

Chapter 8

TOLERANT LOCALLY TESTABLE CODES

In this chapter, we revisit the notion of local testers (as defined in Section 2.3) that was the focus of Chapter 7.

8.1 Introduction

In the definition of LTCs, there is no requirement on the tester for input strings that are very close to a codeword (it has to reject “far” away received words). This “asymmetry” in the way the tester accepts and rejects an input reflects the way Probabilistically Checkable Proofs (or PCPs) [6, 5] are defined, where we only care about accepting perfectly correct proofs with high probability. However, the crux of error-correcting codes is to tolerate and *correct* a few errors that could occur during transmission of the codeword (and not just be able to detect errors). In this context, the fact that a tester can reject received words with few errors is not satisfactory. A more desirable (and stronger) requirement in this scenario would be the following— we would like the tester to make a quick decision on whether or not the purported codeword is close to any codeword. If the tester declares that there is probably a close-by codeword, we then use a decoding algorithm to decode the received word. If on the other hand, the tester rejects, then we assume with high confidence that the received word is far away from all codewords and not run our expensive decoding algorithm.

In this chapter, we introduce the concept of *tolerant testers*. These are testers which reject (w.h.p) received words far from every codeword (like the “standard” local testers) and accept (w.h.p) close-by received words (unlike the “standard” ones which only need to accept codewords). We will refer to codes that admit a tolerant tester as tolerant LTCs. In particular we get tolerant testers that (i) make $O(1)$ queries and work with codes of near constant rate codes and (ii) make sub-linear number of queries and work with codes of constant rate.

8.2 Preliminaries

Recall that for any two vectors $u, v \in [q]^n$, $\delta(u, v)$ denotes the (relative) Hamming distance between them. We will abuse the notation a bit and for any $S \subseteq [q]^n$, use $\delta(u, S)$ to denote $\min_{v \in S} \delta(u, v)$. We now formally define a *tolerant* tester.

Definition 8.1. For any linear code \mathcal{C} over \mathbb{F}_q of block length n and distance d , and $0 \leq c_1 \leq c_2 \leq 1$, a (c_1, c_2) -tolerant tester T for \mathcal{C} with query complexity $p(n)$ (or simply p when

the argument is clear from the context) is a probabilistic polynomial time oracle Turing machine such that for every vector $v \in \mathbb{F}_q^n$:

1. If $\delta(v, \mathcal{C}) \leq \frac{c_1 d}{n}$, T upon oracle access to v accepts with probability at least $\frac{2}{3}$ (tolerance),
2. If $\delta(v, \mathcal{C}) > \frac{c_2 d}{n}$, T rejects with probability at least $\frac{2}{3}$ (soundness),
3. T makes $p(n)$ probes into the string (oracle) v .

A code is said to be (c_1, c_2, p) -testable if it admits a (c_1, c_2) -tolerant tester of query complexity $p(\cdot)$.

A tester has *perfect completeness* if it accepts any codeword with probability 1. As pointed out earlier, local testers are just $(0, c_2)$ -tolerant testers with perfect completeness. We will refer to these as *standard* testers henceforth. Note that our definition of tolerant testers is per se not a generalization of standard testers since we do not require perfect completeness for the case when the input v is a codeword. However, all our constructions will inherit this property from the standard testers we obtain them from.

Recall one of the applications of tolerant testers mentioned earlier: a tolerant tester is used to decide if the expensive decoding algorithm should be used. In this scenario, one would like to set the parameters c_1 and c_2 such that the tester is tolerant up to the decoding radius. For example, if we have an unique decoding algorithm which can correct up to $\frac{d}{2}$ errors, a particularly appealing setting of parameters would be $c_1 = \frac{1}{2}$ and c_2 as close to $\frac{1}{2}$ as possible. However, we would not be able to achieve such large c_1 . In general we will aim for positive constants c_1 and c_2 with $\frac{c_2}{c_1}$ being as small as possible while minimizing $p(n)$.

One might hope that the existing standard testers could also be tolerant testers. We give a simple example to illustrate the fact that this is not the case in general. Consider the tester for the Reed-Solomon (RS) codes of dimension $k+1$: pick $k+2$ points uniformly at random and check if the degree k univariate polynomial obtained by interpolating on the first $k+1$ points agrees with the input on the last point. It is well known that this is a standard tester [96]. However, this is not a tolerant tester. Assume we have an input which differs from a degree k polynomial in only one point. Thus, for $\binom{n-1}{k+1}$ choices of $k+2$ points, the tester would reject, that is, the rejection probability is $\frac{\binom{n-1}{k+1}}{\binom{n}{k+2}} = \frac{k+2}{n}$ which is greater than $\frac{1}{3}$ for high rate RS codes.

Another pointer towards the inherent difficulty in coming up with a tolerant tester is the work of Fischer and Fortnow [39] which shows that there are certain boolean properties which have a standard tester with constant number of queries but for which every tolerant tester requires at least $n^{\Omega(1)}$ queries.

In this chapter, we examine existing standard testers and convert some standard testers into tolerant ones. In Section 8.3 we record a few general facts which will be useful in

performing this conversion. The ultimate goal, if this can be realized at all, would be to construct tolerant LTCs of constant rate which can be tested using $O(1)$ queries (we remark that such a construction has not been obtained even without the requirement of tolerance). In this work, we show that we can achieve either constant number of queries with slightly sub-constant rate (Section 8.4) as well as constant rate with sub-linear number of queries (Section 8.5.1). That is, something non-trivial is possible in both the domains: (a) constant rate, and (b) constant number of queries. Specifically, in Section 8.4 we discuss binary codes which encode k bits into codewords of length $n = k \cdot \exp(\log^\varepsilon k)$ for any $\varepsilon > 0$, and can be tolerant tested using $O(1/\varepsilon)$ queries. In Section 8.5.1, following [14], we will study the simple construction of LTCs using products of codes — this yields asymptotically good codes which are tolerant testable using a sub-linear number n^γ of queries for any desired $\gamma > 0$. An interesting common feature of the codes in Section 8.4 and 8.5.1 is that they can be constructed from any code that has good distance properties and which in particular need not admit a local tester with sub-linear query complexity. In Section 8.6 we discuss the tolerant testability of Reed-Muller codes, which were considered in Chapter 7.

The overall message from this chapter is that a lot of the work on locally testable code constructions extends fairly easily to also yield tolerant locally testable codes. However, there does not seem to be a generic way to “compile” a standard tester to a tolerant tester for an arbitrary code.

8.3 General Observations

In this section we will spell out some general properties of tolerant testers and subsequently use them to design tolerant testers for some existing codes. All the testers we refer to are *non-adaptive testers* which decide on the locations to query all at once based only on the random choices. The motivation for the definition below will be clear in Section 8.4.

Definition 8.2. Let $0 < \alpha \leq 1$. A tester T is $(\langle s_1, q_1 \rangle, \langle s_2, q_2 \rangle, \alpha)$ -smooth if there exists a set $A \subseteq [n]$ where $|A| = \alpha n$ with the following properties:

- T queries at most q_1 points in A , and for every $x \in A$, the probability that each of these queries equals location x is at most $\frac{s_1}{|A|}$, and
- T queries at most q_2 points in $[n] \setminus A$, and for every $x \in [n] \setminus A$, the probability that each of these queries equals location x is at most $\frac{s_2}{n - |A|}$.

As a special case a $(\langle 1, q \rangle, \langle 0, 0 \rangle, 1)$ -smooth tester makes a total of q queries each of them distributed uniformly among the n possible probe points. The following lemma follows easily by an application of the union bound.

Lemma 8.1. For any $0 < \alpha < 1$, a $(\langle s_1, q_1 \rangle, \langle s_2, q_2 \rangle, \alpha)$ -smooth $(0, c_2)$ -tolerant tester T with perfect completeness is a (c_1, c_2) -tolerant tester T' , where $c_1 = \frac{n\alpha(1-\alpha)}{3d \max\{q_1 s_1(1-\alpha), q_2 s_2 \alpha\}}$.

Proof. The soundness follows from the assumption on T . Assume $\delta(v, \mathcal{C}) \leq \frac{c_1 d}{n}$ and let $f \in \mathcal{C}$ be the closest codeword to v . Suppose that f differs from v in a set A' of yd places among locations in A , and a set B' of $(\beta - y)d$ places among locations in $[n] \setminus A$, where we have $\beta \leq c_1$ and $0 \leq y \leq \beta$. The probability that any of the at most q_1 (resp. q_2) queries of T into A (resp. $[n] \setminus A$) falls in A' (resp. B') is at most $\frac{s_1 y d}{\alpha n}$ (resp. $\frac{s_2 (\beta - y) d}{(1 - \alpha) n}$). Clearly, whenever T does not query a location in $A' \cup B'$, it accepts (since T has perfect completeness). Thus, an easy calculation shows that the probability that T rejects v is at most

$$\frac{c_1 d}{n} \max\left\{\frac{s_1 q_1}{\alpha}, \frac{s_2 q_2}{1 - \alpha}\right\}$$

which is $1/3$ for the choice of c_1 stated in the lemma. \square

The above lemma is not useful for us unless the relative distance and the number of queries are constants. Next we sketch how to design tolerant testers from existing *robust* testers with certain properties. We first recall the definition of robust testers from [14].

A standard tester T has two inputs: an oracle for the received word v and a random string s . Depending on s , T generates q query positions i_1, \dots, i_q , fixes a circuit C_s and then accepts if $C_s(v_f(s)) = 1$ where $v_f(s) = \langle v_{i_1}, \dots, v_{i_q} \rangle$. The robustness of T on inputs v and s , denoted by $\rho^T(v, s)$, is defined to be the minimum, over all strings y such that $C_s(y) = 1$, of $\delta(v_f(s), y)$. The expected robustness of T on v is the expected value of $\rho^T(v, s)$ over the random choices of s and would be denoted by $\rho^T(v)$.

A standard tester T is said to be c -robust for \mathcal{C} if for every $v \in \mathcal{C}$, the tester accepts with probability 1, and for every $v \in \mathbb{F}_q^n$, $\delta(v, \mathcal{C}) \leq c \cdot \rho^T(v)$.

The tolerant version T' of the standard c -robust tester T is obtained by accepting an oracle v on random input s , if $\rho^T(v, s) \leq \tau$ for some threshold τ . (Throughout the chapter τ will denote the threshold.) We will sometimes refer to such a tester as one with threshold τ . Recall that a standard tester T accepts if $\rho^T(v, s) = 0$. We next show that T' is sound.

The following lemma follows from the fact that T is c -robust:

Lemma 8.2. *Let $0 \leq \tau \leq 1$, and let $c_2 = \frac{(\tau+2)cn}{3d}$. For any $v \in \mathbb{F}_q^n$, if $\delta(v, \mathcal{C}) > \frac{c_2 d}{n}$, then the tolerant tester T' with threshold τ rejects v with probability at least $\frac{2}{3}$.*

Proof. Let $v \in \mathbb{F}_q^n$ be such that $\delta(v, \mathcal{C}) > \frac{c_2 d}{n}$. By the definition of robustness, the expected robustness, $\rho^T(v)$ is at least $\frac{c_2 d}{nc}$, and thus at least $(\tau + 2)/3$ by the choice of c_2 . By the standard averaging argument, we can have $\rho^T(v, s) \leq \tau$ on at most a fraction $1/3$ of the of the random choices of s for T (and hence T'). Therefore, $\rho^T(v, s) > \tau$ with probability at least $2/3$ over the choice of s and thus T' rejects v with probability at least $2/3$. \square

We next mention a property of the query pattern of T which would make T' tolerant. Let S be the set of all possible choices for the random string s . Further for each s , let $p^T(s)$ be the set of positions queried by T .

Definition 8.3. *A tester T has a partitioned query pattern if there exists a partition $s_1 \cup \dots \cup S_m$ of the random choices of T for some m , such that for every i ,*

- $\cup_{s \in S_i} p^T(s) = \{1, 2, \dots, n\}$, and
- For all $s, s' \in S_i$, $p^T(s) \cap p^T(s') = \emptyset$ if $s \neq s'$.

Lemma 8.3. *Let T have a partitioned query pattern. For any $v \in \mathbb{F}_q^n$, if $\delta(v, \mathcal{C}) \leq \frac{c_1 d}{n}$, where $c_1 = \frac{nr}{3d}$, then the tolerant test T' with threshold τ rejects with probability at most $\frac{1}{3}$.*

Proof. Let S_1, \dots, S_m be the partition of S , the set of all random choices of the tester T . For each j , by the properties of S_j , $\sum_{s \in S_j} \rho^T(v, s) \leq \delta(v, \mathcal{C})$. By an averaging argument and by the assumption on $\delta(v, \mathcal{C})$ and the value of c_1 , at least $\frac{2}{3}$ fraction of the choices of s in S_j have $\rho^T(v, s) \leq \tau$ and thus, T' accepts. Recalling that S_1, \dots, S_m was a partition of S , for at least $\frac{2}{3}$ of the choices of s in S , T' accepts. This completes the proof. \square

8.4 Tolerant Testers for Binary Codes

One of the natural goals in the study of tolerant codes is to design explicit tolerant binary codes with constant relative distance and as large a rate as possible. In the case of standard testers, Ben-Sasson et al [11] give binary locally testable codes which map k bits to $k \cdot \exp(\log^\varepsilon k)$ bits for any $\varepsilon > 0$ and which are testable with $O(1/\varepsilon)$ queries. Their construction uses objects called PCPs of Proximity (PCPP) which they also introduce in [11]. In this section, we show that a simple modification to their construction yields tolerant testable binary codes which map k bits to $k \cdot \exp(\log^\varepsilon k)$ bits for any $\varepsilon > 0$. We note that a similar modification is used by Ben-Sasson et al to give a relaxed locally decodable codes [11] but with worse parameters (specifically they give codes with block length $k^{1+\varepsilon}$).

8.4.1 PCP of Proximity

We start with the definition¹ of a Probabilistic Checkable proof of Proximity (PCPP). A pair language is simply a language whose elements are naturally a pair of strings, i.e., it is some collection of strings (x, y) . A notable example is $\text{CIRCUITVAL} = \{\langle C, a \rangle \mid \text{Boolean circuit } C \text{ evaluates to 1 on assignment } a\}$.

Definition 8.4. *Fix $0 \leq \gamma \leq 1$. A probabilistic verifier V is a PCPP for a pair language L with proximity parameter γ and query complexity $q(\cdot)$ if the following conditions hold:*

- (Completeness) *If $(x, y) \in L$ then there exists a proof π such that V accepts by accessing the oracle $y \circ \pi$ with probability 1.*
- (Soundness) *If y is γ -far from $L(x) = \{y \mid (x, y) \in L\}$, then for all proofs π , V accepts by accessing the oracle $y \circ \pi$ with probability strictly less than $\frac{1}{4}$.*

¹The definition here is a special case of the general PCPP defined in [11] which would be sufficient for our purposes.

- (Query complexity) For any input x and proof π , V makes at most $q(|x|)$ queries in $y \circ \pi$.

Note that a PCPP differs from a standard PCP in that it has a more relaxed soundness condition but its queries into part of the input y are also counted in its query complexity.

Ben-Sasson et. al. give constructions of PCPPs with the following guarantees:

Lemma 8.4 ([11]). *Let $\varepsilon > 0$ be arbitrary. There exists a PCP of proximity for the pair language $\text{CIRCUITVAL} = \{(C, x) \mid C \text{ is a boolean circuit and } C(x) = 1\}$ whose proof length, for inputs circuits of size s , is at most $s \cdot \exp(\log^{\varepsilon/2} s)$ and for $t = \frac{2 \log \log s}{\log \log \log s}$ the verifier of proximity has query complexity $O(\max\{\frac{1}{\gamma}, \frac{1}{\varepsilon}\})$ for any proximity parameter γ that satisfies $\gamma \geq \frac{1}{t}$. Furthermore, the queries of the verifier are non-adaptive and each of the queries which lie in the input part x are uniformly distributed among the locations of x .*

The fact that the queries to the input part are uniformly distributed follows by an examination of the verifier construction in [11]. In fact, in the extended version of that paper, the authors make this fact explicit and use it in their construction of relaxed locally decodable codes (LDCs). To achieve a tolerant LTC using the PCPP, we will need all queries of the verifier to be somewhat uniformly or smoothly distributed. We will now proceed to make the queries of the PCPP verifier that fall into the “proof part” π near-uniform. This will follow a fairly general method suggested in [11] to smoothen out the query distribution, which the authors used to obtain relaxed locally decodable codes from the PCPP. We will obtain tolerant LTCs instead, and in fact will manage to do so without a substantial increase in the encoding length (i.e., the encoding length will remain $k \cdot 2^{\log^\varepsilon k}$). On the other hand, the best encoding length achieved for relaxed LDCs in [11] is $k^{1+\varepsilon}$ for constant $\varepsilon > 0$. We begin with the definition of a mapping that helps smoothen out the query distribution.

Definition 8.5. *Given any $v \in \mathbb{F}_q^n$ and $\vec{p} = \langle p_i \rangle_{i=1}^n$ with $p_i \geq 0$ for all $i \in [n]$ and $\sum_{i=1}^n p_i = 1$, we define the mapping $\text{Repeat}(\cdot, \cdot)$ as follows: $\text{Repeat}(v, \vec{p}) \in \mathbb{F}_q^{n'}$ such that v_i is repeated $\lfloor 4np_i \rfloor$ times in $\text{Repeat}(v, \vec{p})$ and $n' = \sum_{i=1}^n \lfloor 4np_i \rfloor$.*

We now show why the mapping is useful. A similar fact appears in [11], but for the sake of completeness we present its proof here.

Lemma 8.5. *For any $v \in \mathbb{F}_q^n$ let a non-adaptive verifier T (with oracle access to v) make $q(n)$ queries and let p_i be the probability that each of these queries probes location $i \in [n]$. Let $c_i = \frac{1}{2n} + \frac{p_i}{2}$ and $\vec{c} = \langle c_i \rangle_{i=1}^n$. Consider the map $\text{Repeat}(v, \vec{c}) : \mathbb{F}_q^n \rightarrow \mathbb{F}_q^{n'}$. Then there exists another tester T' for strings of length n' with the following properties:*

1. T' makes $2q(n)$ queries on $v' = \text{Repeat}(v, \vec{c})$ each of which probes location j , for any $j \in [n']$, with probability at most $\frac{2}{n}$, and
2. for every $v \in \mathbb{F}_q^n$, the decision of T' on v' is identical to that of T on v . Further, $3n < n' \leq 4n$.

Proof. We first add q dummy queries to T each of which are uniformly distributed, and then permute the $2q$ queries in a random order. Note that each of the $2q$ queries is now identically distributed. Moreover, any position in v is probed with probability at least $\frac{1}{2n}$ for each of the $2q$ queries. For the rest of the proof we will assume that T makes $2q$ queries for each of which any $i \in [n]$ is probed with probability $c_i = \frac{p_i}{2} + \frac{1}{2n}$. Let $r_i = \lfloor 4nc_i \rfloor$. Note that $r_i \leq 4nc_i$ and $r_i > 4nc_i - 1$. Recalling that $n' = \sum_{i=1}^n r_i$ and $\sum_{i=1}^n c_i = 1$, we have $3n < n' \leq 4n$.

T' just simulates T in the following manner: if T queries v_i for any $i \in [n]$, T' queries one of the r_i copies of v_i in v' uniformly at random. It is clear that the decision of T' on $v' = \text{Repeat}(v, \vec{c})$ is identical to that of T on v . We now look at the query distribution of T' . T' queries any $j \in [n']$, where $v'_j = v_i$, with probability $p'_j = c_i \cdot \frac{1}{r_i}$. Recalling the lower bound on r_i , we have $p'_j \leq \frac{c_i}{4nc_i - 1}$ which is at most $\frac{1}{2n}$ since clearly $c_i \geq \frac{1}{2n}$. We showed earlier that $n' \leq 4n$ which implies $p'_j \leq \frac{2}{n'}$ as required. \square

One might wonder if we can use Lemma 8.5 to smoothen out the queries made by the verifier of an arbitrary LTC to obtain a tolerant LTC. That is, whether the above allows one to compile the verifier for any LTC in a black-box manner to obtain a tolerant verifier. We will now argue (informally) that this technique alone will not work. Let C_1 be an $[n, k, d]_q$ LTC with a standard tester T_1 that makes q identically distributed queries with distribution p_i , $1 \leq i \leq n$, such that $p_i \geq 1/2n$ for each i . Create a new $[n+1, k, d]_q$ code C_2 whose $(n+1)$ 'th coordinate is just a copy of the n 'th coordinate, i.e., corresponding to each codeword $(c_1, c_2, \dots, c_n) \in \mathbb{F}_q^n$ of C_1 , we will have a codeword $(c_1, c_2, \dots, c_n, c_n) \in \mathbb{F}_q^{n+1}$ of C_2 . Consider the following tester T_2 for C_2 : Given oracle access to $v \in \mathbb{F}_q^{n+1}$, with probability $1/2$ check whether $v_n = v_{n+1}$, and with probability $1/2$ run the tester T_1 on the first n coordinates of v . Clearly, T_2 is a standard tester for C_2 .

Now, consider what happens in the conversion procedure of Lemma 8.5 to get (C', T') from (C_2, T_2) . Note that by Lemmas 8.5 and 8.3, T' is tolerant. Let $\vec{q} = (q_1, \dots, q_{n+1})$ be the query distribution of T_2 . Since T_2 queries (v_n, v_{n+1}) with probability $1/2$, the combined number of locations of $v' = \text{Repeat}(v, \vec{q})$ corresponding to v_n, v_{n+1} will be about $1/2$ of the total length n' . Now let v' be obtained from a codeword of C' by corrupting just these locations. The tester T' will accept such a v' with probability at least $1/2$, which contradicts the soundness requirement since v' is $1/2$ -far from C' . Therefore, using the behavior of the original tester T_2 as just a black-box, we cannot in general argue that the construction of Lemma 8.5 maintains good soundness.

Applying the transformation of Lemma 8.5 to the proximity verifier and proof of proximity of Lemma 8.4, we conclude the following.

Proposition 8.6. *Let $\varepsilon > 0$ be arbitrary. There exists a PCP of proximity for the pair language $\text{CIRCUITVAL} = \{(C, x) \mid C \text{ is a boolean circuit and } C(x) = 1\}$ with the following properties:*

1. *The proof length, for inputs circuits of size s , is at most $s \cdot \exp(\log^{\varepsilon/2} s)$, and*

2. for $t = \frac{2 \log \log s}{\log \log \log s}$ the verifier of proximity has query complexity $O(\max\{\frac{1}{\gamma}, \frac{1}{\varepsilon}\})$ for any proximity parameter γ that satisfies $\gamma \geq \frac{1}{t}$.

Furthermore, the queries of the verifier are non-adaptive with the following properties:

1. Each query made to one of the locations of the input x is uniformly distributed among the locations of x , and
2. each query to one of the locations in the proof of proximity π probes each location with probability at most $2/|\pi|$ (and thus is distributed nearly uniformly among the locations of π).

8.4.2 The Code

We now outline the construction of the locally testable code from [11]. The idea behind the construction is to make use of a PCPP to aid in checking if the received word is a codeword is far away from being one. Details follow.

Suppose we have a binary code $C_0 : \{0, 1\}^k \rightarrow \{0, 1\}^m$ of distance d defined by a parity check matrix $H \in \{0, 1\}^{(m-k) \times m}$ that is sparse, i.e., each of whose rows has only an absolute constant number of 1's. Such a code is referred to as a low-density parity check code (LDPC). For the construction below, we will use any such code which is asymptotically good (i.e., has rate k/m and relative distance d/m both positive as $m \rightarrow \infty$). Explicit constructions of such codes are known using expander graphs [95]. Let V be a verifier of a PCP of proximity for membership in C_0 ; more precisely, the proof of proximity of an input string $w \in \{0, 1\}^m$ will be a proof that $\tilde{C}_0(w) = 1$ where \tilde{C}_0 is a linear-sized circuit which performs the parity checks required by H on w (the circuit will have size $O(m) = O(k)$ since H is sparse and C_0 has positive rate). Denote by $\pi(x)$ be the proof of proximity guaranteed by Proposition 8.6 for the claim that the input $C_0(x)$ is a member of C_0 (i.e., satisfies the circuit \tilde{C}_0). By Proposition 8.6 and fact that the size of \tilde{C}_0 is $O(k)$, the length of $\pi(x)$ can be made at most $k \exp(\log^{\varepsilon/2} k)$.

The final code is defined as $\mathcal{C}_1(x) = (C_0(x)^t, \pi(x))$ where $t = \frac{(\log k - 1)|\pi(x)|}{|C_0(x)|}$. The repetition of the code part $C_0(x)$ is required in order to ensure good distance, since the length of the proof part $\pi(x)$ typically dominates and we have no guarantee on how far apart $\pi(x_1)$ and $\pi(x_2)$ for $x_1 \neq x_2$ are.

For the rest of this section let ℓ denote the proof length. The tester T_1 for \mathcal{C}_1 on an input $w = (w_1, \dots, w_t, \pi) \in \{0, 1\}^{tm+\ell}$ picks $i \in [t]$ at random and runs the PCPP verifier V on $w_i \circ \pi$. It also performs a few rounds of the following consistency checks: pick $i_1, i_2 \in [t]$ and $j_1, j_2 \in [m]$ at random and check if $w_{i_1}(j_1) = w_{i_2}(j_2)$. Ben-Sasson et al in [11] show that T_1 is a standard tester. However, T_1 need not be a tolerant tester. To see this, note that the proof part of \mathcal{C}_1 forms a $\frac{1}{\log k}$ fraction of the total length. Now consider a received word $w_{rec} = (w_0, \dots, w_0, \pi')$ where $w_0 \in C_0$ but π' is not a correct proof for w_0 being a valid

codeword in c_0 . Note that w_{rec} is close to \mathcal{C}_1 . However, T_1 is not guaranteed to accept w_{rec} with high probability.

The problem with the construction above was that the proof part was too small: a natural fix is to make the proof part a constant fraction of the codeword. We will show that this is sufficient to make the code tolerant testable. We also remark that a similar idea was used by Ben-Sasson et. al. to give efficient constructions for relaxed locally decodable codes [11].

Construction 8.1. Let $0 < \beta < 1$ be a parameter, $C_0 : \{0, 1\}^k \rightarrow \{0, 1\}^m$ be a good ² binary code and V be a PCP of proximity verifier for membership in C_0 . Finally let $\pi(x)$ be the proof corresponding to the claim that $C_0(x)$ is a codeword in C_0 . The final code is defined as $\mathcal{C}_2(x) = (C_0(x)^{r_1}, \pi(x)^{r_2})$ with $r_1 = \frac{(1-\beta)\log k|\pi(x)|}{m}$ and $r_2 = \beta \log k$.³

For the rest of the section the proof length $|\pi(x)|$ will be denoted by ℓ . Further the proximity parameter and the number of queries made by the PCPP verifier V would be denoted by γ_p and q_p respectively. Finally let ρ_0 denote the relative distance of the code C_0 .

The tester T_2 for \mathcal{C}_2 is also the natural generalization of T_1 . For a parameter q_r (to be instantiated later) and input $w = (w_1, \dots, w_{r_1}, \pi_1, \dots, \pi_{r_2}) \in \{0, 1\}^{r_1 m + r_2 \ell}$, T_2 does the following:

1. Repeat the next two steps twice.
2. Pick $i \in [r_1]$ and $j \in [r_2]$ randomly and run V on $w_i \circ \pi_j$.
3. Do q_r repetitions of the following: pick $i_1, i_2 \in [r_1]$ and $j_1, j_2 \in [m]$ randomly and check if $w_{i_1}(j_1) = w_{i_2}(j_2)$.

The following lemma captures the properties of the code \mathcal{C}_2 and its tester T_2 .

Lemma 8.7. *The code \mathcal{C}_2 in Construction 8.1 and the tester T_2 (with parameters β and q_r respectively) above have the following properties:*

1. The code \mathcal{C}_2 has block length $n = \log k \cdot \ell$ with minimum distance d lower bounded by $(1 - \beta)\rho_0 n$.
2. T_2 makes a total of $q = 2q_p + 4q_r$ queries.
3. T_2 is $(\langle 1, q \rangle, \langle 2, 2q_p \rangle, 1 - \beta)$ -smooth.

²This means that $m = O(k)$ and the encoding can be done by circuits of nearly linear size $s_0 = \tilde{O}(k)$.

³The factor $\log k$ overhead is overkill, and a suitably large constant will do, but since the proof length $|\pi(x)|$ will anyway be larger than $|x|$ by more than a polylogarithmic factor in the constructions we use, we can afford this additional $\log k$ factor and this eases the presentation somewhat.

4. T_2 is a (c_1, c_2) -tolerant tester with $c_1 = \frac{n\beta(1-\beta)}{6d \max\{(2q_r+q_p)\beta, 2(1-\beta)q_p\}}$ and $c_2 = \frac{n}{d}(\gamma_p + \frac{4}{q_r} + \beta)$.

Proof. From the definition of \mathcal{C}_2 , it has block length $n = r_1m + r_2\ell = \frac{(1-\beta)\ell \log k}{m} \cdot m + \beta \log k \cdot \ell = \log k \cdot \ell$. Further as \mathcal{C}_0 has relative distance ρ_0 , \mathcal{C}_2 has relative distance at least $\frac{r_1\rho_0m}{\ell \log k} = (1-\beta)\rho_0$.

T_2 makes the same number of queries as V which is q_p in Step 2. In Step 3, T_2 makes $2q_r$ queries. As T_2 repeats Steps 2 and 3 twice, we get the desired query complexity.

To show the smoothness of T_2 we need to define the appropriate subset $A \subset [n]$ such that $|A| = (1-\beta)n$. Let A be the set of indices with the code part: i.e. $A = [r_1m]$. T_2 makes $2q_r$ queries in A in Step 3 each of which is uniformly distributed. Further by Proposition 8.6, T_2 in step 2 makes at most q_p queries in A which are uniformly distributed and at most q_p queries in $[n] \setminus A$ each of which are within a factor 2 of being queried uniformly at random. To complete the proof of property 3 note that T_2 repeats step 2 and 3 twice.

The tolerance of T_2 follows from property 3 and Lemma 8.1. For the soundness part note that if $w = (w_1, \dots, w_{r_1}, \pi_1, \dots, \pi_{r_2}) \in \{0, 1\}^{r_1m+r_2\ell}$ is γ -far from \mathcal{C}_2 then $w' = (w_1, \dots, w_{r_1})$ is at least $\frac{\gamma m - r_2\ell}{n} = \frac{\gamma m - \beta n}{n} = \gamma - \beta$ far from the repetition code $\mathcal{C}' = \{C_0(x)^{r_1} | x \in \{0, 1\}^k\}$. For $\gamma = c_2d/n$ with the choice of c_2 in the lemma, we have $\gamma - \beta \geq \gamma_p + 4/q_r$. The rest of the proof just follows the proof in [11] (also see [68, Chap. 12]) of the soundness of the tester T_1 for the code \mathcal{C}_1 — for the sake of completeness we complete the proof here. We will show that one invocation of Steps 2 and 3 results in T_2 accepting w with probability strictly less than $\frac{1}{2}$. The two repetitions of Steps 2 and 3 reduces this error to at most $\frac{1}{4}$.

Let $u \in \{0, 1\}^m$ be the string such that u^t is the “repetition sequence” that is closest to w' , that is one that minimizes $\Delta(w', u^t) = \sum_{i=1}^{r_1} \Delta(w_{i_1}, u)$. We now consider two cases:

- **Case 1:** $\Delta(w', u^t) \geq r_1m/q_r$. In this case, a single execution of the test in Step 3 rejects with probability

$$\begin{aligned} \mathbb{E}_{i_1, i_2 \in [r_1]} [\Delta(w_{i_1}, w_{i_2})/m] &= \frac{1}{m(r_1)^2} \sum_{i_2} \sum_{i_1} \Delta(w_{i_1}, w_{i_2}) \\ &\geq \frac{1}{m(r_1)^2} \sum_{i_2} \sum_{i_1} \Delta(w_{i_1}, u) \\ &= \frac{1}{mr_1} \sum_{i_1=1}^{r_1} \Delta(w_{i_1}, u) \\ &= \Delta(w', u^t)/(mr_1) \\ &\geq 1/q_r, \end{aligned}$$

where the first inequality follows from the choice of u and the second inequality

follows from the case hypothesis. Thus, after q_r repetitions the test will accept with probability $(1 - 1/q_r)^{q_r} < 1/e < 1/2$.

- **Case 2:** $\Delta(w', u^t) < r_1 m/q_r$. In this case we have the following (where for any subset S of vectors and a vector u , we will use $\Delta(u, S) = \min_{v \in S} \Delta(u, v)$):

$$\frac{\Delta(u, C_0)}{r_1} = \frac{\Delta(u^t, C')}{r_1 m} \geq \frac{\Delta(w', C') - \Delta(w', u^t)}{r_1 m} \geq \gamma_p + 4/q_r - 1/q_r = \gamma_p + 3/q_r, \quad (8.1)$$

where the first inequality follows from the triangle inequality and the last inequality follows from the case hypothesis (and the fact that w' is $\gamma_p + 4/q_r$ -far from C'). Now by the case hypothesis, for an average i , $\Delta(w_i, u) \leq m/q_r$. Thus, by a Markov argument, at most one thirds of the w_i 's are $3/q_r$ -far from u . Since u is $\gamma_p + 3/q_r$ -far from C_0 (by (8.1)), this implies (along with triangle inequality) that for at least two thirds of the w_i 's are γ_p -far from C_0 . Thus, by the property of the PCPP, for each such w_i the test in Step 2 should accept with probability at most $1/4$. Thus the total acceptance probability in this case is at most $\frac{1}{3} \cdot 1 + \frac{2}{4} \cdot \frac{1}{4} = \frac{1}{2}$, as desired.

Thus, in both cases the tester accepts w with probability at most $1/2$, as required. \square

Fix any $0 < \gamma < 1$ and let $\beta = \frac{\gamma}{2}$, $\gamma_p = \frac{\gamma}{6}$, $q_r = \frac{12}{\gamma}$. With these settings we get $\gamma_p + \frac{4}{q_r} + \beta = \gamma$ and $q_p = O(\frac{1}{\gamma})$ from Proposition 8.6 with the choice $\varepsilon = 2\gamma$. Finally, $q = 2q_p + 4q_r = O(\frac{1}{\gamma})$. Substituting the parameters in c_2 and c_1 , we get $c_2 = \frac{\gamma m}{d}$ and

$$\frac{c_1 d}{n} = \frac{\gamma}{24 \max\{\gamma(q_r + q_p/2), (2 - \gamma)q_p\}} = \Omega(\gamma^2).$$

Also note that the minimum distance $d \geq (1 - \beta)\rho_0 n = (1 - \frac{\gamma}{2})\rho_0 n \geq \frac{\rho_0}{2}n$. Thus, we have the following result for tolerant testable binary codes.

Theorem 8.1. *There exists an absolute constant $\alpha_0 > 0$ such that for every γ , $0 < \gamma < 1$, there exists an explicit binary linear code $\mathcal{C} : \{0, 1\}^k \rightarrow \{0, 1\}^n$ where $n = k \cdot \exp(\log^\gamma k)$ with minimum distance $d \geq \alpha_0 n$ which admits a (c_1, c_2) -tolerant tester with $c_2 = O(\gamma)$, $c_1 = \Omega(\gamma^2)$ and query complexity $O(\frac{1}{\gamma})$.*

The claim about explicitness follows from the fact that the PCPP of Lemma 8.4 and hence Proposition 8.6 has an explicit construction. The claim about linearity follows from the fact that the PCPP for CIRCUITVAL is a linear function of the input when the circuit computes linear functions — this aspect of the construction is discussed in detail in Chapter 9 in [68].

8.5 Product of Codes

Tensor product of codes (or just *product of codes*) is simple way to construct new codes from existing codes such that the constructed codes have testers with sub-linear query complexity even though the original code need not admit a sub-linear complexity tester [14]. We start with the definition of product of codes.

Definition 8.6 (Tensor Product of Codes). *Given \mathcal{C}_1 and \mathcal{C}_2 that are $[k_1, n_1, d_1]$ and $[k_2, n_2, d_2]$ codes, their tensor product, denoted by $\mathcal{C}_1 \otimes \mathcal{C}_2$, consists of $n_2 \times n_1$ matrices such that every row of the matrix is a codeword in \mathcal{C}_1 and every column is a codeword in \mathcal{C}_2 .*

It is well known that $\mathcal{C}_1 \otimes \mathcal{C}_2$ is an $[n_1 n_2, k_1 k_2, d_1 d_2]$ code.

A special case in which we will be interested is when $\mathcal{C}_1 = \mathcal{C}_2 = \mathcal{C}$. In such a case, given an $[n, k, d]_q$ code \mathcal{C} , the product of \mathcal{C} with itself, denoted by \mathcal{C}^2 , is a $[n^2, k^2, d^2]_q$ code such that a codeword (viewed as a $n \times n$ matrix) restricted to any row or column is a codeword in \mathcal{C} . It can be shown that this is equivalent to the following [100]. Given the $k \times n$ generator matrix M of \mathcal{C} , \mathcal{C}^2 is precisely the set of matrices in the set $\{M^T \cdot X \cdot M \mid X \in F_q^{k \times k}\}$.

8.5.1 Tolerant Testers for Tensor Products of Codes

A very natural test for \mathcal{C}^2 is to randomly choose a row or a column and then check if the restriction of the received word on that row or column is a codeword in \mathcal{C} (which can be done for example by querying all the n points in the row or column). Unfortunately, as we will see in Section 8.5.2, this test is not robust in general.

Ben-Sasson and Sudan in [14] considered the more general product of codes \mathcal{C}^t for $t \geq 3$ (where \mathcal{C}^t denotes \mathcal{C} tensored with itself $t - 1$ times) along with the following general tester: Choose at random $b \in \{1, \dots, t\}$ and $i \in \{1, \dots, n\}$ and check if b^{th} coordinate of the received word (which is an element of $\mathbb{F}_q^{n^t}$) when restricted⁴ to i is a codeword in \mathcal{C}^{t-1} . It is shown in [14] that this test is robust, in that if a received word is far from \mathcal{C}^t , then many of the tested substrings will be far from \mathcal{C}^{t-1} . This tester lends itself to recursion: the test for \mathcal{C}^{t-1} can be reduced to a test for \mathcal{C}^{t-2} and so on till we need to check whether a word in $\mathbb{F}_q^{n^2}$ is a codeword of \mathcal{C}^2 . This last check can be done by querying all the n^2 points, out of the n^t points in the original received word, thus leading to a sub-linear query complexity. As shown in [14], the reduction can be done in $\log t$ stages by the standard halving technique.

Thus, even though \mathcal{C} might not have a tester with a small query complexity, we can test \mathcal{C}^t with a polylogarithmic number of queries.

We now give a tolerant version of the test for product of codes given by Ben-Sasson and Sudan [14]. In what follows $t \geq 4$ will be a power of two. As mentioned above the tester T for the tensor product \mathcal{C}^t reduces the test to checking if some restriction of the given string belong to \mathcal{C}^2 . For the rest of this section, with a slight abuse of notation let $v_f \in \mathbb{F}_q^{n^2}$ denote

⁴For the $t = 2$ case b signifies either row or column and i denotes the row/column index.

the final restriction being tested. In what follows we assume that by looking at all points in any $v \in \mathbb{F}_q^{n^2}$ one can determine if $\delta(v, \mathcal{C}^2) \leq \tau$ in time polynomial in n^2 .

The tolerant version of the test of [14] is a simple modification as mentioned in Section 8.3: reduce the test on \mathcal{C}^t to \mathcal{C}^2 as in [14] and then accept if v_f is τ -close to \mathcal{C}^2 .

First we make the following observation about the test in [14]. The test recurses $\log t$ times to reduce the test to \mathcal{C}^2 . At step l , the test chooses a random coordinate b_l (this will just be a random bit) and fixes the value of the b_l^{th} coordinate of the current \mathcal{C}^{2^l} to an index i_l (where i_l takes values in the range $1 \leq i_l \leq n^{2^{l-1}}$). The key observation here is that for each fixed choice of $b_1, \dots, b_{\log t}$, distinct choices of $i_1, \dots, i_{\log t}$ correspond to querying disjoint sets n^2 points in the original $v \in \mathbb{F}_q^{n^2}$ string, which together form a partition of all coordinates of v . In other words, T has a *partitioned* query pattern, which will be useful to argue tolerance. For soundness, we use the results in [14], which show that their tester is $C^{\log t}$ -robust for $C = 2^{32}$.

Applying Lemmas 8.2 and 8.3, therefore, we have the following result:

Theorem 8.2. *Let $t \geq 4$ be a power of two and $0 < \tau \leq 1$. There exist $0 < c_1 < c_2 \leq 1$ with $\frac{c_2}{c_1} = C^{\log t}(1 + 2/\tau)$ such that the proposed tolerant tester for \mathcal{C}^t is a (c_1, c_2) -tolerant tester with query complexity $N^{2/t}$ where N is the block length of \mathcal{C}^t . Further, c_1 and c_2 are constants (independent of N) if t is a constant and \mathcal{C} has constant relative distance.*

Corollary 8.3. *For every $\gamma > 0$, there is an explicit family of asymptotically good binary linear codes which are tolerant testable using n^γ queries, where n is the block length of the concerned code. (The rate, relative distance and thresholds c_1, c_2 for the tolerant testing depend on γ .)*

8.5.2 Robust Testability of Product of Codes

Recall that a standard tester for a code is robust if for every received word which is far from being a codeword, the tester not only rejects the codeword with high probability but also with high probability the tester's local view of the received word is far from any accepting view (see Section 8.3 for a more formal definition).

As was mentioned before for the product code $\mathcal{C}_1 \otimes \mathcal{C}_2$, there is a natural tester (which we call $T_{\mathcal{C}_1 \otimes \mathcal{C}_2}$)— flip a coin; if it is heads check if a random row is a codeword in \mathcal{C}_1 ; if it is tails, check if a random column is a codeword in \mathcal{C}_2 . This test is indeed robust in a couple of special cases— for example, when both \mathcal{C}_1 and \mathcal{C}_2 are Reed-Solomon codes (see Section 8.6.1 for more details) and when both \mathcal{C}_1 and \mathcal{C}_2 are themselves tensor product of a code [14].

P. Valiant showed that there are linear codes \mathcal{C}_1 and \mathcal{C}_2 such that $\mathcal{C}_1 \otimes \mathcal{C}_2$ is not robustly testable [103]. Valiant constructs linear codes $\mathcal{C}_1, \mathcal{C}_2$ and a matrix v such that every row (and column) of v is “close” to some codeword in \mathcal{C}_1 (and \mathcal{C}_2) while v is “far” from every codeword in $\mathcal{C}_1 \otimes \mathcal{C}_2$ (where close and far are in the sense of hamming distance).

However, Valiant's construction *does not* work when \mathcal{C}_1 and \mathcal{C}_2 are the same code. In this section, we show a reduction from Valiant's construction to exhibit a code \mathcal{C} such that \mathcal{C}^2 is not robustly testable.

Preliminaries and Known Results

\mathcal{C} is said to be robustly testable if it has a $\Omega(1)$ -robust tester. For a given code \mathcal{C} of block length n over \mathbb{F}_q and a vector $v \in \mathbb{F}_q^n$, the (relative) Hamming distance of v to the closest codeword in \mathcal{C} is denoted by $\delta_{\mathcal{C}}(v)$.

Asking whether $T_{\mathcal{C}_1 \otimes \mathcal{C}_2}$ is a robust tester has the following nice interpretation. The q queries i_1, \dots, i_q are either rows or columns of the received word v . Let the row or column corresponding to the random seed s be denoted by v^s . Then the robustness of $T_{\mathcal{C}_1 \otimes \mathcal{C}_2}$ on inputs (v, s) , $\rho^{T_{\mathcal{C}_1 \otimes \mathcal{C}_2}}(v, s)$ is just $\delta_{\mathcal{C}_1}(v^s)$ when i_s corresponds to a row and $\delta_{\mathcal{C}_2}(v^s)$ when i_s corresponds to a column. Therefore the expected robustness of $T_{\mathcal{C}_1 \otimes \mathcal{C}_2}$ on v is the average of the following two quantities: the average relative distance of the rows of v from \mathcal{C}_1 and the average relative distance of the columns of v from \mathcal{C}_2 .

In particular, if $T_{\mathcal{C}_1 \otimes \mathcal{C}_2}$ is $\Omega(1)$ -robust then it implies that for every received word v such that all rows (and columns) of v are $o(1)$ -close to \mathcal{C}_1 (and \mathcal{C}_2), v is $o(1)$ -close to $\mathcal{C}_1 \otimes \mathcal{C}_2$. P. Valiant proved the following result.

Theorem 8.4 ([103]). *There exist linear codes $[n_1, k_1, d_1 = n_1/10]$ and $[n_2 = n_1^2, k_2, d_2 = n_2/10]$ (call them \mathcal{C}_1 and \mathcal{C}_2) and a $n_2 \times n_1$ received word v such that every row of v is $1/n_1$ -close to \mathcal{C}_1 and every column of v is a codeword \mathcal{C}_2 but v is $1/20$ -far from $\mathcal{C}_1 \otimes \mathcal{C}_2$.*

Note that in the above construction, $n_2 \neq n_1$ and in particular \mathcal{C}_1 and \mathcal{C}_2 are not the same code.

Reduction from the Construction of Valiant

In this section, we prove the following result.

Theorem 8.5. *Let $\mathcal{C}_1 \neq \mathcal{C}_2$ be $[n_1, k_1, d_1 = \Omega(n_1)]$ and $[n_2, k_2, d_2 = \Omega(n_2)]$ codes respectively (with $n_2 > n_1$) and let v be a $n_2 \times n_1$ matrix such that every row (and column) of v is $g(n_1)$ -close to \mathcal{C}_1 ($g(n_2)$ -close to \mathcal{C}_2) but v is ρ -far from $\mathcal{C}_1 \otimes \mathcal{C}_2$. Then there exists a linear code \mathcal{C} with parameters $n, k, d = \Omega(n)$ and a received word v' such that every row (and column) of v' is $g(n_1)/2$ -close (and $g(n_2)/2$ -close) to \mathcal{C} but v' is $\rho/4$ -far from \mathcal{C}^2 .*

Proof. We will first assume that n_1 divides n_2 and let $m = \frac{n_2}{n_1}$. For any $x \in \Sigma^{k_1}$ and $y \in \Sigma^{k_2}$, let

$$\mathcal{C}(\langle x, y \rangle) = \langle (\mathcal{C}_1(x))^m, \mathcal{C}_2(y) \rangle$$

Thus, $k = k_1 + k_2$ and $n = mn_1 + n_2$. Also as $d_1 = \Omega(n_1)$ and $d_2 = \Omega(n_2)$, $d = \Omega(n)$.

We now construct the $n \times n$ matrix v' from v . The lower left $n_2 \times mn_1$ sub-matrix of v' contains the matrix v^m where v^m is the horizontal concatenation of m copies of v (which is a $n_2 \times n_1$ matrix). Every other entry in v' is 0. See figure 8.1 for an example with $m = 2$.

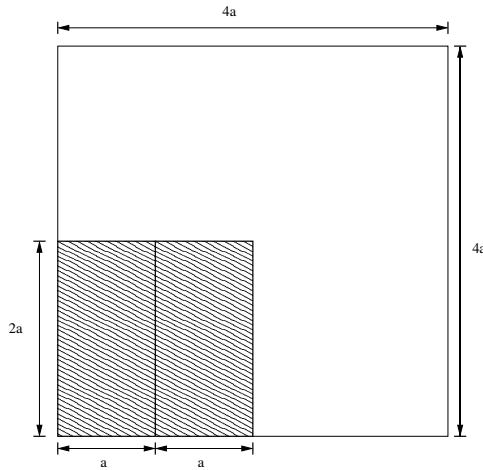


Figure 8.1: The construction of the new received word v' from v for the case when $n_1 = a$, $n_2 = 2a$ and $m = 2$. The shaded boxes represent v and the unshaded regions has all 0s.

Let w be the codeword in $\mathcal{C}_1 \otimes \mathcal{C}_2$ closest to v and construct w' in the same manner as v' was constructed from v . We first claim that w' is the codeword in⁵ \mathcal{C}^2 closest to v' . For the sake of contradiction, assume that there is some other codeword w'' in \mathcal{C}^2 such that $\Delta(v', w'') < \Delta(v', w')$. For any $2n' \times 2n'$ matrix u let u_{lb} denote the lower left $n' \times n'$ sub-matrix of u . Note that by definition of \mathcal{C} , $w''_{lb} = x^m$ where $x \in \mathcal{C}_1 \otimes \mathcal{C}_2$. Further, as v' (necessarily) has 0 everywhere other than v'_{lb} and $\Delta(v', w'') < \Delta(v', w')$, it holds that $\Delta(v, w) > \Delta(v, x)$ which contradicts the definition of w .

Finally, it is easy to see that

$$\delta_{\mathcal{C}^2}(v') = \Delta(v', w')/n^2 = \Delta(v, w)m/(mn_1 + n_2)^2 = \Delta(v, w)/(4n_1n_2) = \frac{\rho}{4}$$

and if for any row (or column), the relative distance of v restricted to that row (or column) from \mathcal{C}_1 (\mathcal{C}_2) is at most α then for every row (or column), the relative distance of v' restricted to that row (or column) from \mathcal{C} is at most $\alpha/2$.

This completes the proof for the case when n_1 divides n_2 . For the case when n_1 does not divide n_2 a similar construction works if one defines \mathcal{C} in the following manner (for any $x \in \Sigma^{k_1}$ and $x_2 \in \Sigma^{k_2}$)

$$\mathcal{C}(\langle x, y \rangle) = \langle (\mathcal{C}_1(x))^{\ell/n_1}, (\mathcal{C}_2(y))^{\ell/n_2} \rangle$$

where $\ell = \text{lcm}(n_1, n_2)$. The received word v' in this case would have its lower left $\ell \times \ell$ matrix as $v^{(\ell/n_1, \ell/n_2)}$ (where $v^{(m_1, m_2)}$ is the matrix obtained by vertically concatenating m_2 copies of v^{m_1}) and it has 0s everywhere else. \square

⁵Note that $w' \in \mathcal{C}^2$ as the all zeros vector is a codeword in both \mathcal{C}_1 and \mathcal{C}_2 and $w \in \mathcal{C}_1 \otimes \mathcal{C}_2$.

Theorem 8.4 and 8.5 imply the following result.

Corollary 8.6. *There exist a linear code \mathcal{C} with linear distance such that the tester $T_{\mathcal{C}^2}$ is not $\Omega(1)$ -robust for \mathcal{C}^2 .*

8.6 Tolerant Testing of Reed-Muller Codes

In this section, we discuss testers for codes based on multivariate polynomials.

8.6.1 Bivariate Polynomial Codes

As we saw in Section 8.5.2, one cannot have a robust standard testers for \mathcal{C}^2 in general. In this subsection, we consider a special case when $\mathcal{C} = \text{RS}[n, k + 1, d = n - k]_q$, that is, the Reed–Solomon code based on evaluation of degree k polynomials over \mathbb{F}_q at n distinct points in the field. We show that the tester for \mathcal{C}^2 considered in Section 8.5.2 is tolerant for this special case. It is well-known (see, for example, Proposition 2 in [88]) that in this case \mathcal{C}^2 is the code with codewords being the evaluations of bivariate polynomials over \mathbb{F}_q of degree k in each variable. The problem of low-degree testing for bivariate polynomials is a well-studied one: in particular we use the work of Polishchuk and Spielman [88] who analyze a tester using axis parallel lines. Call a bivariate polynomial to be one of degree (k_1, k_2) if the maximum degrees of the two variables are k_1 and k_2 respectively. In what follows, we denote by $Q' \in \mathbb{F}_q^{n \times n}$ the received word to be tested (thought of as an $n \times n$ matrix), and let $Q(x, y)$ be the degree (k, k) polynomial whose encoding is closest to Q' .

We now specify the tolerant tester T' . The upper bound of $1 - \sqrt{1 - d/n}$ on τ comes from the fact that this is largest radius for which decoding an $\text{RS}[n, k + 1, d]$ code is known to be solvable in polynomial time [63].

1. Fix τ where $0 \leq \tau \leq 1 - \sqrt{1 - d/n}$.
2. With probability $\frac{1}{2}$ choose $b = 0$ or $b = 1$.
 - If $b = 0$, choose a row r randomly and reject if $\delta(Q'(r, \cdot), P(\cdot)) > \tau$ for every univariate polynomial P of degree k and accept otherwise.
 - If $b = 1$, choose a column c randomly and reject if $\delta(Q'(\cdot, c), P(\cdot)) > \tau$ for every univariate polynomial P of degree k and accept otherwise.

The following theorem shows that T' is a tolerant tester.

Theorem 8.7. *There exists an absolute constant $c_0 > 0$ such that for $\tau \leq 1 - \sqrt{1 - d/n}$, the tester T' with threshold τ is a (c_1, c_2, \sqrt{N}) -tolerant tester for \mathcal{C}^2 (where $\mathcal{C} = \text{RS}[n, k + 1, d]$) where $c_1 = \frac{n\tau}{3d}$, $c_2 = \frac{2nc_0(\tau+2)}{3d}$ and N is the block length of \mathcal{C}^2 .*

Proof. To analyze T' let $R^*(r, \cdot)$ be the closest degree k univariate polynomial (breaking ties arbitrarily) for each row r . Similarly construct $C^*(\cdot, c)$. We will use the following refinement of the Bivariate testing lemma of [88]:

Lemma 8.8 ([88, 13]). *There exists an universal constant $c_0 \leq 128$ such that the following holds. If $8k \leq n$ then $\delta(Q', \mathcal{C}^2) = \delta(Q', Q) \leq c_0 \cdot (\delta(R^*, Q') + \delta(C^*, Q'))$.*

The following proposition shows that the standard tester version of T' (that is T' with $\tau = 0$) is a robust tester–

Proposition 8.9. *T' with $\tau = 0$ is a $2c_0$ robust tester, where c_0 is the constant from Lemma 8.8.*

Proof. By the definition of the row polynomial R , for any row index r , the robustness of the tester with $b = 0$ and r , $\rho(Q', \langle b, r \rangle) = \delta(Q'(r, \cdot), R^*(r, \cdot))$. Similarly for $b = 1$, we have $\rho(Q', \langle b, c \rangle) = \delta(Q'(\cdot, c), C^*(\cdot, c))$. Now the expected robustness of the test is given by

$$\begin{aligned} \rho(Q') &= \Pr[b = 0] \sum_{i=1}^n \Pr[r = i] \cdot \delta(Q'(r, \cdot), R^*(r, \cdot)) + \\ &\quad \Pr[b = 1] \sum_{j=1}^n \Pr[c = j] \cdot \delta(Q'(\cdot, c), C^*(\cdot, c)) \\ &= \frac{1}{2}(\delta(Q', R^*) + \delta(Q', C^*)). \end{aligned}$$

Using Lemma 8.8, we get $\delta(Q', Q) \leq 2c_0\rho(Q')$, as required. \square

From the description of T' , it is clear that it has a *partitioned* query pattern. There are two partitions: one for the rows (corresponding to the choice $b = 0$) and one for the columns (corresponding to the choice $b = 1$).

Lemmas 8.2 and 8.3 prove Theorem 8.7 where c_0 is the constant from Lemma 8.8. \square

8.6.2 General Reed-Muller Codes

We now turn our attention to testing of general Reed-Muller codes. Recall that $\text{RM}_q(k, m)$ the linear code consisting of evaluations of m -variate polynomials over \mathbb{F}_q of *total degree* at most k at all points in \mathbb{F}_q^m .⁶ To test codewords of $\text{RM}_q(k, m)$, we need to, given a function $f : \mathbb{F}_q^m \rightarrow \mathbb{F}_q$ as a table of values, test if f is close to an m -variate polynomial of total degree k . We will do this using the following natural and by now well-studied *low-degree test* which we call the *lines test*: pick a random line in \mathbb{F}_q^m and check if the restriction of f on the line is a univariate polynomial of degree at most k . In order to achieve tolerance,

⁶The results of the previous section were for polynomials which had degree in each *individual* variable bounded by some value; here we study the total degree case.

we will modify the above test to accept if the restriction of f on the picked line is within distance τ from some degree k univariate polynomial, for a threshold τ . Using the analysis of the low-degree test from [42], we can show the following.

Theorem 8.8. *For $0 \leq \tau \leq 1 - \sqrt{k/q}$ and $q = \Omega(k)$, $\text{RM}_q(k, m)$ is (c_1, c_2, p) testable with $c_1 = \frac{n\tau}{3d}$, $c_2 = \frac{3(\tau+2)n}{d}$ and $p = n^{1-1/m}$ where $n = q^m$ and d are the block length and the distance of the code.*

Proof. Recall that our goal is to test if a given function $f : \mathbb{F}_q^m \rightarrow \mathbb{F}_q$ is close to an m -variate polynomial of total degree k . For any $x, h \in \mathbb{F}_q^m$, a line passing through x in direction h is given by the set $L_{x,h} = \{x + th | t \in \mathbb{F}_q\}$. Further define $P_{x,h}^f(\cdot)$ to be the univariate polynomial of degree at most k which is closest (in Hamming distance) from the restriction of f on $L_{x,h}$. We will use the following result.

Theorem 8.9 ([42]). *There exists a constant c such that for all k , if q is a prime power that is at least ck , then given a function $f : \mathbb{F}_q^m \rightarrow \mathbb{F}_q$ with*

$$\rho \stackrel{\text{def}}{=} E_{x,h \in \mathbb{F}_q^m} \Pr_{t \in \mathbb{F}_q} [P_{x,h}^f(t) \neq f(x + th)] \leq \frac{1}{9},$$

there exists an m -variate polynomial g of total degree at most k such that $\text{dist}(f, g) \leq 2\rho$.

The above result clearly implies that the line test is robust which we record in the following corollary.

Corollary 8.10. *There exists a constant c such that the line test for $\text{RM}_q(k, m)$ with $q \geq ck$ is 9 -robust.*

The line test picks a random line by choosing x and h randomly. Consider the case when h is fixed. It is not hard to check that for there is a partition of $\mathbb{F}_q^m = X_1 \cup \dots \cup X_q$ where each X_i has size q^{m-1} such that $\cup_{x \in X_i} L_{x,h} = \mathbb{F}_q^m$. In other words:

Proposition 8.10. *The line test has a partitioned query pattern.*

The proposed tolerant tester for $\text{RM}_q(k, m)$ is as follows: pick $x, h \in \mathbb{F}_q^m$ uniformly at random and check if the input restricted to $L_{x,h}$ is τ -close to some univariate polynomial of degree k . If so accept, otherwise reject. When the threshold τ satisfies $\tau \leq 1 - \sqrt{k/q}$, the test can be implemented in polynomial time [63]. From Corollary 8.10, Proposition 8.10, Lemmas 8.2 and 8.3, the above is indeed a tolerant tester for $\text{RM}_q(k, m)$, and Theorem 8.8 follows. \square

8.7 Bibliographic Notes and Open Questions

The results in this chapter (other than those in Section 8.5.2) were presented in [57]. Results in Section 8.5.2 appear in [26].

In the general context of property testing, the notion of tolerant testing was introduced by Parnas *et al.* [83] along with the related notion of distance approximation. Parnas *et al.* also give tolerant testers for clustering. We feel that codeword-testing is a particularly natural instance to study tolerant testing. (In fact, if LTCs were defined solely from a coding-theoretic viewpoint, without their relevance and applications to PCPs in mind, we feel that it is likely that the original definition itself would have required tolerant testers.)

The question of whether the natural tester for $C_1 \otimes C_2$ is a robust one was first explicitly asked by Ben-Sasson and Sudan [14]. P. Valiant showed that in general, the answer to the question is no. Dinur, Sudan and Wigderson [30] further show that the answer is positive if at least one of C_1 or C_2 is a *smooth* code, for a certain notion of smoothness. They also show that any non-adaptive and *bi-regular* binary linear LTC is a smooth code. A bi-regular LTC has a tester that in every query probes the same number of positions and every bit in the received word is queried by the same number of queries. The latter requirement (in the terminology of this chapter) is that the tester is $(\langle 1, q \rangle, \langle 0, 0 \rangle, 1)$ -smooth, where the tester makes q queries. The result of [30] however only works with constant query complexity. Note that for such an LTC, Lemma 8.1 implies that the code is also tolerant testable.

Obtaining non-trivial lower bounds on the the block length of codes that are locally testable with very few (even 3) queries is an extremely interesting question. This problem has remained open and resisted even moderate progress despite all the advancements in constructions of LTCs. The requirement of having a tolerant local tester is a stronger requirement. While we have seen that we can get tolerance with similar parameters to the best known LTCs, it remains an interesting question whether the added requirement of tolerance makes the task of proving lower bounds more tractable. In particular,

Open Question 8.1. *Does there exist a code with constant rate and linear distance that has a tolerant tester that makes constant number of queries ?*

This seems like a good first step in making progress towards understanding whether locally testable codes with constant rate and linear distance exist, a question which is arguably one of the grand challenges in this area. For interesting work in this direction which proves that such codes, if they exist, cannot also be *cyclic*, see [10].

The standard testers for Reed-Muller codes considered in Section 8.6 (and hence, the tolerant testers derived from them) work only for the case when the size of the field is larger than the degree of the polynomial being tested. Results in Chapter 7 and those in [74] give a standard tester for Reed-Muller codes which works for all fields. These testers do have a partitioned query pattern— however, it is not clear if the testers are robust. Thus, our techniques to convert it into a tolerant tester fail. It will be interesting to show the following result.

Open Question 8.2. *Design tolerant testers for RM codes over any finite field.*