exists $y \in \{0,1\}^{\text{poly}(|x|)}$ such that the following two conditions hold:

1. Machine $M$ accepts $x$, on input $x$ and non-deterministic choices $y$, when for every $i$ it holds that the $i^{\text{th}}$ oracle query made by $M$ is answered by the value $a^{(i)}(x,y)$. We stress that we do not assume here that these answers are consistent with $S'$; we merely refer to the decision of $M$ on a given input,

---

[5]Do not get confused by the fact that the class of oracles may *not* be closed under complementation. From the point of view of the oracle machine, the oracle is merely a function, and the machine may do with its answer whatever it pleases (and in particular negate it).

when it makes a specific sequence of non-deterministic choices, and is given specific oracle answers.

2. Each bit $a^{(i)}(x,y)$ is consistent with $S'$; that is, for every $i$, it holds that $a^{(i)}(x,y)=1$ if and only if $q^{(i)}(x,y)\in S'$.

Denoting the first event by $A(x,y)$ and letting $q(x,y)\leq \mathrm{poly}(|x|)$ denote the number of queries made by $M$, it follows that $x\in S$ if and only if

$$\exists y\left(A(x,y)\ \wedge\ \bigwedge_{i=1}^{q(x,y)}\left((a^{(i)}(x,y)=1)\Leftrightarrow(q^{(i)}(x,y)\in S')\right)\right).$$

Denoting the verification algorithm of $S'$ by $V'$, it holds that $x\in S$ if and only if

$$\exists y\left(A(x,y)\ \wedge\ \bigwedge_{i=1}^{q(x,y)}\left((a^{(i)}(x,y)=1)\Leftrightarrow\exists y_1^{(i)}\forall y_2^{(i)}\cdots Q_k y_k^{(i)}\,V'(q^{(i)}(x,y),y_1^{(i)},...,y_k^{(i)})=1\right)\right).$$

The proof is completed by observing that the foregoing expression can be rearranged to fit the definition of $\Sigma_{k+1}$. Details follow.

Starting with the foregoing expression, we first pull all quantifiers outside, and obtain a quantified expression with $k+1$ alternations, starting with an existential quantifier.[6] (We get $k+1$ alternations rather than $k$, because $a^{(i)}(x,y)=0$ introduces an expression of the form $\neg\exists y_1^{(i)}\forall y_2^{(i)}\cdots Q_k y_k^{(i)}\,V'(q^{(i)}(x,y),y_1^{(i)},...,y_k^{(i)})=1$, which in turn is equivalent to the expression $\forall y_1^{(i)}\exists y_2^{(i)}\cdots\overline{Q}_k y_k^{(i)}\,\neg V'(q^{(i)}(x,y),y_1^{(i)},...,y_k^{(i)})=1$).) Once this is done, we may incorporate the computation of all the $q^{(i)}(x,y)$'s (and $a^{(i)}(x,y)$'s) as well as the polynomial number of invocations of $V'$ (and other logical operations) into the new verification algorithm $V$. It follows that $S\in\Sigma_{k+1}$.

∎

**A general perspective – what does $\mathcal{C}_1^{\mathcal{C}_2}$ mean?** By the foregoing discussion it should be clear that the class $\mathcal{C}_1^{\mathcal{C}_2}$ can be defined for two complexity classes $\mathcal{C}_1$ and $\mathcal{C}_2$, *provided that $\mathcal{C}_1$ is associated with a class of machines that extends naturally in a way that allows for oracle access.* Actually, the class $\mathcal{C}_1^{\mathcal{C}_2}$ *is not defined based on the class $\mathcal{C}_1$ but rather by analogy to it.* Specifically, suppose that $\mathcal{C}_1$ is the class of sets that are recognizable (or rather accepted) by machines of certain type (e.g., deterministic or non-deterministic) with certain resource bounds (e.g., time and/or space bounds). Then, we consider analogous oracle machines (i.e., of the same type and with the same resource bounds), and say that $S\in\mathcal{C}_1^{\mathcal{C}_2}$ if there exists an adequate oracle machine $M_1$ (i.e., of this type and resource bounds) and a set $S_2\in\mathcal{C}_2$ such that $M_1^{S_2}$ accepts the set $S$.

---

[6]For example, note that for predicates $P_1$ and $P_2$, the expression $\exists y\,(P_1(y)\Leftrightarrow\exists z\,P_2(y,z))$ is equivalent to the expression $\exists y\,((P_1(y)\wedge\exists z\,P_2(y,z))\vee((\neg P_1(y)\wedge\neg\exists z\,P_2(y,z)))$, which in turn is equivalent to the expression $\exists y\exists z'\forall z''\,((P_1(y)\wedge P_2(y,z'))\vee((\neg P_1(y)\wedge\neg P_2(y,z'')))$. Note that pulling the quantifiers outside in $\wedge_{i=1}^t\exists y^{(i)}\forall z^{(i)}P(y^{(i)},z^{(i)})$ yields an expression of the type $\exists y^{(1)},...,y^{(t)}\forall z^{(1)},...,z^{(t)}\wedge_{i=1}^t P(y^{(i)},z^{(i)})$.

**Decision problems that are Cook-reductions to NP, revisited.**   Using the foregoing notation, the class of decision problems that are Cook-reductions to $\mathcal{NP}$ is denoted $\mathcal{P}^{\mathcal{NP}}$, and thus is a subset of $\mathcal{NP}^{\mathcal{NP}} = \Sigma_2$ (see Exercise 3.9).   In contrast, recall that the class of decision problems that are Karp-reductions to $\mathcal{NP}$ equals $\mathcal{NP}$.

### 3.2.3   The P/poly-versus-NP Question and PH

As stated in Section 3.1, a main motivation for the definition of $\mathcal{P}/\text{poly}$ is the hope that it can serve to separate $\mathcal{P}$ from $\mathcal{NP}$ (by showing that $\mathcal{NP}$ is not even contained in $\mathcal{P}/\text{poly}$, which is a (strict) superset of $\mathcal{P}$).  In light of the fact that $\mathcal{P}/\text{poly}$ extends far beyond $\mathcal{P}$ (and in particular contains undecidable problems), one may wonder if this approach does not run the risk of asking too much (since it may be that $\mathcal{NP}$ is in $\mathcal{P}/\text{poly}$ even if $\mathcal{P} \neq \mathcal{NP}$).  The common feeling is that the added power of non-uniformity is irrelevant with respect to the P-vs-NP Question. Ideally, we would like to know that $\mathcal{NP} \subset \mathcal{P}/\text{poly}$ may occur only if $\mathcal{P} = \mathcal{NP}$ (which may be phrased as saying that the Polynomial-time Hierarchy collapses to its zero level).  The following result seems to get close to such an implication, showing that $\mathcal{NP} \subset \mathcal{P}/\text{poly}$ may occur only if the Polynomial-time Hierarchy collapses to its second level.

**Theorem 3.12** *If $\mathcal{NP} \subset \mathcal{P}/\text{poly}$ then $\Sigma_2 = \Pi_2$.*

Recall that $\Sigma_2 = \Pi_2$ implies $\mathcal{PH} = \Sigma_2$ (see Proposition 3.10).  Thus, an unexpected behavior of the non-uniform complexity class $\mathcal{P}/\text{poly}$ implies an unexpected behavior in the world of uniform complexity (i.e., the ability to reduce any constant number of quantifier alternations to two quantifier alternations).

**Proof:**   Using the hypothesis (i.e., $\mathcal{NP} \subset \mathcal{P}/\text{poly}$) and starting with an arbitrary set $S \in \Pi_2$, we shall show that $S \in \Sigma_2$.  Loosely speaking, $S \in \Pi_2$ means that $x \in S$ if and only if for all $y$ there exists a $z$ such that some (fixed) polynomial-time verifiable condition regarding $(x, y, z)$ holds.  Note that the residual condition regarding $(x, y)$ is of the NP-type, and thus (by the hypothesis) it can be verified by a polynomial-size circuit.  This suggests saying that $x \in S$ if and only if there exists an *adequate* circuit $C$ such that for all $y$ it holds that $C(x, y) = 1$.  Thus, we managed to switch the order of the universal and existential quantifiers.  Specifically, the resulting assertion is of the desired $\Sigma_2$-type provided that we can either verify the *adequacy condition* in co$\mathcal{NP}$ (or even in $\Sigma_2$) or keep out of trouble even in the case that $x \notin S$ and $C$ is inadequate.  In the following proof we implement the latter option by observing that the hypothesis yields small circuits for NP-search problems (and not only for NP-decision problems).  Specifically, we obtain (small) circuits that given $(x, y)$ find an NP-witness for $(x, y)$ (whenever such a witness exists) and relying on the fact that we can efficiently verify the correctness of NP-witnesses.  (The alternative approach of providing a coNP-type procedure for verifying the adequacy of the circuit is pursued in Exercise 3.11.)

   Let $S$ be an arbitrary set in $\Pi_2$.  Then, by Proposition 3.9, there exists a polynomial $p$ and a set $S' \in \mathcal{NP}$ such that $S = \{x : \forall y \in \{0, 1\}^{p(|x|)} \ (x, y) \in S'\}$.

Let $R' \in \mathcal{PC}$ be the witness-relation corresponding to $S'$; that is, there exists a polynomial $p'$, such that $x' = \langle x, y \rangle \in S'$ if and only if there exists $z \in \{0, 1\}^{p'(|x'|)}$ such that $(x', z) \in R'$. It follows that

$$ S = \{ x : \forall y \in \{0, 1\}^{p(|x|)} \exists z \in \{0, 1\}^{p'(|\langle x, y \rangle|)} \ (\langle x, y \rangle, z) \in R' \}. $$

By the reduction of $\mathcal{PC}$ to $\mathcal{NP}$ (see the proof of Theorem 2.6 and further discussion in Section 2.2.1), the theorem's hypothesis (i.e., $\mathcal{NP} \subseteq \mathcal{P}/\text{poly}$) implies the existence of polynomial-size circuits for solving the search problem of $R'$. Using the existence of these circuits, it follows that for any $x \in S$ there exists a small circuit $C'$ such that for every $y$ it holds that $C'(x, y) \in R'(x, y)$, whereas for any $x \notin S$ there exists a $y$ such that $\langle x, y \rangle \notin S'$ and hence $C'(x, y) \notin R'(x, y)$ for any circuit $C'$ (for the trivial reason that $R'(x, y) = \emptyset$). But let us first spell-out what we mean by polynomial-size circuits for solving a search problem as well as further justify their existence for the search problem of $R'$.

In Section 3.1, we have focused on polynomial-size circuits that solve decision problems. However, the definition sketched in Section 3.1.1 also applies to solving search problems, provided that an appropriate encoding is used for allowing solutions of possibly varying lengths (for instances of fixed length) to be presented as strings of fixed length. Next observe that combining the Cook-reduction of $\mathcal{PC}$ to $\mathcal{NP}$ with the hypothesis $\mathcal{NP} \subseteq \mathcal{P}/\text{poly}$, implies that $\mathcal{PC}$ is Cook-reducible to $\mathcal{P}/\text{poly}$. In particular, this implies that any search problem in $\mathcal{PC}$ can be solved by a family of polynomial-size circuits. Note that the resulting circuit that solves $n$-bit long instances of such a problem may incorporate polynomially (in $n$) many circuits, each solving a decision problem for $m$-bit long instances, where $m \in [\text{poly}(n)]$. Needless to say, the size of the resulting circuit that solves the search problem of the aforementioned $R' \in \mathcal{PC}$ (for instances of length $n$) is upper-bounded by $\text{poly}(n) \cdot \sum_{m=1}^{\text{poly}(n)} \text{poly}(m)$.

It follows that $x \in S$ if and only if *there exists a* $\text{poly}(|x| + p(|x|))$*-size circuit* $C'$ *such that for all* $y \in \{0, 1\}^{p(|x|)}$ *it holds that* $(\langle x, y \rangle, C'(x, y)) \in R'$. Note that in the case that $x \in S$ we use the circuit $C'$ that is guaranteed for inputs of length $|x| + p(|x|)$ by the foregoing discussion[7], whereas in the case that $x \notin S$ it does not matter which circuit $C'$ is used (because in that case there exists a $y$ such that for all $z$ it holds that $(\langle x, y \rangle, z)) \notin R'$).

The key observation regarding the foregoing condition (i.e., $\exists C' \forall y \ (\langle x, y \rangle, C'(x, y)) \in R'$) is that it is of the desired form (of a $\Sigma_2$ statement). Specifically, consider the polynomial-time verification procedure $V$ that given $x, y$ and the description of the circuit $C'$, first computes $z \leftarrow C'(x, y)$ and accepts if and only if $(\langle x, y \rangle, z) \in R'$, where the latter condition can be verified in polynomial-time (because $R' \in \mathcal{PC}$). Denoting the description of a potential circuit by $\langle C' \rangle$, the aforementioned (polynomial-time) computation of $V$ is denoted $V(x, \langle C' \rangle, y)$, and indeed $x \in S$ if and only if

$$ \exists \langle C' \rangle \in \{0, 1\}^{\text{poly}(|x| + p(|x|))} \forall y \in \{0, 1\}^{p(|x|)} \ V(x, \langle C' \rangle, y) = 1. $$

---

[7] Thus, $C'$ may actually depend only on $|x|$, which in turn determines $p(|x|)$.

Having established that $S \in \Sigma_2$ for an arbitrary $S \in \Pi_2$, we conclude that $\Pi_2 \subseteq \Sigma_2$. The theorem follows (by applying Exercise 3.9.4). ∎

# Chapter Notes

The class $\mathcal{P}/\text{poly}$ was defined by Karp and Lipton [130] as part of a general formulation of "machines which take advice" [130]. They also noted the equivalence to the traditional formulation of polynomial-size circuits as well as the effect of uniformity (Proposition 3.4).

The Polynomial-Time Hierarchy ($\mathcal{PH}$) was introduced by Stockmeyer [201]. A third equivalent formulation of $\mathcal{PH}$ (via so-called "alternating machines") can be found in [48].

The implication of the failure of the conjecture that $\mathcal{NP}$ is not contained in $\mathcal{P}/\text{poly}$ on the Polynomial-time Hierarchy (i.e., Theorem 3.12) was discovered by Karp and Lipton [130]. This interesting connection between non-uniform and uniform complexity provides the main motivation for presenting $\mathcal{P}/\text{poly}$ and $\mathcal{PH}$ in the same chapter.

# Exercises

**Exercise 3.1 (a small variation on the definitions of $\mathcal{P}/\text{poly}$)** Using an adequate encoding of strings of length smaller than $n$ as $n$-bit strings (e.g., $x \in \cup_{i<n}\{0,1\}^i$ is encoded as $x10^{n-|x|-1}$), define circuits (resp., machines that take advice) as devices that can handle inputs of various lengths up to a given bound (rather than as devices that can handle inputs of a fixed length). Show that the class $\mathcal{P}/\text{poly}$ remains invariant under this change (and Theorem 3.6 remains valid).

**Exercise 3.2 (sparse sets)** A set $S \subset \{0,1\}^*$ is called sparse if there exists a polynomial $p$ such that $|S \cap \{0,1\}^n| \leq p(n)$ for every $n$.

1. Prove that any sparse set is in $\mathcal{P}/\text{poly}$. Note that a sparse set may be undecidable.

2. Prove that a set is in $\mathcal{P}/\text{poly}$ if and only if it is Cook-reducible to some sparse set.

**Guideline:** For the forward direction of Part 2, encode the advice sequence $(a_n)_{n\in\mathbb{N}}$ as a sparse set $\{(1^n, i, \sigma_{n,i}) : n \in \mathbb{N}, i \leq |a_n|\}$, where $\sigma_{n,i}$ is the $i^{\text{th}}$ bit of $a_n$. For the opposite direction, note that on input $x$ the Cook-reduction makes queries of length at most $\text{poly}(|x|)$, and so emulating the reduction on an input of length $n$ only requires knowledge of all the strings that are in the sparse set and have length at most $\text{poly}(n)$.

**Exercise 3.3 (advice hierarchy)** Prove that for any two functions $\ell, \delta : \mathbb{N} \to \mathbb{N}$ such that $\ell(n) < 2^{n-1}$ and $\delta$ is unbounded, it holds that $\mathcal{P}/\ell$ is strictly contained in $\mathcal{P}/(\ell + \delta)$.

**Guideline:** For every sequence $\overline{a} = (a_n)_{n\in\mathbb{N}}$ such that $|a_n| = \ell(n) + \delta(n)$, consider the set $S_{\overline{a}}$ that encodes $\overline{a}$ such that $x \in S_{\overline{a}} \cap \{0,1\}^n$ if and only if the $\text{idx}(x)^{\text{th}}$ bit in $a_n$

equals 1 (and $\text{idx}(x) \leq |a_n|$), where $\text{idx}(x)$ denotes the index of $x$ in $\{0,1\}^n$. For more details see Section 4.1.

**Exercise 3.4** Prove that $\Sigma_2$ contains all sets that are Cook-reducible to $\mathcal{NP}$.

**Guideline:** This is quite obvious when using the definition of $\Sigma_2$ as presented in Section 3.2.2; see Exercise 3.9. Alternatively, the fact can be proved by using *some* of the ideas that underlie the proof of Theorem 2.33, while noting that a conjunction of NP and coNP assertions forms an assertion of type $\Sigma_2$ (see also the second part of the proof of Proposition 3.11).

**Exercise 3.5** Let $\Delta = \mathcal{NP} \cap \text{co}\mathcal{NP}$. Prove that $\Delta$ equals the class of decision problems that are Cook-reducible to $\Delta$ (i.e., $\Delta = \mathcal{P}^\Delta$).

**Guideline:** See proof of Theorem 2.33.

**Exercise 3.6 (the class $\Pi_i$)** Recall that $\Pi_k$ is defined to equal $\text{co}\Sigma_k$, which in turn is defined to equal $\{\{0,1\}^* \setminus S : S \in \Sigma_k\}$. Prove that for any natural number $k$, a decision problem $S \subseteq \{0,1\}^*$ is in $\Pi_k$ if there exists a polynomial $p$ and a polynomial time algorithm $V$ such that $x \in S$ if and only if

$$\forall y_1 \in \{0,1\}^{p(|x|)} \exists y_2 \in \{0,1\}^{p(|x|)} \forall y_3 \in \{0,1\}^{p(|x|)} \cdots Q_k y_k \in \{0,1\}^{p(|x|)}$$
$$\text{s.t. } V(x, y_1, ..., y_k) = 1$$

where $Q_k$ is a universal quantifier if $k$ is odd and is an existential quantifier otherwise.

**Exercise 3.7 (complete problems for the various levels of $\mathcal{PH}$)** A $k$-alternating quantified Boolean formula is a quantified Boolean formula with up to $k$ alternations between existential and universal quantifiers, starting with an existential quantifier. For example, $\exists z_1 \exists z_2 \forall z_3 \phi(z_1, z_2, z_3)$ (where the $z_i$'s are Boolean variables) is a 2-alternating quantified Boolean formula. Prove that the problem of *deciding whether or not a k-alternating quantified Boolean formula is valid* is $\Sigma_k$-complete under Karp-reductions. That is, denoting the aforementioned problem by kQBF, prove that kQBF is in $\Sigma_k$ and that every problem in $\Sigma_k$ is Karp-reducible to kQBF.

**Exercise 3.8 (on the relation between $\mathcal{PH}$ and $\mathcal{AC}^0$)** Note that there is an obvious analogy between $\mathcal{PH}$ and constant-depth polynomial-size circuits of unbounded fan-in, where existential (resp., universal) quantifiers are represented by "large" $\bigvee$ (resp., $\bigwedge$) gates. To articulate this relationship, consider the following definitions.

- A family of circuits $\{C_N\}$ is called highly uniform if there exists a polynomial-time algorithm that answers local queries regarding the structure of the relevant circuit. Specifically, on input $(N, u, v)$, the algorithm determines the type of gates represented by the vertices $u$ and $v$ in $C_N$ as well as whether there exists a directed edge from $u$ to $v$. Note that this algorithm operates in time that polylogarithmic in the size of $C_N$.

  We focus on family of polynomial-size circuits, meaning that the size of $C_N$ is polynomial in $N$, which in turn represents the number of inputs to $C_N$.

- Fixing a polynomial $p$, a $p$-succinctly represented input $Z \in \{0,1\}^N$ is a circuit $c_Z$ of size at most $p(\log_2 N)$ such that for every $i \in [N]$ it holds that $c_Z(i)$ equals the $i^{\text{th}}$ bit of $Z$.

- For a fixed family of highly uniform circuits $\{C_N\}$ and a fixed polynomial $p$, the problem of evaluating a succinctly represented input is defined as follows. *Given $p$-succinct representation of an input $Z \in \{0,1\}^N$, determine whether or not $C_N(Z) = 1$.*

For every $k$ and every $S \in \Sigma_k$, show that there exists a family of highly uniform unbounded fan-in circuits of depth $k$ and polynomial-size such that $S$ is Karp-reducible to evaluating a succinctly represented input (with respect to that family of circuits). That is, the reduction should map an instance $x \in \{0,1\}^n$ to a $p$-succinct representation of some $Z \in \{0,1\}^N$ such that $x \in S$ if and only if $C_N(Z) = 1$. (Note that $Z$ is represented by a circuit $c_Z$ of size at most $p(\log_2 N)$, and that it follows that $|c_Z| \leq \text{poly}(n)$ and thus $N \leq \exp(\text{poly}(n))$.)[8]

**Guideline:** Let $S \in \Sigma_k$ and let $V$ be the corresponding verification algorithm as in Definition 3.8. That is, $x \in S$ if and only if $\exists y_1 \forall y_2 \cdots Q_k y_k$, where each $y_i \in \{0,1\}^{\text{poly}(|x|)}$ such that $V(x, y_1, ..., y_k) = 1$. Then, for $m = \text{poly}(|x|)$ and $N = 2^{k \cdot m}$, consider the fixed circuit $C_N(Z) = \bigvee_{i_1 \in [2^m]} \bigwedge_{i_2 \in [2^m]} \cdots Q'_{i_k \in [2^m]} Z_{i_1,i_2,...,i_k}$, and the problem of evaluating $C_N$ at an input consisting of the truth-table of $V(x, \cdots)$ (i.e., when setting $Z_{i_1,i_2,...,i_k} = V(x, i_1, ..., i_k)$, where $[2^m] \equiv \{0,1\}^m$). Note that the size of $C_N$ is $O(N)$.[9]

**Exercise 3.9** Verify the following facts:

1. For every $k \geq 1$, it holds that $\Sigma_k \subseteq \mathcal{P}^{\Sigma_k} \subseteq \Sigma_{k+1}$.

   (Note that, for any complexity class $\mathcal{C}$, the class $\mathcal{P}^{\mathcal{C}}$ is the class of sets that are Cook-reducible to some set in $\mathcal{C}$. In particular, $\mathcal{P}^{\mathcal{P}} = \mathcal{P}$.)

2. For every $k \geq 1$, $\Pi_k \subseteq \mathcal{P}^{\Pi_k} \subseteq \Pi_{k+1}$.

   (Hint: For any complexity class $\mathcal{C}$, it holds that $\mathcal{P}^{\mathcal{C}} = \mathcal{P}^{\text{co}\mathcal{C}}$ and $\mathcal{P}^{\mathcal{C}} = \text{co}\mathcal{P}^{\mathcal{C}}$.)

3. For every $k \geq 1$, it holds that $\Sigma_k \subseteq \Pi_{k+1}$ and $\Pi_k \subseteq \Sigma_{k+1}$. Thus, $\mathcal{PH} = \cup_k \Pi_k$.

4. For every $k \geq 1$, if $\Sigma_k \subseteq \Pi_k$ (resp., $\Pi_k \subseteq \Sigma_k$) then $\Sigma_k = \Pi_k$.

   (Hint: For any $S \in \Pi_k$ (resp., $S \in \Sigma_k$), apply the hypothesis to $\{0,1\}^* \setminus S$.)

**Exercise 3.10** In continuation to Exercise 3.7, prove that following claims:

1. SAT is computationally equivalent to 1QBF.

---

[8] Assuming $\mathcal{P} \neq \mathcal{NP}$, it cannot be that $N \leq \text{poly}(n)$ (because circuit evaluation can be performed in time polynomial in the size of the circuit).

[9] Advanced comment: the computational limitations of $\mathcal{AC}^0$ circuits (see, e.g., [78, 110]) imply limitations on the functions of a *generic* input $Z$ that the aforementioned circuits $C_N$ can compute. Unfortunately, these computational limitations do not seem to provide useful information on the limitations of functions of inputs $Z$ that have succinct representation (as obtained by setting $Z_{i_1,i_2,...,i_k} = V(x, i_1, ..., i_k)$, where $V$ is a polynomial-time algorithm). This fundamental problem is "resolved" in the context of "relativization" by providing $V$ with oracle access to an arbitrary input of length $N$ (or so); cf. [78].

2. For every $k \geq 1$, it holds that $\mathcal{P}^{\Sigma_k} = \mathcal{P}^{\mathtt{kQBF}}$ and $\Sigma_{k+1} = \mathcal{NP}^{\mathtt{kQBF}}$.

   **Guideline:** Prove that if $S$ is $\mathcal{C}$-complete then $\mathcal{P}^{\mathcal{C}} = \mathcal{P}^S$. Note that $\mathcal{P}^{\mathcal{C}} \subseteq \mathcal{P}^S$ uses the polynomial-time reductions of $\mathcal{C}$ to $S$, whereas $\mathcal{P}^S \subseteq \mathcal{P}^{\mathcal{C}}$ uses $S \in \mathcal{C}$.

**Exercise 3.11 (an alternative proof of Theorem 3.12)** In continuation to the discussion in the proof of Theorem 3.12, use the following guidelines to provide an alternative proof of Theorem 3.12.

1. First, prove that if $S'$ is downwards self-reducible (as defined in Exercise 2.13) then the correctness of circuits deciding $S'$ can be decided in $\mathrm{co}\mathcal{NP}$. Specifically, denoting by $\chi$ the characteristic function of $S'$, show that the set

$$\mathrm{ckt}_\chi \stackrel{\mathrm{def}}{=} \{(1^n, \langle C \rangle) : \forall w \in \{0,1\}^n \; C(w) = \chi(w)\}$$

   is in $\mathrm{co}\mathcal{NP}$.

   **Guideline:** Using the more flexible formulation suggested in Exercise 3.1, it suffices to verify that, for every $i < n$ and every $i$-bit string $w$, the value $C(w)$ equals the output of the downwards self-reduction on input $w$ when obtaining answers according to $C$. Thus, for every $i < n$, the correctness of $C$ on inputs of length $i$ follows from its correctness on inputs of length less than $i$. Needless to say, the correctness of $C$ on the empty string (or on all inputs of constant length) can be verified by comparison to the fixed value of $\chi$ on the empty string (resp., the values of $\chi$ on a constant number of strings).

2. Recalling that $\mathtt{SAT}$ is downwards self-reducible and that $\mathcal{NP}$ reduces to $\mathtt{SAT}$, derive Theorem 3.12 as a corollary of Part 1.

**Exercise 3.12** In continuation to Part 2 of Exercise 3.2, we consider the class of sets that are Karp-reducible to a sparse set. It can be proved that this class contains $\mathtt{SAT}$ if and only if $\mathcal{P} = \mathcal{NP}$ (see [76]). Here, we only consider the special case in which the sparse set is contained in a polynomial-time decidable set that is itself sparse (e.g., the latter set may be $\{1\}^*$, in which case the former set may be an arbitrary unary set). Actually, prove the following seemingly stronger claim:

> *if $\mathtt{SAT}$ is Karp-reducible to a set $S \subseteq G$ such that $G \in \mathcal{P}$ and $G \setminus S$ is sparse then $\mathtt{SAT} \in \mathcal{P}$.*

Using the hypothesis, we outline a polynomial-time procedure for solving the search problem of SAT, and leave the task of providing the details as an exercise. The procedure conducts a DFS on the tree of all possible partial truth assignment to the input formula, while truncating the search at nodes that are roots of sub-trees that were already demonstrated to contain no satisfying assignment (at the leaves).[10]

**Guideline:** The key observation is that each internal node (which yields a formula derived from the initial formulae by instantiating the corresponding partial truth assignment) is mapped by the Karp-reduction either to a string not in $G$ (in which case we conclude

---

[10] For an $n$-variable formulae, the leaves of the tree correspond to all possible $n$-bit long strings, and an internal node corresponding to $\tau$ is the parent of nodes corresponding to $\tau 0$ and $\tau 1$.

that the sub-tree contains no satisfying assignments and backtrack from this node) or to a string in $G$. In the latter case, unless we already know that this string is not in $S$, we *start a scan of the sub-tree rooted at this node*. However, once we backtrack from this internal node, we know that the corresponding element of $G$ is not in $S$, and we will never scan again a sub-tree rooted at a node that mapped to this element. Also note that once we reach a leaf, we can check by ourselves whether or not it corresponds to a satisfying assignment to the initial formula.

(Hint: When analyzing the forgoing procedure, note that on input an $n$-variable formulae $\phi$ the number of times we start to scan a sub-tree is at most $n \cdot |\cup_{i=1}^{\mathrm{poly}(|\phi|)} \{0,1\}^i \cap (G \setminus S)|$.)

# Chapter 4

# More Resources, More Power?

*More electricity, less toil.*

The Israeli Electricity Company, 1960s

*Is it indeed the case that the more resources one has, the more one can achieve?*
The answer may seem obvious, but the obvious answer (of yes) actually presumes
that the worker knows how much resources are at his/her disposal. In this case,
when allocated more resources, the worker (or computation) can indeed achieve
more. But otherwise, nothing may be gained by adding resources.

In the context of computational complexity, an algorithm knows the amount of
resources that it is allocated if it can determine this amount without exceeding the
corresponding resources. This condition is satisfies in all "reasonable" cases, but it
may not hold in general. The latter fact should not be that surprising: we already
know that some functions are not computable and if these functions are used to
determine resources then the algorithm may be in trouble. Needless to say, this
discussion requires some formalization, which is provided in the current chapter.

**Summary:** When using "nice" functions to determine the algorithm's
resources, it is indeed the case that more resources allow for more tasks
to be performed. However, when "ugly" functions are used for the same
purpose, increasing the resources may have no effect. By nice functions
we mean functions that can be computed without exceeding the amount
of resources that they specify (e.g., $t(n) = n^2$ or $t(n) = 2^n$). Naturally,
"ugly" functions do not allow to present themselves in such nice forms.

The forgoing discussion refers to a uniform model of computation and
to (natural) resources such as time and space complexities. Thus, we
get results asserting, for example, that there are functions computable
in cubic-time but not in quadratic-time. In case of non-uniform models

of computation, the issue of "nicety" does not arise, and it is easy to establish separations between levels of circuit complexity that differ by any unbounded amount.

Results that *separate* the class of problems solvable within one resource bound from the class of problems solvable within a larger resource bound are called hierarchy theorems. Results that indicate the non-existence of such separations, hence indicating a "gap" in the growth of computing power (or a "gap" in the existence of algorithms that utilize the added resources), are called gap theorems. A somewhat related phenomenon, called speed-up theorems, refers to the inability to define the complexity of some problems.

**Caveat:** Uniform complexity classes based on specific resource bounds (e.g., cubic-time) are model dependent. Furthermore, the tightness of separation results (i.e., how much more time is required to solve an additional computational problem) is also model dependent. Still the existence of such separations is a phenomenon common to all reasonable and general models of computation (as referred to in the Cobham-Edmonds Thesis). In the following presentation, we will explicitly differentiate model-specific effects from generic ones.

**Organization:** We will first demonstrate the "more resources yield more power" phenomenon in the context of non-uniform complexity. In this case the issue of "knowing" the amount of resources allocated to the computing device does not arise, because each device is tailored to the amount of resources allowed for the input length that it handles (see Section 4.1). We then turn to the time complexity of uniform algorithms; indeed, hierarchy and gap theorems for time-complexity, presented in Section 4.2, constitute the main part of the current chapter. We end by mentioning analogous results for space-complexity (see Section 4.3, which may also be read after Section 5.1).

## 4.1   Non-uniform complexity hierarchies

The model of machines that use advice (cf. §1.2.4.2 and Section 3.1.2) offers a very convenient setting for separation results. We refer specifically, to classes of the form $\mathcal{P}/\ell$, where $\ell : \mathbb{N} \to \mathbb{N}$ is an arbitrary function (see Definition 3.5). Recall that every Boolean function is in $\mathcal{P}/2^n$, by virtue of a trivial algorithm that is given as advice the truth-table of the function restricted to the relevant input length. An analogous algorithm underlies the following separation result.

**Theorem 4.1** *For any two functions $\ell', \delta : \mathbb{N} \to \mathbb{N}$ such that $\ell'(n) + \delta(n) \leq 2^n$ and $\delta$ is unbounded, it holds that $\mathcal{P}/\ell'$ is strictly contained in $\mathcal{P}/(\ell' + \delta)$.*

**Proof:** Let $\ell \stackrel{\text{def}}{=} \ell' + \delta$, and consider the algorithm $A$ that given advice $a_n \in \{0,1\}^{\ell(n)}$ and input $i \in \{1, ..., 2^n\}$ (viewed as an $n$-bit long string), outputs the $i^{\text{th}}$ bit of $a_n$ if $i \leq |a_n|$ and zero otherwise. Clearly, for any $\overline{a} = (a_n)_{n \in \mathbb{N}}$ such that

$|a_n| = \ell(n)$, it holds that the function $f_{\overline{a}}(x) \stackrel{\text{def}}{=} A(a_{|x|}, x)$ is in $\mathcal{P}/\ell$. Furthermore, different sequences $\overline{a}$ yield different functions $f_{\overline{a}}$. We claim that some of these functions $f_{\overline{a}}$ are not in $\mathcal{P}/\ell'$, thus obtaining a separation.

The claim is proved by considering all possible (polynomial-time) algorithms $A'$ and all possible sequences $\overline{a}' = (a'_n)_{n \in \mathbb{N}}$ such that $|a'_n| = \ell'(n)$. Fixing any algorithm $A'$, we consider the number of $n$-bit long functions that are correctly computed by $A'(a'_n, \cdot)$. Clearly, the number of these functions is at most $2^{\ell'(n)}$, and thus $A'$ may account for at most $2^{-\delta(n)}$ fraction of the functions $f_{\overline{a}}$ (even when restricted to $n$-bit strings). This consideration holds for every $n$ and every possible $A'$, and thus the measure of the set of functions that are computable by algorithms that take advice of length $\ell'$ is zero.[1]  ■

A somewhat less tight bound can be obtained by using the model of Boolean circuits. In this case some slackness is needed in order to account for the gap between the upper and lower bounds regarding the number of Boolean functions over $\{0,1\}^n$ that are computed by Boolean circuits of size $s < 2^n$. Specifically (see Exercise 4.1), an obvious lower-bound on this number is $2^{s/O(\log s)}$ whereas an obvious upper-bound is $s^{2s} = 2^{2s \log_2 s}$. (Compare these bounds to the lower-bound $2^{\ell'(n)}$ and the upper-bound $2^{\ell'(n) + ((\delta(n)-2)/2)}$, which were used in the proof of Theorem 4.1.)

## 4.2 Time Hierarchies and Gaps

In this section we show that in the "reasonable cases" increasing time-complexity allows for more problems to be solved, whereas in "pathological cases" it may happen that even a dramatic increase in the time-complexity provides no additional computing power. As hinted in the introductory comments to the current chapter, the "reasonable cases" correspond to time bounds that can be determined by the algorithm itself within the specified time complexity.

We stress that also in the aforementioned "reasonable cases", the added power does not necessarily refer to natural computational problems. That is, like in the case of non-uniform complexity (i.e., Theorem 4.1), the hierarchy theorems are proved by introducing artificial computational problems. Needless to say, we do not know of natural problems in $\mathcal{P}$ that are provably unsolvable in cubic (or some other fixed polynomial) time (on, say, a two-tape Turing machine). Thus, although $\mathcal{P}$ contains an infinite hierarchy of computational problems, each requiring significantly more time than the other, we know of no such hierarchy of natural computational problems. In contrast, it is our experience that any natural problem shown to be solvable in polynomial-time is followed by a sequence of results that eventually establish algorithms having running-time that is bounded by a moderate polynomial.

---

[1]It suffices to show that this measure is strictly less than one. This is easily done by considering, for every $n$, the performance of any algorithm $A'$ having description of length shorter than $(\delta(n) - 2)/2$ on all inputs of length $n$.

### 4.2.1 Time Hierarchies

Note that the non-uniform computing devices, considered in Section 4.1, were explicitly given the relevant resource bounds (e.g., the length of advice). Actually, they were given the resources themselves (e.g., the advice itself) and did not need to monitor their usage of these resources. In contrast, when designing algorithms of arbitrary time-complexity $t : \mathbb{N} \to \mathbb{N}$, we need to make sure that the algorithm does not exceed the time bound. Furthermore, when invoked on input $x$, the algorithm is not given the time bound $t(|x|)$ explicitly, and a reasonable design methodology is to have the algorithm compute this bound (i.e., $t(|x|)$) before doing anything else. This, in turn, requires the algorithm to read the entire input (see Exercise 4.3) as well as to compute $t(n)$ using $O(t(n))$ (or so) time. The latter requirement motivates the following definition (which is related to the standard definition of "fully time constructibility" (cf. [117, Sec. 12.3])).

**Definition 4.2** (time constructible functions): *A function $t : \mathbb{N} \to \mathbb{N}$ is called* time constructible *if there exists an algorithm that on input $n$ outputs $t(n)$ using at most $t(n)$ steps.*

Equivalently, we may require that the mapping $1^n \mapsto t(n)$ be computable within time complexity $t$. We warn that the foregoing definition is model dependent; however, typically nice functions are computable even faster (e.g., in $\mathrm{poly}(\log t(n))$ steps), in which case the model-dependency is irrelevant (for reasonable and general models of computation, as referred to in the Cobham-Edmonds Thesis). For example, in any reasonable and general model, functions like $t_1(n) = n^2$, $t_2(n) = 2^n$, and $t_3(n) = 2^{2^n}$ are computable in $\mathrm{poly}(\log t_i(n))$ steps.

Likewise, for a fixed model of computation (to be understood from the context) and for any function $t : \mathbb{N} \to \mathbb{N}$, we *denote by* $\mathrm{DTIME}(t)$ *the class of decision problems that are solvable in time complexity $t$*. We call the reader's attention to Exercise 4.7 that asserts that in many cases $\mathrm{DTIME}(t) = \mathrm{DTIME}(t/2)$.

#### 4.2.1.1 The Time Hierarchy Theorem

In the following theorem, we refer to the model of two-tape Turing machines. In this case we obtain quite a tight hierarchy in terms of the relation between $t_1$ and $t_2$. We stress that, using the Cobham-Edmonds Thesis, this results yields (possibly less tight) hierarchy theorems for any reasonable and general model of computation.

> **Teaching note:** The standard statement of Theorem 4.3 asserts that *for any time constructible function $t_2$ and every function $t_1$ such that $t_2 = \omega(t_1 \log t_1)$ and $t_1(n) > n$ it holds that* $\mathrm{DTIME}(t_1)$ *is strictly contained in* $\mathrm{DTIME}(t_2)$. The current version is only slightly weaker, but it allows a somewhat simpler and more intuitive proof. We comment on the proof of the standard version of Theorem 4.3 after proving the current version.

**Theorem 4.3** (time hierarchy for two-tape Turing machines): *For any time constructible function $t_1$ and every function $t_2$ such that $t_2(n) \geq (\log t_1(n))^2 \cdot t_1(n)$ and $t_1(n) > n$ it holds that $\mathrm{DTIME}(t_1)$ is strictly contained in $\mathrm{DTIME}(t_2)$.*

As will become clear from the proof, an analogous result holds for any model in which a universal machine can emulate $t$ steps of another machine in $O(t \log t)$ time, where the constant in the O-notation depends on the emulated machine. Before proving Theorem 4.3, we derive the following corollary.

**Corollary 4.4** (time hierarchy for any reasonable and general model): *For any reasonable and general model of computation there exists a positive polynomial $p$ such that for any time-computable function $t_1$ and every function $t_2$ such that $t_2 > p(t_1)$ and $t_1(n) > n$ it holds that* DTIME$(t_1)$ *is strictly contained in* DTIME$(t_2)$.

It follows that, for every such model and every polynomial $t$ (such that $t(n) > n$), there exist problems in $\mathcal{P}$ that are not in DTIME$(t)$. It also follows that $\mathcal{P}$ is a strict subset of $\mathcal{E}$ or even of "quasi-polynomial time"; moreover, $\mathcal{P}$ is a strict subset of DTIME$(q)$, where $q(n) = n^{\log_2 n}$ (or even $q(n) = n^{\log_2 \log_2 n}$).

**Proof of Corollary 4.4:** Letting DTIME$_2$ denote the classes that correspond to two-tape Turing machines, we have DTIME$(t_1) \subseteq$ DTIME$_2(t_1')$ and DTIME$(t_2) \supseteq$ DTIME$_2(t_2')$, where $t_1' = \text{poly}(t_1)$ and $t_2'$ is defined such that $t_2(n) = \text{poly}(t_2'(n))$. The latter unspecified polynomials, hereafter denoted $p_1$ and $p_2$ respectively, are the ones guaranteed by the Cobham-Edmonds Thesis. Also, the hypothesis that $t_1$ is time-computable implies that $t_1' = p_1(t_1)$ is time-constructible with respect to the two-tape Turing machine model. Thus, for a suitable choice of the polynomial $p$, it holds that

$$t_2'(n) = p_2^{-1}(t_2(n)) > p_2^{-1}(p(t_1(n))) > p_2^{-1}(p(p_1^{-1}(t_1'(n)))) > t_1'(n)^2 \,.$$

Invoking Theorem 4.3, we have DTIME$_2(t_2') \supset$ DTIME$_2(t_1')$, and the corollary follows. $\blacksquare$

**Proof of Theorem 4.3:** The idea is to construct a Boolean function $f$ such that all machines having time complexity $t_1$ fail to compute $f$. This is done by associating each possible machine $M$ a different input $x_M$ (e.g., $x_M = \langle M \rangle$), and making sure that $f(x_M) \neq M'(x)$, where $M'(x)$ denotes an emulation of $M(x)$ that is suspended after $t_1(|x|)$ steps. Actually, we are going to use a mapping $\mu$ of inputs to machines (i.e., $\mu(x_M) = M$), such that each machine is in the range of $\mu$ and $\mu$ is very easy to compute. Thus, by construction, $f \notin$ DTIME$(t_1)$.

The issue is presenting an algorithm for computing $f$. This algorithm is straightforward: On input $x$, it computes $t = t_1(|x|)$, determines the machine $M = \mu(x)$ that corresponds to $x$ (outputting a default value of no such machine exists), emulates $M(x)$ for $t$ steps, and returns the value $1 - M'(x)$. The question is how much time is required for this emulation. We should bear in mind that the time complexity of our algorithm needs to be analyzed in the two-tape Turing machine model, whereas $M$ itself is a two-tape Turing machine. We start by implementing our algorithm on a three-tape Turing-machine, and next emulate this machine on a two-tape Turing-machine.

The obvious implementation of our algorithm on a three-tape Turing-machine uses two tapes for the emulation itself and the third tape for the emulation procedure. Thus, each step of the the two-tape machine $M$ is emulated using $O(|\langle M \rangle|)$

steps (on the three-tape machine).[2] This includes also the amortized complexity of maintaining a step-counter for the emulation (see Exercise 4.4). Next, we need to emulate the foregoing three-tape machine on a two-tape machine. This is done by using the fact (cf., e.g., [117, Thm. 12.6]) that $t'$ steps of a three-tape machine can be emulated on a two-tape machine in $O(t' \log t')$ steps. Thus, the complexity of computing $f$ on input $x$ is upper-bounded by $O(T_{\mu(x)}(|x|) \log T_{\mu(x)}(|x|))$, where $T_M(n) = O(|\langle M \rangle| \cdot t_1(n))$ denotes the cost of emulating $t_1(n)$ steps of the two-tape machine $M$ on a three-tape machine (as in the foregoing discussion).

It turns out that the quality of the result we obtain depends on the mapping $\mu$ of inputs to machines. Using the naive (identity) mapping (i.e., $\mu(x) = x$) we can only establish the theorem for $t_2(n) = \omega((n \cdot t_1(n)) \cdot \log(n \cdot t_1(n)))$, because in this case $T_{\mu(x)}(|x|) = O(|x| \cdot t_1(|x|))$. (Note that in this case $x_M = \langle M \rangle$ is a description of $M$.) The theorem follows by associating with machine $M$ the input $x_M = 0^m 1 \langle M \rangle$, where $m = 2^{|\langle M \rangle|}$; that is, we may use the mapping $\mu$ such that $\mu(x) = M$ if $x = 0^{2^{|\langle M \rangle|}} 1 \langle M \rangle$ and $\mu(x)$ equals some fixed machine otherwise. In this case $|\mu(x)| < \log_2 |x| < \log t_1(|x|)$ and so $T_{\mu(x)}(|x|) = O((\log t_1(|x|)) \cdot t_1(|x|))$.
∎

---

**Teaching note:** Proving the standard version of Theorem 4.3 cannot be done by associating a sufficiently long input $x_M$ with each machine $M$, because this does not allow to get rid from an additional unbounded factor in $T_{\mu(x)}(|x|)$ (i.e., the $|\mu(x)|$ factor that multiplies $t_1(|x|)$). Note that the latter factor needs to be computable (at the very least) and thus cannot be accounted for by the generic $\omega$-notation that appears in the standard version (cf. [117, Thm. 12.9]). Instead, a different approach is taken (see Footnote 3).

---

**Technical Comments.** The proof of Theorem 4.3 associates with each potential machine an input and makes this machine err on this input. The aforementioned association is rather flexible: it should merely be efficiently computed (in the direction from the input to a possible machine) and should be sufficiently shrinking (in that direction). Specifically, we used the mapping $\mu$ such that $\mu(x) = M$ if $x = 0^{2^{|\langle M \rangle|}} 1 \langle M \rangle$ and $\mu(x)$ equals some fixed machine otherwise. We comment that each machine can be made to err on infinitely many inputs by redefining $\mu$ such that $\mu(x) = M$ if $0^{2^{|\langle M \rangle|}} 1 \langle M \rangle$ is a prefix of $x$ (and $\mu(x)$ equals some fixed machine otherwise). We also comment that, in contrast to the proof of Theorem 4.3, the

---

[2] This overhead accounts both for searching the code of $M$ for the adequate action and for the effecting of this action (which may refer to a larger alphabet than the one used by the emulator).

[3] In the standard proof the function $f$ is not defined with reference to $t_1(|x_M|)$ steps of $M(x_M)$, but rather with reference to the result of emulating $M(x_M)$ while using a total of $t_2(|x_M|)$ steps in the emulation process (i.e., in the algorithm used to compute $f$). This guarantees that $f$ is in DTIME$(t_2)$, and "pushes the problem" to showing that $f$ is not in DTIME$(t_1)$. It also explains why $t_2$ (rather than $t_1$) is assumed to be time constructible. As for the foregoing problem, it is resolved by observing that for each relevant machine (i.e., having time complexity $t_1$) the executions on any sufficiently long input will be fully emulated. Thus, we merely need to associate with each $M$ a disjoint set of infinitely many inputs and make sure that $M$ errs on each of these inputs.

proof of Theorem 1.5 utilizes a rigid mapping of inputs to machines (i.e., there $\mu(x) = M$ if $x = \langle M \rangle$).

**Digest: Diagonalization.** The last comment highlights the fact that the proof of Theorem 4.3 is merely a sophisticated version of the proof of Theorem 1.5. Both proofs refer to versions of the universal function, which in the case of the proof of Theorem 4.3 is (implicitly) defined such that its value at $(\langle M \rangle, x)$ equals $M'(x)$, where $M'(x)$ denotes an emulation of $M(x)$ that is suspended after $t_1(|x|)$ steps.[4] Actually, both proofs refers to the "diagonal" of the aforementioned function, denoted $d$, which in the case of the proof of Theorem 4.3 is only defined implicitly. Denoting the former function by $U$, the diagonal function is defined such that $d(x) = U(\langle \mu(x) \rangle, x)$. This is actually a definitional schema, as the choice of the function $\mu$ remains unspecified. Indeed, setting $\mu(x) = x$ corresponds to a "real" diagonal in the matrix depicting $U$, but other settings of $\mu$ as a 1-1 function can be viewed as "kind of diagonals" too. Either way, $f$ is defined such that $f(x) \neq d(x)$ on every $x$. This guarantees that no machine of time-complexity $t_1$ can compute $f$, and the focus is on presenting an algorithm that computes $f$ (which, needless to say, has time-complexity greater than $t_1$). Part of the proof of Theorem 4.3 is devoted to selecting $\mu$ in a way that minimizes the time-complexity of computing $f$, whereas in the proof of Theorem 1.5 we merely need to guarantee that $f$ is computable.

#### 4.2.1.2 Impossibility of speed-up for universal computation

The Time Hierarchy Theorem (Theorem 4.3) implies that the computation of a universal machine cannot be significantly sped-up. That is, consider the function $\mathbf{u}'(\langle M \rangle, x, t) \overset{\text{def}}{=} y$ if on input $x$ machine $M$ halts within $t$ steps and outputs the string $y$, and $\mathbf{u}'(\langle M \rangle, x, t) \overset{\text{def}}{=} \bot$ if on input $x$ machine $M$ makes more than $t$ steps. Recall that the value of $\mathbf{u}'(\langle M \rangle, x, t)$ can be computed in $\widetilde{O}(|x| + |\langle M \rangle| \cdot t)$ steps. Theorem 4.3 implies that this value cannot be computed with significantly less steps.

**Theorem 4.5** *There exists no two-tape Turing machine that, on input $\langle M \rangle, x$ and $t$, computes $\mathbf{u}'(\langle M \rangle, x, t)$ in $o((t + |x|) \cdot f(M)/\log^2(t + |x|))$ steps, where $f$ is an arbitrary function.*

A similar result holds for any reasonable and general model of computation (cf., Corollary 4.4). In particular, it follows that $\mathbf{u}'$ is not computable in polynomial time (because the input $t$ is presented in binary). In fact, one can show that the set $\{(\langle M \rangle, x, t) : \mathbf{u}'(\langle M \rangle, x, t) \neq \bot\}$ is not in $\mathcal{P}$; see Exercise 4.5.

**Proof:** Suppose (towards the contradiction) that, for every fixed $M$, given $x$ and $t > |x|$, the value of $\mathbf{u}'(\langle M \rangle, x, t)$ can be computed in $o(t/\log^2 t)$ steps, where the o-notation hides a constant that may depend on $M$. Consider an arbitrary time constructible $t_1$ (s.t. $t_1(n) > n$) and an arbitrary set $S \in \text{DTIME}(t_2)$, where

---
[4] Needless to say, in the proof of Theorem 1.5, $M' = M$.

$t_2(n) = t_1(n) \cdot \log^2 t_1(n)$. Let $M$ be a machine of time complexity $t_2$ that decides membership in $S$, and consider an algorithm that, on input $x$, first computes $t = t_1(|x|)$, and then computes (and outputs) the value $\mathbf{u}'(\langle M \rangle, x, t \log^2 t)$. By the time constructibility of $t_1$, the first computation can be implemented in $t$ steps, and by the contradiction hypothesis the same holds for the second computation. Thus, $S$ can be decided in $\mathrm{DTIME}(2t_1) = \mathrm{DTIME}(t_1)$, implying that $\mathrm{DTIME}(t_2) = \mathrm{DTIME}(t_1)$, which in turn contradicts Theorem 4.3.   ■

### 4.2.1.3   Hierarchy theorem for non-deterministic time

Analogously to $\mathrm{DTIME}$, for a fixed model of computation (to be understood from the context) and for any function $t : \mathbb{N} \rightarrow \mathbb{N}$, we *denote by* $\mathrm{NTIME}(t)$ *the class of sets that are accepted by some non-deterministic machine of time complexity $t$.* Alternatively, analogously to the definition of $\mathcal{NP}$, a set $S \subseteq \{0,1\}^*$ is in $\mathrm{NTIME}(t)$ if there exists a *linear-time* algorithm $V$ such that the two conditions hold

1. For every $x \in S$ there exists $y \in \{0,1\}^{t(|x|)}$ such that $V(x, y) = 1$.

2. For every $x \notin S$ and every $y \in \{0,1\}^*$ it holds that $V(x, y) = 0$.

We warn that the two formulations are not identical, but in sufficiently strong models (e.g., two-tape Turing machines) they are related up to logarithmic factors (see Exercise 4.6). The hierarchy theorem itself is similar to the one for deterministic time, except that here we require that $t_2(n) \geq (\log t_1(n+1))^2 \cdot t_1(n+1)$ (rather than $t_2(n) \geq (\log t_1(n))^2 \cdot t_1(n)$). That is:

**Theorem 4.6** (non-deterministic time hierarchy for two-tape Turing machines): *For any time-constructible and monotonicly non-decreasing function $t_1$ and every function $t_2$ such that $t_2(n) \geq (\log t_1(n+1))^2 \cdot t_1(n+1)$ and $t_1(n) > n$ it holds that $\mathrm{NTIME}(t_1)$ is strictly contained in $\mathrm{NTIME}(t_2)$.*

**Proof:** We cannot just apply the proof of Theorem 4.3, because the Boolean function $f$ defined there requires the ability to determine whether $M$ accepts the input $x_M$ in $t_1(|x_M|)$ steps. In the current context, $M$ is a non-deterministic machine and so the only way we know how to determine this question (both for a yes and no answers) is to try all the $(2^{t_1(|x_M|)})$ relevant executions. But this would put $f$ in $\mathrm{DTIME}(2^{t_1})$, rather than in $\mathrm{NTIME}(\widetilde{O}(t_1))$, and so a different approach is needed.

    We associate with each machine $M$, a large interval of strings (viewed as integers), denoted $I_M = [\alpha_M, \beta_M]$, such that the various intervals do not intersect and such that it is easy to determine for each string $x$ in which interval it resides. For each $x \in [\alpha_M, \beta_M - 1]$, we define $f(x) = 1$ if and only if there exists a non-deterministic computation of $M$ that accepts the input $x' \stackrel{\text{def}}{=} x + 1$ in $t_1(|x'|) \leq t_1(|x| + 1)$ steps. Thus, unless either $M$ accepts each string in the interval $I_M$ or rejects each such string, it (i.e., $M$) fails to accept $\{x : f(x) = 1\}$. So it is left to deal with the case that $M$ is invariant on $I_M$, which is where the definition of the value of $f(\beta_M)$ comes into play: We define $f(\beta_M)$ to equal *zero* if and only

if there exists a non-deterministic computation of $M$ that accepts the input $\alpha_M$ in $t_1(|\alpha_M|)$ steps. We shall select $\beta_M$ to be large enough relative to $\alpha_M$ such that we can afford to try all possible computations of $M$ on input $\alpha_M$. Details follow.

We present the following non-deterministic machine for accepting the set $\{x : f(x) = 1\}$. We assume that on input $x$ it is easy to determine the machine $M$ as well as the interval $n[\alpha_M, \beta_M]$ in which $x$ reside. On input $x \in [\alpha_M, \beta_M - 1]$, this non-deterministic machine emulates a (single) non-deterministic computation of $M$ on input $x' = x + 1$, and decides accordingly. Indeed, this emulation can be performed in time $(\log t_1(|x + 1|))^2 \cdot t_1(|x + 1|) \leq t_2(|x|)$. On input $x = \beta_M$, our machine just tries all $2^{t_1(|\alpha_M|)}$ executions of $M$ on input $\alpha_M$ and decides in a suitable manner; that is, our machine emulates all $2^{t_1(|\alpha_M|)}$ possible executions of $M(\alpha_M)$ and accepts $\alpha_M$ if and only if all the emulated executions ended rejecting $\alpha_M$. Note that this part of the emulation is deterministic, and it amounts to emulating $T_M \overset{\text{def}}{=} 2^{t_1(|\alpha_M|)} \cdot t_1(|\alpha_M|)$ steps of $M$. By a suitable choice of the interval $[\alpha_M, \beta_M]$, this number (i.e., $T_M$) is smaller than $t_1(|\beta_M|)$ (e.g., $|\beta_M| \geq T_M$ implies $T_M \leq t_1(|\beta_M|)$), and it follows that $T_M$ steps of $M$ can be emulated in time $(\log_2 t_1(|\beta_M|))^2 \cdot t_1(|\beta_M|) \leq t_2(|\beta_M|)$. Thus, $f$ is in $\mathrm{NTIME}(t_2)$.

Finally, we show that defining $f$ as in the foregoing indeed guarantees that it is not in $\mathrm{NTIME}(t_1)$. Suppose on the contrary, that some non-deterministic machine $M$ of time complexity $t_1$ accepts the set $\{x : f(x) = 1\}$. We define a Boolean function $A_M$ such that $A_M(x) = 1$ if and only if there exists a non-deterministic computation of $M$ that accepts the input $x$, and note that by the contradiction hypothesis $A_M(x) = f(x)$. Focusing on the interval $[\alpha_M, \beta_M]$, we have $A_M(x) = f(x)$ for every $x \in [\alpha_M, \beta_M]$, which (combined with the definition of $f$) implies that $A_M(x) = f(x) = A_M(x + 1)$ for every $x \in [\alpha_M, \beta_M - 1]$ and $A_M(\beta_M) = f(\beta_M) = 1 - A_M(\alpha_M)$. Thus, we reached a contraction (because we got $A_M(\alpha_M) = \cdots = A_M(\beta_M) = 1 - A_M(\alpha_M)$). ■

## 4.2.2 Time Gaps and Speed-Up

In contrast to Theorem 4.3, there exists functions $t : \mathbb{N} \to \mathbb{N}$ such that $\mathrm{DTIME}(t) = \mathrm{DTIME}(t^2)$ (or even $\mathrm{DTIME}(t) = \mathrm{DTIME}(2^t)$). Needless to say, these functions are not time-constructible (and thus the aforementioned fact does not contradict Theorem 4.3). The reason for this phenomenon is that, for such functions $t$, there exists not algorithms that have time complexity above $t$ but below $t^2$ (resp., $2^t$).

**Theorem 4.7** (the time gap theorem): *For every non-decreasing computable function $g : \mathbb{N} \to \mathbb{N}$ there exists a non-decreasing computable function $t : \mathbb{N} \to \mathbb{N}$ such that $\mathrm{DTIME}(t) = \mathrm{DTIME}(g(t))$.*

The forgoing examples referred to $g(m) = m^2$ and $g(m) = 2^m$. Since we are mainly interested in dramatic gaps (i.e., super-polynomial functions $g$), the model of computation does not matter here (as long as it is reasonable and general).

**Proof:** Consider an enumeration of all possible algorithms (or machines), which also includes machines that do not halt on some inputs. (Recall that we cannot

enumerate only all machines that halt on every input.)  Let $t_i$ denote the time complexity of the $i^{\text{th}}$ algorithm; that is, $t_i(n) = \infty$ if the $i^{\text{th}}$ machine does not halt on some $n$-bit long input and otherwise $t_i(n) = \max_{x \in \{0,1\}^n} \{T_i(x)\}$, where $T_i(x)$ denotes the number of steps taken by the $i^{\text{th}}$ machine on input $x$.

The basic idea is to define $t$ such that no $t_i$ is "sandwiched" between $t$ and $g(t)$, and thus no algorithm will have time complexity between $t$ and $g(t)$. Intuitively, if $t_i(n)$ is finite, then we may define $t$ such that $t(n) > t_i(n)$ and thus guarantee that $t_i(n) \notin [t(n), g(t(n))]$, whereas if $t_i(n) = \infty$ then any finite value of $t(n)$ will do (because then $t_i(n) > g(t(n))$). Thus, for every $m$ and $n$, we can define $t(n)$ such that $t_i(n) \notin [t(n), g(t(n))]$ for every $i \in [m]$ (e.g., $t(n) = \max_{i \in [m]:t_i(n) \neq \infty} \{t_i(n)\} + 1$).[5] This yields a weaker version of the theorem, in which the function $t$ is not computable.

The problem is that we want $t$ to be computable, whereas given $n$ we cannot tell whether or not $t_i(n)$ is finite. However, we do not really need to make the latter decision: for each candidate value $v$ of $t(n)$, we should just determine whether or not $t_i(n) \in [v, g(v)]$, which can be decided by running the $i^{\text{th}}$ machine for at most $g(v) + 1$ steps (on each $n$-bit long string). That is, as far as the $i^{\text{th}}$ machine is concerned, we should just find a value $v$ such that either $v > t_i(n)$ or $g(v) < t_i(n)$ (which includes the case $t_i(n) = \infty$). This can be done by starting with $v = v_0$ (where, say, $v_0 = n + 1$), and increasing $v$ until either $v > t_i(n)$ or $g(v) < t_i(n)$. The point is that if $t_i(n)$ is finite then we output $v = t_i(n) + 1$ after performing $\sum_{j=v_0}^{t_i(n)} 2^n \cdot j$ emulation steps and otherwise we output $v = v_0$ after emulating $2^n \cdot (g(v_0) + 1)$ steps. Bearing in mind that we should deal with all possible machines, we obtain the following procedure for setting $t(n)$.

Let $\mu : \mathbb{N} \to \mathbb{N}$ be any unbounded and computable function (e.g., $\mu(n) = n$ will do). Starting with $v = n + 1$, we keep incrementing $v$ until $v$ satisfies, for every $i \in \{1, ..., \mu(n)\}$, either $t_i(n) < v$ or $t_i(n) > g(v)$. This condition can be verified by computing $\mu(n)$ and $g(v)$, and emulating the execution of of the $\mu(n)$ machines on each of the $n$-bit long strings for $g(v)$ steps. The procedure sets $t(n)$ to equal the first value $v$ satisfying the aforementioned condition, and halts. To show that the procedure halts on every $n$, consider the set $H_n \subseteq \{1, ..., \mu(n)\}$ of indices of the relevant machines that halt on all inputs of length $n$. Then the procedure definitely halts before reaching the value $v = T_n + 2$, where $T_n = \max_{i \in H_n} \{t_i(n)\}$. (Indeed, the procedure may halt with a value $v \leq T_n$, but this will happen only if $g(v) < T_n$.)

For the foregoing function $t$, we claim that $\text{DTIME}(t) = \text{DTIME}(g(t))$. Indeed, let $S \in \text{DTIME}(g(t))$ and suppose that the $i^{\text{th}}$ algorithm decides $S$ in time at most $g(t)$; that is, for every $n$, it holds that $t_i(n) \leq g(t(n))$. Then (by the construction), for every $n$ satisfying $\mu(n) \geq i$, it holds that $t_i(n) < t(n)$, and it follows that the $i^{\text{th}}$ algorithm decides $S$ in time at most $t$ on all but finitely many inputs. Combining this algorithm with a "look-up table" machine that handles the exceptional inputs, the theorem follows.   ∎

---

[5]We may assume, without loss of generality, that $t_1(n) = 1$ for every $n$; e.g., by letting the machine that always halts after a single step be the first machine in our enumeration.

**Comment:** The function $t$ defined by the foregoing proof is computable in time that exceeds $g(t)$. Specifically, the presented procedure computes $t(n)$ (as well as $g(f(n))$) in time $\widetilde{O}(2^n \cdot g(t(n)) + T_g(t(n)))$, where $T_g(m)$ denotes the number of steps required to compute $g(m)$ on input $m$.

**Speed-up Theorems.** Theorem 4.7 can be viewed as asserting that some time complexity classes (i.e., $\mathrm{DTIME}(g(t))$ in the theorem) collapse to lower classes (i.e., to $\mathrm{DTIME}(t)$). A conceptually related phenomenon is of problems that have no optimal algorithm (not even in a very mild sense); that is, every algorithm for these ("pathological") problems can be drastically sped-up. It follows that the complexity of these problems can not be defined (i.e., as the complexity of the best algorithm solving this problem). The following drastic speed-up theorem should not be confused with the linear speed-up that is an artifact of the definition of a Turing machine (see Exercise 4.7).[6]

**Theorem 4.8** (the time speed-up theorem): *For every computable* (and super-linear) *function $g$ there exists a decidable set $S$ such that if $S \in \mathrm{DTIME}(t)$ then $S \in \mathrm{DTIME}(t')$ for $t'$ satisfying $g(t'(n)) < t(n)$.*

Taking $g(n) = n^2$ (or $g(n) = 2^n$), the theorem asserts that, for every $t$, if $S \in \mathrm{DTIME}(t)$ then $S \in \mathrm{DTIME}(\sqrt{t})$ (resp., $S \in \mathrm{DTIME}(\log t)$). Note that Theorem 4.8 can be applied any (constant) number of times, which means that we cannot give a reasonable estimate to the complexity of deciding membership in $S$. In contrast, recall that in some important cases, optimal algorithms for solving computational problems do exist. Specifically, algorithms solving (candid) search problems in NP cannot be speed-up (see Theorem 2.31), nor can the computation of a universal machine (see Theorem 4.5).

We refrain from presenting a proof of Theorem 4.8, but comment on the complexity of the sets involved in this proof. The proof (presented in [117, Sec. 12.6]) provides a construction of a set $S$ in $\mathrm{DTIME}(t') \setminus \mathrm{DTIME}(t'')$ for $t'(n) = h(n - O(1))$ and $t''(n) = h(n - \omega(1))$, where $h(n)$ denoted $g$ iterated $n$ times on 2 (i.e., $h(n) = g^{(n)}(2)$, where $g^{(i+1)}(m) = g(g^{(i)}(m))$ and $g^{(1)} = g$). The set $S$ is constructed such that for every $i > 0$ there exists a $j > i$ and an algorithm that decides $S$ in time $t_i$ but not in time $t_j$, where $t_k(n) = h(n - k)$.

# 4.3 Space Hierarchies and Gaps

Hierarchy and Gap Theorems analogous to Theorem 4.3 and Theorem 4.7, respectively, are known for space complexity. In fact, since space-efficient emulation of space-bounded machines is simpler than time-efficient emulations of time-bounded machines, the results tend to be sharper. This is most conspicuous in the case of

---

[6]We note that the linear speed-up phenomenon was implicitly addressed in the proof of Theorem 4.3, by allowing an emulation overhead that depends on the length of the description of the emulated machine.

the separation result (stated next), which is optimal (in light of linear speed-up results; see Exercise 4.7).

Before stating the result, we need a few preliminaries. We refer the reader to §1.2.3.4 for a definition of space complexity (and to Chapter 5 for further discussion). As in case of time complexity, we consider a specific model of computation, but the results hold for any other reasonable and general model. Specifically, we consider three-tape Turing machines, because we designate two special tapes for input and output. For any function $s : \mathbb{N} \to \mathbb{N}$, we *denote by* DSPACE($s$) *the class of decision problems that are solvable in space complexity $s$*. Analogously to Definition 4.2, we call a function $s : \mathbb{N} \to \mathbb{N}$ space constructible if there exists an algorithm that on input $n$ outputs $s(n)$ using at most $s(n)$ cells of the work-tape. Actually, functions like $s_1(n) = \log n$, $s_2(n) = (\log n)^2$, and $s_3(n) = 2^n$ are computable using $\log s_i(n)$ space.

**Theorem 4.9** (space hierarchy for three-tape Turing machines): *For any space constructible function $s_2$ and every function $s_1$ such that $s_2 = \omega(s_1)$ and $s_1(n) > \log n$ it holds that* DSPACE($s_1$) *is strictly contained in* DSPACE($s_2$).

Theorem 4.9 is analogous to the traditional version of Theorem 4.3 (rather to the one we presented), and is proven using the alternative approach sketched in Footnote 3. The details are left as an exercise (see Exercise 4.9).

# Chapter Notes

The material presented in this chapter predates the theory of NP-completeness and the dominant stature of the P-vs-NP Question. At these early days, the field (to be known as complexity theory) did not yet develop an independent identity and its perspectives were dominated by two classical theories: the theory of computability (and recursive function) and the theory of formal languages. Nevertheless, we believe that the results presented in this chapter are interesting for two reasons. Firstly, as stated up-front, these results address the natural question of under what conditions is it the case that more resources help. Secondly, the inapplicability of these results to questions of the type of P-vs-NP illustrates the non-generic flavor of the latter questions as referring to specific (and natural) aspects of the "complexity of computation".

The hierarchy theorems (e.g., Theorem 4.3) were proved by Hartmanis and Stearns [109]. Gap theorems (e.g., Theorem 4.7, often referred to as Borodin's Gap Theorem) were proven by Borodin [43]. A axiomatic treatment of complexity measures and corresponding speed-up theorems (e.g., Theorem 4.8, often referred to as Blum's Speed-up Theorem) are due to Blum [35].

# Exercises

**Exercise 4.1** Let $F_n(s)$ denote the number of different Boolean functions over $\{0,1\}^n$ that are computed by Boolean circuits of size $s$. Prove that, for any $s < 2^n$, it holds that $F_n(s) \geq 2^{s/O(\log s)}$ and $F_n(s) \leq s^{2s}$.

**Guideline:** Any Boolean function $f : \{0,1\}^\ell \to \{0,1\}$ can be computed by a circuit of size $s_\ell = O(\ell \cdot 2^\ell)$. Thus, for every $\ell \leq n$, it holds that $F_n(s_\ell) \geq 2^{2^\ell} > 2^{s_\ell/O(\log s_\ell)}$. On the other hand, the number of circuits of size $s$ is less than $2^s \cdot \binom{s^2}{s}$, where the second factor represents the number of possible choices of pair of gates that feed any gate in the circuit.

**Exercise 4.2 (advice can speed-up computation)** For every time constructible function $t$, show that there exists a set $S$ in $\mathrm{D{\small TIME}}(t^2) \setminus \mathrm{D{\small TIME}}(t)$ that can be decided in linear-time using an advice of linear length (i.e., $S \in \mathrm{D{\small TIME}}(\ell)/\ell$ where $\ell(n) = O(n)$).

**Guideline:** Starting with a set $S' \in \mathrm{D{\small TIME}}(T^2) \setminus \mathrm{D{\small TIME}}(T)$, where $T(m) = t(2^m)$, consider the set $S = \{x0^{2^{|x|}-|x|} : x \in S'\}$.

**Exercise 4.3** Referring to a reasonable model of computation (and assuming that the input length is not given explicitly (e.g., as in Definition 10.10)), prove that any algorithm that has sub-linear time-complexity actually has constant time-complexity.

**Guideline:** Consider the question of whether or not there exists an infinite set of strings $S$ such that when invoked on any input $x \in S$ the algorithm reads all of $x$. Note that if $S$ is infinite then the algorithm cannot have sub-linear time-complexity, and prove that if $S$ is finite then the algorithm has constant time-complexity.

**Exercise 4.4 (constant amortized time step-counter)** A step-counter is an algorithm that runs for a number of steps that is specified in its input. Actually, such an algorithm may run for a somewhat larger number of steps but halt after issuing a number of "signals" as specified in its input, where these signals are defined as entering (and leaving) a designated state (of the algorithm). A step-counter may be run in parallel to another procedure in order to suspend the execution after a desired number of steps (of the other procedure) has elapsed. Show that there exists a simple deterministic machine that, on input $n$, halts after issuing $n$ signals while making $O(n)$ steps.

**Guideline:** A slightly careful implementation of the straightforward algorithm will do, when coupled with an "amortized" time complexity.

**Exercise 4.5 (a natural set in $\mathcal{E} \setminus \mathcal{P}$)** In continuation to the proof of Theorem 4.5, prove that the set $\{(\langle M \rangle, x, t) : \mathtt{u}'(\langle M \rangle, x, t) \neq \perp\}$ is in $\mathcal{E} \setminus \mathcal{P}$, where $\mathcal{E} \overset{\text{def}}{=} \cup_c \mathrm{D{\small TIME}}(e_c)$ and $e_c(n) = 2^{cn}$.

**Exercise 4.6** Prove that the two definitions of N{\small TIME}, presented in §4.2.1.3, are related up to logarithmic factors. Note the importance of condition that $V$ has linear (rather than polynomial) time-complexity.

**Guideline:** When emulating a non-deterministic machine by the verification procedure $V$, encode the non-deterministic choices in $y$ such that $|y|$ is slightly larger than the number of steps taken by the original machine. Specifically, having $|y| = O(t \log t)$, where $t$ denotes the number of steps taken by the original machine, allows to emulate the latter in linear time (i.e., linear in $|y|$).

**Exercise 4.7 (linear speed-up of Turing machine)** Prove that any problem that can be solved by a two-tape Turing machine that has time-complexity $t$ can be solved by another two-tape Turing machine having time-complexity $t'$, where $t'(n) = O(n) + (t(n)/2)$.

**Guideline:** Consider a machine that uses a larger alphabet, capable of encoding a constant (denoted $c$) number of symbols of the original machine, and thus capable of emulating $c$ steps of the original machine in $O(1)$ steps, where the constant in the notation $O(1)$ is a universal constant (independent of $c$). Note that the $O(n)$ term accounts to a preprocessing required to encode the input in the work-alphabet of the new machine, and that a similar result for one-tape Turing machine seems to require a $O(n^2)$ term.

**Exercise 4.8** In continuation to Exercise 4.7, state and prove an analogous result for space complexity, when using the standard definition of space as recalled in Section 4.3. (Note that this result does not hold with respect to "binary space complexity" as defined in Section 5.1.)

**Exercise 4.9** Prove Theorem 4.9. As a warm-up, assume that $s_1$ (rather than $s_2$) is space constructible.

**Guideline:** Note that providing a space-efficient emulation of one machine by another machine is easier than providing an analogous time-efficient emulation.

# Appendix E

# Explicit Constructions

> *It is easier for a camel to go through the eye of a needle, than for a rich man to enter into the kingdom of God.*
>
> Matthew, 19:24.

Complexity theory provides a clear definition of the intuitive notion of an explicit construction. Furthermore, it also suggests a hierarchy of different levels of explicitness, referring to the ease of constructing the said object. The basic levels of explicitness are provided by considering the complexity of fully constructing the object (e.g., the time it takes to print the truth-table of a finite function). In this context, explicitness often means outputting a full description of the object in time that is polynomial in the length of that description. Stronger levels of explicitness emerge when considering the complexity of answering natural queries regarding the object (e.g., the time it takes to evaluate a fixed function at a given input). In this context, (strong) explicitness often means answering such queries in polynomial-time. The aforementioned themes are demonstrated in our brief overview of explicit constructions of *error correcting codes* and *expander graphs*. These constructions are, in turn, used in various parts of the main text.

**Summary:** We review several popular constructions of error correcting codes, culminating with the construction of a concatenated code that combines a Reed-Solomon code with a "mildly explicit" construction of a small code. We also review briefly the notions of locally testable and locally decodable codes, and a useful "list decoding bound" (i.e., bounding the number of codewords that are close to any single sequence).

We review the two standard definitions of expanders, two levels of explicitness, and two properties of expanders that are related to (single-step and multi-step) random walks on them. We then review two explicit constructions of expander graphs.

# E.1    Error Correcting Codes

In this section we highlight some issues and aspects of coding theory that are most relevant to the current book. The interested reader is referred to [205] for a more comprehensive treatment of the computational aspects of coding theory. Structural aspects of coding theory, which are at the traditional focus of that field, are covered in standard textbook such as [153].

Loosely speaking, an error correcting code is a mapping of strings to longer strings such that any two different strings are mapped to a corresponding pair of strings that are far apart (and not merely different). Specifically, $C : \{0,1\}^k \to \{0,1\}^n$ is a (binary) code of distance $d$ if for every $x \neq y \in \{0,1\}^k$ it holds that $C(x)$ and $C(y)$ differ on at least $d$ bit positions.

It will be useful to extend this definition to sequences over an arbitrary alphabet $\Sigma$, and to use some notations. Specifically, for $x \in \Sigma^m$, we denote the $i^{\text{th}}$ symbol of $x$ by $x_i$ (i.e., $x = x_1 \cdots x_m$), and consider codes over $\Sigma$ (i.e., mappings of $\Sigma$-sequences to $\Sigma$-sequences). The mapping (code) $C : \Sigma^k \to \Sigma^n$ has distance $d$ if for every $x \neq y \in \Sigma^k$ it holds that $|\{i : C(x)_i \neq C(y)_i\}| \geq d$. The members of $\{C(x) : x \in \Sigma^k\}$ are called codewords (and in some texts this set itself is called a code).

In general, we define a metric, called Hamming distance, over the set of $n$-long sequences over $\Sigma$. The Hamming distance between $y$ and $z$, where $y, z \in \Sigma^n$, is defined as the number of locations on which they disagree (i.e., $|\{i : y_i \neq z_i\}|$). The Hamming weight of such sequences is defined as the number of non-zero elements (assuming that one element of $\Sigma$ is viewed as zero). Typically, $\Sigma$ is associated with an additive group, and in this case the distance between $y$ and $z$ equals the Hamming weight of $w = y - z$, where $w_i = y_i - z_i$ (for every $i$).

**Asymptotics.**   We will actually consider infinite families of codes; that is, $\{C_k : \Sigma_k^k \to \Sigma_k^{n(k)}\}_{k \in S}$, where $S \subseteq \mathbb{N}$ (and typically $S = \mathbb{N}$). (N.B., we allow $\Sigma_k$ to depend on $k$.) We say that such a family has distance $d : \mathbb{N} \to \mathbb{N}$ if for every $k \in S$ it holds that $C_k$ has distance $d(k)$. Needless to say, both $n = n(k)$ (called the block-length) and $d(k)$ depend on $k$, and the aim is to have a linear dependence (i.e., $n(k) = O(k)$ and $d(k) = \Omega(n(k))$). In such a case, one talks of the relative rate of the code (i.e., the constant $k/n(k)$) and its relative distance (i.e., the constant $d(k)/n(k)$).

In general, we will often refer to *relative distances* between sequences. For example, for $y, z \in \Sigma^n$, we say that $y$ and $z$ are $\varepsilon$-close (resp., $\varepsilon$-far) if $|\{i : y_i \neq z_i\}| \leq \varepsilon \cdot n$ (resp., $|\{i : y_i \neq z_i\}| \geq \varepsilon \cdot n$).

**Explicitness.**   A mild notion of explicitness refers to constructing the list of all codewords in time that is polynomial in its length (which is exponential in $k$). A more standard notion of explicitness refers to generating a specific codeword (i.e., producing $C(x)$ when given $x$), which coincides with the encoding task mentioned next. Stronger notions of explicitness refer to other computational problems concerning codes (see next).

**Computational problems.** The most basic computational tasks associated with codes are encoding and decoding (under noise). The definition of the encoding task is straightforward (i.e., map $x \in \Sigma_k^k$ to $C_k(x)$), and an efficient algorithm is required to compute each symbol in $C_k(x)$ in $\mathrm{poly}(k, \log |\Sigma_k|)$-time.[1] When defining the decoding task we note that "minimum distance decoding" (i.e., given $w \in \Sigma_k^{n(k)}$, find $x$ such that $C_k(x)$ is closest to $y$ (in Hamming distance)) is just one natural possibility. Two related variants, regarding a code of distance $d$, are:

Unique decoding: Given $w \in \Sigma_k^{n(k)}$ that is at Hamming distance less than $d(k)/2$ from some codeword $C_k(x)$, retrieve the corresponding decoding of $C_k(x)$ (i.e., retrieve $x$).

Needless to say, this task is well-defined because there cannot be two different codewords that are each at Hamming distance less than $d(k)/2$ from $w$.

List decoding: Given $w \in \Sigma_k^{n(k)}$ and a parameter $d' \geq d(k)/2$, output a list of all $x \in \Sigma_k^k$ that are at Hamming distance at most $d'$ from $w$.

Typically, one considers the case that $d' < d(k)$. See Section E.1.3 for discussion of upper-bounds on the number of codewords that are within a certain distance from a generic sequence.

Two additional computational tasks are considered in Section E.1.2.

**Linear codes.** Associating $\Sigma_k$ with some finite field, we call a code $C_k : \Sigma_k^k \rightarrow \Sigma_k^{n(k)}$ linear if it satisfies $C_k(x + y) = C_k(x) + C_k(y)$, where $x$ and $y$ (resp., $C_k(x)$ and $C_k(y)$) are viewed as $k$-dimensional (resp., $n(k)$-dimensional) vectors over $\Sigma_k$, and the arithmetic is of the corresponding vector space. A useful property of linear codes is that their distance equals the Hamming weight of the lightest codeword other than $C_k(0^k)$; that is, $\min_{x \neq y}\{|\{i : C_k(x)_i \neq C_k(y)_i\}|\}$ equals $\min_{x \neq 0^k}\{|\{i : C_k(x)_i \neq 0\}|\}$. Another useful property is that the code is fully specified by a $k$-by-$n(k)$ matrix, called the generating matrix, that consists of the codewords of some fixed basis of $\Sigma_k^k$. That is, the set of all codewords is obtained by taking all $|\Sigma_k|^k$ different linear combination of the rows of the generating matrix.

## E.1.1 A few popular codes

Our focus will be on explicitly constructible codes; that is, (families of) codes of the form $\{C_k : \Sigma_k^k \rightarrow \Sigma_k^{n(k)}\}_{k \in S}$ that are coupled with efficient encoding and decoding algorithms. But before presenting a few such codes, let us consider a non-explicit construction.

**Proposition E.1** (random linear codes): *Let $c > 1$ and $n, d : \mathbb{N} \rightarrow \mathbb{N}$ be such that, for all sufficiently large $k$, it holds that $n(k) > \max(c \cdot k/(1 - H_2(d(k)/n(k))), 2d(k))$,*

---

[1]This formulation is not the one common in coding theory, but it is the most natural one for our applications. On one hand, this formulation is applicable also to codes with super-polynomial block-length. On the other hand, this formulation does not support a discussion of practical algorithms that compute the codeword faster than by computing each of its bits separately.

where $H_2(\alpha) \overset{\text{def}}{=} \alpha \log_2(1/\alpha) + (1 - \alpha) \log_2(1/(1 - \alpha))$. *Then, for all sufficiently large $k$, with high probability, a random linear transformation of $\{0,1\}^k$ to $\{0,1\}^{n(k)}$ constitutes a code of distance $d(k)$.*

Thus, *for every constant $\delta \in (0, 0.5)$ there exists a constant $\rho > 0$ and an infinite family of codes $\{C_k : \{0,1\}^k \to \{0,1\}^{k/\rho}\}_{k \in \mathbb{N}}$ of relative distance $\delta$.* Specifically, $\rho = (1 - H_2(\delta))/c$ will do.

**Proof:** We consider a uniformly selected $k$-by-$n(k)$ generating matrix over GF(2), and upper-bound the probability that it yields a linear code of distance less than $d(k)$. We use a union bound on all possible $2^k - 1$ linear combinations of the rows of the generating matrix, where for each such combination we compute the probability that it yields a vector of Hamming weight less than $d(k)$. Observe that the result of each such linear combination is uniformly distributed over $\{0,1\}^{n(k)}$, and thus has Hamming weight less than $d(k)$ with probability $\sum_{i=0}^{d(k)-1} \binom{n(k)}{i} \cdot 2^{-n(k)} < 2^{-(1-H_2(d(k)/n(k))) \cdot n(k)}$. Using $(1 - H_2(d(k)/n(k))) \cdot n(k) > c \cdot k$, the proposition follows. $\blacksquare$

### E.1.1.1  A mildly explicit version of Proposition E.1

Note that Proposition E.1 yields a (deterministic) $\exp(k \cdot n(k))$-time algorithm that finds a linear code of distance $d(k)$. The time bound can be improved to $\exp(k + n(k))$, by observing that we may choose the rows of the generating matrix one by one, making sure that all non-empty linear combinations of the current rows have weight at least $d(k)$. Note that the proof of Proposition E.1 can be adapted to assert that as long as we have less than $k$ rows a random choice of the next row will do with high probability. Note that in the case that $n(k) = O(k)$, this yields an algorithm that runs in time that is polynomial in the size of the code (i.e., the number of codewords). Needless to say, this mild level of explicitness is inadequate for most coding applications; however, it will be useful to us in §E.1.1.5.

### E.1.1.2  The Hadamard Code

The Hadamard code is the longest (non-repetitive) *linear* code over $\{0,1\} \equiv$ GF(2). That is, $x \in \{0,1\}^k$ is mapped to the sequence of all $n(k) = 2^k$ possible linear combinations of its bits (i.e., bit locations in the codewords are associated with $k$-bit strings, and location $\alpha \in \{0,1\}^k$ in the codeword of $x$ holds the value $\sum_{i=1}^{k} \alpha_i x_i$). It can be verified that each non-zero codeword has weight $2^{k-1}$, and thus this code has relative distance $d(k)/n(k) = 1/2$ (albeit its block-length $n(k)$ is exponential in $k$).

Turning to the computational aspects, we note that encoding is very easy. As for decoding, the warm-up discussion at the beginning of the proof of Theorem 7.7 provides a very fast probabilistic algorithm for unique decoding, whereas Theorem 7.8 provides a very fast probabilistic algorithm for list decoding.

We mention that the Hadamard code has played a key role in the proof of the PCP Theorem (Theorem 9.16); see §9.3.2.1.

**A propos long codes.** We note that the longest (non-repetitive) binary code (called the Long-Code and introduced in [26]) is extensively used in the design of "advanced" PCP systems (see, e.g., [111, 112]). In this code, a $k$-bit long string $x$ is mapped to the sequence of $n(k) = 2^{2^k}$ values, each corresponding to the evaluation of a different Boolean function at $x$; that is, bit locations in the codewords are associated with Boolean functions such that the location associated with $f : \{0, 1\}^k \to \{0, 1\}$ in the codeword of $x$ holds the value $f(x)$.

### E.1.1.3 The Reed–Solomon Code

A Reed-Solomon code is defined for a non-binary alphabet, which is associated with a finite field of $n$ elements, denoted $\mathrm{GF}(n)$. For any $k < n$, we consider the mapping of univariate degree $k - 1$ polynomials over $\mathrm{GF}(n)$ to their evaluation at all field elements. That is, $p \in \mathrm{GF}(n)^k$ (viewed as such a polynomial), is mapped to the sequence $(p(\alpha_1), ..., p(\alpha_n))$, where $\alpha_1, ..., \alpha_n$ is a canonical enumeration of the elements of $\mathrm{GF}(n)$.[2]

The Reed-Solomon code offers infinite families of codes with constant rate and constant relative distance (e.g., by taking $n(k) = 3k$ and $d(k) = 2k$), but the alphabet size grows with $k$ (or rather with $n(k) > k$). Efficient algorithms for unique decoding and list decoding are known (see [204] and references therein). These computational tasks correspond to the extrapolation of polynomials based on a noisy version of their values at all possible evaluation points.

### E.1.1.4 The Reed–Muller Code

Reed-Muller codes generalize Reed-Solomon codes by considering multi-variate polynomials rather than univariate polynomials. Consecutively, the alphabet may be any finite field, and in particular the two-element field $\mathrm{GF}(2)$. Reed-Muller codes (and variants of them) are extensively used in complexity theory; for example, they underly Construction 7.11 and the PCP constructed at the end of §9.3.2.2. The relevant property of these codes is that, under a suitable setting of parameters that satisfies $n(k) = \mathrm{poly}(k)$, they allow super fast "codeword testing" and "self-correction" (see discussion in Section E.1.2).

For any prime power $q$ and parameters $m$ and $r$, we consider the set, denoted $P_{m,r}$, of all $m$-variate polynomials of total degree at most $r$ over $\mathrm{GF}(q)$. Each polynomial in $P_{m,r}$ is represented by the $k = \log_q |P_{m,r}|$ coefficients of all relevant monomials, where in the case that $r < q$ it holds that $k = \binom{m+r}{m}$. We consider the code $C : \mathrm{GF}(q)^k \to \mathrm{GF}(q)^n$, where $n = q^m$, mapping $m$-variate polynomials of total degree at most $r$ to their values at all $q^m$ evaluation points. That is, the $m$-variate polynomial $p$ of total degree at most $r$ is mapped to the sequence of values $(p(\overline{\alpha}_1), ..., p(\overline{\alpha}_n))$, where $\overline{\alpha}_1, ..., \overline{\alpha}_n$ is a canonical enumeration of all the $m$-tuples of $\mathrm{GF}(q)$. The relative distance of this code is lower-bounded by $(q - r)/q$.

---

[2]Alternatively, we may map $(v_1, ..., v_k) \in \mathrm{GF}(n)^k$ to $(p(\alpha_1), ..., p(\alpha_n))$, where $p$ is the unique univariate polynomial of degree $k - 1$ that satisfies $p(\alpha_i) = v_i$ for $i = 1, ..., k$. Note that this modification amounts to a linear transformation of the generating matrix.

In typical applications one sets $r = \Theta(m^2 \log m)$ and $q = \text{poly}(r)$, which yields $k > m^m$ and $n = \text{poly}(r)^m = \text{poly}(m^m)$. Thus we have $n(k) = \text{poly}(k)$ but not $n(k) = O(k)$. As we shall see in Section E.1.2, the advantage (in comparison to the Reed-Solomon code) is that codeword testing and self-correction can be performed at complexity related to $q = \text{poly}(\log n)$. Actually, in most complexity applications, a variant in which only $m$-variate polynomials of individual degree $r' = r/m$ are used. In this case, an alternative presentation analogous to the one presented in Footnote 2 is preferred: The information is viewed as a function $f : H^m \to \text{GF}(q)$, where $H \subset \text{GF}(q)$ is of size $r' + 1$, and is encoded by the evaluation at all points in $\text{GF}(q)^m$ of the $m$-variate polynomial of individual degree $r'$ that extends the function $f$.

### E.1.1.5    Binary codes of constant relative distance and constant rate

Recall that we seek binary codes of constant relative distance and constant rate. Proposition E.1 asserts that such codes exists, but does not provide an explicit construction. The Hadamard code is explicit but does not have a constant rate (to say the least (since $n(k) = 2^k$)).[3] The Reed-Solomon code has constant relative distance and constant rate but uses a non-binary alphabet (which grows at least linearly with $k$). We achieve the desired construction by using the paradigm of concatenated codes [73], which is of independent interest. (Indeed, concatenated codes may be viewed as a simple version of the proof composition paradigm presented in §9.3.2.2.)

Intuitively, concatenated codes are obtained by first encoding information, viewed as a sequence over a large alphabet, by some code and next encoding each resulting symbol, which is viewed as a sequence of over a smaller alphabet, by a second code. Formally, consider $\Sigma_1 \equiv \Sigma_2^{k_2}$ and two codes, $C_1 : \Sigma_1^{k_1} \to \Sigma_1^{n_1}$ and $C_2 : \Sigma_2^{k_2} \to \Sigma_2^{n_2}$. Then, the concatenated code of $C_1$ and $C_2$, maps $(x_1, ..., x_{k_1}) \in \Sigma_1^{k_1} \equiv \Sigma_2^{k_1 k_2}$ to $(C_2(y_1), ..., C_2(y_{n_1}))$, where $(y_1, ..., y_{n_1}) = C_1(x_1, ..., x_{k_1})$.

Note that the resulting code $C : \Sigma_2^{k_1 k_2} \to \Sigma_2^{n_1 n_2}$ has constant rate and constant relative distance if both $C_1$ and $C_2$ have these properties. Encoding in the concatenated code is straightforward. To decode a corrupted codeword of $C$, we view the input as an $n_1$-long sequence of blocks, where each block is an $n_2$-long sequence over $\Sigma_2$. Applying the decoder of $C_2$ to each block, we obtain $n_1$ sequences (each of length $k_2$) over $\Sigma_2$, and interpret each such sequence as a symbol of $\Sigma_1$. Finally, we apply the decoder of $C_1$ to the resulting $n_1$-long sequence (over $\Sigma_1$), and interpret the resulting $k_1$-long sequence (over $\Sigma_1$) as a $k_1 k_2$-long sequence over $\Sigma_2$. The key observation is that *if $w \in \Sigma_2^{n_1 n_2}$ is $\varepsilon_1 \varepsilon_2$-close to $C(x_1, ..., x_{k_1}) = (C_2(y_1), ..., C_2(y_{n_1}))$ then at least $(1 - \varepsilon_1) \cdot n_1$ of the blocks of $w$ are $\varepsilon_2$-close to the corresponding $C_2(y_i)$.*[4]

We are going to consider the concatenated code obtained by using the Reed-

---

[3]Binary Reed-Muller codes also fail to simultaneously provide constant relative distance and constant rate.

[4]This observation offers unique decoding from a fraction of errors that is the product of the fractions (of error) associated with the two original codes. Stronger statements regarding unique decoding of the concatenated code can be made based on more refined analysis (cf. [73]).

Solomon Code $C_1 : \mathrm{GF}(n_1)^{k_1} \to \mathrm{GF}(n_1)^{n_1}$ as the large code, setting $k_2 = \log_2 n_1$, and using the mildly explicit version of Proposition E.1, $C_2 : \{0,1\}^{k_2} \to \{0,1\}^{n_2}$ as the small code. We use $n_1 = 3k_1$ and $n_2 = O(k_2)$, and so the concatenated code is $C : \{0,1\}^k \to \{0,1\}^n$, where $k = k_1 k_2$ and $n = n_1 n_2 = O(k)$. The key observation is that $C_2$ can be constructed in $\exp(k_2)$-time, whereas here $\exp(k_2) = \mathrm{poly}(k)$. Furthermore, both encoding and decoding with respect to $C_2$ can be performed in time $\exp(k_2) = \mathrm{poly}(k)$. Thus, we get:

**Theorem E.2** (an explicit good code): *There exists constants $\delta, \rho > 0$ and an explicit family of binary codes of rate $\rho$ and relative distance at least $\delta$. That is, there exists a polynomial-time* (encoding) *algorithm $C$ such that $|C(x)| = |x|/\rho$ (for every $x$) and a polynomial-time* (decoding) *algorithm $D$ such that for every $y$ that is $\delta/2$-close to some $C(x)$ it holds that $D(y) = x$. Furthermore, $C$ is a linear code.*

The linearity of $C$ is justified by using a Reed-Solomon code over the extension field $F = \mathrm{GF}(2^{k_2})$, and noting that this code induces a linear transformation over $\mathrm{GF}(2)$. Specifically, the value of a polynomial $p$ over $F$ at a point $\alpha \in F$ can be obtained as a linear transformation of the coefficient of $p$, when viewed as $k_2$-dimensional vectors over $\mathrm{GF}(2)$.

**Relative distance approaching one half.** Starting with a Reed-Solomon code of relative distance $\delta_1$ and a smaller code $C_2$ of relative distance $\delta_2$, we obtain a concatenated code of relative distance $\delta_1 \delta_2$. Note that, for any constant $\delta_1 < 1$, there exists a Reed-Solomon code $C_1 : \mathrm{GF}(n_1)^{k_1} \to \mathrm{GF}(n_1)^{n_1}$ of relative distance $\delta_1$ and constant rate (i.e., $1 - \delta_1$). Giving up on constant rate, we may start with a Reed-Solomon code of block-length $n_1(k_1) = \mathrm{poly}(k_1)$ and distance $n_1(k_1) - k_1$ over $[n_1(k_1)]$, and use a Hadamard code (encoding $[n_1(k_1)]$ by $\{0,1\}^{n_1(k_1)}$) in the role of the small code $C_2$. This yields a (concatenated) binary code of block length $n(k) = n_1(k)^2$ and distance $(n_1(k) - k) \cdot n_1(k)/2$. Thus, *the resulting explicit code has relative distance approximately* $(1/2) - (k/\sqrt{n(k)})$.

## E.1.2  Two additional computational problems

In this section we briefly review relaxations of two traditional coding theoretic tasks. The purpose of these relaxations is enabling super-fast (randomized) algorithms that provide meaningful information. Specifically, these algorithms may run in sub-linear (e.g., poly-logarithmic) time, and thus cannot possibly solve the unrelaxed version of the problem.

**Local testability.** This task refers to testing whether a given word is a codeword (in a predetermine code), based on (randomly) inspecting few locations in the word. Needless to say, we can only hope to make an approximately correct decision; that is, accept each codeword and reject with high probability each word that is *far* from the code. (Indeed, this task is within the framework of property testing; see Section 10.1.2.)

**Local decodability.** Here the task is to recover a specified bit in the plaintext by (randomly) inspecting few locations in a mildly corrupted codeword. This task is somewhat related to the task of self-correction (i.e., recovering a specified bit in the codeword itself, by inspecting few locations in the mildly corrupted codeword).

Note that the Hadamard code is both locally testable and locally decodable as well as self-correctable (based on a constant number of queries into the word); these facts were demonstrated and extensively used in §9.3.2.1. However, the Hadamard code has an exponential block-length (i.e., $n(k) = 2^k$), and the question is whether one can achieve analogous results with respect to a shorter code (e.g., $n(k) = \text{poly}(k)$). As hinted in §E.1.1.4, the answer is positive (when we refer to performing these operations in time that is poly-logarithmic in $k$):

**Theorem E.3** *For some constant $\delta > 0$ and polynomials $n, q : \mathbb{N} \to \mathbb{N}$, there exists an explicit family of codes $\{C_k : [q(k)]^k \to [q(k)]^{n(k)}\}_{k \in \mathbb{N}}$ of relative distance $\delta$ that can be locally testable and locally decodable in $\text{poly}(\log k)$-time. That is, the following three conditions hold.*

1. Encoding: *There exists a polynomial time algorithm that on input $x \in [q(k)]^k$ returns $C_k(x)$.*

2. Local Testing: *There exists a probabilistic polynomial-time oracle machine $T$ that given $k$ (in binary)[5] and oracle access to $w \in [q(k)]^{n(k)}$ distinguishes the case that $w$ is a codeword from the case that $w$ is $\delta/2$-far from any codeword. Specifically:*

   (a) *For every $x \in [q(k)]^k$ it holds that $\Pr[T^{C_k(x)}(k) = 1] = 1$.*
   (b) *For every $w \in [q(k)]^{n(k)}$ that is $\delta/2$-far from any codeword of $C_k$ it holds that $\Pr[T^w(k) = 1] \leq 1/2$.*

   *As usual, the error probability can be reduced by repetitions.*

3. Local Decoding: *There exists a probabilistic polynomial-time oracle machine $D$ that given $k$ and $i \in [k]$ (in binary) and oracle access to any $w \in [q(k)]^{n(k)}$ that is $\delta/2$-close to $C_k(x)$ returns $x_i$; that is, $\Pr[D^w(k, i) = x_i] \geq 2/3$.*

   Self correction *holds too: there exists a probabilistic polynomial-time oracle machine $M$ that given $k$ and $i \in [n(k)]$ (in binary) and oracle access to any $w \in [q(k)]^{n(k)}$ that is $\delta/2$-close to $C_k(x)$ returns $C_k(x)_i$; that is, $\Pr[D^w(k, i) = C_k(x)_i] \geq 2/3$.*

We stress that all these oracle machines work in time that is polynomial in the binary representation of $k$, which means that they run in time that is poly-logarithmic in $k$. The code asserted in Theorem E.3 is a (small modification of a) Reed-Muller code, for $r = m^2 \log m < q(k) = \text{poly}(r)$ and $[n(k)] \equiv \text{GF}(q(k))^m$ (see §E.1.1.4).[6]

---

[5]Thus, the running time of $T$ is $\text{poly}(|k|) = \text{poly}(\log k)$.

[6]The modification is analogous to the one presented in Footnote 2: For a suitable choice of $k$ points $\overline{\alpha}_1, ..., \overline{\alpha}_k \in \text{GF}(q(k))^m$, we map $v_1, ..., v_k$ to $(p(\overline{\alpha}_1), ..., p(\overline{\alpha}_n))$, where $p$ is the unique $m$-variate polynomial of degree at most $r$ that satisfies $p(\overline{\alpha}_i) = v_i$ for $i = 1, ..., k$.

The aforementioned oracle machines query the oracle $w : [n(k)] \rightarrow \mathrm{GF}(q(k))$ at a non-constant number of locations. Specifically, self-correction for location $i \in \mathrm{GF}(q(k))^m$ is performed by selecting a random line (over $\mathrm{GF}(q(k))^m$) that passes through $i$, recovering the values assigned by $w$ to all $q(k)$ points on this line, and performing univariate polynomial extrapolation (under mild noise). Local testability is easily reduced to self-correction, and (under the aforementioned modification) local decodability is a special case of self-correction.

**Constant number of queries.** The local testing and decoding algorithms asserted in Theorem E.3 make a polylogarithmic number of queries into the oracle. In contrast, the Hadamard code supports these operation using a *constant number of queries. Can this be obtained with much shorter codewords?* For local testability the answer is definitely positive. One can obtain such locally testable codes with length that is nearly linear (i.e., linear up to polylogarithmic factors; see [33, 62]). For local decodability based on a constant number of queries, the shortest known code has super-polynomial length (see [227]). In light of this state of affairs, we advocate a relaxation of the local decodability task (e.g., the one studied in [32]).

The interested reader is referred to [89], which includes more details on locally testable and decodable codes as well as a wider perspective. (Note, however, that this survey was written prior to [62] and [227], which address two major open problems discussed in [89].)

## E.1.3 A list decoding bound

A necessary condition for the feasibility of the list decoding task is that the list of codewords that are close to the given word is short. In this section we present an upper-bound on the length of such lists, noting that this bound has found several applications in complexity theory (and specifically to studies related to the contents of this book). In contrast, we do not present far more famous bounds (which typically refer to the relation among the main parameters of codes (i.e., $k, n$ and $d$)), because they seem irrelevant to the contents of this book.

We start with a general statement that refers to any alphabet $\Sigma \equiv [q]$, and later specialize it to the case that $q = 2$. Especially in the general case, it is natural and convenient to consider the agreement (rather than the distance) between sequences over $[q]$. Furthermore, it is natural to focus on agreement rate of at least $1/q$, and it is convenient to state the following result in terms of the "excessive agreement rate" (i.e., the excess beyond $1/q$).[7]

**Lemma E.4** (Part 2 of [101, Thm. 15]): *Let $C : [q]^k \rightarrow [q]^n$ be an arbitrary code of distance $d \leq n - (n/q)$, and let $\eta_{\mathrm{C}} \stackrel{\mathrm{def}}{=} (1 - (d/n)) - (1/q) \geq 0$ denote the corresponding upper-bound on the excessive agreement rate between codewords.*

---

[7]Indeed, we only consider codes with distance $d \leq (1 - 1/q) \cdot n$ and words that are at distance at most $d$ from the code. Note that $1/q$ is a natural threshold for an upper-bound on the relative agreement between sequences over $[q]$, because a random sequence is expected to agree with any fixed sequence on a $1/q$ fraction of the locations.

*Suppose that $\eta \in (0, 1)$ satisfies*

$$\eta > \sqrt{\left(1 - \frac{1}{q}\right) \cdot \eta_{\mathrm{c}}} \tag{E.1}$$

*Then, for any $w \in [q]^n$, the number of codewords that agree with $w$ on at least $((1/q) + \eta) \cdot n$ positions* (i.e., are at distance at most $(1 - ((1/q) + \eta)) \cdot n$ from $w$) *is upper-bounded by*

$$\frac{(1 - (1/q))^2 - (1 - (1/q)) \cdot \eta_{\mathrm{c}}}{\eta^2 - (1 - (1/q)) \cdot \eta_{\mathrm{c}}} \tag{E.2}$$

In the binary case (i.e., $q = 2$), Eq. (E.1) requires $\eta > \sqrt{\eta_{\mathrm{c}}/2}$ and Eq. (E.2) yields the upper-bound $(1 - 2\eta_{\mathrm{c}})/(4\eta^2 - 2\eta_{\mathrm{c}})$. We highlight two specific cases:

1. At the end of §D.4.2.2, we refer to this bound (for the binary case) while setting $\eta_{\mathrm{c}} = (1/k)^2$ and $\eta = 1/k$. Indeed, in this case $(1 - 2\eta_{\mathrm{c}})/(4\eta^2 - 2\eta_{\mathrm{c}}) = O(k^2)$.

2. In the case of the Hadamard code, we have $\eta_{\mathrm{c}} = 0$. Thus, for every $w \in \{0, 1\}^n$ and every $\eta > 0$, the number of codewords that are $(0.5 - \eta)$-close to $w$ is at most $1/(4\eta^2)$.

In the general case (and specifically for $q \gg 2$) it is useful to simplify Eq. (E.1) by $\eta > \min\{\sqrt{\eta_{\mathrm{c}}}, (1/q) + \sqrt{\eta_{\mathrm{c}} - (1/q)}\}$ and Eq. (E.2) by $\frac{1}{\eta^2 - \eta_{\mathrm{c}}}$.

## E.2   Expander Graphs

Loosely speaking, expander graphs are graphs of small degree that exhibit various properties of cliques. In particular, we refer to properties such as the relative sizes of cuts in the graph, and the rate at which a random walk converges to the uniform distribution (relative to the logarithm of the graph size to the base of its degree).

**Some technicalities.**   Typical presentations of expander graphs refer to one of several variants. For example, in some sources, expanders are presented as bipartite graphs, whereas in others they are presented as ordinary graphs (and are in fact very far from being bipartite). We shall follow the latter convention. Furthermore, at times we implicitly consider an augmentation of these graphs where self-loops are added to each vertex. For simplicity, we also allow parallel edges.

   We often talk of expander graphs while we actually mean an infinite collection of graphs such that each graph in this collection satisfies the same property (which is informally attributed to the collection). For example, when talking of a $d$-regular expander (graph) we actually refer to an infinite collection of graphs such that each of these graphs is $d$-regular. Typically, such a collection (or family) contains a single $N$-vertex graph for every $N \in \mathbb{S}$, where $\mathbb{S}$ is an infinite subset of $\mathbb{N}$. Throughout this section, we denote such a collection by $\{G_N\}_{N \in \mathbb{S}}$, with the understanding that $G_N$ is a graph with $N$ vertices and $\mathbb{S}$ is an infinite set of natural numbers.

## E.2.1 Definitions and Properties

We consider two definitions of expander graphs, two different notions of explicit constructions, and two useful properties of expanders.

### E.2.1.1 Two Mathematical Definitions

We start with two different definitions of expander graphs. These definitions are qualitatively equivalent and even quantitatively related. We start with an algebraic definition, which seems technical in nature but is actually the definition typically used in complexity theoretic applications, since it directly implies various "mixing properties" (see §E.2.1.3). We later present a very natural combinatorial definition (which is the source of the term "expander").

**The algebraic definition (spectral gap).** Identifying graphs with their adjacency matrix, we consider the eigenvalues (and eigenvectors) of a graph (or rather of its adjacency matrix). Any $d$-regular graph $G = (V, E)$ has the uniform vector as an eigenvector corresponding to the eigenvalue $d$, and if $G$ is connected and not bipartite then (the absolute values of) all other eigenvalues are strictly smaller than $d$. The second eigenvalue, denoted $\lambda_2(G) < d$, of such a graph $G$ is thus a tight upper-bound on the *absolute value* of all the other eigenvalues. Using the connection to the combinatorial definition, it follows that $\lambda_2(G) < d - \Omega(1/|V|^2)$ holds (for every connected non-bipartite $d$-regular graph $G$). The algebraic definition of expanders refers to an infinite family of $d$-regular graphs and requires the existence of a *constant* eigenvalue bound that holds for all the graphs in the family.

**Definition E.5** *An infinite family of $d$-regular graphs, $\{G_N\}_{N \in \mathbb{S}}$, where $\mathbb{S} \subseteq \mathbb{N}$, satisfies the eigenvalue bound $\lambda$ if for every $N \in \mathbb{S}$ it holds that $\lambda_2(G_N) \leq \lambda$.*

In such a case we say that the family has spectral gap $d - \lambda$. It will be often convenient to consider relative (or normalized) versions of these quantities, obtained by division by $d$.

**The combinatorial definition (expansion).** Loosely speaking, expansion requires that any (not too big) set of vertices of the graph has a relatively large set of neighbors. Specifically, a graph $G = (V, E)$ is $c$-expanding if, for every set $S \subset V$ of cardinality at most $|V|/2$, it holds that

$$\Gamma_G(S) \stackrel{\text{def}}{=} \{v : \exists u \in S \text{ s.t. } (u, v) \in E\} \tag{E.3}$$

has cardinality at least $(1 + c) \cdot |S|$. Equivalently (assuming the existence of self-loops on all vertices), we may require that $|\Gamma_G(S) \setminus S| \geq c \cdot |S|$. Clearly, every connected graph $G = (V, E)$ is $(1/|V|)$-expanding. The combinatorial definition of expanders refers to an infinite family of $d$-regular graphs and requires the existence of a *constant* expansion bound that holds for all the graphs in the family.

**Definition E.6** *An infinite family of $d$-regular graphs, $\{G_N\}_{N \in \mathbb{S}}$ is $c$-expanding if for every $N \in \mathbb{S}$ it holds that $G_N$ is $c$-expanding.*

The two definitions of expander graphs are related (see [10, Sec. 9.2] or [118, Sec. 4.5]).

**Theorem E.7** *Let $G$ be a non-bipartite $d$-regular graph.*

1. *The graph $G$ is $c$-expanding for $c \geq (d - \lambda_2(G))/2d$.*

2. *If $G$ is $c$-expanding then $d - \lambda_2(G) \geq c^2/(4 + 2c^2)$.*

Thus, any non-zero bound on the combinatorial expansion of a family of $d$-regular graphs yields a non-zero bound on its spectral gap, and vice versa. Note, however, that the back-and-forth translation between these definitions is not tight. The applications presented in the main text refer to the algebraic definition, and the loss incurred in Theorem E.7 is immaterial for them.

**Amplification.** The quality of expander graphs improves by raising them to any power $t > 1$ (i.e., raising their adjacency matrix to the $t^{\text{th}}$ power), which corresponds to considering graphs in which $t$-paths are replaced by edges. Using the algebraic definition, we have $\lambda_2(G^t) = \lambda_2(G)^t$, but indeed the degree also gets raised to the power $t$. Still, the ratio $\lambda_2(G^t)/d^t$ deceases with $t$. An analogous phenomenon occurs also under the combinatorial definition, provided that some suitable modifications are applied. For example, if $G = (V, E)$ is $c$-expanding (i.e., for every $S \subseteq V$ it holds that $|\Gamma_G(S)| \geq \min((1 + c) \cdot |S|, |V|/2)$), then for every $S \subseteq V$ it holds that $|\Gamma_{G^t}(S)| \geq \min((1 + c)^t \cdot |S|, |V|/2)$.

**The optimal eigenvalue bound.** For every $d$-regular graph $G = (V, E)$, it holds that $\lambda_2(G) \geq 2\gamma_G \cdot \sqrt{d-1}$, where $\gamma_G = 1 - O(1/\log_d |V|)$. Thus, $2\sqrt{d-1}$ is a lower-bound on the eigenvalue bound of any infinite family of $d$-regular graphs.

### E.2.1.2   Two levels of explicitness

A mild level of explicit constructiveness refers to the complexity of constructing the entire object (i.e., graph). Thus, an infinite family of graphs $\{G_N\}_{N \in \mathbb{S}}$ is said to be explicitly constructible if there exists a *polynomial-time algorithm that, on input $1^N$ (where $N \in \mathbb{S}$), outputs the list of the edges in the $N$-vertex graph $G_N$.*

The aforementioned level of explicitness suffices when the application requires holding the entire graph and/or when the running-time of the application is lower-bounded by the size of the graph. In contrast, other applications only refer to a huge virtual graph (which is much bigger than their running time), and only require the computation of the neighborhood relations in such a graph. In this case, the following stronger level of explicitness is relevant.

A strongly explicit construction of an infinite family of ($d$-regular) graphs $\{G_N\}_{N \in \mathbb{S}}$ is a *polynomial-time algorithm that on input $N$ (in binary), a vertex $v$ in the $N$-vertex graph $G_N$ and an index $i$ ($i \in \{1, ..., d\}$), returns the $i^{\text{th}}$ neighbor of $v$.* That is, the neighbor is determined in time that is polylogarithmic in the size of the graph. Needless to say, the strong level of explicitness implies the basic level.

An additional requirement, which is often forgotten but is very important, refers to the "tractability" of the set $\mathbb{S}$. Specifically, we require the existence of an *efficient algorithm that given any* $n \in \mathbb{N}$ *finds an* $s \in \mathbb{S}$ *such that* $n \leq s < 2n$. Corresponding to the foregoing definitions, efficient may mean either running in time $\mathrm{poly}(n)$ or running in time $\mathrm{poly}(\log n)$. The requirement that $n \leq s < 2n$ suffices in most applications, but in some cases a smaller interval (e.g., $n \leq s < n + \sqrt{n}$) is required, whereas in other cases a larger interval (e.g., $n \leq s < \mathrm{poly}(n)$) suffices.

**Greater flexibility.** In continuation to the foregoing paragraph, we comment that expanders can be combined in order to obtain expanders for a wider range of sizes. For example, two $d$-regular $c$-expanding graphs, $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$ where $|V_1| \leq |V_2|$ and $c \leq 1$, can be combined into a $(d+1)$-regular $c/2$-expanding graph on $|V_1| + |V_2|$ vertices by connecting the two graphs with a perfect matching of $V_1$ and $|V_1|$ of the vertices of $V_2$ (and adding self-loops to the remaining vertices of $V_2$). More generally, the $d$-regular $c$-expanding graphs, $G_1 = (V_1, E_1)$ through $G_t = (V_t, E_t)$, where $N \stackrel{\mathrm{def}}{=} \sum_{i=1}^{t-1} |V_i| \leq |V_t|$, yield a $(d+1)$-regular $c/2$-expanding graph on $\sum_{i=1}^{t} |V_i|$ vertices by using a perfect matching of $\cup_{i=1}^{t-1} V_i$ and $N$ of the vertices of $V_t$.

### E.2.1.3   Two properties

The following two properties provide a quantitative interpretation to the statement that expanders approximate the complete graph. The deviation from the latter is represented by an error term that is linear in $\lambda/d$.

**The mixing lemma.** The following lemma is folklore and has appeared in many papers. Loosely speaking, the lemma asserts that expander graphs (for which $d \gg \lambda$) have the property that the fraction of edges between two large sets of vertices approximately equals the product of the densities of these sets. This property is called *mixing*.

**Lemma E.8** (Expander Mixing Lemma): *For every $d$-regular graph $G = (V, E)$ and for every two subsets $A, B \subseteq V$ it holds that*

$$\left| \frac{|(A \times B) \cap E_2|}{|E_2|} - \frac{|A|}{|V|} \cdot \frac{|B|}{|V|} \right| \leq \frac{\lambda_2(G)\sqrt{|A| \cdot |B|}}{d \cdot |V|} \leq \frac{\lambda_2(G)}{d} \tag{E.4}$$

*where $E_2$ denotes the set of directed edges that correspond to the undirected edges of $G$ (i.e., $E_2 = \{(u, v) : \{u, v\} \in E\}$ and $|E_2| = d|V|$).*

**Proof:** Let $N \stackrel{\mathrm{def}}{=} |V|$ and $\lambda \stackrel{\mathrm{def}}{=} \lambda_2(G)$. For any subset of the vertices $S \subseteq V$, we denote its density in $V$ by $\rho(S) \stackrel{\mathrm{def}}{=} |S|/N$. Hence, Eq. (E.4) is restated as

$$\left| \frac{|(A \times B) \cap E_2|}{d \cdot N} - \rho(A) \cdot \rho(B) \right| \leq \frac{\lambda \sqrt{\rho(A) \cdot \rho(B)}}{d}.$$

We proceed by providing bounds on the value of $|(A \times B) \cap E_2|$. To this end we let $\overline{a}$ denote the $N$-dimensional Boolean vector having 1 in the $i^{\text{th}}$ component if and only if $i \in A$. The vector $\overline{b}$ is defined similarly. Denoting the adjacency matrix of the graph $G$ by $M = (m_{i,j})$, we note that $|(A \times B) \cap E_2|$ equals $\overline{a}^\top M \overline{b}$ (because $(i, j) \in (A \times B) \cap E_2$ if and only if it holds that $i \in A$, $j \in B$ and $m_{i,j} = 1$). We consider the *orthogonal eigenvector basis*, $\overline{e_1}, ..., \overline{e_N}$, where $\overline{e_1} = (1, ..., 1)^\top$ and $\overline{e_i}^\top \overline{e_i} = N$ for each $i$, and write each vector as a linear combination of the vectors in this basis. Specifically, we denote by $a_i$ the coefficient of $\overline{a}$ in the direction of $\overline{e_i}$; that is, $a_i = (\overline{a}^\top \overline{e_i})/N$ and $\overline{a} = \sum_i a_i \overline{e_i}$. Note that $a_1 = (\overline{a}^\top \overline{e_1})/N = |A|/N = \rho(A)$ and $\sum_{i=1}^N a_i^2 = (\overline{a}^\top \overline{a})/N = |A|/N = \rho(A)$. Similarly for $\overline{b}$. It now follows that

$$
\begin{aligned}
|(A \times B) \cap E_2| &= \overline{a}^\top M \left( b_1 \overline{e_1} + \sum_{i=2}^N b_i \overline{e_i} \right) \\
&= \rho(B) \cdot \overline{a}^\top M \overline{e_1} + \sum_{i=2}^N b_i \cdot \overline{a}^\top M \overline{e_i} \\
&= \rho(B) \cdot d \cdot \overline{a}^\top \overline{e_1} + \sum_{i=2}^N b_i \lambda_i \cdot \overline{a}^\top \overline{e_i}
\end{aligned}
$$

where $\lambda_i$ denotes the $i^{\text{th}}$ eigenvalue of $M$ (and indeed $\lambda_1 = d$). Thus,

$$
\begin{aligned}
\frac{|(A \times B) \cap E_2|}{dN} &= \rho(B)\rho(A) + \sum_{i=2}^N \frac{\lambda_i b_i a_i}{d} \\
&\in \left[ \rho(B)\rho(A) \pm \frac{\lambda}{d} \cdot \sum_{i=2}^N a_i b_i \right]
\end{aligned}
$$

Using $\sum_{i=1}^N a_i^2 = \rho(A)$ and $\sum_{i=1}^N b_i^2 = \rho(B)$, and applying Cauchy-Schwartz Inequality, we bound $\sum_{i=2}^N a_i b_i$ by $\sqrt{\rho(A)\rho(B)}$. The lemma follows.    ∎

**The random walk lemma.**   Loosely speaking, the first part of the following lemma asserts that, as far as remaining trapped in some subset of the vertex set is concerned, a random walk on an expander approximates a random walk on the complete graph.

**Lemma E.9** (Expander Random Walk Lemma): *Let $G = ([N], E)$ be a d-regular graph, and consider walks on $G$ that start from a uniformly chosen vertex and take $\ell - 1$ additional random steps, where in each such step we uniformly selects one out of the $d$ edges incident at the current vertex and traverses it.*

**Theorem 8.28 (restated):** *Let $W$ be a subset of $[N]$ and $\rho \stackrel{\text{def}}{=} |W|/N$. Then the probability that such a random walk stays in $W$ is at most*

$$
\rho \cdot \left( \rho + (1 - \rho) \cdot \frac{\lambda_2(G)}{d} \right)^{\ell - 1} \tag{E.5}
$$

**Exercise 8.37 (restated):** *For any $W_0, ..., W_{\ell-1} \subseteq [N]$, the probability that a random walk of length $\ell$ intersects $W_0 \times W_1 \times \cdots \times W_{\ell-1}$ is at most*

$$\sqrt{\rho_0} \cdot \prod_{i=1}^{\ell-1} \sqrt{\rho_i + (\lambda/d)^2}, \tag{E.6}$$

*where $\rho_i \stackrel{\text{def}}{=} |W_i|/N$.*

The basic principle underlying Lemma E.9 was discovered by Ajtai, Komlos, and Szemerédi [4], who proved a bound as in Eq. (E.6). The better analysis yielding Theorem 8.28 is due to Kahale [127, Cor. 6.1]. More general bounds that refer to the probability of visiting $W$ for a number of times that approximates $|W|/N$ are given in [82], which actually considers an even more general problem (i.e., obtaining Chernoff-type bounds for random variables that are generated by a walk on a Markov Chain).

**Proof of Equation (E.6):** The basic idea is to view the random walk as the evolution of a corresponding probability vector under suitable transformations. The transformations correspond to taking a random step in $G$ and to passing through a "sieve" that keeps only the entries that correspond to the current set $W_i$. The key observation is that the first transformation shrinks the component that is orthogonal to the uniform distribution, whereas the second transformation shrinks the component that is in the direction of the uniform distribution. Details follow.

Let $A$ be a matrix representing the random walk on $G$ (i.e., $A$ is the adjacency matrix of $G$ divided by $d$), and let $\hat{\lambda}$ denote the absolute value of the second largest eigenvalue of $A$ (i.e., $\hat{\lambda} \stackrel{\text{def}}{=} \lambda_2(G)/d$). Note that the uniform distribution, represented by the vector $\overline{u} = (N^{-1}, ..., N^{-1})^\top$, is the eigenvector of $A$ that is associated with the largest eigenvalue (which is 1). Let $P_i$ be a 0-1 matrix that has 1-entries only on its diagonal, and furthermore entry $(j, j)$ is set to 1 if and only if $j \in W_i$. Then, the probability that a random walk of length $\ell$ intersects $W_0 \times W_1 \times \cdots \times W_{\ell-1}$ is the sum of the entries of the vector

$$\overline{v} \stackrel{\text{def}}{=} P_{\ell-1}A \cdots P_2 A P_1 A P_0 \overline{u}. \tag{E.7}$$

We are interested in upper-bounding $\|\overline{v}\|_1$, and use $\|\overline{v}\|_1 \leq \sqrt{N} \cdot \|\overline{v}\|$, where $\|\overline{z}\|_1$ and $\|\overline{z}\|$ denote the $L_1$-norm and $L_2$-norm of $\overline{z}$, respectively (e.g., $\|\overline{u}\|_1 = 1$ and $\|\overline{u}\| = N^{-1/2}$). The key observation is that the linear transformation $P_i A$ shrinks every vector.

**Main Claim.** For every $\overline{z}$, it holds that $\|P_i A \overline{z}\| \leq (\rho_i + \hat{\lambda}^2)^{1/2} \cdot \|\overline{z}\|$.

**Proof.** Intuitively, $A$ shrinks the component of $\overline{z}$ that is orthogonal to $\overline{u}$, whereas $P_i$ shrinks the component of $\overline{z}$ that is in the direction of $\overline{u}$. Specifically, we decompose $\overline{z} = \overline{z_1} + \overline{z_2}$ such that $\overline{z_1}$ is the projection of $\overline{z}$ on $\overline{u}$ and $\overline{z_2}$ is the component orthogonal to $\overline{u}$. Then, using the triangle inequality and other obvious facts (which

imply $\|P_i A \overline{z_1}\| = \|P_i \overline{z_1}\|$ and $\|P_i A \overline{z_2}\| \le \|A \overline{z_2}\|$), we have

$$
\begin{aligned}
\|P_i A \overline{z_1} + P_i A \overline{z_2}\| &\le \|P_i A \overline{z_1}\| + \|P_i A \overline{z_2}\| \\
&\le \|P_i \overline{z_1}\| + \|A \overline{z_2}\| \\
&\le \sqrt{\rho_i} \cdot \|\overline{z_1}\| + \hat{\lambda} \cdot \|\overline{z_2}\|
\end{aligned}
$$

where the last inequality uses the fact that $P_i$ shrinks any uniform vector by eliminating $1 - \rho_i$ of its elements, whereas $A$ shrinks the length of any eigenvector except $\overline{u}$ by a factor of at least $\hat{\lambda}$. Using the Cauchy-Schwartz inequality[8], we get

$$
\begin{aligned}
\|P_i A \overline{z}\| &\le \sqrt{\rho_i + \hat{\lambda}^2} \cdot \sqrt{\|\overline{z_1}\|^2 + \|\overline{z_2}\|^2} \\
&= \sqrt{\rho_i + \hat{\lambda}^2} \cdot \|\overline{z}\|
\end{aligned}
$$

where the equality is due to the fact that $\overline{z_1}$ is orthogonal to $\overline{z_2}$.  $\square$

Recalling Eq. (E.7) and using the Main Claim (and $\|\overline{v}\|_1 \le \sqrt{N} \cdot \|\overline{v}\|$), we get

$$
\begin{aligned}
\|\overline{v}\|_1 &\le \sqrt{N} \cdot \|P_{\ell-1} A \cdots P_2 A P_1 A P_0 \overline{u}\| \\
&\le \sqrt{N} \cdot \left( \prod_{i=1}^{\ell-1} \sqrt{\rho_i + \hat{\lambda}^2} \right) \cdot \|P_0 \overline{u}\|.
\end{aligned}
$$

Finally, using $\|P_0 \overline{u}\| = \sqrt{\rho_0 N \cdot (1/N)^2} = \sqrt{\rho_0 / N}$, we establish Eq. (E.6).  ∎

**Rapid mixing.**   A property related to Lemma E.9 is that a random walk starting at any vertex converges to the uniform distribution on the expander vertices after a logarithmic number of steps. Using notation as in the proof of Eq. (E.6), we *claim that for every starting distribution $\overline{s}$* (including one that assigns all weight to a single vertex), *it holds that* $\|A^\ell \overline{s} - \overline{u}\|_1 \le \sqrt{N} \cdot \hat{\lambda}^\ell$, which is meaningful for any $\ell > 0.5 \cdot \log_{1/\hat{\lambda}} N$. The claim is proved by recalling that $\|A^\ell \overline{s} - \overline{u}\|_1 \le \sqrt{N} \cdot \|A^\ell \overline{s} - \overline{u}\|$ and using the fact that $\overline{s} - \overline{u}$ is orthogonal to $\overline{u}$ (because the former is a zero-sum vector). Thus, $\|A^\ell \overline{s} - \overline{u}\| = \|A^\ell (\overline{s} - \overline{u})\| \le \hat{\lambda}^\ell \|\overline{s} - \overline{u}\|$ and using $\|\overline{s} - \overline{u}\| < 1$ the claim follows.

## E.2.2   Constructions

Many explicit constructions of expanders were discovered, starting in [154] and culminating in the optimal construction of [150] where $\lambda = 2\sqrt{d-1}$. Most of these constructions are quite simple (see, e.g., §E.2.2.1), but their analysis is based on non-elementary results from various branches of mathematics. In contrast, the construction of Reingold, Vadhan, and Wigderson [180], presented in §E.2.2.2,

---

[8]That is, we get $\sqrt{\rho_i}\|z_1\| + \hat{\lambda}\|z_2\| \le \sqrt{\rho_i + \hat{\lambda}^2} \cdot \sqrt{\|z_1\|^2 + \|z_2\|^2}$, by using $\sum_{i=1}^{n} a_i \cdot b_i \le \left( \sum_{i=1}^{n} a_i{}^2 \right)^{1/2} \cdot \left( \sum_{i=1}^{n} b_i{}^2 \right)^{1/2}$, with $n = 2$, $a_1 = \sqrt{\rho_i}$, $b_1 = \|z_1\|$, etc.

is based on an iterative process, and its analysis is based on a relatively simple algebraic fact regarding the eigenvalues of matrices.

Before turning to these explicit constructions we note that it is relatively easy to prove the existence of 3-regular expanders, by using the Probabilistic Method (cf. [10]) and referring to the combinatorial definition of expansion.

**Theorem E.10** *For some constant $\lambda < 3$ there exists a family of $(3, \lambda)$-expanders for any even graph size.*

**Proof Sketch:**[9] As a warm-up, one may establish the existence of $d$-regular expanders, for some constant $d$. In particular, foreseeing the case of $d = 3$, consider a random graph $G$ on the vertex set $V = \{0, ..., n-1\}$ constructed by augmenting the fixed edge set $\{\{i, i+1 \bmod n\} : i = 0, ..., n-1\}$ with $d-2$ uniformly (and independently) chosen perfect matchings of the vertices of $F \stackrel{\text{def}}{=} \{0, ..., (n/2)-1\}$ to the vertices of $L \stackrel{\text{def}}{=} \{n/2, ..., n-1\}$. For a sufficiently small universal constant $\varepsilon > 0$, we upper-bound the probability that such a random graph is not $\varepsilon$-expanding. Noting that for every set $S$ it holds that $|\Gamma_G(S \cap F) \cap F| \geq |S \cap F| - 1$ (and similarly for $L$), we focus on the sizes of $|(\Gamma_G(S \cap F) \cap L) \setminus \Gamma_G(S \cap L)|$ and $|(\Gamma_G(S \cap L) \cap F) \setminus \Gamma_G(S \cap F)|$. Assuming without loss of generality that $|S \cap F| \geq |S \cap L|$, we upper-bound the probability that there exists a set $S \subset V$ of size at most $n/2$ such that $|(\Gamma_G(S \cap F) \cap L) \setminus \Gamma_G(S \cap L)| < \varepsilon|S|$. Fixing a set $S$, the corresponding probability is upper-bounded by $p_S^{d-2}$, where $p_S$ denotes the probability that a uniformly selected matching of $F$ to $L$ matches $S \cap F$ to a set that contains less than $\varepsilon|S|$ elements in $L \setminus \Gamma_G(S \cap L)$. That is,

$$p_S \stackrel{\text{def}}{=} \sum_{i=0}^{\varepsilon|S|-1} \frac{\binom{|L|-\ell}{i} \cdot \binom{\ell}{|S \cap F|-i}}{\binom{|L|}{|S \cap F|}} \leq \frac{\binom{(n/2)-\ell}{\varepsilon|S|} \cdot \binom{\ell+\varepsilon|S|}{|S \cap F|}}{\binom{n/2}{|S \cap F|}}$$

where $\ell = |\Gamma_G(S \cap L) \cap L|$. Indeed, we may focus on the case that $|S \cap F| \leq \ell + \varepsilon|S|$ (because in the other case $p_S = 0$), and observe that for every $\alpha < 1/2$ there exists a sufficiently small $\varepsilon > 0$ such that $p_S < \binom{n}{|S|}^{-\alpha}$. The claim follows for $d \geq 5$, by using a union bound on all sets (and setting $\alpha = 1/3$).

To deal with the case $d = 3$, we use a more sophisticated union bound. Specifically, fixing an adequate constant $t > 6$ (e.g., $t = 1/\sqrt{\varepsilon}$), we decompose $S$ into $S'$ and $S''$, where $S'$ contains the elements of $S$ that reside on $t$-long arithmetic subsequences of $S$ that use an step increment of either 1 or 2, and $S'' = S \setminus S'$. It can be shown that $|\Gamma_G(S'') \setminus S| > |S''|/2t$ (hint: an arithmetic subsequence has neighborhood greater than itself whereas a suitable partition of the elements to such subsequences guarantees that the overall excess is at least half the individual

---

[9] The proof is much simpler in the case that one refers to the alternative definition of combinatorial expansion in which for each relevant set $S$ it holds that $|\Gamma_G(S) \setminus S| \geq \varepsilon \cdot |S|$. In this case, for a sufficiently small $\varepsilon > 0$ and all sufficiently large $n$, a random 3-regular $n$-vertex graph is $\varepsilon$-expanding with overwhelmingly high probability. The proof proceeds by considering a (not necessarily simple) graph $G$ generated by three perfect matchings of the elements of $[n]$. For every $S \subseteq [n]$ of size at most $n/2$ and for every set $T$ of size $\varepsilon|S|$, we consider the probability that $\Gamma_G(S) \subseteq S \cup T$. The argument is concluded by applying a union bound.

count). Thus, if $|S''| > 2|S|/t$ then $|\Gamma_G(S'') \setminus S| > |S|/t^2$. Hence, it suffices to consider the case $|S''| \leq 2|S|/t$ (and $t \geq 6$) and prove that $|\Gamma_G(S')| > (1 + (4/t)) \cdot |S'|$. The gain is that, when applying the union bound, it suffices to consider less than $\sum_{j=1}^{n'/t} 2^j \cdot \binom{n}{2j} < \binom{n}{3n'/t}$ possible sets $S'$ of size $n'$, which are each a union of at most $n'/t$ arithmetic sequences that use an step increment of either 1 or 2.  $\square$

### E.2.2.1   The Margulis–Gabber–Galil Expander

For every natural number $m$, consider the graph with vertex set $\mathbb{Z}_m \times \mathbb{Z}_m$ and edge set in which every $\langle x, y \rangle \in \mathbb{Z}_m \times \mathbb{Z}_m$ is connected to the vertices $\langle x \pm y, y \rangle$, $\langle x \pm (y + 1), y \rangle$, $\langle x, y \pm x \rangle$, and $\langle x, y \pm (x + 1) \rangle$, where the arithmetic is modulo $m$. This yields an extremely simple and explicit 8-regular graph with second eigenvalue that is bounded by a constant $\lambda < 8$ that is independent of $m$. Thus we get:

**Theorem E.11** *For some constant $\lambda < 8$ there exists a strongly explicit construction of a family of $(8, \lambda)$-expanders for graph sizes $\{m^2 : m \in \mathbb{N}\}$. Furthermore, the neighbors of a vertex can be computed in logarithmic-space.*[10]
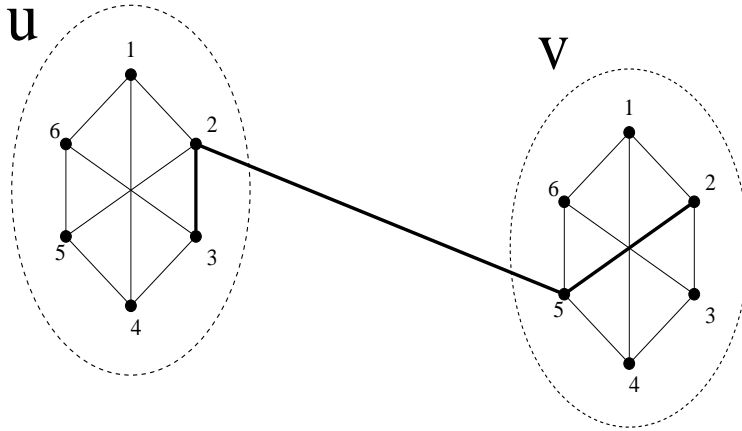
An appealing property of Theorem E.11 is that, for every $n \in \mathbb{N}$, it directly yields expanders with vertex set $\{0,1\}^n$. This is obvious in case $n$ is even, but can be easily achieved also for odd $n$ (e.g., use two copies of the graph for $n - 1$, and connect the two copies by the obvious perfect matching).

Theorem E.11 is due to Gabber and Galil [79], building on the basic approach suggested by Margulis [154]. We mention again that the optimal construction of [150] achieves $\lambda = 2\sqrt{d-1}$, but there are annoying restrictions on the degree $d$ (i.e., $d - 1$ should be a prime congruent to 1 modulo 4) and on the graph sizes for which this construction works.

### E.2.2.2   The Iterated Zig-Zag Construction

The starting point of the following construction is a very good expander $G$ of *constant size*, which may be found by an exhaustive search. The construction of a large expander graph proceeds in iterations, where in the $i^{\text{th}}$ iteration the current graph $G_i$ and the fixed graph $G$ are combined, resulting in a larger graph $G_{i+1}$. The combination step guarantees that the expansion property of $G_{i+1}$ is at least as good as the expansion of $G_i$, while $G_{i+1}$ maintains the degree of $G_i$ and is a constant times larger than $G_i$. The process is initiated with $G_1 = G^2$ and terminates when we obtain a graph $G_t$ of approximately the desired size (which requires a logarithmic number of iterations).

---

[10]In fact, under a suitable encoding of the vertices and for $m$ that is a power of two, the neighbors can be computed by a on-line algorithm that uses a constant amount of space. The same holds also for a variant in which each vertex $\langle x, y \rangle$ is connected to the vertices $\langle x \pm 2y, y \rangle$, $\langle x \pm (2y + 1), y \rangle$, $\langle x, y \pm 2x \rangle$, and $\langle x, y \pm (2x + 1) \rangle$. (This variant yields a better known bound on $\lambda$, i.e., $\lambda \leq 5\sqrt{2} \approx 7.071$.)

*In this example $G'$ is 6-regular and $G$ is a 3-regular graph having six vertices. In the graph $G'$ (not shown), the 2nd edge of vertex $u$ is incident at $v$, as its 5th edge. The wide 3-segment line shows one of the corresponding edges of $G'\textcircled{z}G$, which connects the vertices $\langle u, 3 \rangle$ and $\langle v, 2 \rangle$.*

Figure E.1: Detail of the zig-zag product of $G'$ and $G$.

**The Zig-Zag product.** The heart of the combination step is a new type of "graph product" called *Zig-Zag product*. This operation is applicable to any pair of graphs $G = ([D], E)$ and $G' = ([N], E')$, provided that $G'$ (which is typically larger than $G$) is $D$-regular. For simplicity, we assume that $G$ is $d$-regular (where typically $d \ll D$). The Zig-Zag product of $G'$ and $G$, denoted $G'\textcircled{z}G$, is defined as a graph with vertex set $[N] \times [D]$ and an edge set that includes an edge between $\langle u, i \rangle \in [N] \times [D]$ and $\langle v, j \rangle$ if and only if $(i, k), (\ell, j) \in E$ and the $k^{\text{th}}$ edge incident at $u$ equals the $\ell^{\text{th}}$ edge incident at $v$. See Figure E.1 as well as further clarification that follows.

---

**Teaching note:** The following paragraph, which provides a formal description of the zig-zag product, can be ignored in first reading but is useful for more advanced discussion.

---

It will be convenient to represent graphs like $G'$ by their **edge rotation function**, denoted $R' : [N] \times [D] \to [N] \times [D]$, such that $R'(u, i) = (v, j)$ if $(u, v)$ is the $i^{\text{th}}$ edge incident at $u$ as well as the $j^{\text{th}}$ edge incident at $v$. For simplicity, we assume that $G$ is edge-colorable with $d$ colors, which in turn yields a natural edge rotation function (i.e., $R(i, \alpha) = (j, \alpha)$ if the edge $(i, j)$ is colored $\alpha$). We will denote by $E_\alpha(i)$ the vertex reached from $i \in [D]$ by following the edge colored $\alpha$ (i.e., $E_\alpha(i) = j$ iff $R(i, \alpha) = (j, \alpha)$). The Zig-Zag product of $G'$ and $G$, denoted $G'\textcircled{z}G$, is then defined as a graph with the vertex set $[N] \times [D]$ and the edge rotation function

$$(\langle u, i \rangle, \langle \alpha, \beta \rangle) \mapsto (\langle v, j \rangle, \langle \beta, \alpha \rangle) \quad \text{if } R'(u, E_\alpha(i)) = (v, E_\beta(j)). \quad \text{(E.8)}$$

That is, edges are labeled by pairs over $[d]$, and the $\langle\alpha,\beta\rangle^{\text{th}}$ edge out of vertex $\langle u,i\rangle \in [N]\times[D]$ is incident at the vertex $\langle v,j\rangle$ (as its $\langle\beta,\alpha\rangle^{\text{th}}$ edge) if $R(u,E_\alpha(i)) = (v,E_\beta(j))$. (That is, based on $\langle\alpha,\beta\rangle$, we take a $G$-step from $\langle u,i\rangle$ to $\langle u,E_\alpha(i)\rangle$, then viewing $\langle u,E_\alpha(i)\rangle \equiv (u,E_\alpha(i))$ as an edge of $G'$ we rotate it to $(v,j') \stackrel{\text{def}}{=} R'(u,E_\alpha(i))$, and take a $G$-step from $\langle v,j'\rangle$ to $\langle v,E_\beta(j')\rangle$, while defining $j = E_\beta(j')$ and using $j' = E_\beta(E_\beta(j')) = E_\beta(j)$.)

Clearly, the graph $G'\textcircled{z}G$ is $d^2$-regular and has $D\cdot N$ vertices. The key fact, proved in [180] (using techniques as in §E.2.1.3), is that the relative eigenvalue of the zig-zag product is upper-bounded by the sum of the relative eigenvalues of the two graphs (i.e., $\bar\lambda_2(G'\textcircled{z}G) \le \bar\lambda_2(G') + \bar\lambda_2(G)$, where $\bar\lambda_2(\cdot)$ denotes the relative eigenvalue of the relevant graph). The (qualitative) fact that $G'\textcircled{z}G$ is an expander if both $G'$ and $G$ are expanders is very intuitive (e.g., consider what happens if $G'$ or $G$ is a clique). Things are even more intuitive if one considers the (related) **replacement product** of $G'$ and $G$, denoted $G'\textcircled{r}G$, *where there is an edge between* $\langle u,i\rangle \in [N]\times[D]$ *and* $\langle v,j\rangle$ *if and only if either* $u=v$ *and* $(i,j)\in E$ *or the* $i^{\text{th}}$ *edge incident at* $u$ *equals the* $j^{\text{th}}$ *edge incident at* $v$.[11]

**The iterated construction.**   The iterated expander construction uses the aforementioned zig-zag product as well as graph squaring. Specifically, the construction starts with the $d^2$-regular graph $G_1 = G^2 = ([D],E^2)$, where $D = d^4$ and $\bar\lambda_2(G) < 1/4$, and proceeds in iterations such that $G_{i+1} = G_i^2\textcircled{z}G$ for $i = 1,2,...,t-1$. That is, in each iteration, the current graph is first squared and then composed with the fixed ($d$-regular $D$-vertex) graph $G$ via the zig-zag product. This process maintains the following two invariants:

1. The graph $G_i$ is $d^2$-regular and has $D^i$ vertices.

   (The degree bound follows from the fact that a zig-zag product with a $d$-regular graph always yields a $d^2$-regular graph.)

2. The relative eigenvalue of $G_i$ is smaller than one half.

   (Here we use the fact that $\bar\lambda_2(G_{i-1}^2\textcircled{z}G) \le \bar\lambda_2(G_{i-1}^2) + \bar\lambda_2(G)$, which in turn equals $\bar\lambda_2(G_{i-1})^2 + \bar\lambda_2(G) < (1/2)^2 + (1/4)$. Note that graph squaring is used to reduce the relative eigenvalue of $G_i$ before increasing it by zig-zag product with $G$.)

To ensure that we can construct $G_i$, we should show that we can actually construct the edge rotation function that correspond to its edge set. This boils down to showing that, given the edge rotation function of $G_{i-1}$, we can compute the edge rotation function of $G_{i-1}^2$ as well as of its zig-zag product with $G$. Note that this computation amounts to two recursive calls to computations regarding $G_{i-1}$ (and two computations that correspond to the constant graph $G$). But since the recursion depth is logarithmic in the size of the final graph, the time spend in the recursive computation is polynomial in the size of the final graph. This suffices for the minimal notion of explicitness, but not for the stronger one.

---

[11]As an exercise, the reader is encouraged to show that if both $G'$ and $G$ are expanders according to the combinatorial definition then so is $G'\textcircled{r}G$.

**The strongly explicit version.** To achieve a *strongly explicit construction*, we slightly modify the iterative construction. Rather than letting $G_{i+1} = G_i^2 \textcircled{z} G$, we let $G_{i+1} = (G_i \times G_i)^2 \textcircled{z} G$, where $G' \times G'$ denotes the *tensor product of $G'$ with itself*; that is, if $G' = (V', E')$ then $G' \times G' = (V' \times V', E'')$, where

$$E'' = \{(\langle u_1, u_2 \rangle, \langle v_1, v_2 \rangle) : (u_1, v_1), (u_2, v_2) \in E'\}$$

with an edge rotation function

$$R''(\langle u_1, u_2 \rangle, \langle i_1, i_2 \rangle) = (\langle v_1, v_2 \rangle, \langle j_1, j_2 \rangle)$$

where $R'(u_1, i_1) = (v_1, j_1)$ and $R'(u_2, i_2) = (v_2, j_2)$. (We still use $G_1 = G^2$.) Using the fact that tensor product preserves the relative eigenvalue (while squaring the degree) and using a $d$-regular $G = ([D], E)$ with $D = d^8$, we note that the modified $G_i = (G_{i-1} \times G_{i-1})^2 \textcircled{z} G$ is a $d^2$-regular graph with $(D^{2^{i-1}-1})^2 \cdot D = D^{2^i-1}$ vertices, and $\bar{\lambda}_2(G_i) < 1/2$ (because $\bar{\lambda}_2((G_{i-1} \times G_{i-1})^2 \textcircled{z} G) \leq \bar{\lambda}_2(G_{i-1})^2 + \bar{\lambda}_2(G)$). Computing the neighbor of a vertex in $G_i$ boils down to a constant number of such computations regarding $G_{i-1}$, but due to the tensor product operation the depth of the recursion is only double-logarithmic in the size of the final graph (and hence logarithmic in the length of the description of vertices in it).

**Digest.** In the first construction, the zig-zag product was used both in order to increase the size of the graph and to reduce its degree. However, as indicated by the second construction (where the tensor product of graphs is the main vehicle for increasing the size of the graph), the primary effect of the zig-zag product is to reduce the degree, and the increase in the size of the graph is merely a side-effect (which is actually undesired in Section 5.2.4). In both cases, graph squaring is used in order to compensate for the modest increase in the relative eigenvalue caused by the zig-zag product. In retrospect, the second construction is the "correct" one, because it decouples three different effects, and uses a natural operation to obtain each of them: Increasing the size of the graph is obtained by tensor product of graphs (which in turn increases the degree), a degree reduction is obtained by the zig-zag product (which in turn increases the relative eigenvalue), and graph squaring is used in order to reduce the relative eigenvalue.

**Stronger bound regarding the effect of the zig-zag product.** In the foregoing description we relied on the fact, proved in [180], that the relative eigenvalue of the zig-zag product is upper-bounded by the sum of the relative eigenvalues of the two graphs. Actually, a stronger upper-bound is proved in [180]: For $g(x, y) = (1 - y^2) \cdot x/2$, it holds that

$$\begin{aligned} \bar{\lambda}_2(G' \textcircled{z} G) &\leq g(\bar{\lambda}_2(G'), \bar{\lambda}_2(G)) + \sqrt{g(\bar{\lambda}_2(G'), \bar{\lambda}_2(G))^2 + \bar{\lambda}_2(G)^2} \quad \text{(E.9)} \\ &\leq 2g(\bar{\lambda}_2(G'), \bar{\lambda}_2(G)) + \bar{\lambda}_2(G) \\ &= (1 - \bar{\lambda}_2(G)^2) \cdot \bar{\lambda}_2(G') + \bar{\lambda}_2(G). \end{aligned}$$

Thus, we get $\bar\lambda_2(G' \textcircled{z} G) \leq \bar\lambda_2(G') + \bar\lambda_2(G)$. Furthermore, Eq. (E.9) yields a non-trivial bound for any $\bar\lambda_2(G'), \bar\lambda_2(G) < 1$, even in case $\bar\lambda_2(G')$ is very close to 1 (as in the proof of Theorem 5.6). Specifically, Eq. (E.9) is upper-bounded by

$$g(\bar\lambda_2(G'), \bar\lambda_2(G)) + \sqrt{\left(\frac{1 - \bar\lambda_2(G)^2}{2}\right)^2 + \bar\lambda_2(G)^2}$$

$$= \frac{(1 - \bar\lambda_2(G)^2) \cdot \bar\lambda_2(G')}{2} + \frac{1 + \bar\lambda_2(G)^2}{2}$$

$$= 1 - \frac{(1 - \bar\lambda_2(G)^2) \cdot (1 - \bar\lambda_2(G'))}{2} \tag{E.10}$$

Thus, $1 - \bar\lambda_2(G' \textcircled{z} G) \geq (1 - \bar\lambda_2(G)^2) \cdot (1 - \bar\lambda_2(G'))/2$. In particular, if $\bar\lambda_2(G) < 1/\sqrt{3}$ then $1 - \bar\lambda_2(G' \textcircled{z} G) > (1 - \bar\lambda_2(G'))/3$. This fact plays an important role in the proof of Theorem 5.6.