



# Testing Periodicity

## Extended Abstract

Oded Lachish  
 University of Haifa  
 Haifa, Israel  
 loded@cs.haifa.ac.il

Ilan Newman  
 University of Haifa  
 Haifa, Israel  
 ilan@cs.haifa.ac.il

### Abstract

A string  $\alpha \in \Sigma^n$  is called *p-periodic*, if for every  $i, j \in \{1, \dots, n\}$ , such that  $i \equiv j \pmod p$ ,  $\alpha_i = \alpha_j$ , where  $\alpha_i$  is the  $i$ -th place of  $\alpha$ . A string  $\alpha \in \Sigma^n$  is said to be *period*( $\leq g$ ), if there exists  $p \in \{1, \dots, g\}$  such that  $\alpha$  is *p-periodic*.

An  $\epsilon$  property tester for *period*( $\leq g$ ) is a randomized algorithm, that for an input  $\alpha$  distinguishes between the case that  $\alpha$  is in *period*( $\leq g$ ) and the case that one needs to change at least  $\epsilon$ -fraction of the letters of  $\alpha$ , so that it will become *period*( $\leq g$ ). The complexity of the tester is the number of letter-queries it makes to the input. We study here the complexity of  $\epsilon$  testers for *period*( $\leq g$ ) when  $g$  varies in the range  $1, \dots, \frac{n}{2}$ . We show that there exists a surprising exponential phase transition in the query complexity around  $g = \log n$ . That is, for every  $\delta > 0$  and for each  $g$ , such that  $g \geq (\log n)^{1+\delta}$ , the number of queries required and sufficient for testing *period*( $\leq g$ ) is polynomial in  $g$ . On the other hand, for each  $g \leq \frac{\log n}{4}$ , the number of queries required and sufficient for testing *period*( $\leq g$ ) is only poly-logarithmic in  $g$ .

We also prove an exact asymptotic bound for testing general periodicity. Namely, that 1-sided error, non adaptive  $\epsilon$ -testing of periodicity (*period*( $\leq \frac{n}{2}$ )) is  $\Theta(\sqrt{n \log n})$  queries.

## 1 Introduction

Periodicity in strings plays an important role in several branches of CS and engineering applications. It is being used as a measure of 'self similarity' in many application regarding string algorithms (e.g pattern matching), computational biology, data analysis and planning (e.g analysis of stock prices, communication patterns etc.), signal and image processing and others. On the other hand, sources of very large streams of data are now common inputs for strategy-planning or trend detection algorithms. Typically, such streams of data are either too large to store entirely in the computer memory, or so large that even linear processing time is not feasible. Thus it would be of interest to develop very fast (sub-linear time) algorithms that test whether a long sequence is periodic or approximately periodic, and in particular, that test if it has a very short period. This calls for algorithms in the framework of Combinatorial Property Testing [1]. In this framework, introduced initially by Rubinfeld and Sudan [2] and formalized by Goldreich et al. [1], one uses a randomized algorithm that queries the input at very few locations and based on this, decides whether it has a given property or it is 'far' from having the property. Indeed related questions to periodicity have already been investigated [3, 4, 5], although the focus here was somewhat different.

In [6], the author constructs an algorithm that approximates in a certain sense the DFT (Discrete Fourier Transform) of a finite sequence in sub linear time. This is quite related but not equivalent to testing how close is a sequence to being periodic. In [3], the authors study some alternative

parametric definitions of periodicity that intend to 'capture the distance' of a sequence to being periodic. They mainly relate the different definitions of periodicity. They also show that there is a tolerant tester for periodicity. That is, they show a simple algorithm, that given  $0 \leq \epsilon_1 < \epsilon_2 \leq 1$ , decides whether a sequence is  $\epsilon_1$ -close to periodic or  $\epsilon_2$ -far from being periodic, using  $O(\sqrt{n} \cdot \text{poly}(\log n))$  queries.

There are other works on sequences sketching [4] etc. but none of those seems to address directly periodicity testing.

The property of being periodic is formalized here in a very general form: Let  $\Sigma$  be a finite alphabet. A string  $\alpha \in \Sigma^n$  is said to be  $p$ -periodic for an integer  $p \in [n]$ , if for every  $i, j \in [n]$ , where  $i \equiv j \pmod{p}$ ,  $\alpha_i = \alpha_j$ , where  $\alpha_i$  is the  $i$ -th character of  $\alpha$ . We say that  $\alpha$  is in  $\text{period}(\leq g)$ , where  $g \in [\frac{n}{2}]$ , if it is  $p$ -periodic for some  $p \leq g$ . We say a string is periodic if it is in  $\text{period}(\leq \frac{n}{2})$ .

We study here the complexity of  $\epsilon$  testing the property  $\text{period}(\leq g)$  when  $g$  varies in the range  $1, \dots, \frac{n}{2}$ . An  $\epsilon$ -tester for  $\text{period}(\leq g)$  is a randomized algorithm, that for an input  $\alpha \in \Sigma^n$  distinguishes between the case that  $\alpha$  is in  $\text{period}(\leq g)$  and the case that one needs to change at least an  $\epsilon$ -fraction of the letters of  $\alpha$ , so that it will be in  $\text{period}(\leq g)$ . The complexity of the tester is the number of letter-queries it makes to the input.

We show that there exists a surprising exponential phase transition around  $g = \log n$ . That is, for every  $\delta > 0$  and for each  $g$ , such that  $g \geq (\log n)^{1+\delta}$ , the number of queries required and sufficient for testing  $\text{period}(\leq g)$  is polynomial in  $g$ . On the other hand, for each  $g \leq \frac{\log n}{4}$ , the number of queries required and sufficient for testing  $\text{period}(\leq g)$  is only poly-logarithmic in  $g$ . We also settle the exact complexity of non-adaptive 1-sided error test for general periodicity (that is,  $\text{period}(\leq n/2)$ ). We show that the exact complexity in this case is  $\theta(\sqrt{n \log n})$ . The upper bound that we prove is an improvement over the result of [3] and uses a construction of a small random set  $A \subseteq [n]$  of size  $\sqrt{n \log n}$  for which the multi-set  $A - A = \{a - b | a, b \in A\}$  contains at least  $\log n$  copies of each member of  $[n/2]$ . This can be trivially done with a set  $A$  of size  $\sqrt{n \log n}$ . We improve on the natural bound by using non-independence of the samples, but still retain the property that the probability of the  $\log n$  copies of each number in  $A - A$  behave as if being not too far from independent. This problem has occurred in various other situation, (e.g [7, 8]) and thus could be interesting in its own.

The rest of the paper is organized as follows. In section 2 we introduce the necessary notations and some very basic observations. Section 3 contains an  $\epsilon$ -test for  $\text{period} \leq g$  that uses  $\sqrt{g \log g}$  queries. In section 4 we construct an  $\epsilon$ -test for  $\text{period}(\leq g)$ ,  $g \leq \frac{\log n}{4}$  that uses only  $\tilde{O}((\log g)^6)$  queries. Sections 5 and 6 contain the corresponding lower bounds, thus showing the claimed phase transition. Finally, the special case of  $g = \theta(n)$  is treated in Section 7 which contains a proof that any 1-sided error non adaptive tester for periodicity ( $\text{period}(\leq \frac{n}{2})$ ) requires  $\Omega(\sqrt{n \log n})$  queries. This implies that the tester in section 3 is asymptotically optimal when  $g = n/2$ .

## 2 Preliminaries

For a string  $\alpha \in \Sigma^n$  and an integer  $i \in [n]$ , ( $[n] = \{1, \dots, n\}$ ) we denote by  $\alpha_i$  the  $i$ -th symbol of  $\alpha$ , that is  $\alpha = \alpha_1 \dots \alpha_n$ . Given a set  $S \subseteq [n]$  such that  $S = \{i_1, i_2, \dots, i_m\}$  and  $i_1 < i_2 < \dots < i_m$  we define  $\alpha_S = \alpha_{i_1} \alpha_{i_2} \dots \alpha_{i_m}$ . In the following  $\Sigma$  will be fixed. Also, unless otherwise stated, all strings are in  $\Sigma^n$ .

For two strings  $\alpha, \beta$  we denote by  $\text{dist}(\alpha, \beta)$  the Hamming distance between  $\alpha$  and  $\beta$ . Namely,  $\text{dist}(\alpha, \beta) = |\{i | \alpha_i \neq \beta_i\}|$ . For a property  $\mathcal{P} \subseteq \Sigma^n$  and a string  $\alpha \in \Sigma^n$ ,  $\text{dist}(\alpha, \mathcal{P}) = \min\{\text{dist}(\alpha, \beta) | \beta \in \mathcal{P}\}$  denotes the distance from  $\alpha$  to  $\mathcal{P}$ . We say that  $\alpha$  is  $\epsilon$ -far from  $\mathcal{P}$  if  $\text{dist}(\alpha, \mathcal{P}) \geq \epsilon n$ , otherwise we say that  $\alpha$  is  $\epsilon$ -close to  $\mathcal{P}$ .

**Definition 2.1.** For a string  $\alpha$  and a subset  $S \subseteq [n]$  we say that  $\alpha_S$  is homogeneous if for all  $i, j \in S$ ,  $\alpha_i = \alpha_j$ .

### Property Testing

The type of algorithms that we consider here are 'property-testers' [1, 9, 10]. These are randomized algorithms that access the input string via a 'location-oracle' which it can query: A query is done by specifying one place in the string to which the answer is the value of the string in the queried location. The complexity of the algorithm is the amount of queries it makes to the string in the worst case. Such an algorithm is said to be an  $\epsilon$ -test for a property  $\mathcal{P}$  of strings of length  $n$ , if it distinguishes with success probability at least  $2/3$  between the case that the input string belongs to  $\mathcal{P}$ , and the case that it is  $\epsilon$ -far from  $\mathcal{P}$ .

### Periodicity

**Definition 2.2.** A string  $\alpha \in \Sigma^n$  has period  $p$  (denoted  $p$ -periodic), if  $\alpha_i = \alpha_j$  for every  $i, j \in [n]$ , such that  $i \equiv j \pmod{p}$ .

Note that a string  $\alpha$  is homogeneous if and only if  $\alpha$  is 1-periodic.

**Definition 2.3.** We say that  $\alpha$  has the property  $\text{period}(\leq g)$  if it is  $p$ -periodic for some  $p \leq g$ . We say that  $\alpha$  is periodic if it has the property  $\text{period}(\leq \frac{n}{2})$ .

**Definition 2.4.** A witness that a string  $\alpha$  is not  $p$ -periodic, denoted as  $p$ -witness, is an unordered pair  $\{i, j\} \subseteq [n]$ , such that  $i \equiv j \pmod{p}$  and  $\alpha_i \neq \alpha_j$ .

According to the definition of periodicity a string  $\alpha$  has a period  $p \leq \frac{n}{2}$  if and only if there does not exist a  $p$ -witness. In the same manner a witness for not having  $\text{period}(\leq g)$  is defined as follows:

**Definition 2.5.** A witness that a string  $\alpha$  is not in  $\text{period}(\leq g)$  is a set of integers  $Q \subseteq [n]$  such that for every  $p \leq g$  there are two integers  $i, j \in Q$  that form a  $p$ -witness.

**Fact 2.6.** A string  $\alpha$  has  $\text{period}(\leq g)$ , if and only if there does not exist a witness that the string is not in  $\text{period}(\leq g)$ .

**Fact 2.7.** If a string  $\alpha \in \Sigma^n$  does not have a period in  $(\frac{g}{2}, g]$ , where  $g \leq \frac{n}{2}$ , then it has no period of at most  $\frac{g}{2}$ . If a string  $\alpha \in \Sigma^n$  is  $\epsilon$ -far from having a period in  $(\frac{g}{2}, g]$ , where  $g \leq \frac{n}{2}$ , then it is  $\epsilon$ -far from having the property  $\text{period}(\leq g)$ .

*Proof.* Observe that if a string  $\alpha \in \Sigma^n$  has period  $p \leq \frac{g}{2}$ , then it also has period  $q$  for every  $q$  that is a multiple of  $p$ . Note also that there must exist such  $q \in (\frac{g}{2}, g]$ . ■

**Definition 2.8.** For  $\alpha \in \Sigma^n$ ,  $p \in [n]$  and  $0 \leq i \leq p-1$ , let  $Z(p, i) = \{j \mid j \equiv i \pmod{p}\}$ . We call  $\alpha_{Z(p, i)}$  the  $i$ -th  $p$ -section of  $\alpha$ .

The following obvious fact relates the distance of a string to  $p$ -periodic and the homogeneity of its  $p$ -sections.

**Fact 2.9.** For each  $\alpha$ ,  $\text{dist}(\alpha, p\text{-periodic}) = \sum_{i=0}^{p-1} \text{dist}(\alpha_{Z(p, i)}, \text{homogeneous})$ .

In further sections we use the following basic  $\epsilon$ -test for  $p$ -period.

### Algorithm p-test

**Input:** a string  $\alpha \in \Sigma^n$ , the string length  $n$ , a period  $p \in [\frac{n}{2}]$  and a distance parameter  $0 < \epsilon < 1$ ;

1. Select  $\frac{1}{\epsilon}$  random unordered pairs (with repetitions)  $\{i, j\} \subset [n]$  such that  $i \equiv j \pmod{p}$ .
2. Reject if one of the selected pairs is a  $p$ -witness (namely a witness for being *non- $p$ -periodic*). Otherwise accept.

**Proposition 2.10.** *Algorithm  $p$ -test is a 1-sided error, non-adaptive  $\epsilon$ -test for  $p$ -periodic. Its query complexity is  $\frac{2}{\epsilon}$ .*

*Proof.* The query complexity is obvious. The test is 1-sided error since it rejects only if it finds a  $p$ -witness. To estimate its error probability let  $\alpha$  be such that  $\text{dist}(\alpha, p\text{-periodic}) \geq \epsilon n$ .

Assume that  $p$  divides  $n$  (the general case is the same except for some  $\lceil \cdot \rceil$  addition which make no essential difference). With this assumption  $|Z(p, i)| = n/p = m$ . For every  $0 \leq i \leq p-1$  set  $\text{dist}(\alpha_{Z(p, i)}, \text{homogeneous}) = d_i$ . For fixed  $i$  and every  $\sigma \in \Sigma$  let  $n_\sigma$  be the number of occurrences of  $\sigma$  in  $Z(p, i)$ . Assume also that we have renumbered the letters in  $\Sigma$  so that  $n_1 \geq n_2, \dots \geq n_k$ . Then  $d_i = n - n_1$  as it is easy to see that the closest homogeneous string is obtained by changing all letters different from  $\sigma_1$  to  $\sigma$  (exactly  $m - n_1$  letters). Hence the number of  $p$ -witnesses in  $Z(p, i)$ ,  $W_i$ , is  $W_i = \frac{1}{2}(\sum n_j)^2 - \sum n_j^2$ . It is easy to see that  $W_i \geq m \cdot d_i/2$ . Now, according to Fact 2.9 we get that,  $\text{dist}(\alpha, p\text{-periodic}) = \sum_{i=0}^{p-1} \text{dist}(\alpha_{Z(p, i)}, \text{homogeneous}) = \sum_{i=0}^{p-1} d_i \leq \frac{2}{m} \sum W_i \leq \frac{2}{m} W$  where  $W$  is the total number of  $p$ -witnesses.

It follows that  $W \geq \frac{m}{2} \text{dist}(\alpha, p\text{-periodic}) \geq \frac{\epsilon n^2}{2p}$ . Note however that the total number of unordered pairs  $i, j$  such that  $j \equiv i \pmod{p}$  is  $\frac{n^2}{2p}$ . Thus we conclude that the probability that a random such pair is a  $p$ -witness is at least  $\epsilon$  and the result follows. ■

We will need the following proposition in a number of lower bound proofs. Let  $m < n/8$ ,  $S \subset [n]$ ,  $|S| = m$  and an assignment  $a : S \rightarrow \{0, 1\}$ . We define the following distribution  $U/S$  on strings of length  $n$ . We choose every letter  $\alpha_i = a_i$  if  $i \in S$ . For all places  $i \notin S$  we choose  $\alpha_i$  to be ‘1’ with probability  $1/2$  and ‘0’ with probability  $1/2$ , independently between different  $i$ ’s. Note that  $U/S$  is just the uniform distribution on all binary strings conditioned on the event that the projection on  $S$  is  $a$ . We have,

**Proposition 2.11.** *Let  $m, S, a$  be as above and let  $\mathcal{G}(g)$  be the event that a string selected according to  $U/S$  is  $\frac{1}{16}$ -far from having  $\text{period}(\leq g)$ . Then  $\text{Prob}_{U/S}(\mathcal{G}(g)) \geq 1 - \frac{1}{g}$ .*

**Proof:** Let  $U$  be the uniform distribution over  $\Sigma^n$ . Note that if a string is  $\frac{1}{16}$ -far from having  $\text{period}(\leq \frac{n}{2})$  then for every  $g \leq \frac{n}{2}$  the string is  $\frac{1}{16}$ -far from having  $\text{period}(\leq g)$ . Hence it is sufficient to prove the proposition only in the case that  $g = \frac{n}{2}$ . According to Fact 2.7, it is enough to prove that the probability that a string selected according to  $U$  is  $\frac{1}{16}$ -far from having a period in  $(\frac{n}{4}, \frac{n}{2}]$ , is at least  $1 - \frac{1}{n}$ . We prove that for each  $p \in (\frac{n}{4}, \frac{n}{2}]$  the probability that a string selected according to  $U$ , is  $\frac{1}{16}$ -close to being  $p$ -periodic, is at most  $\frac{4}{n^2}$ , and therefore by the union bound we are done.

Indeed let  $p \in (\frac{n}{4}, \frac{n}{2}]$ , and let  $\alpha$  be a string selected according to  $U$ . Then for every  $i \in [p] - m$ ,  $\text{Prob}_U(\text{dist}(\alpha_{Z(p, i)}, \text{homogeneous}) \geq 1) \geq \frac{1}{2}$ . Thus, since the  $Z(p, i)$ ’s are mutually disjoint and  $p - m > \frac{n}{8}$ , Chernoff bound implies that,

$$\text{Prob}_U \left( \sum_{i \in [p]} \text{dist}(\alpha_{Z(p, i)}, \text{homogeneous}) \leq \frac{n}{16} \right) \leq e^{-\frac{n}{16}} \leq \frac{4}{n^2}. \quad \blacksquare$$

### 3 Upper Bound for Testing $\text{period}(\leq g)$

In this section we construct a 1-sided error, non adaptive  $\epsilon$ -test for  $\text{period}(\leq g)$  that uses  $O(\sqrt{g \log g}/\epsilon^2)$  queries. The construction consists of two stages. We first show a 1-sided error,

non adaptive algorithm,  $PT$ , for testing periodicity ( $period(\leq \frac{n}{2})$ ), that uses  $O(\frac{\sqrt{n \log n}}{\epsilon^2})$  queries. Then we show how to use algorithm  $PT$  in order to construct the claimed algorithm for testing  $period(\leq g)$ .

The intuition behind the algorithm  $PT$  is simple. We first explain it by hinting to a test for  $period(\leq \frac{n}{2})$  that uses  $O(\frac{\sqrt{n \log n}}{\epsilon^2})$  queries. By Fact 2.7  $\alpha$  is  $\epsilon$ -far from  $period(\leq n)$  if and only if  $\alpha$  is  $\epsilon$ -far from being  $p$ -periodic for every  $p \in (\frac{n}{4}, \frac{n}{2}]$ . We show that if  $\alpha$  is  $\epsilon$ -far from  $period(\leq \frac{n}{2})$ , we can find a  $p$ -witness for each  $p \in (\frac{n}{4}, \frac{n}{2}]$  with probability at least  $\frac{1}{n}$ . This is sufficient since by the union bound we will be done.

Suppose we can construct a set  $Q \subseteq [n]$  of size  $O(\sqrt{n \log n})$ , such that  $Q$  contains at least  $\log n$   $p$ -witnesses for each  $p \in (\frac{n}{4}, \frac{n}{2}]$ . Then obviously we achieve the above goal. Indeed such a set  $Q$  can be constructed (e.g by choosing for every  $p \in (\frac{n}{4}, \frac{n}{2}]$   $\log n$  random pairs  $i, j$  such  $i \equiv j \pmod{p}$ ).

To reduce the size of  $Q$ , we construct such a set  $Q$  of size  $\sqrt{n \log n}$  with the above properties while giving up the independence between the  $\log n$  pairs for each  $p$ . In turn we will have to show that the probability that none will be a  $p$ -witness is still low enough.

We next present the  $\epsilon$ -test for  $g = \frac{n}{2}$  (periodicity).

**Definition 3.1.** Let  $\ell = \sqrt{n \log n}$  and let  $J = \{I_i\}_{i=1}^{\ell}$  be the set of pairwise disjoint intervals  $I_i = \left( (i-1)\sqrt{\frac{n}{\log n}}, i\sqrt{\frac{n}{\log n}} \right]$ .

We now present the  $\epsilon$ -test for  $period(\leq \frac{n}{2})$ .

## Algorithm PT

**Input:** a string  $\alpha \in \Sigma^n$ , the string length  $n$  and a distance parameter  $0 < \epsilon < 1$ ;

Let  $\ell, J$  be as in Definition 3.1.

1. Select a set of integers  $T = \{t_1, \dots, t_{\ell}\}$  by choosing one random point from each  $I \in J$ .
2. Repeat the following independently  $m = \frac{2^{10} \cdot \log n}{\epsilon^2}$  times: uniformly select an interval  $I \in J$ . Let  $J^*$  be the set of intervals that were selected and let  $H = \cup_{I \in J^*} I$ .
3. Reject if for every  $p \in (\frac{n}{4}, \frac{n}{2}]$  the set  $H \cup T$  contains a  $p$ -witness. Otherwise accept.

**Theorem 3.2.** Algorithm  $PT$  is a 1-sided error, non-adaptive  $\epsilon$ -test for periodicity. Its query complexity is  $O(\sqrt{n \log n}/\epsilon^2)$ .

**Proof:** The query complexity is obvious. The test is 1-sided error since it rejects only if it finds a  $p$ -witness for every  $p \in (\frac{n}{4}, \frac{n}{2}]$  and thus by Facts 2.6 and 2.7 it can't reject a periodic string. The error probability estimate follows directly using the union bound and Lemma 3.3 below. ■

**Lemma 3.3.** Let  $p \in (\frac{n}{4}, \frac{n}{2}]$  and  $\alpha \in \Sigma^n$  be  $\epsilon$ -far from  $p$ -periodic. Then with probability at most  $\frac{1}{n}$  the set  $H \cup T$  of algorithm  $PT$  does not contain a  $p$ -witness for  $\alpha$ .

In the proof of Lemma 3.3 we use the following simple facts about the interval set  $J$ .

**Fact 3.4.**

1. The intervals in  $J$  are pairwise disjoint, each of size  $\sqrt{\frac{n}{\log n}}$  and  $\cup_{I \in J} I = [n]$ .
2. For every  $p \in (\frac{n}{4}, \frac{n}{2}]$  and  $t \in [n]$  there is an interval  $I \in J$  such that there exists a unique  $s \in I$  for which  $s \equiv t \pmod{p}$ .
3. Fix  $p \in (\frac{n}{4}, \frac{n}{2}]$  and let  $T \subset [n]$  be any set of  $\ell$  points such that  $|T \cap I| = 1$  for every  $I \in J$ . Then for any fixed interval  $I \in J$  there are at most 8 points  $s$  for which there exists  $t \in T$  such that  $s \equiv t \pmod{p}$ .

The proof of Lemma 3.3 is in the Appendix at subsection 8.1.

**Theorem 3.5.** *For any  $g \leq \frac{n}{2}$  there is a 1-sided error, non-adaptive  $\epsilon$ -test for  $\text{period}(\leq g)$ . Its query complexity is  $O(\sqrt{g \log g}/\epsilon^2)$ .*

**Proof (idea):** Let  $\alpha \in \Sigma^n$ . We think of  $\alpha$  as being composed of  $\frac{n}{2g}$  pieces of length  $n' = 2g$  each. We now run the  $\epsilon$ -test for periodicity for strings of length  $n'$  and for each query  $q \in [n']$  we query  $q$  in one the pieces that is chosen randomly and independently for each query. We avoid further details here. ■

## 4 Upper Bound for Testing $\text{period}(\leq g)$ , where $g \leq \frac{\log n}{4}$

In this section we describe an algorithm for testing whether a string has  $\text{period}(\leq g)$ , where  $g \leq \frac{\log n}{4}$  that has query complexity  $\text{poly}(\log g)$ . The technicalities are somewhat involved, however, the intuition is simple and motivated by the following reasoning: Our goal is to find a witness for not being  $r$ -periodic for every  $r \in [g]$ . One thing that we could do easily is to check whether  $\alpha$  is  $q$ -periodic where  $q$  is the product of all numbers in  $[g]$  (this number would be a actually too big, but for understanding the following intuition this should suffice). Now if  $\alpha$  is not  $q$ -periodic then it is certainly not in  $\text{period}(\leq g)$  and we are done. Yet as far as we know  $\alpha$  maybe far from  $\text{period}(\leq g)$  and  $q$ -periodic. Our goal is to show that if a  $\alpha$  is far from  $p$ -periodic, where  $p \leq g$ , then there exists  $q' \gg g$ , such that,  $\alpha$  is far from  $q'$ -periodic and  $p$  divides  $q'$ . Indeed it follows from Lemma 4.6 that if  $\alpha$  is far from  $p$  then there exists such a  $q'$  as above. This together with the following definition enables us to construct our  $\epsilon$ -test.

**Definition 4.1.** *For integers  $n, g \in [n]$ , a gcd-cover of  $[g]$  is a set  $E \subseteq [n]$  such that for each  $\ell \leq g$ , there exists a subset  $I \subseteq E$ , that satisfies  $\ell = \text{gcd}(I)$ .*

The importance of a gcd-cover of  $[g]$  is the following: We prove that if  $\ell = \text{gcd}(I)$  for an integer  $\ell$  and a subset  $I$ , then, the assumption that  $\alpha$  is far from being  $\ell$ -periodic implies that it is also sufficiently far from being  $t$ -periodic for some  $t \in I$ . Using this, let  $E$  be a gcd-cover of  $[g]$ . Then we are going to query for each  $t \in E$  enough pairs  $i, j \in [n]$  such that  $i \equiv j \pmod t$ . We say that such pairs  $i, j$  cover  $t$ . Now for each  $\ell \leq g$  there is a set  $I_\ell \subseteq E$  such that  $\ell = \text{gcd}(I_\ell)$ . The pairs of queries that cover  $I_\ell$  cover  $\ell$ . Thus if  $\alpha$  is ‘far’ from being  $\ell$ -periodic then, as  $\text{gcd}(I_\ell) = \ell$ , there exists a period  $t$ ,  $t \in I_\ell$ , such that the string is far from from  $t$ -periodic. We thus expect that this  $t$  will distinguish between strings that are  $\ell$ -periodic and those that are  $k\ell$ -periodic but far from being  $\ell$ -periodic.

As the complexity of the test will depend crucially on the size  $E$  (the gcd-cover of  $E$ , we need to guarantee the existence of a small gcd cover. This is done in the following Lemma, which also brings in an additional technical requirement.

**Lemma 4.2.** *For every integers  $n$  and  $g \in \left[\frac{\log n}{4}\right]$ , there exists a set  $E \subseteq [\sqrt{n}]$ , of size  $|E| = O((\log g)^3)$ , that is a gcd-cover of  $[g]$ .*

The proof of Lemma 4.2 is in the Appendix at subsection 8.1.

A gcd-cover of  $[g]$  as in the Lemma is called an *efficient gcd-cover of  $[g]$* . We prove lemma 4.2 in the Appendix at subsection 8.2. We next describe our algorithm.

### Algorithm SPT

**Input:** a string  $\alpha \in \Sigma^n$ , a distance parameter  $0 < \epsilon < 1$  and a threshold period  $g \leq \frac{\log n}{4}$ ;

1. Set  $E$  to be an *efficient gcd-cover* of  $[g]$ .
2. Let  $M = \frac{|E| \cdot \log(8|E|)}{2\epsilon}$ . For each  $\ell \in E$  select  $M$  pairs of integers uniformly and with repetitions from the set  $\{\{x, y\} \mid x \equiv y \pmod{\ell} \text{ and } x, y \in [n]\}$ . Let  $Q$  be the resulting (multi)set.
3. Reject if for each  $\ell \in [g]$  the set  $Q$  contains a witness that  $\alpha$  is not in  $\ell$ -periodic. Otherwise accept.

**Theorem 4.3.** *Algorithm SPT is a 1-sided error, non-adaptive  $\epsilon$ -test for  $\text{period}(\leq g)$ , where  $g \leq \frac{\log n}{4}$ . Its query complexity is  $\tilde{O}\left(\frac{(\log g)^6}{\epsilon}\right)$ .*

**Proof:** Since for every member of  $E$  we used  $M$  queries and  $|E| = O((\log g)^3)$  the estimate on the query complexity follows. The fact that the algorithm never rejects a  $\ell$ -periodic string, for  $\ell \leq g$  is immediate, as the algorithm rejects only when it finds a witness the input string is not in  $\text{period}(\leq g)$ , where  $g \leq \frac{\log n}{4}$ .

In order to compute the success probability of algorithm *SPT* we first need the following definition.

**Definition 4.4.** *A string  $\alpha \in \Sigma^n$  is called  $\epsilon$ -bad if for every  $\ell \leq g$  and every  $S \subseteq E$ , where  $\ell = \text{gcd}(S)$ , there exists  $s \in S$  for which  $\alpha$  is  $\frac{\epsilon}{2|E|}$ -far from being  $s$ -periodic.*

Note: if  $\alpha$  is  $\epsilon$ -bad and  $s, \ell$  as in the definition, then a witness for it showing that it is not  $s$ -periodic is also a witness that it is not  $\ell$ -periodic.

The proof of the Theorem now follows from Lemma 4.5 and Lemma 4.6. ■

**Lemma 4.5.** *Algorithm SPT rejects every  $\alpha$  that is  $\epsilon$ -bad with probability at least  $\frac{3}{4}$ .*

*Proof.* Let  $\alpha \in \Sigma^n$  be  $\epsilon$ -bad. Let  $B \subseteq E$  be the set that contains every  $b \in E$  such that  $\alpha$  is  $\frac{\epsilon}{2|E|}$ -far from being  $b$ -periodic. By Definition 4.4, for every  $\ell \leq g$  there exists a  $b \in B$  such that  $\ell \mid b$ , and hence if we find a  $b$ -witness for each  $b \in B$  - this is also a witness that  $\alpha$  is not in  $\text{period}(\leq g)$ . Thus, it is enough to prove that for each  $b \in B$ , the probability that there is no  $b$ -witness in  $Q$  is at most  $\frac{1}{8|E|}$ , as then by the union bound we are done.

Let  $\delta = \frac{\epsilon}{2|E|}$  and let  $b$  be a member of  $B$ , namely  $\alpha$  is  $\delta$ -far from being  $b$ -periodic. By Corollary 2.10, for a random  $x, y \in [n]$ , such that  $x \equiv y \pmod{b}$  we have that  $\text{Prob}(\alpha_x \neq \alpha_y) \geq \delta$ . Since  $M$  random pairs are being queried for each  $\ell \in E$  and in particular for  $b$ , the probability that a  $b$ -witness is not found is at most  $(1 - \delta)^M \leq \frac{1}{(8|E|)}$ . ■

**Lemma 4.6.** *If a string  $\alpha \in \Sigma^n$  is  $\epsilon$ -far from having short  $\text{period}(\leq g)$ , where  $g \leq (\log n)/4$  then it is a  $\epsilon$ -bad string.*

The proof of Lemma 4.6 is in the Appendix at subsection 8.3.

## 5 Lower Bound for Testing $\text{period}(\leq g)$

In this section we prove that every adaptive, 2-sided error,  $\frac{1}{32}$ -test for  $\text{period}(\leq g)$  uses  $\Omega(\sqrt{g/(\log g \cdot \log n)})$  queries. We prove this for  $\Sigma = \{0, 1\}$ . This implies the same bound for every alphabet that contains at least two symbols.

**Theorem 5.1.** *Any adaptive 2-sided, error  $\frac{1}{32}$ -test for  $\text{period}(\leq g)$ , uses  $\Omega(\sqrt{g/(\log g \cdot \log n)})$  queries.*

*Proof.* Fix  $g \leq n/2$ . We prove the theorem by using Yao's principle. That is, we construct a distribution  $\mathcal{D}$  over legitimate instances (strings that are either in  $period(\leq g)$ , or strings that are  $\frac{1}{32}$ -far from  $period(\leq g)$ ) and prove that any adaptive deterministic tester, that uses  $m = o(\sqrt{g/(\log g \cdot \log n)})$  queries, gives an incorrect answer with probability greater than  $\frac{1}{3}$ .

In order to define  $\mathcal{D}$  we use auxiliary distributions  $\mathcal{D}_P, \mathcal{D}_N$  and the following notation. Let  $Primes = \{p \mid p \text{ is a prime such that } p \leq g\}$  and let  $r = |Primes|$ . According to the *Prime Number Theorem* [11, 12, 13],  $|Primes| = \theta(\frac{g}{\log g})$ .

We now define distributions  $\mathcal{D}_P, \mathcal{D}_N$ .

- $\mathcal{D}_N$  is simply the uniform distribution over  $\Sigma^n$ .
- An instance  $\alpha$  of length  $n$  is selected according to distribution  $\mathcal{D}_P$  as follows. Uniformly select  $p \in Primes$ , then uniformly select  $\omega \in \Sigma^p$  and finally set  $\alpha = (\omega)^{\frac{n}{p}}$ .

We next define distribution  $\mathcal{D}$ . Let  $\mathcal{G}$  be the event that  $\alpha$  is  $\frac{1}{16}$ -far from having  $period(\leq g)$ . Let  $\mathcal{D}_{N/\mathcal{G}}$  be the distribution  $\mathcal{D}_N$  given that event  $\mathcal{G}$  is true. A string  $\alpha$  is selected according to distribution  $\mathcal{D}$  by choosing one of the distributions  $\mathcal{D}_P, \mathcal{D}_{N/\mathcal{G}}$  with equal probability and then selecting  $\alpha$  according to the distribution chosen. Namely  $\mathcal{D} = \frac{1}{2}\mathcal{D}_P + \frac{1}{2}\mathcal{D}_{N/\mathcal{G}}$ .

We don't work directly with  $\mathcal{D}$  but rather with a simpler distribution  $\mathcal{D}'$  which approximates  $\mathcal{D}$  well enough and is defined as follows  $\mathcal{D}' = \frac{1}{2}\mathcal{D}_P + \frac{1}{2}\mathcal{D}_N$ . Let  $\mathcal{B}$  be the event that the tester gives an incorrect answer. We prove that  $\text{Prob}_{\mathcal{D}'}(\mathcal{B}) \geq \frac{2}{5}$ . This is indeed sufficient as  $\text{Prob}_{\mathcal{D}}(\mathcal{B}) \geq \text{Prob}_{\mathcal{D}'}(\mathcal{B}) - \text{Prob}_{\mathcal{D}_N}(\overline{\mathcal{G}})$ . Using Proposition 2.11,  $\text{Prob}_{\mathcal{D}}(\mathcal{B}) \geq \frac{1-o(1)}{2} - \frac{1}{g} > \frac{1}{3}$ .

We assume with out loss of generality that for any string of length  $n$  the tester uses the same number of queries. Hence we can view the tester as a full binary decision tree of depth  $m$  which is labeled as follows. Each node of the tree represents a query location, for each internal node one of the outgoing edges is labeled by 1 and the other by 0, where 0, 1 represent the answers to the query, and each leaf is labeled either by "accept" or by "reject" according to the decision of the algorithm.

For each leaf  $l$  in the tree we associate a pair  $Q_l, f_l$ , where  $Q_l$  is the set of queries on the path from the root to the leaf  $l$  and  $f_l : Q_l \rightarrow \{0, 1\}$  is a mapping between each query and its answer, that is the labellings on the edges of the path from the root to the leaf  $l$ . Let  $L$  be the set of all leaves,  $L_0$  be the set of all leaves that are labeled by "reject" and let  $L_1$  be the set of all leaves that are labeled by "accept".

Let  $h_\ell : \Sigma^n \rightarrow \{0, 1\}$  be a function that is 1 only on strings  $\alpha \in \Sigma^n$  such that for every  $q \in Q_\ell$  we have  $\alpha_q = f_\ell(q)$ . That is  $h_\ell(\alpha) = 1$  if and only if  $\alpha$  coincides with  $Q_\ell, f_\ell$ . Let  $far : \Sigma^n \rightarrow \{0, 1\}$  be a function that is 1 only on strings  $\alpha \in \Sigma^n$  such that are  $\frac{1}{16}$ -far from  $period(\leq g)$ . Thus the  $\text{Prob}_{\mathcal{D}'}(\mathcal{B})$  is at least

$$\frac{1}{2}\sum_{\ell \in L_0} \text{prob}_{\mathcal{D}_P}[h_\ell(\alpha) = 1] + \frac{1}{2}\sum_{\ell \in L_1} \text{prob}_{\mathcal{D}_N}[far(\alpha) = 1 \wedge h_\ell(\alpha) = 1]$$

We will prove that for each  $\ell \in L_0$ ,  $\text{prob}_{\mathcal{D}_P}[h_\ell(\alpha) = 1] \geq \frac{1-o(1)}{2^m}$ , and for each  $\ell \in L_1$ ,  $\text{prob}_{\mathcal{D}_N}[far(\alpha) = 1 \wedge h_\ell(\alpha) = 1] \geq \frac{1-o(1)}{2^m}$ . This implies that  $\text{Prob}_{\mathcal{D}'}(\mathcal{B}) \geq \frac{1}{2}\sum_{\ell \in L} \frac{1-o(1)}{2^m} \geq \frac{1-o(1)}{2}$ .

Indeed, recall that a string is selected according to  $\mathcal{D}_P$  by first selecting  $z \in Primes$  and then selecting a  $z$ -periodic string. For  $Q \subset [n]$ ,  $|Q| = m$  let  $\mathcal{A}(Q)$  be the event that for  $\alpha$  selected according to  $\mathcal{D}_P$ , there exists no  $j, k \in Q$  such that  $j \equiv k \pmod{z}$ . For any fixed  $i < j \in Q$  there are at most  $\log n$  prime divisors of  $j - i$ . Hence for each  $\ell \in L_0$ ,  $\text{prob}_{\mathcal{D}_P}[\mathcal{A}(Q_\ell)] \geq 1 - \frac{m^2 \cdot \log n}{r} \geq 1 - o(1)$ . Observe that for any fixed  $\ell$  if  $\mathcal{A}(Q_\ell)$  occurs then according to the definition of  $\mathcal{D}_P$ , for each  $q \in Q_\ell$ ,  $\alpha_q$  is selected uniformly and independently of any other  $q'$ . Hence,  $\text{prob}_{\mathcal{D}_P}[h_\ell(\alpha) = 1] \geq \text{prob}_{\mathcal{D}_P}[h_\ell(\alpha) = 1 \wedge \mathcal{A}(Q_\ell)] \geq \text{prob}_{\mathcal{D}_P}[h_\ell(\alpha) = 1 \mid \mathcal{A}(Q_\ell)] \cdot \text{Prob}_{\mathcal{D}_P}[\mathcal{A}(Q_\ell)] \geq \frac{1}{2^m}(1 - o(1))$ .



Observe that for each  $\ell \in L_1$ ,  $\text{Prob}_{\mathcal{D}_N}[far(\alpha) = 1 \wedge h_\ell(\alpha) = 1] \geq \text{Prob}_{\mathcal{D}_N}[far(\alpha) = 1 \mid h_\ell(\alpha) = 1] \cdot \text{Prob}_{\mathcal{D}_N}[h_\ell(\alpha) = 1]$ . By the definition of  $\mathcal{D}_N$ ,  $\text{Prob}_{\mathcal{D}_N}[h_\ell(\alpha) = 1] = \frac{1}{2^m}$ . Also, using the definition just before Proposition 2.11 we see that  $\text{Prob}_{\mathcal{D}_N}[far(\alpha) = 1 \mid h_\ell(\alpha) = 1] = \text{prob}_{U(Q_\ell)}[far(\alpha) = 1] \geq 1 - \frac{1}{g}$  (where the last inequality is by Proposition 2.11). ■

## 6 Lower Bound for Testing $period(\leq g)$ , where $g \leq \frac{\log n}{4}$

In this section we prove that every 2-sided error,  $\frac{1}{16}$ -test for  $period(\leq g)$ , where  $g \leq \frac{\log n}{4}$ , uses  $\Omega((\log g)^{\frac{1}{4}})$  queries. We prove this for  $\Sigma = \{0, 1\}$ . This implies the same bound for every alphabet that contains at least two symbols. This also shows that the test presented in Section 4 cannot be dramatically improved.

**Theorem 6.1.** *Any 2-sided error  $\frac{1}{16}$ -test for  $period(\leq g)$ , where  $g \leq \frac{\log n}{4}$ , requires  $\Omega((\log g)^{\frac{1}{4}})$  queries.*

*Proof.* The proof is very similar to the proof of Theorem 5.1. We just describe here the two probabilities  $\mathcal{D}_{\mathcal{P}}$  and  $\mathcal{D}_{\mathcal{N}}$  that are concentrated on positive inputs (those that have  $period(\leq g)$ ) and negative inputs (those that are  $\frac{1}{16}$ -far from being  $period(\leq g)$ ) respectively.

Let  $S = \{p_1, \dots, p_k\}$  where  $k = 1 + \sqrt{\log g}$  and each  $p_i$  is a prime such that  $2^{\sqrt{\log g} - 0.9} \leq p_i \leq 2^{\sqrt{\log g}}$ . According to the *Prime Number Theorem* [11, 12, 13] there exists at least  $\frac{1}{8\sqrt{\log g}} \cdot 2^{\sqrt{\log g}} > 2 + \sqrt{\log g}$  such primes. Hence the set  $S$  is well defined. Let  $t = \prod_{p \in S} p$  and let  $Z = \{\frac{t}{p} \mid p \in S\}$ . Note that  $t > g$  while  $z \leq g$  for every  $z \in Z$ .

We now define distributions  $\mathcal{D}_{\mathcal{P}}$ ,  $\mathcal{D}_{\mathcal{N}}$  (we assume w.l.o.g that  $t$  divides  $n$ ).

- We use an auxiliary distribution  $\mathcal{U}$  in order to define  $\mathcal{D}_{\mathcal{N}}$ . To select an instance according to  $\mathcal{U}$ , select a random string  $\omega \in \Sigma^t$  and then set  $\alpha = (\omega)^{\frac{n}{t}}$ . Let  $G$  be the event that  $\alpha$  is  $\frac{1}{16}$ -far from being  $period(\leq g)$ . Then  $\mathcal{D}_{\mathcal{N}}$  is defined as the distribution  $\mathcal{U}|G$ , namely  $\mathcal{U}$  given  $G$ .
- An instance  $\alpha$  of length  $n$  is selected according to distribution  $\mathcal{D}_{\mathcal{P}}$  as follows. Uniformly a select  $z \in Z$ , select uniformly a string  $\omega \in \Sigma^z$  and then set  $\alpha = (\omega)^{\frac{n}{z}}$ .

We omit further details from this extended abstract. ■

## 7 Lower Bound for Testing Periodicity

Let  $\Sigma = \{0, 1\}$ , we prove the following lower bound on the number of queries that is needed for testing  $period(\leq \frac{n}{2})$  over  $\Sigma$ . This clearly implies the same lower bound for any alphabet that contains at least two letter. It also shows that the test presented in Section 3 is asymptotically optimal.

**Theorem 7.1.** *Any non-adaptive 1-sided error  $\frac{1}{16}$ -test for periodicity requires  $\Omega(\sqrt{n \log n})$  queries.*

**Proof:** A 1-sided error test rejects only when the input string is not periodic. Therefore according to corollary 2.6 such a test rejects only if the set of queries it uses contains a witness that the string is not periodic. To prove the Theorem we use Yao's principle (the easy direction). Namely, we construct a distribution on  $1/16$ -far inputs and show that any *deterministic* non-adaptive test that queries at most  $\frac{\sqrt{n \log n}}{100}$  queries finds a witness for non-periodicity with probability at most  $1/3$ .

Let  $U$  be the uniform distribution over  $\Sigma^n$ , and let  $\mathcal{G}$  be the event that a string selected according to  $U$  is  $\frac{1}{16}$ -far from being periodic. Let  $D$  be the distribution  $U/\mathcal{G}$  ( $U$  given  $\mathcal{G}$ ). Namely,  $D$  is uniform over strings in  $\Sigma^n$  that are  $1/16$ -far from being periodic. Let  $\mathcal{B}$  be the event that the set of queries  $Q$  contains a witness. Proposition 2.11  $\text{Prob}_U(\mathcal{G}) \geq 1 - \frac{1}{n}$ . Thus, it is enough to prove that  $\text{prob}_U(\mathcal{B}) < \frac{1}{4}$ , since  $\text{prob}_D(\mathcal{B}) \leq \text{Prob}_U(\mathcal{B}) + \text{Prob}_U(\overline{\mathcal{G}}) \leq \frac{1}{4} + \frac{1}{n}$ . The proof follows from Lemma 7.2 below. ■

**Lemma 7.2.** *Let  $Q \subseteq [n]$  be a set of  $\frac{\sqrt{n \log n}}{100}$  queries. Then for  $\alpha$  chosen according to  $U$ , the probability that  $Q$  contains a witness that  $\alpha$  is not periodic is at most  $\frac{1}{4}$ .*

*Proof.* Let  $Q, U$  be as in the Lemma. For a given input  $\alpha \in \Sigma^n$  let  $\alpha_q \in \Sigma$  be the answer to query  $q \in Q$ . Let  $G$  be the complete graph on  $Q$ . We call an edge  $(x, y)$  a  $t$ -edge if  $p = |x - y|$ . For every  $p \in (\frac{n}{4}, \frac{n}{2}]$ , let  $L_p$  be the event that all the vertices (queries) in  $G$  that are connected to  $t$ -edges,  $t \equiv 0 \pmod{p}$  have the same answer  $\sigma \in \Sigma$ . Note that if  $L_p$  occurs for some  $p$  then  $Q$  does not contain a witness that  $\alpha$  is not  $p$ -periodic. Thus it is enough to prove that  $\text{Prob}_U(\cup_{p=n/4}^{n/2} L_p) \geq \frac{3}{4}$ .

The events  $L_p$  are generally dependent. We are going to construct a set  $S^* \subseteq (\frac{n}{4}, \frac{n}{2}]$ , of size  $\Omega(\frac{\sqrt{n}}{(\log n)^{\frac{3}{2}}})$ , for which: (a) The events  $L_p, p \in S^*$  are independent. (b) For each  $p \in S^*$  there exists at most  $\frac{\log n}{25}$   $t$ -edges in  $G$  such that  $t \equiv 0 \pmod{p}$ .

Given such set  $S^*$ , the proof is completed as for any fixed  $p \in S^*$  all queries that are connected with  $p$ -edges are labeled by '1' with probability at least  $\frac{1}{2^{\log n/25}} = n^{-1/25}$ . Namely  $L_p$  occurs. Hence,  $\text{Prob}_U(\forall p \in S^*, \bar{L}_p) \leq \left(1 - \frac{1}{n^{1/25}}\right)^{|S^*|}$ , which for a large enough  $n$  is smaller than  $\frac{1}{n}$ .

For each  $p \in (\frac{n}{4}, \frac{n}{2}]$  let  $G_p$  be the subgraph of  $G$  that is spanned by all  $t$ -edges with  $t \equiv 0 \pmod{p}$ . Note that if for two different  $p, p' \in (\frac{n}{4}, \frac{n}{2}]$ ,  $G_p$  and  $G_{p'}$  are vertex disjoint then  $L_p$  and  $L_{p'}$  are independent.

By our assumptions  $G$  contains at most  $|Q|^2/2 \leq n \log n/200$  edges. Note also that for each  $p \in (\frac{n}{4}, \frac{n}{2}]$  there are at most 4 multiples of  $p$ . Hence, there is a set  $P \subseteq (\frac{n}{4}, \frac{n}{2}]$ , of size  $n/8$ , such that for every  $q \in P$ ,  $G$  contains at most  $(\log n)/25$   $t$ -edges for which  $t$  is a multiple of  $p$ . The set  $S^*$  will be a subset of  $P$  and thus the second requirement above automatically holds. Furthermore, the way we construct  $S^*$  will automatically ensure that  $G_p, G_q$  are vertex disjoint for each  $q, p \in S^*$ . Hence by the remark above, this will automatically ensure the first requirement.

We specify now the actual construction of  $S^*$ : We start with  $S^* = \phi$  and  $P$  as above (with  $n/8$  elements). We add to  $S^*$  an arbitrary  $p \in P$  and delete  $p$  from  $P$ . We set  $G' = G_p$ . Next we delete from  $P$  all  $q \in P$  such that there is a  $t$ -edge whose end point is in  $G'$  and  $q$  divides  $t$ . Namely, we avoid using all  $q \in P$  for which  $G_q$  vertex-intersects  $G'$ . This ensures that for the next chosen  $p \in P$ , that is picked,  $G_p$  will be disjoint from  $G'$ . We now iterate the process each time picking a new  $p$  from the new set  $P$ , set  $G' := G' \cup G_p$ ,  $S^* := S^* \cup \{p\}$  and deleting  $p$  and the corresponding  $q$ 's as above from  $P$ .

We claim that we can construct this way  $S^*$  of size  $|S^*| \geq \frac{\sqrt{n}}{(\log n)^{\frac{3}{2}}}$ . Indeed, each time we add a new  $p \in P$  to  $S^*$ , as  $G_p$  contains at most  $\log n/25$  edges, we add at most  $2 \log n/25$  new vertices to  $G'$ . We then delete all  $q$ 's for which  $G_q$  intersects the current  $G'$ . This results in deleting at most  $|Q| \cdot 2 \log n/25$  additional members from  $P$ . Hence, this process can continue as long as  $P$  is not empty, that is, for at least  $|P|/(|Q| \cdot 2 \log n/25) = \Omega(\frac{\sqrt{n}}{(\log n)^{\frac{3}{2}}})$  times. ■

## References

- [1] S. Goldwasser O. Goldreich and D. Ron. Property testing and its connection to learning and approximation. *Journal of the ACM*, 45:653–750, 1998.
- [2] R. Rubinfeld and M. Sudan. Robust characterization of polynomials with applications to program testing. *SIAM Journal of Computing*, 25:252–271, 1996.
- [3] F. Ergun, S. Muthukrishnan, and C. Sahinalp. Sublinear methods for detecting periodic trends in data streams. In *LATIN 2004, Proc. of the 6th Latin American Symposium on Theoretical Informatics (LATIN'04)*, pages 16–28, 2004.
- [4] P. Indyk, N. Koudas, and S. Muthukrishnan. Identifying representative trends in massive time series data sets using sketches. In *VLDB 2000, Proceedings of 26th International Conference on Very Large Data Bases, September 10-14, 2000, Cairo, Egypt*, pages 363–372. Morgan Kaufmann, 2000.
- [5] R. Krauthgamer O. Sasson. Property testing of data dimensionality. In *Proceedings of the fourteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 18–27. Society for Industrial and Applied Mathematics, 2003.
- [6] A. C. Gilbert, S. Guha, P. Indyk, S. Muthukrishnan, and M. Strauss. Near-optimal sparse fourier representations via sampling. In *STOC 2002, Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 152–161, 2002.
- [7] Alex Samorodnitsky and Luca Trevisan. A PCP characterization of NP with optimal amortized query complexity. pages 191–199, 2000.
- [8] Johan Hästad and Avi Wigderson. Simple analysis of graph tests for linearity and pcp. *Random Struct. Algorithms*, 22(2):139–160, 2003.
- [9] E. Fischer. The art of uninformed decisions: A primer to property testing. *The computational complexity column of The Bulletin of the European Association for Theoretical Computer Science*, 75:97–126, 2001.
- [10] D. Ron. Property testing (a tutorial). In *Handbook of Randomized computing*. Kluwer Press, 2001.
- [11] J. Hadamard. Sur la distribution des zéros de la fonction  $\zeta(s)$  et ses conséquences arithmétiques. *Bull. Soc. Math. France*, 24:199–220, 1896.
- [12] V. Poussin. Recherces analytiques sur la théorie des nombres premiers. *Ann. Soc. Sci. Bruxelles*, 1897.
- [13] D.J. Newman. Simple analytic proof of the prime number theorem. *Amer. Math. Monthly*, 87:693–696, 1980.

## 8 Appendix

### 8.1 Proof of Lemma 3.3

Fix  $p$ ,  $\alpha$  as in the Lemma. Using Fact 2.9 we get  $\text{dist}(\alpha, p\text{-periodic}) = \sum_{i=0}^{p-1} \text{dist}(\alpha_{Z(p,i)}, \text{homogeneous}) \geq \epsilon n$ . Note, however, that as  $p > \frac{n}{4}$  each  $p$ -section of  $\alpha$  is of size at most 4. Hence, there exists a set  $S \subseteq \{0, \dots, p-1\}$  of size at least  $\frac{\epsilon n}{4}$ , such that for every  $i \in S$ ,  $\text{dist}(\alpha_{Z(p,i)}, \text{homogeneous}) \geq 1$ . In particular, for each  $i \in S$  there is an unordered pair of indexes  $\{j, k\} \subset [n]$ , such that  $j, k \equiv i \pmod{p}$  and  $\alpha_j \neq \alpha_k$ . Let  $W$  be the set of size  $|S|$  that contains exactly one such pair for each  $i \in S$ . Let  $W_R$  be the set of the 'right-ends' of pairs in  $W$ , namely  $W_R = \{j \mid (i, j) \in W, i < j\}$ . Note that  $|W_R| = |S|$ .

Recall that  $\cup_{I \in J} I = [n]$  (Fact 3.4) and since  $|W_R| \geq \epsilon \frac{n}{4}$ , there exists a set  $G_R \subseteq J$  of size  $|G_R| \geq \frac{\epsilon |J|}{8} = \frac{\epsilon \sqrt{n} \log n}{8}$  such that for every  $I \in G_R$ ,  $|I \cap W_R| \geq \frac{\epsilon |I|}{8}$ .

In step 1 of algorithm *PT* the set  $T$  is formed by choosing a random member  $t_i$  in each  $I \in J$ . By the definition of  $G_R$ ,  $\text{Prob}(t_i \in W_R \mid I_i \in G_R) \geq \frac{\epsilon}{8}$ . Thus by the Chernoff bound  $\text{Prob}\left(|T \cap W_R| < \frac{\epsilon |G_R|}{16}\right) \leq e^{-\frac{\epsilon |G_R|}{64}}$ . Let  $R = T \cap W_R$  and let  $\mathcal{A}$  be the event that  $|R| \geq \frac{\epsilon |G_R|}{16}$ .

Let  $L = \{j \mid k \in R, (j, k) \in W\}$ , namely  $L$  contains the 'left-ends' of the pairs of  $W$  that corresponds to the right-end points in  $R$ . Let  $G_L = \{I \in J \mid I \cap L \neq \emptyset\}$ . Namely, an interval  $I$  is  $G_L$  if it contains a point from  $L$ . In step 2 algorithm *PT* uniformly selects a set of intervals  $J^*$  by repeatedly selecting a random interval  $I \in J$ . If  $I \in J^* \cap G_L$ , then  $I \cup T$  contains a  $p$ -witness as  $I$  contains a left-end of a  $p$ -witness whose right-end must be (by definition) in  $R \subseteq T$ . Thus it is enough to show that with probability at most  $\frac{1}{n}$ ,  $J^* \cap G_L = \emptyset$ .

Indeed, assuming that the event  $\mathcal{A}$  is satisfied.  $|G_L| \geq \frac{\epsilon^2 |J|}{2^{10}}$ . This is true as each member of  $L$  is a left-end of a  $p$ -witness for different  $t \in T$ , thus by Fact 3.4 each interval in  $G_L$  can contain at most 8 points from  $L$ . However,  $|L| = |R| \geq \frac{\epsilon |G_R|}{16} \geq \frac{\epsilon^2 n}{2^6}$  hence  $|G_L| \geq \frac{|L|}{8}$  as claimed.

Finally, we conclude that for one randomly chosen interval  $I \in_R J$   $\text{Prob}(I \notin G_L) \leq 1 - \frac{\epsilon^2}{2^{10}}$ . Thus by selecting  $m$  independent  $I$ 's the probability that none of the  $I$ 's is in  $G_L$  is at most  $(1 - \frac{\epsilon^2}{2^{10}})^m < \frac{1}{n}$ . ■

### 8.2 Proof of Lemma 4.2

Let  $n, g$  be fixed with  $g \leq \frac{\log n}{4}$ . We use the following notations.

Let  $\text{Primes} = \{p \mid p \text{ is a prime such that } p \leq g\}$  and let  $m = |\text{Primes}|$ . According to the prime number theorem  $m \leq g/\log g$ . For an integer  $t$  let the *unique prime decomposition* of  $t$ ,  $\text{UPD}(t)$  be the unique set of pairs  $\text{UPD}(t) = \{(p_i, \eta_i) \mid (p_i, \eta_i) \mid t, i = 1, \dots, m\}$ , such that  $p_i, i = 1, \dots, m$  are distinct members of  $\text{Primes}$  and  $t = \prod_{p_i \in \text{Primes}} p_i^{\eta_i}$ . Let  $P(t)$  be the set of prime divisors of  $t$ . Note that in  $\text{UPD}(t)$  we include also the primes that do not divide  $t$ .

The following gives a characterization of a *gcd cover* of  $[g]$ , from which we need only the sufficient part.

**Proposition 8.1.** *A set  $E \subseteq [n]$  is a gcd-cover of  $[g]$  if the following condition is satisfied.*

*For every  $t \in [g]$  and every  $(p_i, \eta_i) \in \text{UPD}(t)$  (including those  $p_i$ 's for which  $\eta_i = 0$ ), there exists  $e \in E$  such that  $(p_i, \eta_i) \in \text{UPD}(e)$  and  $t$  divides  $e$ .*

*Proof.* Assume that the premises of the Proposition hold. Fix  $t$  with unique prime decomposition  $\text{UPD}(t) = \{(p_i, \eta_i) \mid i = 1, \dots, m\}$ . For each  $(p_i, \eta_i) \in \text{UPD}(t)$  let  $e_i$  denote the appropriate member of  $E$  for which  $(p_i, \eta_i) \in \text{UPD}(e_i)$  and such that  $t$  divides  $e_i$ . Then it is easy to see that  $\text{gcd}(\{e_1, \dots, e_m\}) = t$ . ■

The way to construct a *gcd*-cover will be via the construction of the following set which we call *prime cover* of  $[g]$ .

**Definition 8.2.** [*Prime Cover of  $[g]$* ] We say that a collection of subsets of Primes,  $\mathcal{R} \subseteq 2^{\text{Primes}}$  is a *prime cover* of  $[g]$  if for every  $t \in [g]$  and every  $p \in \text{Primes}$  there exists an  $R \in \mathcal{R}$  such that  $p \in R$  and  $R \cap P(t) \subseteq \{p\}$ .

Given a small size collection, in which all sets have small cardinality gives an immediate construction of an efficient *gcd* cover. This is asserted by the following Claim.

**Claim 8.3.** Let  $\mathcal{R} \subseteq 2^{\text{Primes}}$  be a *prime cover* of  $[g]$  then there is a set  $E \subseteq [n]$  that is a *gcd*-cover of  $[g]$  of size at most  $|\mathcal{R}| \log g$ . Moreover, for every  $e \in E$ ,  $e < g^{(c_0 \cdot g)/\log g}$  where  $c_0$  is a universal constant.

**Proof:** For each  $p \in \text{Primes}$  let  $\kappa(p)$  be the maximum integer such that  $p^{\kappa(p)} \leq g$ .

Assume that  $\mathcal{R}$  is a *prime cover* of  $[g]$ . For each  $R \in \mathcal{R}$  we define the following set of at most  $\log g$  numbers.

$$E(R) = \left\{ \prod_{r \in R} r^{\min\{i, \kappa(r)\}} \prod_{q \in \text{Primes} \setminus R} q^{\kappa(q)} \mid i = 0, \dots, \log g \right\}.$$

Note that each such number is of size  $g^{(c_0 \cdot g)/\log g}$  as each multiplier is at most  $g$  and there are at most  $(c_0 \cdot g)/\log g$  primes up to  $g$ . Now let  $E = \cup_{R \in \mathcal{R}} E(R)$ . We claim that  $E$  is a *gcd* cover of  $[g]$ . To see this fix any  $t \in [g]$  and let  $(p, \eta) \in UPD(t)$ . Let  $R$  be such that  $p \in R$  and  $R \cap P(t) \subseteq \{p\}$ . Then  $e = p^\eta \prod_{q \in \text{Primes} - \{p\}} r^{\kappa(r)}$  is in  $E(R)$  and meets the conditions of Proposition 8.1 for  $t$  and  $(p, \eta)$ . Note that the upper bound on the size of  $E$  is as asserted. ■

We are left now with the task of constructing an ‘efficient’ *prime cover* of  $[g]$  in terms of the collection size and the cardinality of its members.

**Proposition 8.4.** For every  $g$  there exists a set of size  $2(\log g)^2$ , that is a *prime cover* of  $[g]$ .

*Proof.* We prove the proposition as follows. We show a probabilistic construction of a family of  $2(\log g)^2$  sets, and prove that with strictly positive probability it is a *prime cover* of  $[g]$ . This implies that there must exist a set of size  $2(\log g)^2$ , that is a *prime cover* of  $[g]$ .

We choose a set  $R \subseteq \text{Primes}$  in the following way. For each  $p \in \text{Primes}$  we put  $p$  in  $R$  with probability  $\frac{1}{\log g}$ . For specific  $t \in [g]$  and  $p \in \text{Primes}$ , the probability that  $p \in R$  and  $R \cap P(t) \subseteq \{p\}$ , is at least  $\frac{1}{\log g} \left(1 - \frac{1}{\log g}\right)^{\log g} \geq \frac{1}{2 \log g}$ . If we repeat this process for  $2(\log g)^2$  independently random  $R$ ’s the probability that every one of them fails to have  $p \in R$  and  $R \cap P(t) \subseteq \{p\}$ , is at most  $\left(1 - \frac{1}{2 \log g}\right)^{2(\log g)^2} \leq \frac{1}{g^2}$ . As there are at most  $m \cdot g$  pairs of  $p \in \text{Primes}$  and  $t \in [g]$ , then with probability at least  $1 - m/g > 1 - 1/\log g$  there must be a collection of  $2(\log g)^2$  sets that succeeds for every pair. ■

### 8.3 Proof Of Lemma 4.6

In order to prove the Lemma 4.6 we first need two claims that relate  $\text{dist}(\alpha, \ell - \text{periodic})$  and  $\text{dist}(\alpha, q - \text{periodic})$ , where  $\ell \neq q$ .

**Claim 8.5.** For every  $\alpha \in \Sigma^n$ , and every  $\ell, q \in [n]$ , such that  $\text{gcd}(\ell, q) = 1$  and  $\ell q \leq n$  the following inequality is satisfied.

$$\text{dist}(\alpha, \text{homogeneous}) \leq 2 \text{dist}(\alpha, \ell - \text{periodic}) + \text{dist}(\alpha, q - \text{periodic}).$$

**Proof:** For simplicity of the calculations we assume in the proof that  $\ell q|n$ . Note first that for  $\ell, q$  as above and every  $i \leq \ell - 1, i' \leq q - 1$  the intersection of the  $i$ 'th  $\ell$ -section and the  $i'$   $q$ -section satisfies  $|Z(\ell, i) \cap Z(q, i')| = \frac{n}{\ell q}$ . This is obvious as for any interval of integers  $[s, s + \ell q - 1]$  those two section intersect in exactly one place (Chinese remainder).

Assume that  $\alpha$  is  $\epsilon_1$ -close to being  $\ell$ -periodic and  $\epsilon_2$ -close to being  $q$ -periodic. We show that  $\text{dist}(\alpha, \text{homogeneous}) \leq (2\epsilon_1 + \epsilon_2)$ . As  $\text{dist}(\alpha, \ell\text{-periodic}) \leq \epsilon_1 n$  there is a  $\ell$ -periodic string  $\beta \in \Sigma^n$  for which  $\text{dist}(\alpha, \beta) \leq \epsilon_1 n$  and hence, by the triangle inequality,  $\text{dist}(\beta, q\text{-periodic}) \leq (\epsilon_1 + \epsilon_2)n$ . By Fact 2.9

$\text{dist}(\beta, q\text{-periodic}) = \sum_{i=0}^{q-1} \text{dist}(Z(q, i), \text{homogeneous}) \leq (\epsilon_1 + \epsilon_2)n$ . Hence there must be an  $i \in \{0, \dots, q-1\}$  for which  $\text{dist}(\alpha_{Z(q, i)}, \text{homogeneous}) \leq (\epsilon_1 + \epsilon_2) \frac{n}{q}$ . Namely,  $Z(q, i)$  contains a sub-string  $\gamma \subseteq Z(q, i)$  that is homogeneous, of length  $|\gamma| \geq (1 - (\epsilon_1 + \epsilon_2)) \frac{n}{q}$ . As the intersection of  $Z(q, i)$  and every  $\ell$ -section is of size  $\frac{n}{\ell q}$ , it follows that  $\gamma$  intersects at least  $s = (1 - (\epsilon_1 + \epsilon_2)) \frac{n}{q} \cdot (\frac{n}{\ell q})^{-1}$   $\ell$ -sections. However, since  $\beta$  is  $\ell$ -periodic all of its  $\ell$ -sections are homogeneous. As  $\gamma$  is homogeneous too, it follows that  $\beta$  is homogeneous on all of its  $\ell$ -sections that intersect  $\gamma$ . Since the union of at least  $s$   $\ell$ -sections that intersect  $\gamma$  is of size at least  $s \cdot \frac{n}{\ell} \geq (1 - (\epsilon_1 + \epsilon_2))n$ , it follows that  $\beta$  is  $(\epsilon_1 + \epsilon_2)$ -close to homogeneous and hence by the triangle inequality  $\alpha$  is  $(2\epsilon_1 + \epsilon_2)$ -close to homogeneous. ■

**Claim 8.6.** Let  $\alpha \in \Sigma^n$  and let  $t_1, t_2 \in [\sqrt{n}]$  with  $r = \gcd(t_1, t_2)$ . Then

$$\text{dist}(\alpha, r\text{-periodic}) \leq 2\text{dist}(\alpha, t_1\text{-periodic}) + \text{dist}(\alpha, t_2\text{-periodic}).$$

**Proof:** For  $t_1, t_2$  co-primes this follows directly from Claim 8.5, as  $\text{dist}(\alpha, \text{homogeneous}) = \text{dist}(\alpha, 1\text{-periodic})$ . Assume then that  $\gcd(t_1, t_2) = r$  and  $t_1 = s_1 \cdot r, t_2 = s_2 \cdot r$ . For each  $i \in [0, r-1]$  let  $\beta(i) = \alpha_{Z(r, i)}$ . By Fact 2.9 we have

$$\text{dist}(\alpha, r\text{-periodic}) = \sum_{i=0}^{r-1} \text{dist}(\alpha_{Z(r, i)}, \text{homogeneous}) = \sum_{i=0}^{r-1} \text{dist}(\beta(i), \text{homogeneous}) \quad (1)$$

As  $s_1, s_2$  are co-primes, Claim 8.5 implies that for each  $i \in [0, r-1]$ ,  $\text{dist}(\beta(i), \text{homogeneous}) \leq 2\text{dist}(\beta(i), s_1\text{-periodic}) + \text{dist}(\beta(i), s_2\text{-periodic})$ .

Plugging this into equation (1) we get,

$$\text{dist}(\alpha, r\text{-periodic}) \leq \sum_{i=0}^{r-1} (2\text{dist}(\beta(i), s_1\text{-periodic}) + \text{dist}(\beta(i), s_2\text{-periodic})) \quad (2)$$

By applying Fact 2.9 twice we get,  $\text{dist}(\alpha, t_1\text{-periodic}) = \sum_{i=0}^{r-1} \sum_{j=0}^{s_1-1} \text{dist}(\alpha_{Z(t_1, i \cdot r + j)}, \text{homogeneous}) = \sum_{i=0}^{r-1} \text{dist}(\alpha_{Z(r, i)}, s_1\text{-periodic})$ .

The last equality is true as for each fixed  $i \in [0, r-1]$  the  $(i \cdot r + j)$ -section of  $\alpha$  with respect to  $t_1$  is exactly the  $j$ -section of  $\alpha_{Z(r, i)}$  with respect to  $s_1$ . Plugging this (and the analogue expression for  $t_2$ ) into equation (2) we get:  $\text{dist}(\alpha, r\text{-periodic}) \leq 2\text{dist}(\alpha, t_1\text{-periodic}) + \text{dist}(\alpha, t_2\text{-periodic})$  as required. ■

We next prove Lemma 4.6. Let  $\alpha \in \Sigma^n$  be  $\epsilon$ -far from  $\text{period}(\leq g)$ . In order to prove that  $\alpha$  is a  $\epsilon$ -bad string we need to prove that for every  $\ell \in [g]$  and every  $S \subseteq E$ , where  $\ell = \gcd(S)$ , there exists  $s \in S$  for which  $\alpha$  is  $\frac{\epsilon}{2|E|}$ -far from being  $s$ -periodic. Let  $\ell$  be a period in  $[g]$  and  $S \subseteq E$ , be a set such that  $\ell = \gcd(S)$ . By induction, using Claim 8.6, we have,  $\text{dist}(\alpha, \ell\text{-periodic}) \leq \sum_{s \in S} 2\text{dist}(\alpha, s\text{-periodic})$ . Hence by averaging there must be a  $s \in S$  for which  $\text{dist}(\alpha, s\text{-periodic}) \geq \frac{\epsilon n}{2|E|}$  as  $|S| \leq |E|$ . ■