

The Trichotomy of Regular Property Testing

Gabriel Bathie ✉ 🏠

LaBRI, Université de Bordeaux

DIENS, Paris, France

Nathanaël Fijalkow ✉ 🏠

LaBRI, CNRS, Université de Bordeaux, France

Corto Mascle ✉ 🏠

LaBRI, Université de Bordeaux, France

Abstract

Property testing is concerned with the design of algorithms making sublinear number of queries to distinguish whether the input satisfies a given property or is far from having this property. A seminal paper of Alon, Krivelevich, Newman, and Szegedy in 2001 introduced property testing of formal languages: the goal is to determine whether an input word belongs to a given language, or is far from any word in that language. They constructed the first property testing algorithm for the class of all regular languages. This opened a line of work with improved complexity results and applications to streaming algorithms. In this work, we show a trichotomy result: the class of regular languages can be divided into three classes, each associated with a query complexity. Our analysis yields effective characterizations for all three classes using so-called minimal blocking sequences, reasoning directly and combinatorially on automata.

2012 ACM Subject Classification Theory of computation → Regular languages

Keywords and phrases property testing, regular languages

1 Introduction

Property testing was defined by Goldreich, Goldwasser, and Ron [14] in 1998: it is the study of very fast randomized approximate decision procedures on huge objects, where very fast typically means sublinear; the algorithm does not even scan the whole input. A very active branch of property testing focuses on graph properties, for instance one can test whether a given graph appears as a subgraph [3] or as an induced subgraph [4], and more generally every monotone graph property can be tested with one-sided error [6]. Other families of objects heavily studied under this algorithmic paradigm include probabilistic distributions [19, 9] combined with privacy constraints [2], numerical functions [8, 20], and programs [11, 10]. We refer to the book of Goldreich [13] for an overview of the field of property testing.

In this paper we continue the line of work initiated by Alon, Krivelevich, Newman, and Szegedy [5] which studies property testing of formal languages: given a language L (a set of finite words), the goal is to determine whether an input word u belongs to the language or is ε -far from it, where ε is the precision parameter. We assume random access to the input word: a query specifies a position in the word and asks for the letter in this position. To measure the distance of a word to a language we assume a metric over words; two natural choices include the Hamming distance (the number of positions at which two words differ) or the edit distance (the number of edits to transform one word into the other one). The seminal paper [5] showed a surprising result: all regular languages (meaning, languages recognised by deterministic finite automata) are testable with $\mathcal{O}(\log^3(\varepsilon^{-1})/\varepsilon)$ queries, where the $\mathcal{O}(\cdot)$ notation hides constants that depend on the language, but, crucially, not on the length of the input word.

A series of papers built upon this work, improving the query complexity (i.e. the number of queries). The original paper [5] identified the class of *trivial* regular languages, those

for which the answer is always *yes* or always *no* for large enough n , and showed that testing membership in a non-trivial regular language requires $\Omega(1/\varepsilon)$ queries. Building upon their work, Magniez and de Rougemont [18] extended their result by giving a tester using $\mathcal{O}(\log^2(\varepsilon^{-1})/\varepsilon)$ queries for the edit distance with moves, and François et al. [12] gave a tester using $\mathcal{O}(1/\varepsilon^2)$ queries for the case of the weighted edit distance. More recently, Bathie and Starikovskaya [7] gave a tester for the edit distance using $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries, and showed that there exists a *hard* regular language that cannot be tested with asymptotically fewer queries. However, there exist *easy* regular languages that can be tested with $\mathcal{O}(1/\varepsilon)$ queries. These results raise the following questions:

1. are there regular languages with a query complexity different from asymptotically 0, $\Theta(1/\varepsilon)$ and $\Theta(\log(\varepsilon^{-1})/\varepsilon)$?
2. is there a combinatorial and effective characterization of the languages in each class?

In this work, we answer both questions almost completely: we show a trichotomy theorem that classifies all regular languages in one of the three classes: *trivial* (asymptotically 0 queries¹), *easy* ($\Theta(1/\varepsilon)$ queries), and *hard* ($\Omega(\log(\varepsilon^{-1})/\varepsilon)$ queries). In the case of languages recognised by strongly connected NFAs, we even provide a matching upper bound of $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ for all languages. Our characterization of the three classes relies on the combinatorial notion of minimal blocking factors (and sequences).

We can therefore ask the meta-question: can we determine whether a given regular language is trivial, easy, or hard? Answering this question has practical motivations: determining to which class the language belongs enables choosing the appropriate most efficient property testing algorithm. We show that the meta-question is complete for the complexity class PSPACE (Turing machines working in polynomial space).

2 Overview of the paper

In this overview we assume familiarity with classical notions; all definitions can be found in Section 3. Let us start with the notion of a property tester for a language L : the goal is to determine whether an input word u belongs to the language L , or whether it is ε -far from it. We say that u of length n is ε -far from L with respect to a metric d over words if all words $v \in L$ satisfy $d(u, v) \geq \varepsilon n$, written $d(u, L) \geq \varepsilon n$. Throughout this work and unless explicitly stated otherwise, we will consider the case where d is the Hamming distance, defined for two words u and v as the number of positions at which they differ if they have the same length, and as $+\infty$ otherwise. In that case, $d(u, L) \geq \varepsilon n$ means that one cannot change a proportion ε of the letters in u to obtain a word in L . We assume random access to the input word: a query specifies a position in the word and asks for the letter in this position.

► **Definition 2.1.** *A property tester for the language L and precision ε is a randomized algorithm T that, for any input u of length n , given random access to u , satisfies the following properties:*

$$\begin{aligned} & \text{if } u \in L, \text{ then } T(u) = 1, & (1) \\ & \text{if } u \text{ is } \varepsilon\text{-far from } L, \text{ then } \mathbb{P}(T(u) = 0) \geq 2/3. & (2) \end{aligned}$$

¹ By *asymptotically 0 queries*, we mean that for every small enough $\varepsilon > 0$, for large enough n , the answer is either *yes* for all words of length n or *no* for all, and only depends on n , thus the algorithm does not need to query the input word.

85 *The query complexity of T is a function of n and ε that counts the maximum number of letters*
 86 *of the input that T reads over all inputs of length n and over all possible random choices. In*
 87 *this paper we are interested in property testers whose query complexity is independent of n ,*
 88 *and only depends on ε .*

89 More precisely, the definition of property testers given above is called “property testers
 90 with perfect completeness”: they always accept positive instances. Because they are based
 91 on the notion of blocking factors that we will discuss below, all known testers for regular
 92 languages [5, 18, 12, 7] have perfect completeness.

93 We say that a tester is *non-adaptive* if the index of a query does not depend on the result
 94 of previous queries. Alternatively, a non-adaptive tester can be understood as an algorithm
 95 that first sends the index of all of its queries, receives the result of all the queries, and then
 96 returns its output.

97 **Infinite languages**

98 Let us make a trivial observation: if L is finite, meaning that it contains finitely many words,
 99 then it is *trivial*. Let N denote the maximum length of a word in L : as L is finite, N is also
 100 finite. We can test L as follows: if the input has length less than N , query all of it and check
 101 whether it is one of the finitely many words of L and answer accordingly. Otherwise, if the
 102 length of the input is greater than N , answer *no*. This tester makes no queries for $n > N$,
 103 hence L is *trivial*. For this reason, we only consider infinite languages L ; this will make some
 104 technical statements nicer.

105 **Easy languages**

106 Let us consider the language $L_1 = a^*$ consisting of words containing only a 's, over the
 107 alphabet $\{a, b\}$. For a word $u \in \{a, b\}^*$, the distance $d(u, L)$ is the number of b 's in u . Here is
 108 a very simple property tester for L_1 : given a word of length n , sample $1/\varepsilon$ letters at random
 109 and answer no if we find a b , and yes otherwise. If $u \in L_1$, it contains no b to the algorithm
 110 returns yes, and if u is ε -far from L_1 , then each sample has probability at least ε to be a b ,
 111 and thus we will find a b with constant probability. One can easily show that $1/\varepsilon$ is a lower
 112 bound on the number of samples to get a property tester for L_1 ; we say that L_1 is *easy*:

113 ► **Definition 2.2.** *We say that L is **easy** if for small enough $\varepsilon > 0$, the optimal query*
 114 *complexity for a property tester for L is $\Theta(1/\varepsilon)$.*

115 **Blocking factors**

116 Extrapolating from the example L_1 , let us introduce the notion of *blocking factors* (also
 117 known as killing words [17]): a word v is a blocking factor for L if it cannot appear as a
 118 factor of a word in L . For instance, b is a blocking factor for L_1 . Note that bb and bbb are
 119 also blocking factors, but b is a minimal blocking factor (there are no strict factors of b
 120 that are blocking factors). Blocking factors were introduced in the original work giving a
 121 property tester for all regular languages [5]. A key insight of our work is to focus on *minimal*
 122 *blocking factors*. One important although simple property we will use is that if L is a regular
 123 language, then the set of minimal blocking factors of L is also a regular language.

124 It turns out that all property testers will be based on extensions of this very simple idea:
 125 we sample a number of positions in the word looking for blocking factors and answer no if
 126 we find a blocking factor, and yes otherwise. To be more precise, the analysis above for L_1
 127 rests on the following property:

4 The Trichotomy of Regular Property Testing

128 If u is ε -far from L_1 , then it contains at least εn disjoint minimal blocking factors.

129 We will show later that this property can be extended to all regular languages.

130 Trivial languages

131 At this point we can revisit the class of trivial languages identified in [5]:

132 ► **Definition 2.3.** *We say that L is **trivial** if for all small enough $\varepsilon > 0$, there exists a*
133 *property tester for L that makes 0 queries for all large enough n .*

134 An example of a trivial language is L_2 consisting of words containing at least one a over the
135 alphabet $\{a, b\}$. For any word u , replacing any letter by a yields a word in L_2 , so $d(u, L_2) \leq 1$.
136 A trivial property tester for L_2 simply answers yes all the time. One of our contributions in
137 this work is a characterization of the trivial languages identified by Alon et al. [5].

138 ► **Lemma 2.4.** *A regular language L is trivial if and only if it has no (minimal) blocking*
139 *factors.*

140 Period and positional words

141 Let us now consider the language $L_3 = (ab)^*$ consisting of words of the form $ab \cdot ab \cdot ab \cdots ab$
142 over the alphabet $\{a, b\}$. Generalizing the ideas used for the analysis of L_1 , a very simple
143 property tester for L_3 goes as follows: given a word of length n , sample $1/\varepsilon$ pairs of letters
144 from a random *even* position and answer no if we find anything else than ab , and yes otherwise.
145 There are two new difficulties: we need to consider factors of length 2, and we want them
146 to start at even positions. The arguments above are naturally extended to prove that this
147 property tester has query complexity $\mathcal{O}(1/\varepsilon)$, and that this is asymptotically tight.

148 What this example shows is that instead of consider words we will need to consider
149 *positional words*, which additionally encode information about the position. In the case
150 of L_3 , we need to distinguish between even and odd positions, so the word $abab$ is better
151 represented as $(0, a)(1, b)(0, a)(1, b)$, where the first index denotes the parity of the position.
152 More generally, we can associate to each regular language a period, and work with positional
153 words encoding the position modulo this period. The notion of blocking factors is naturally
154 extended to positional words, for instance $(0, a)(1, a)$ is a blocking factor, but $(1, b)(0, a)$ is
155 not.

156 Almost characterizing easy languages

157 Generalizing the ideas presented above, one can prove the following lemma:

158 ► **Lemma 2.5.** *Let L be a regular language. If there are finitely many minimal blocking*
159 *factors for L , then L is easy.*

160 Indeed, in that case there is an upper bound ℓ on the length of minimal blocking factors,
161 which depends only on L . We construct a property tester as follows: we sample $1/\varepsilon$ factors
162 of length ℓ and answer no if we find a blocking factor, and yes otherwise. One can prove that
163 this yields a property tester with query complexity $\mathcal{O}(1/\varepsilon)$. Unfortunately, this is not quite
164 a characterization: the converse implication does not hold, let us explain why using another
165 example.

166 Blocking sequences

167 Let us now consider the language L_4 consisting of words such that there are no c after a
 168 b over the alphabet $\{a, b, c\}$. The minimal blocking factors are of the form $ba^n c$ for $n \geq 0$,
 169 so there are infinitely many, the above argument above does not apply to this language.
 170 However, L_4 is *easy*: let us construct a property tester. We sample $1/\varepsilon$ letters at random
 171 and answer no if the sample contains a c after a b , and yes otherwise. To prove that this
 172 yields a property tester, we rely on the following property:

173 If u is ε -far from L_4 , then it can be decomposed $u = u_1 u_2$ where u_1 contains at least
 174 εn letters b and u_2 contains at least εn letters c .

175 What this example shows is that blocking factors are not enough: we need to consider
 176 sequences of factors, yielding the notion of blocking sequences. Intuitively, a blocking
 177 sequence for L is a sequence of (positional) words such that if the sequence appears as factors
 178 of some word u then $u \notin L$. For L_4 , the minimal blocking sequence is (b, c) .

179 Getting back to the almost characterization of *easy* languages sketched above, we will
 180 prove that L is *easy* if and only if there are finitely many minimal blocking sequences for L .
 181 The structure of the proof follows the original paper [5], considering first the case where L is
 182 recognised by a strongly connected automaton, and then extending it to the general case.
 183 Along the way, we will show that if L is recognised by a strongly connected automaton,
 184 then the characterization above holds: L is *easy* if and only there are finitely many minimal
 185 blocking factors for L . Introducing blocking sequences is necessary to deal with automata
 186 with more than one strongly connected component.

187 Hard languages

188 The remaining case is languages L which have infinitely many minimal blocking sequences.
 189 Let us illustrate this case on an example. We start from the parity language P consisting of
 190 words such that there is an even number of b 's, over the alphabet $\{a, b\}$. If the goal would
 191 be to distinguish between $u \in P$ and $u \notin P$, any property tester would require scanning
 192 the whole input word. However, relaxing with the Hamming distance makes the question
 193 different: every word is at distance at most 1 from P by swapping at most one letter, so the
 194 language is *trivial*. Now, consider L_5 consisting of words such that inbetween each letter
 195 \sharp , there is an even number of b 's, over the alphabet $\{a, b, \sharp\}$. Intuitively, L_5 encodes an
 196 arbitrary number of parity instances. Bathie and Starikovskaya [7] proved a lower bound of
 197 $\Omega(\log(\varepsilon^{-1})/\varepsilon)$ on the query complexity of (non-adaptive) property testers for L_5 , matching
 198 the property testing algorithm they constructed for all regular languages.

199 ► **Definition 2.6.** *We say that L is **hard** if for all small enough $\varepsilon > 0$, the optimal query*
 200 *complexity of a tester for L is $\Omega(\log(\varepsilon^{-1})/\varepsilon)$.*

201 Inspecting the minimal blocking sequences for L_5 , we find infinitely many: this is no
 202 coincidence, we will extend Bathie and Starikovskaya's proof to show that any regular
 203 language with infinitely many minimal blocking sequences is *hard*.

204 The trichotomy theorem

205 Our main technical result is stated below. Recall that the case of finite languages is *easy*, so
 206 we focus on infinite languages.

207 ► **Theorem 2.7.** *Let L be an infinite regular language, let us write $MBS(L)$ for the set of*
 208 *minimal blocking sequences of L .*

- 209 ■ *L is trivial if and only if $MBS(L)$ is empty;*
- 210 ■ *L is easy if and only if $MBS(L)$ is finite and nonempty;*
- 211 ■ *L is hard if and only if $MBS(L)$ is infinite.*

212 This trichotomy theorem closes a line of work on improving query complexity for property
 213 testers and identifying easier subclasses of regular languages. As mentioned above, the proof
 214 considers first the case where L is recognised by a strongly connected automaton, and then
 215 extends the results to the general case (following [5]):

- 216 ■ For the strongly connected case, we extend the ideas from [5] using the framework
 217 of minimal blocking factors, thereby simplifying the exposition, and obtain optimal
 218 property testers for trivial and easy languages, together with matching lower bounds.
 219 Our novel contributions here concern hard languages. First, we construct a property
 220 tester with query complexity $\Theta(\log(\varepsilon^{-1})/\varepsilon)$ for all regular languages recognised by a
 221 strongly connected automaton. This is an improvement over the similar result of Bathie
 222 and Starikovskaya [7], which works under the edit distance, while ours is designed for the
 223 Hamming distance. As the edit distance never exceeds the Hamming distance, the set
 224 of words that are ε -far with respect to the former is contained in the set of words ε -far
 225 for the latter. Therefore, an ε -tester for the Hamming distance is also an ε -tester for
 226 the edit distance, and our result supersedes and generalizes theirs in the case of strongly
 227 connected automata. Second, we prove a matching lower bound, again inspired by but
 228 strongly generalizing a result from Bathie and Starikovskaya [7], which was for a single
 229 language (L_5 discussed above), to all regular languages with infinitely many minimal
 230 blocking factors. We use Yao’s minmax principle [21]: this involves constructing a hard
 231 distribution over inputs, and showing that any deterministic property testing algorithms
 232 cannot distinguish between positive and negative instances against this distribution.
- 233 ■ The general case follows a similar outline and builds upon the results in the strongly
 234 connected case. The notion of (minimal) blocking sequences enables smooth yet technical
 235 generalization of most of the results based on blocking factors for strongly connected
 236 automata. The main difficulty is the case of hard languages, and more precisely the lower
 237 bound. The complication here is that it is not enough to consider strongly connected
 238 components in isolation: there exists finite automata which contain a strongly connected
 239 component that induces a hard language, yet the language of the whole automaton is
 240 easy. The case where both are hard also occurs. Our proof defines a notion of “portals”
 241 which allows us to extract “crucial” connected components and show that hardness of
 242 these components imply hardness of the whole language.

243 The meta-question

244 Once the trichotomy theorem is established, the natural pending question is whether it
 245 can be made effective: the meta-question is then, given a regular language L , to determine
 246 whether it is trivial, easy, or hard. One could use this procedure to run the appropriate
 247 algorithm. On the positive front, our characterizations are indeed effective, in particular
 248 since for a regular language L , the set of minimal blocking sequences of L is another regular
 249 language, which can be effectively computed. However, we prove that the complexity of
 250 checking whether this set is empty, finite, or infinite (corresponding to the trichotomy), is
 251 PSPACE-complete.

252 **Outline**

253 The missing definitions are given in Section 3. We treat strongly connected automata
 254 in Section 4, and the general case in Section 5. The complexity of the meta-question is
 255 established in Section 6.

256 **3 Preliminaries**257 **Words and languages.**

258 In this work, we consider a finite set Σ , the *alphabet*, whose elements are called *letters*. Words
 259 over Σ are finite sequences of letters, and Σ^* (resp. Σ^+) denotes the set of all words (resp.
 260 nonempty words) over Σ . A subset of Σ^* is called a language over Σ . The length of a word
 261 $u \in \Sigma^*$, denoted $|u|$, is the number of letters that it contains, and for $i \in [0, |u| - 1]$, we
 262 use $u[i]$ to denote the i -th letter of u . Given two words $u, v \in \Sigma^*$, the *concatenation* $u \cdot v$
 263 (more concisely denoted uv) of u and v is the word composed of the letters of u followed by
 264 the letters of v . This operation is associative, hence (Σ^*, \cdot) is a monoid. Its unique neutral
 265 element, the empty word, is denoted γ .

266 Given a word $u \in \Sigma^*$ and two integers $0 \leq i, j \leq |u| - 1$, define $u[i..j]$ as the word
 267 $u[i]u[i+1] \dots u[j]$ if $i \leq j$ and γ otherwise. Further, we let $u[i..j)$ denote the word $u[i..j-1]$.
 268 A word w is a *factor* (resp. *prefix*, *suffix*) of u if there exist indices i, j such that $w = u[i..j]$
 269 (resp. with $i = 0, j = |u| - 1$). We use $w \preceq u$ to denote “ w is a factor of u ”. Furthermore, if
 270 w is a factor of u and $w \neq u$, we say that w is a *proper factor* of u .

271 **Finite automata.**

272 **► Definition 3.1** (Nondeterministic Finite automaton). A nondeterministic finite automaton
 273 (NFA) \mathcal{A} is a transition system defined by a tuple $(Q, \Sigma, \delta, q_0, F)$, with Q a finite set of states,
 274 Σ a finite alphabet, $\delta : Q \times \Sigma \rightarrow 2^Q$ the transition function, q_0 is the initial state and F is
 275 the set of final states.

276 We say that \mathcal{A} is deterministic (resp. complete) if $|\delta(q, a)| \leq 1$ (resp. ≥ 1) for all $q \in Q, a \in \Sigma$.
 277 We say that there is a transition from $p \in Q$ to $q \in Q$ labeled by $w \in \Sigma^*$, denoted
 278 $p \xrightarrow{w} q$, if there exists states $p_0 = p, p_1, \dots, p_{|w|} = q$ such that for every $i = 0, \dots, |w| - 1$,
 279 $p_{i+1} \in \delta(p_i, w[i])$. In this case, we say that q is reachable from p and that p is co-reachable
 280 from q . The *language recognized by \mathcal{A}* , denoted $L(\mathcal{A})$, is the set of words that label a transition
 281 from the initial state to a final state, i.e.

$$282 \quad L(\mathcal{A}) = \{w \in \Sigma^* \mid \exists q_f \in F : q_0 \xrightarrow{w} q_f\}.$$

283 We say that an NFA is *trim* if every state is reachable from the initial state and co-reachable
 284 from some final state. An NFA \mathcal{A} can always be converted into a trim NFA that recognizes
 285 the same language by removing the states of \mathcal{A} that are either not reachable from q_0 or not
 286 co-reachable from any final state.

287 **Property testing.**

288 **► Definition 3.2.** Let L be a language, let u be a word of length n , let $\varepsilon > 0$ be a precision
 289 parameter and let $d : \Sigma^* \times \Sigma^* \rightarrow \mathbb{N} \cup \{+\infty\}$ be a metric. We say that the word u is ε -far
 290 from L w.r.t. d if $d(u, L) \geq \varepsilon n$, where

$$291 \quad d(u, L) := \inf_{v \in L} d(u, v).$$

292 Throughout this work and unless explicitly stated otherwise, we will consider the case where
 293 d is the Hamming distance, defined for two words u and v as the number of positions at
 294 which they differ if they have the same length, and as $+\infty$ otherwise.

295 Graphs and periodicity.

296 We now recall tools introduced by Alon et al [5] to deal with periodicity in finite automata.

297 The period $\lambda = \lambda(G)$ of a graph G is the greatest common divisor of the length of the
 298 cycles in G . If G is acyclic, we set $\lambda(G) = \infty$. Following the work of Alon et al [5], we will
 299 use the following property of directed graphs.

300 ► **Fact 3.3** (From [5, Lemma 2.3]). *Let $G = (V, E)$ be a nonempty, strongly connected directed*
 301 *graph with finite period $\lambda = \lambda(G)$. Then there exists a partition $V = Q_0 \sqcup \dots \sqcup Q_{\lambda-1}$ and a*
 302 *reachability constant $\rho = \rho(G)$ that does not exceed $3|V|^2$ such that:*

303 1. *For every $0 \leq i, j \leq \lambda - 1$ and for every $u \in Q_i, v \in Q_j$, the length of any directed path*
 304 *from u to v in G is equal to $(j - i) \pmod{\lambda}$.*

305 2. *For every $0 \leq i, j \leq \lambda - 1$, for every $u \in Q_i, v \in Q_j$ and for every integer $r \geq \rho$, if*
 306 *$r = (j - i) \pmod{\lambda}$, then there exists a directed path from u to v in G of length r .*

307 The sets $(Q_i : i = 0, \dots, \lambda - 1)$ are the *periodicity classes* of G . In what follows, we will
 308 slightly abuse notation and use Q_i for arbitrary non-negative integers i to mean $Q_{i \pmod{\lambda}}$
 309 when $i \geq \lambda$.

310 Given a finite automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$, we can naturally obtain the underlying
 311 directed graph by removing the labels from the transitions: it is the graph $G = (Q, E)$ where
 312 $E = \{(p, q) \in Q^2 \mid \exists a \in \Sigma : q \in \delta(p, a)\}$. In what follows, we naturally extend the notions of
 313 period², reachability constant and periodicity classes to finite automata through this graph G .
 314 The numbering of the periodicity classes is defined up to a shift mod λ : we say that Q_0
 315 is the class that contains the initial state q_0 . Similarly, we say that a finite automaton is
 316 strongly connected if the underlying graph is strongly connected.

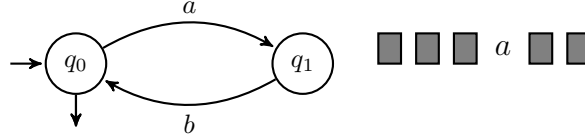
317 A *strongly connected component* S of an automaton \mathcal{A} is a maximal subset of states such
 318 that every state of S is reachable from every other one. Its *periodicity* is the periodicity of
 319 the subgraph induced by \mathcal{A} over S .

320 Positional words and positional languages.

321 To motivate the following definitions, let us recall the example discussed in the overview.
 322 Consider a simple deterministic automaton with two states for the language $(ab)^*$: there is a
 323 transition labeled by a starting from one state but not from the other. The parity of the
 324 position of the factor ab in a word carries an important piece of information: if the position
 325 is odd, then we know that the word containing the factor is not in $(ab)^*$. Furthermore, while
 326 b appears at position 1 in ab , if this ab appears at an odd position in u then b appears at an
 327 even position in u . This leads to the definition of *positional words*.

328 ► **Definition 3.4** (Positional words). *Let p be a positive integer. A p -positional word is a*
 329 *word over the alphabet $\mathbb{Z}/p\mathbb{Z} \times \Sigma$ of the form $(n \pmod{p}, a_0)((n+1) \pmod{p}, a_1) \dots ((n+\ell)$
 330 $\pmod{p}, a_\ell)$. If $u = a_0 \dots a_\ell$, we write $(n : u)$ to denote this word.*

² Note that in this context, an *aperiodic automaton* means an automaton with an aperiodic underlying graph, which is not the same thing as a counter-free automaton, which are sometimes called aperiodic automata.



■ **Figure 1** An automaton recognising the language $(ab)^*$. A witness that a word is not in this language is an a on an odd position or a b on an even position.

331 With this definition, starting with the 2-positional word $(0 : u)$, the factor ab at an odd
 332 position in u is $(1, a)(0, b)$, and the positional factor corresponding to the b is $(0, b)$. In this
 333 case, even when taking factors of factors of u , we still retain the (congruence classes of the)
 334 indices in the original word.

335 Any strongly connected finite automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ can naturally be extended
 336 into an automaton over $\lambda(\mathcal{A})$ -positional words as follows. Let $Q_0, \dots, Q_{\lambda-1}$ be the partition
 337 of the states of \mathcal{A} given by Fact 3.3, where $\lambda = \lambda(\mathcal{A})$ is the periodicity of \mathcal{A} . The *positional*
 338 *extension* of \mathcal{A} is the finite automaton $\widehat{\mathcal{A}}$ defined by:

$$339 \quad \widehat{\mathcal{A}} = (Q, \mathbb{Z}/\lambda\mathbb{Z} \times \Sigma, \delta', q_0, F) \quad \text{where } \delta'(q, (i, a)) = \begin{cases} \delta(q, a) & \text{if } q \in Q_i, \\ \emptyset & \text{otherwise.} \end{cases}$$

340 By fact Fact 3.3, any transition from a state of Q_i goes to a state in Q_{i+1} , hence $\widehat{\mathcal{A}}$ recognized
 341 well-formed λ -positional words. We call the language recognized by $\widehat{\mathcal{A}}$ the *positional language*
 342 *of* \mathcal{A} , and denote it $\mathcal{TL}(\mathcal{A})$. This definition is motivated by the following property:

343 ► **Property 3.5.** *For any word $u \in \Sigma^*$, we have $u \in \mathcal{L}(\mathcal{A})$ if and only if $(0 : u) \in \mathcal{TL}(\mathcal{A})$.*

344 For the reasons that we exposed earlier, positional words make it easier to manipulate
 345 factors with positional information, hence we phrase our property testing results in terms of
 346 positional languages. Notice that a property tester for $\mathcal{TL}(\mathcal{A})$ immediately gives a property
 347 tester for $\mathcal{L}(\mathcal{A})$, as one can simulate queries to $(0 : u)$ with queries to u by simply pairing
 348 the index of the query modulo $\lambda(\mathcal{A})$ with its result.

349 **4 Strongly Connected NFAs**

350 We first study the case of strongly connected NFAs, which are NFAs such that for any pair of
 351 states $p, q \in Q$, there exists a word w such that $p \xrightarrow{w} q$. We show that the query complexity
 352 of the language of such an NFA \mathcal{A} can be characterized by the cardinality of the set of
 353 *minimal blocking factors* of \mathcal{A} , which are factor-minimal $\lambda(\mathcal{A})$ -positional words that witness
 354 the fact that a word does not belong to $\mathcal{TL}(\mathcal{A})$. In this section, we consider a fixed NFA
 355 \mathcal{A} and simply use “positional words” to refer to λ -positional words, where $\lambda = \lambda(\mathcal{A})$ is the
 356 period of \mathcal{A} .

357 ► **Definition 4.1** (Blocking factors). *Let \mathcal{A} be a strongly connected NFA. A positional word τ*
 358 *is a blocking factor of \mathcal{A} if for any other positional word μ we have $\tau \preceq \mu \Rightarrow \mu \notin \mathcal{TL}(\mathcal{A})$.*

359 *Further, we say that τ is a minimal blocking factor of \mathcal{A} if no proper factor of τ is blocking*
 360 *a blocking factor of \mathcal{A} . We use $MBF(\mathcal{A})$ to denote the set of all minimal blocking words of \mathcal{A} .*

361 Intuitively and in terms of automata, $(i : u)$ is blocking for \mathcal{A} if it does not label any transition
 362 in \mathcal{A} labeled by u starting from a state of Q_i . This property is formally established later in
 363 Lemma 4.6.

364 The main result of this section is the following:

- 365 ► **Theorem 4.2.** *Let L be an infinite language recognised by a strongly connected NFA \mathcal{A} .*
 366 **1.** *L is hard if and only if $MBF(\mathcal{A})$ is infinite.*
 367 **2.** *L is easy if and only if $MBF(\mathcal{A})$ is finite and nonempty.*
 368 **3.** *L is trivial if and only if $MBF(\mathcal{A})$ is empty.*

369 Our approach

370 The definition of blocking factors gives a simple but powerful framework to design property
 371 testers for $L(\mathcal{A})$: using random sampling, attempt to find a blocking factor in $(0 : u)$; if one
 372 is found, reject u , otherwise accept u . If $u \in L(\mathcal{A})$, then $(0 : u)$ does not contain a blocking
 373 factor, and we always accept u . The other case is where insight is required: one needs to find
 374 a sampling strategy that had a good probability of finding a blocking factor in $(0 : u)$ for
 375 any u ε -far from $L(\mathcal{A})$. A central tool for building such a sampling strategy is Lemma 4.8,
 376 which shows that any word ε -far from $L(\mathcal{A})$ contains many blocking factors. This approach,
 377 introduced by Alon et al. [5], is used by all existing property testing algorithms for regular
 378 languages.

379 This section is organized as follows. First, in Section 4.1, we give a few tools that help
 380 us deal with positional words and blocking factors in strongly connected NFA. Next, in
 381 Section 4.3 we tackle the case of trivial and easy languages (i.e. items Item 2 and Item 3
 382 of Theorem 4.2). In Section 4.4, we prove that there exists an ε tester using $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$
 383 queries for any language $\mathcal{L}(\mathcal{A})$. Finally, in Section 4.5, we show that any language not in the
 384 “easy” class requires $\Omega(\log(\varepsilon^{-1})/\varepsilon)$ queries, thereby proving that there is no intermediate
 385 query complexity class between easy and hard, and completing the trichotomy.

386 4.1 Positional words, blocking factors and strongly connected NFAs

387 Before diving into the proof of Theorem 4.2, we establish a few properties of positional
 388 words that will help us ensuring that we are creating well-formed positional words. In
 389 Section 4.2, we highlight the connection between property testing and blocking factors in
 390 strongly connected NFAs.

391 We start with the following facts, which are consequences of Fact 3.3.

- 392 ► **Fact 4.3.** *Let n be a nonnegative integer, let w be a word of length n . If for some states*
 393 *$p \in Q_i, q \in Q_j$ of \mathcal{A} we have $p \xrightarrow{w} q$, then the indices i, j satisfy the equation*

$$394 \quad j - i = |w| \pmod{\lambda}$$

- 395 ► **Fact 4.4.** *Let $\tau = (i : u)$ and $\mu = (j : v)$ be positional words. If $\tau \preceq \mu$, then there exists*
 396 *positional words η, η' with $|\eta| = i - j \pmod{\lambda}$ such that $\mu = \eta\tau\eta'$. In particular, this implies*
 397 *that there exists words w, w' with $|w| = i - j \pmod{\lambda}$ such that $v = ww'$.*

398 The next property shows that chaining positional words in the automaton $\widehat{\mathcal{A}}$ results
 399 in well-formed positional words, in the sense that its letters are numbered by consecutive
 400 numbers modulo λ .

- 401 ► **Property 4.5.** *Let p, q, r be states of $\widehat{\mathcal{A}}$ and let τ, μ be two positional words such that $p \xrightarrow{\tau} q$*
 402 *and $q \xrightarrow{\mu} r$. Then $\tau\mu$ is a well-formed positional word, i.e. there exists a word w and an*
 403 *integer $i \in \mathbb{Z}/\lambda\mathbb{Z}$ such that $\tau\mu = (i : w)$.*

404 **Proof.** Let i, j be the respective indices of the periodicity classes of p and q , i.e. we have
 405 $p \in Q_i$ and $q \in Q_j$. Then there exist words u, v such that $\tau = (i : u)$ and $\mu = (j : v)$.
 406 Furthermore, by Fact 3.3, the length of any path from p to q is equal to $j - i \pmod{\lambda}$, hence
 407 the last letter of τ is $(j - 1, a)$ for some $a \in \Sigma$ and the words can be chained correctly, i.e.
 408 $\tau\mu = (i : uv)$. \blacktriangleleft

409 These properties allows us to formalize the intuition we gave earlier about blocking factors.

410 **► Lemma 4.6.** *A positional word $\tau = (i : u)$ is a blocking factor for \mathcal{A} iff for every states*
 411 *$p \in Q_i, q \in Q$, we have $p \not\stackrel{u}{\rightarrow} q$.*

412 **Proof.** We first show that if there exists states $p \in Q_i, q \in Q$ such that $p \stackrel{u}{\rightarrow} q$, then τ is not
 413 blocking, i.e. there exists $\mu \in \mathcal{TL}(\mathcal{A})$ such that $\tau \preceq \mu$. As \mathcal{A} is strongly connected, there
 414 exist positional words η, η' such that $q_0 \stackrel{\eta}{\rightarrow} p$ and $q \stackrel{\eta'}{\rightarrow} q_f$ for some $q_f \in F$. By Property 4.5,
 415 the positional word $\mu = \eta\tau\eta'$ is well formed. Furthermore, it labels a transition from q_0 to
 416 q_f , hence it is in $\mathcal{TL}(\mathcal{A})$, and τ is not blocking.

417 For the converse, assume that τ is non-blocking: we show that there exists two states
 418 $p \in Q_i, q \in Q$ such that $p \stackrel{u}{\rightarrow} q$. As τ is non-blocking, there exists a positional word $\mu = (0 : w)$
 419 such that $\tau \preceq \mu$ and there exists a final state r such that $q_0 \stackrel{\mu}{\rightarrow} r$, and equivalently, $q_0 \stackrel{w}{\rightarrow} r$. By
 420 Fact 4.4, since $\tau \preceq \mu$, there exists words v, v' such that $w = vv'$ and the length of v is equal
 421 to i modulo λ . In particular, the path $q_0 \stackrel{w}{\rightarrow} r$ can be decomposed into $q_0 \stackrel{v}{\rightarrow} p \stackrel{u}{\rightarrow} q \stackrel{w'}{\rightarrow} r$: in
 422 particular, we have $p \stackrel{u}{\rightarrow} q$. It only remains to show that p is in Q_i : this follows by Fact 4.3
 423 since $|v| = i \pmod{\lambda}$. \blacktriangleleft

424 Finally, the Hamming distance between u and $\mathcal{L}(\mathcal{A})$ is the same as the distance between
 425 $(0 : u)$ and $\mathcal{TL}(\mathcal{A})$.

426 \triangleright **Claim 4.7.** For any word $u \in \Sigma^*$, we have $d(u, \mathcal{L}(\mathcal{A})) = d((0 : u), \mathcal{TL}(\mathcal{A}))$.

427 **Proof.** The \leq part is straightforward. For the reverse inequality, it suffices to see that in
 428 any minimal substitution sequence from $(0 : u)$ to a positional word in $\mathcal{TL}(\mathcal{A})$, no operation
 429 changes only an index in an (index, letter) pair. \triangleleft

430 This allows us to interchangeably use the statements “ u is ε -far from $\mathcal{L}(\mathcal{A})$ ” and “ $(0 : u)$ is
 431 ε -far from $\mathcal{TL}(\mathcal{A})$ ”.

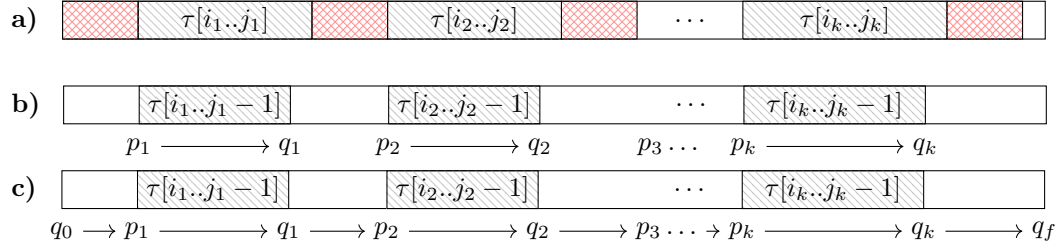
432 4.2 Strongly connected NFAs and blocking factors

433 Alon et al. [5, Lemma 2.6] first noticed that if a word u is ε -far from $L(\mathcal{A})$, then it contains
 434 $\Omega(\varepsilon n)$ short factors that witness the fact that u is not in $L(\mathcal{A})$. We start by translating
 435 the lemma of Alon et al. on “short witnesses” to the framework of blocking factors. More
 436 precisely, we show that if u is ε -far from $L(\mathcal{A})$, then $(0 : u)$ contains many disjoint blocking
 437 factors (Lemma 4.8).

438 **► Lemma 4.8.** *Let $\varepsilon > 0$, let u be a word of length $n \geq 6m^2/\varepsilon$ and assume that $L(\mathcal{A})$*
 439 *contains at least one word of length n . If $\tau = (0 : u)$ is ε -far from $\mathcal{TL}(\mathcal{A})$, then τ contains at*
 440 *least $\varepsilon n/(6m^2)$ disjoint blocking factors.*

441 **Proof.** We build a set \mathcal{P} of disjoint blocking factors of τ as follows: we process u from left to
 442 right, starting at index $i_1 = \rho$. Next, at iteration t , set j_t to be the smallest integer greater
 443 than or equal to i_t and smaller than $n - \rho$ such that $\tau[i_t..j_t]$ is a blocking factor. If there is
 444 no such integer, we stop the process. Otherwise, we add $\tau[i_t..j_t + \rho - 1]$ to the set \mathcal{P} , and
 445 iterate starting from the index $i_{t+1} = j_t + \rho$.

446 Let k denote the size of \mathcal{P} . We will show that we can substitute at most $3(k+1)m^2$
 447 positions in τ to obtain a word in $\mathcal{TL}(\mathcal{A})$. (See Figure 2 for an illustration of this construction.)
 448 Using the assumption that τ is ε -far from $\mathcal{TL}(\mathcal{A})$ (which follows from Claim 4.7) will give us
 449 the desired bound on k .



■ **Figure 2 a)** The decomposition process returns k factors $\tau[i_1, j_t], \dots, \tau[i_k, j_k]$ (represented as diagonally hatched in gray regions), separated together and with the start of the text by padding regions of $\rho - 1$ letters (red crosshatched regions). **b)** After removing the last letter, each previously blocking factor now labels a transition between some pair of states p_t, q_t . **c)** We use the padding regions to bridge between consecutive factors as well as the start and end of the word.

450 For every t , we chose j_t to be minimal so that $\tau[i_t..j_t]$ is blocking, hence $\tau[i_t..j_t - 1]$ is not
 451 blocking, and therefore $\tau[i_t..j_t - 1]$ labels a run from some $p_t \in Q_{i_t}$ to a $q_t \in Q_{j_t}$. Therefore,
 452 using the strong connectivity of \mathcal{A} and Fact 3.3, we can edit the last ρ letters of the block
 453 $\tau[i_t..j_t + \rho - 1]$ to obtain a non-blocking factor that labels a transition from p_t to p_{t+1} . Using
 454 the ρ letters at the start and the end of the word, we add transitions from an initial state
 455 to p_1 and from q_k to a final state: the assumption that $\mathcal{L}(\mathcal{A})$ contains a word of length n
 456 ensures that Q_n contains a final state, hence this is always possible. The resulting word is in
 457 $\mathcal{TL}(\mathcal{A})$ and was obtained from τ using $(k+1)\rho \leq 3(k+1)m^2$ substitutions. As τ is ε -far
 458 from $\mathcal{TL}(\mathcal{A})$, we obtain the following bound on k :

$$\begin{aligned}
 459 \quad 3(k+1)m^2 \geq \varepsilon n &\Rightarrow k+1 \geq \frac{\varepsilon n}{3m^2} \\
 &\Rightarrow k \geq \frac{\varepsilon n}{3m^2} - 1 \\
 460 \quad & \\
 461 \quad &\Rightarrow k \geq \frac{\varepsilon n}{6m^2}
 \end{aligned}$$

462 The last implication uses the assumption that $n \geq 6m^2/\varepsilon$. ◀

463 Lemma 4.8 allows us to handle three cases of Theorem 4.2, namely we use it to construct
 464 a tester with $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries for any regular language, to construct a tester with
 465 $\mathcal{O}(1/\varepsilon)$ queries for regular languages with a finite number of blocking factors and to show
 466 the triviality of languages with no blocking factors.

467 4.3 The finite case

468 Using the framework described in the previous subsection, we show that when $\text{MBF}(\mathcal{A})$ is
 469 finite then $\mathcal{L}(\mathcal{A})$ can be tested with $\mathcal{O}(1/\varepsilon)$ queries, and furthermore if $\text{MBF}(\mathcal{A})$ is empty,
 470 then $\mathcal{L}(\mathcal{A})$ is trivial.

471 **Automata with no blocking factors.**

472 First, observe that if $\text{MBF}(\mathcal{A})$ is empty, there are no blocking factors, and no word can contain
 473 a blocking factor. Hence, the decomposition procedure used in the proof of Lemma 4.8
 474 terminates with $k = 0$, which shows that, if $\mathcal{L}(\mathcal{A}) \cap \Sigma^n$ is nonempty, then any word of length
 475 n is at distance at most $3m^2$ of $\mathcal{L}(\mathcal{A})$. Therefore, for any $\varepsilon > 0$ and for $n \geq 3m^2/\varepsilon$, no word
 476 of length n is ε -far from $\mathcal{L}(\mathcal{A})$, and the tester that always accepts without queries is correct.

477 **A tester for the finite-but-nonempty case.**

478 To design a property tester with $\mathcal{O}(1/\varepsilon)$ queries, recall that, from Lemma 4.8, if u is ε -far
 479 from $L(\mathcal{A})$, then $(0 : u)$ contains many disjoint blocking factors. We then extract from each
 480 of these blocking factor a minimum blocking factor: because $\text{MBF}(\mathcal{A})$ is finite, the length of
 481 each of these minimal factors is bounded by a constant C independent of u , hence a constant
 482 number of queries is enough to read one such factor. Finally, we show in Lemma 4.9 that
 483 sampling $\mathcal{O}(1/\varepsilon)$ factors is enough; the result follows.

484 ► **Lemma 4.9.** *Let \mathcal{A} be a trim strongly connected NFA. If $\text{MBF}(\mathcal{A})$ is finite, then the*
 485 *language $L = L(\mathcal{A})$ can be tested with $\mathcal{O}(1/\varepsilon)$ queries.*

486 **Proof.** If $\text{MBF}(\mathcal{A})$ is finite, then there exists a constant C such that every minimal blocking
 487 factor of \mathcal{A} has length at most C .

488 Let m denote the number of states of \mathcal{A} . Given a word u of length n , we first check the
 489 following:

- 490 ■ If $n < 6m^2/\varepsilon$, read all of u , run the automaton \mathcal{A} on u and accept if and only if \mathcal{A} accepts.
- 491 ■ If $L(\mathcal{A})$ does not contain words of length n , reject. This can be checked efficiently using
 492 a simple dynamic programming algorithm.

493 The above procedure uses at most $\mathcal{O}(1/\varepsilon)$ queries, and if both checks fail, then u satisfies
 494 the hypotheses of Lemma 4.8. We then use the following procedure:

- 495 ■ sample independently $K = 6m^2 \ln(3)/\varepsilon$ random factors of length C in $(0 : u)$. To sample
 496 a factor, choose an index i uniformly at random in $\{1, \dots, n\}$, and return $(i : u[i..i + C])$.
- 497 ■ rejects if at least one of these factors is blocking for \mathcal{A} .

498 We show that this algorithm is an ε -tester for L .

499 First, if $u \in L$, then no factor of $(0 : u)$ is blocking, and the algorithm accepts with
 500 probability 1.

501 Now, assume that u is ε -far from L . By Lemma 4.8, $(0 : u)$ contains at least $N = \varepsilon n/(6m^2)$
 502 disjoint blocking factors. Each of these blocking factors induces at least one minimal blocking
 503 factor, i.e. $(0 : u)$ contains at least N disjoint minimal blocking factors. Each of these
 504 factors has length at most C , therefore the probability that the sampling procedure returns
 505 a factor that contains one of them is at least $N/n = \varepsilon/(6m^2)$. By repeating independently
 506 $K = 6m^2 \ln(3)/\varepsilon$ times, the probability of *not* finding any of the blocking factors is at most
 507 $(1 - N/n)^K \leq e^{-KN/n} = e^{-\ln 3} = 1/3$, therefore the algorithm rejects u with probability at
 508 least $2/3$ and satisfies Definition 2.1.

509 This tester uses $6Cm^2/\varepsilon = \mathcal{O}(1/\varepsilon)$ queries. ◀

510 **Lower bound in the nonempty case.**

511 It remains to show that if $\text{MBF}(\mathcal{A})$ is nonempty, then testing $\mathcal{L}(\mathcal{A})$ requires $\Omega(1/\varepsilon)$ queries.
 512 Alon et al. [5] showed that “non-trivial” regular languages require $\Omega(1/\varepsilon)$ queries, using a
 513 notion of triviality defined differently from ours. They define non-trivial languages as follows:

514

515 **► Definition 4.10** ([5, Definition 3.1]). *A language L is non-trivial if there exists a constant*
 516 *$\varepsilon_0 > 0$, so that for infinitely many values of n the set $L \cap \Sigma^n$ is nonempty, and there exists a*
 517 *word $w \in \Sigma^n$ so that $d(w, L) \geq \varepsilon_0 n$.*

518 Their lower bound is the following:

519 **► Fact 4.11** ([5, Proposition 2]). *Let L be a non-trivial (in the sense of Alon et. al) regular*
 520 *language. Then for all sufficiently small $\varepsilon > 0$, any ε -tester for L requires $\Omega(1/\varepsilon)$ queries.*

521 To prove the lower bound in item 2) of Theorem 4.2, we show that if a language is
 522 non-trivial in our sense, i.e. $\text{MBF}(\mathcal{A})$ is nonempty, then it is non-trivial in the sense of Alon
 523 et al.: we then get our lower bound by applying theirs.

524 **► Lemma 4.12.** *Let \mathcal{A} be a strongly connected NFA such that $\text{MBF}(\mathcal{A})$ is nonempty and*
 525 *denote $L = \mathcal{L}(\mathcal{A})$. Then there exists a constant $\varepsilon_0 > 0$ such that for infinitely many values*
 526 *of n the set $L \cap \Sigma^n$ is nonempty and there exists a word $w \in \Sigma^n$ so that $d(w, L) \geq \varepsilon_0 n$.*

527 **Proof.** As \mathcal{A} is strongly connected, L is infinite, hence there are infinitely many integers n
 528 such that $L \cap \Sigma^n$ is nonempty. We show that there exists a constant ε_0 such that for large
 529 enough n such that $L \cap \Sigma^n$ is nonempty, there is a word of length n that is ε_0 -far from L .

530 Since $\text{MBF}(L)$ is nonempty, it contains at least one blocking factor, which is of the form
 531 $(i : u)$ for some $i \in \mathbb{Z}/\lambda\mathbb{Z}$. Let C denote the smallest multiple of λ greater than the length of
 532 u , let x denote an arbitrary word of length C with u as a prefix, and let $\varepsilon_0 = 1/(2C)$. We
 533 proceed to show that for any sufficiently large $n \geq 2(C + \lambda)$ such that $L \cap \Sigma^n$ is nonempty,
 534 there exists a word $w \in \Sigma^n$ such that $d(w, L) \geq \varepsilon_0 n$. We construct the word w by replacing
 535 a portion of v with disjoint copies of x , where x is an arbitrary word of length C that has
 536 u as a prefix. More precisely, we define w as $w = v[..i]x^k v[i + k \cdot C + 1..]$ where $k = \lceil \varepsilon_0 n \rceil$
 537 disjoint copies of x . This word has length n as $i + k \cdot C \leq \lambda + C \lceil \varepsilon_0 n \rceil \leq n$.

538 We now claim that $d(w, L) \geq k \geq \varepsilon_0 n$. First, notice that as C is a multiple of λ , all
 539 k copies of x (and therefore of u) in w start at position equal to i modulo λ . Therefore,
 540 any such occurrence of u induces an occurrence of $(i : u)$ in $(0 : w)$. Next, consider a word
 541 w' obtained by performing less than k substitutions on w . Some copy of u in w' will be
 542 untouched, hence $(i : u) \preceq (0 : w')$, and therefore $w' \notin L$. Overall, we have

$$543 \quad d(w, L) = d((0 : w), \mathcal{L}(\mathcal{A})) \geq k \geq \varepsilon_0 n.$$

544 We have shown that there exists ε_0 such that for infinitely many n , $L \cap \Sigma^n$ is nonempty
 545 and there exists a word $w \in \Sigma^n$ so that $d(w, L) \geq \varepsilon_0 n$, hence L is non-trivial in the sense of
 546 Alon et al, and their lower bound applies. \blacktriangleleft

547 4.4 An efficient generic property tester for regular languages.

548 In this section, we show that for any strongly connected NFA \mathcal{A} , there exists ε -property
 549 tester for $L(\mathcal{A})$ that uses $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.

550 **► Theorem 4.13.** *Let \mathcal{A} be a strongly connected NFA. For any $\varepsilon > 0$, there exists an*
 551 *ε -property tester for $L(\mathcal{A})$ that uses $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.*

552 Note that this result is an improvement over the similar result of Bathie and Starikovskaya [7]:
 553 while both testers use the same number of queries, theirs works under the edit distance,
 554 while ours is designed for the Hamming distance. As the edit distance never exceeds the
 555 Hamming distance, the set of words that are ε -far with respect to the former is contained in

556 the set of words ε -far for the latter. Therefore, an ε -tester for the Hamming distance is also
 557 an ε -tester for the edit distance, and our result supersedes and generalizes theirs.

558 The algorithm for Theorem 4.13 is given in Algorithm 1. The procedure is fairly simple:
 559 the algorithm samples at random factors of various lengths in u , and rejects if and only if
 560 at least one of these factors is blocking. On the other hand, the correctness of the tester is
 561 far from trivial. The lengths and the number of factors of each lengths are chosen so that
 562 the number of queries is $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ and the probability of finding a blocking factor is
 563 maximized, regardless of their repartition in u .

■ **Algorithm 1** Generic ε -property tester using $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries

```

1: function SAMPLE( $u, \ell$ )
2:    $i \leftarrow \text{UNIFORM}(0, n - 1)$ 
3:    $l \leftarrow \max(i - \ell, 0), r \leftarrow \min(i + \ell, n - 1)$ 
4:    $\eta \leftarrow (l : u[l..r])$ 
5:   return  $v$ 
6: function TESTER( $u, \varepsilon$ )
7:    $\beta \leftarrow 12m^2/\varepsilon$ 
8:   if  $L(\mathcal{A}) \cap \Sigma^n = \emptyset$  then
9:     Reject
10:  else if  $n < \beta$  then
11:    Query all of  $u$  and run  $\mathcal{A}$  on it
12:    Accept if and only if  $\mathcal{A}$  accepts
13:  else
14:     $\mathcal{F} \leftarrow \emptyset$ 
15:     $T \leftarrow \lceil \log(\beta) \rceil$ 
16:    for  $t = 0$  to  $T$  do
17:       $\ell_t \leftarrow 2^t, r_t \leftarrow \lceil 2 \ln(3)\beta/\ell_t \rceil$ 
18:      for  $i = 1$  to  $r_t$  do
19:         $\mathcal{F} \leftarrow \mathcal{F} \cup \{\text{SAMPLE}(u, \ell_t)\}$ 
20:    Reject if and only if  $\mathcal{F}$  contains a factor blocking for  $\mathcal{A}$ .

```

564 We now turn to proving these properties formally.

565 \triangleright **Claim 4.14.** The tester given in Algorithm 1 makes $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries to u .

566 **Proof.** If $n \leq \beta$, then the tester makes $|u| \leq \beta = \mathcal{O}(1/\varepsilon)$ queries, and the claim holds. The
 567 SAMPLE function with parameter ℓ makes at most 2ℓ queries to u . Therefore, if $n \geq \beta$,
 568 the tester it makes at most $\ell_t \cdot r_t = \mathcal{O}(\beta)$ queries for every $t = 0, \dots, T$, which adds up to
 569 $\mathcal{O}(T \cdot \beta) = \mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries. \blacktriangleleft

570 Next, we show an extension of Lemma 4.8 that shows that if u is ε -far from $L(\mathcal{A})$, then
 571 $(0 : u)$ contains $\Omega(\varepsilon n)$ blocking factors of length $\mathcal{O}(1/\varepsilon)$.

572 \blacktriangleright **Lemma 4.15.** Let $\varepsilon > 0$, let u be a word of length $n \geq 6m^2/\varepsilon$ and assume that $L(\mathcal{A})$
 573 contains at least one word of length n . If u is ε -far from $L(\mathcal{A})$, then the positional word
 574 $(0 : u)$ contains at least $\varepsilon n/(12m^2)$ disjoint blocking factors of length at most $12m^2/\varepsilon$.

575 **Proof.** Let u, \mathcal{A} be a word and an automaton satisfying the above hypotheses. By Lemma 4.8,
 576 $(0 : u)$ contains at least $\varepsilon n/(6m^2)$ disjoint blocking factors. As these factors are disjoint, at
 577 most half of them (that is, $\varepsilon n/(12m^2)$ of them) can have length greater than $12m^2/\varepsilon$, as the
 578 sum of their lengths cannot exceed n . \blacktriangleleft

579 For the correctness analysis, we assume that u is ε -far from $L(\mathcal{A})$, and show that
 580 Algorithm 1 finds a blocking factor of $(0 : u)$ with probability at least $2/3$.

581 ► **Lemma 4.16.** *In the last Else block, if u is ε -far from $L(\mathcal{A})$, then Algorithm 1 rejects
 582 with probability at least $2/3$.*

583 **Proof.** Assume that u is ε -far from $L(\mathcal{A})$. As we are in the last Else block of Algorithm 1,
 584 $L(\mathcal{A}) \cap \Sigma^n$ is not empty (i.e. $L(\mathcal{A})$ contains a word of length n) and $n \geq \beta$, therefore the
 585 conditions of Lemma 4.15 are satisfied. Let \mathcal{B} denote the set of minimal blocking factors in
 586 $(0 : u)$ given by Lemma 4.15: we have $|\mathcal{B}| \geq n/\beta$. We conceptually divide the blocking factors
 587 in \mathcal{B} into different categories depending on their length: for $t = 0, \dots, T$, let B_t denote the
 588 subset of \mathcal{B} of blocking factors of length at most $\ell_t = 2^t$. We then carefully analyze the
 589 probability that randomly sampled factors of length $2\ell_t$ contains a blocking factor from B_t ,
 590 and show that over all t , at least one blocking factor is found with probability at least $2/3$.

591 ▷ **Claim 4.17.** If in a call to SAMPLE, the value i is such that there exists indices $l, r, l \leq i \leq r$,
 592 such that $(0 : u)[l, r]$ is a blocking factor of \mathcal{A} of length at most ℓ , then the factor η returned
 593 by the function is blocking for \mathcal{A} .

594 As the factors given by Lemma 4.15 are disjoint, the probability p_t that the factor returned
 595 by SAMPLE is blocking is lower bounded by

$$596 \quad p_t \geq \frac{1}{n} \sum_{\tau \in B_t} |\tau|$$

597 The SAMPLE function is called $r_t = 2 \ln(3)\beta/\ell_t$ times independently for each t , hence the
 598 probability p that the algorithm samples a blocking factor satisfies the following:

$$\begin{aligned} 599 \quad (1 - p) &= \prod_{t=0}^T (1 - p_t)^{r_t} \leq \exp\left(-\sum_{t=0}^T p_t r_t\right) \\ 600 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{t=0}^T \frac{1}{\ell_t} \sum_{\tau \in B_t} |\tau|\right) \\ 601 \quad &= \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{\tau \in \mathcal{B}} |\tau| \sum_{t=\lceil \log |\tau| \rceil}^T 2^{-t}\right) \\ 602 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{\tau \in \mathcal{B}} |\tau| \cdot 2^{-\lceil \log |\tau| \rceil}\right) \\ 603 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{\tau \in \mathcal{B}} |\tau| \frac{1}{2^{|\tau|}}\right) \\ 604 \quad &= \exp\left(-\frac{2 \ln(3)\beta}{n} \cdot \frac{|\mathcal{B}|}{2}\right) \\ 605 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \cdot \frac{n}{2\beta}\right) \\ 606 \quad &\leq \exp(-\ln(3)) = 1/3 \end{aligned}$$

607 It follows that $p \geq 2/3$, and Algorithm 1 satisfies Definition 2.1. ◀

4.5 Lower bound when there are infinitely many minimal blocking words

We now show that languages with infinitely many blocking factors are hard, i.e. any tester for such a language requires $\Omega(\log(\varepsilon^{-1})/\varepsilon)$ queries.

► **Theorem 4.18.** *Let \mathcal{A} be a trim strongly connected automaton. If $MBF(\mathcal{A})$ is infinite, then there exists a constant ε_0 such that for any $\varepsilon < \varepsilon_0$, any ε -property tester for $L = \mathcal{L}(\mathcal{A})$ uses $\Omega(\log(\varepsilon^{-1})/\varepsilon)$ queries.*

Our proof of this result will look familiar to readers acquainted with the lower bound of Bathie and Starikovskaya [7, Theorem 15]: our proof extends theirs to any language with arbitrarily long minimal blocking words. One difference is that our lower bound applies to ε -testers for the Hamming distance, instead of the edit distance. This is a weakening assumption as the edit distance never exceeds the Hamming distance, but it appears to be needed in the proof of Lemma 4.23.

Our proof is based on (a consequence of) Yao's minmax principle, which we recall here.

► **Fact 4.19** (From Yao's Minmax Principle [21]). *Let $f : \mathbb{R} \rightarrow \mathbb{N}$ be a nondecreasing function. Let \mathcal{T} denote the set of all algorithms using less than $f(\varepsilon)$ queries, and let \mathcal{T}_D denote the subset of deterministic algorithms. Let \mathcal{D} be a probability distribution over Σ^* . Then, we have*

$$\inf_{T \in \mathcal{T}} \sup_{x \in \Sigma^*} \mathbb{P}_T(T \text{ errs on } x) \geq \inf_{T \in \mathcal{T}_D} \mathbb{P}_{x \sim \mathcal{D}}(T \text{ errs on } x).$$

Therefore, to show that any randomized algorithm with less than $\log(\varepsilon^{-1})/\varepsilon$ queries errs with large probability, it suffices to exhibit a probability distribution over inputs such that any *deterministic* tester errs with large probability on this distribution.

We will construct a hard distribution using long minimal blocking factors, and show that with large probability, any deterministic algorithm using less than $\log(\varepsilon^{-1})/\varepsilon$ queries has the same query results for many pairs of positive and ε -far instances. As the tester is deterministic, it must answer the same on all these pairs, and therefore make an error with large probability.

Our proof of Theorem 4.18 goes through the following steps:

1. first, show that with high probability, an input u sampled w.r.t. \mathcal{D} is either in or ε -far from L (Lemma 4.23),
2. show that with high probability, any deterministic tester that makes fewer than $c \cdot \log(\varepsilon^{-1})/\varepsilon$ queries (for a suitable constant c) cannot distinguish whether the instance u is positive or ε -far,
3. combine the above to prove Theorem 4.18 via Fact 4.19.

4.5.1 Constructing a hard distribution

Let $\varepsilon > 0$ be sufficiently small and let n be a large enough integer. In what follows, m denotes the number of states of \mathcal{A} . To construct the hard distribution \mathcal{D} , we will use an infinite family of blocking factors that share a common structure, given by the following lemma.

► **Lemma 4.20.** *If $MBF(\mathcal{A})$ is infinite, then there exist positional words ϕ, ν_+, ν_-, χ such that:*

1. the words ν_+ and ν_- have the same length,
2. there exists a constant $S = 2^{\mathcal{O}(m)}$ such that $|\phi|, |\nu_+|, |\nu_-|, |\chi| \leq S$,

649 **3.** *there exists an index $i_* \in \mathbb{Z}/\lambda\mathbb{Z}$ and a state $q_* \in Q_{i_*}$ such that for every integer $r \geq 1$,*
 650 *$\tau_{-,r} = \phi(\nu_-)^r z$ is blocking for \mathcal{A} , and for every $s < r$, we have*

651
$$q_* \xrightarrow{\tau_{+,r,s}} q_* \text{ where } \tau_{+,r,s} = \phi(\nu_-)^j \nu_+ (\nu_-)^{r-1-s} \chi.$$

652 *In particular, $\tau_{+,r,s}$ is not blocking for \mathcal{A} .*

653 The crucial property here is that $\tau_{-,r}$ and $\tau_{+,r,s}$ are very similar: they have the same
 654 length, differ in at most S letters, yet one of them is blocking and the other is not. The
 655 proof of this lemma is deferred to Appendix A.

656 We now use the words $\tau_{-,r}$ and $\tau_{+,r,s}$ and the constant S to describe how to sample an
 657 input $\mu = (0 : u)$ of length n w.r.t. \mathcal{D} .

658 Let π be a uniformly random bit. If $\pi = 1$, we will construct a positive instance
 659 $\mu \in \mathcal{TL}(\mathcal{A})$, and otherwise the instance will be ε -far from $\mathcal{TL}(\mathcal{A})$ with high probability.
 660 We divide the interval $[1..n]$ into $k = \varepsilon n$ intervals of length $\ell = 1/\varepsilon$, plus small initial and
 661 final segments μ_i and μ_f of length $\mathcal{O}(\rho)$ to be specified later. For the sake of simplicity, we
 662 assume that k and ℓ are integers and that λ divides ℓ . For $j = 1, \dots, k$, let a_j, b_j denote
 663 the endpoints of the j -th interval. For each interval, we sample independently at random a
 664 variable τ_j with the following distribution:

665
$$\tau_j = \begin{cases} t, & \text{w.p. } p_t = 3 \cdot 2^t S \varepsilon / \log((S\varepsilon)^{-1}) \text{ for } t = 1, 2, \dots, \log((S\varepsilon)^{-1}), \\ 0, & \text{w.p. } p_0 = 1 - \sum_{t=1}^{\log((S\varepsilon)^{-1})} p_t. \end{cases} \quad (3)$$

666 The event $\tau_j > 0$ means that the j -th interval is filled with with $N \approx 2^{-\tau_j}/\varepsilon$ “special” factors.
 667 When $\pi = 0$, these “special” factors will be minimal blocking factors $\tau_{-,r}$ for $r = 2^{\tau_j}$, whereas
 668 when $\pi = 1$, they will instead be similar non-blocking factors $\tau_{+,r,s}$ for a uniformly random
 669 s : they will be hard to distinguish with few queries. On the other hand, the event $\tau_j = 0$
 670 means that the j -th interval contains no specific information. More precisely, we choose a
 671 positional word η_* of length ℓ such that $q_* \xrightarrow{w_*} q_*$: by Fact 3.3, this is possible as $\ell = 0$
 672 (mod λ). Then, if $\tau_j = 0$, we set $\mu[a_j..b_j] = \eta_*$, regardless of the value of π .

673 Formally, if $\tau_j > 0$, let $r = 2^{\tau_j}$, $N = 2^{-\tau_j}/(S\varepsilon)$ and let η be a word of length $\ell - N \cdot |\tau_{-,r}|$
 674 such that $q_* \xrightarrow{\eta} q_*$: such a word exists as λ divides ℓ and $|\tau_{-,r}|$. We construct the j -th
 675 interval as follows:

- 676 ■ if $\pi = 0$, we set $\mu[a_j..b_j] = (\tau_{-,r})^N \eta$,
- 677 ■ if $\pi = 1$, we select $s \in [0..r-1]$ uniformly at random, and set $\mu[a_j..b_j] = (\tau_{+,r,s})^N \eta$.

678 Finally, the initial and final fragments μ_i and μ_f of μ are chosen to be the shortest words
 679 that label a transition from q_0 to q_* and q_* to a final state, respectively.

680 4.5.2 Properties of the distribution \mathcal{D}

681 We now conclude the proof of Theorem 4.2 by studying properties of the distribution \mathcal{D} .

682 ► **Observation 4.21.** *If ε is small enough, \mathcal{D} is well-defined, i.e. for every t between 0 and*
 683 *$\log((S\varepsilon)^{-1})$, we have $0 \leq p_t \leq 1$.*

684 ► **Observation 4.22.** *If $\pi = 1$, then $\mu \in \mathcal{TL}(\mathcal{A})$.*

685 Next, we show that when $\pi = 0$, the resulting instance is ε -far from L with high probability.

686 ► **Lemma 4.23.** *Conditioned on $\pi = 0$, the probability of the event $\mathcal{F} = \{\mu \text{ is } \varepsilon\text{-far from}$*
 687 *$\mathcal{TL}(\mathcal{A})\}$ goes to 1 as n goes to infinity.*

688 **Proof.** When $\pi = 0$, the procedure for sampling μ puts blocking factors of the form $(i_* : x)$
 689 at positions equal to $i_* \bmod \lambda$. Any word containing such a factor at such a position is not
 690 in $\mathcal{TL}(\mathcal{A})$, therefore any sequence of substitutions that transforms μ into a word of $\mathcal{TL}(\mathcal{A})$
 691 must make at least one substitution in every such factor. Consequently, the distance between
 692 μ and $\mathcal{TL}(\mathcal{A})$ is at least the number of blocking factors in μ . To prove the lemma, we show
 693 that this number is at least εn with high probability, by showing that it is larger than εn by
 694 a constant factor in expectation and using a concentration argument.

695 Let B_j denote the number of blocking factors in the j -th interval: it is equal to $2^{-\tau_j}/(S\varepsilon)$
 696 when $\tau_j > 0$ and to 0 otherwise.

697 \triangleright **Claim 4.24.** Let $B = \sum_{j=1}^k B_j$, and let $E = \mathbb{E}[B]$. We have $E \geq 2\varepsilon n$.

698 **Claim proof.** By direct calculation:

$$\begin{aligned}
 699 \quad E &= \sum_{j=1}^k \mathbb{E}[B_j] && \text{by linearity} \\
 700 \quad &= \sum_{j=1}^k \sum_{t=1}^{\log(S/\varepsilon)} 2^{-t}/(S\varepsilon) \cdot p_t && \text{def. of expectation} \\
 701 \quad &= \sum_{j=1}^k \sum_{t=1}^{\log(S/\varepsilon)} 2^{-t}/(S\varepsilon) \cdot 3 \cdot 2^t \varepsilon S / \log(S/\varepsilon) && \text{def. of } p_t \\
 702 \quad &= \sum_{j=1}^k \sum_{t=1}^{\log(S/\varepsilon)} 3 / \log(S/\varepsilon) \\
 703 \quad &= 3k \geq 2\varepsilon n
 \end{aligned}$$

704 ◀

705 We will now show that $\mathbb{P}(B < \varepsilon n)$ goes to 0 as n goes to infinity. By Claim 4.24, we have
 706 $B < \varepsilon n \Rightarrow E - B \geq \varepsilon n$, and therefore $\mathbb{P}(B < \varepsilon n) \leq \mathbb{P}(E - B \geq \varepsilon n)$. The random variable
 707 B is the sum of k independent random variables, each taking values between 0 and $1/(S\varepsilon)$.
 708 Therefore, by Hoeffding's Inequality (Lemma B.1), we have

$$\begin{aligned}
 709 \quad \mathbb{P}(E - B < \varepsilon n) &\leq \exp\left(-\frac{2\varepsilon^2 n^2}{k/(S\varepsilon)^2}\right) \\
 710 &\leq \exp\left(-\frac{2S^2 \varepsilon^4 n^2}{\varepsilon n}\right) \text{ as } k \leq \varepsilon n \\
 711 &\leq \exp(-2S^2 \varepsilon^3 n)
 \end{aligned}$$

712 This probability goes to 0 as n goes to infinity, which concludes the proof. ◀

713 \blacktriangleright **Corollary 4.25.** For large enough n , we have $\mathbb{P}(\mathcal{F}) \geq 5/12$.

714 Intuitively, our distribution is hard to test because positive and negative instance are
 715 very similar. Therefore, a tester with few queries will likely not be able to tell them apart:
 716 the perfect completeness constraint forces the tester to accept in that case. Below, we prove
 717 this last part formally.

718 \blacktriangleright **Lemma 4.26.** Let T be a deterministic tester with perfect completeness (i.e. one sided
 719 error, always accepts $\tau \in \mathcal{TL}(\mathcal{A})$) and let q_j denote the number of queries that it makes in
 720 the j -th interval. Conditioned on the event $\mathcal{M} = \{\forall j \text{ s.t. } \tau_j > 0, q_j < 2^{\tau_j}\}$, the probability
 721 that T accepts u is 1.

722 **Proof.** We show that if there exists a τ with non-zero probability w.r.t. \mathcal{D} under \mathcal{M} that T
 723 rejects, then there exists a word $\tau' \in \mathcal{TL}(\mathcal{A})$ that T rejects that also has non-zero probability,
 724 contradicting the fact that T has perfect completeness.

725 Let τ be the word rejected by T : as T has perfect completeness, hence $\tau \notin \mathcal{TL}(\mathcal{A})$, and
 726 there must be at least one interval with $\tau_j > 0$. Consider every interval j such that $\tau_j > 0$: it
 727 is of the form $(\tau_{-,r})^N v$ where $r = 2^{\tau_j}$ and $\tau_{-,r} = \phi(\nu_-)^r \chi$. Therefore, if $q_j < 2^{\tau_j}$, then there
 728 is a copy of ν_- that has not been queried by T across all copies of $\tau_{-,r}$. Consider the word
 729 τ' obtained by replacing this copy of ν_- with ν_+ in all N copies of $\tau_{-,r}$ in the block. The
 730 result block is of the form $(\tau_{+,r,s})^N v$ for some $s < r$, and by construction it is not blocking.
 731 Applying this operation to all blocks results in a word τ' that is in $\mathcal{TL}(\mathcal{A})$. Furthermore,
 732 τ' has non-zero probability under \mathcal{D} conditioned on \mathcal{M} : it can be obtained by flipping the
 733 random bit π and choosing the right index s in every block. \blacktriangleleft

734 Next, we show that if a tester makes few queries, then the event \mathcal{M} has large probability.

735 **► Lemma 4.27.** *Let T be a deterministic tester, let q_j denote the number of queries that*
 736 *it makes in the j -th interval, and assume that T makes at most $\frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$ queries, i.e.*
 737 *$\sum_j q_j \leq \frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$. The probability of the event $\mathcal{M} = \{\forall j \text{ s.t. } \tau_j > 0, q_j < 2^{\tau_j}\}$ is at*
 738 *least $11/12$.*

739 **Proof.** We show that the probability of $\overline{\mathcal{M}}$, the complement of \mathcal{M} , is at most $1/12$. We
 740 have:

$$\begin{aligned}
 741 \quad \mathbb{P}(\overline{\mathcal{M}}) &= \mathbb{P}(\exists j : \tau_j > 0 \wedge q_j \geq 2^{\tau_j}) \\
 742 \quad &\leq \sum_j \mathbb{P}(\tau_j > 0 \wedge q_j \geq 2^{\tau_j}) && \text{by union bound} \\
 743 \quad &\leq \sum_j \sum_{t=1}^{\lfloor \log q_j \rfloor} p_t \\
 744 \quad &= \sum_j \sum_{t=1}^{\lfloor \log q_j \rfloor} \frac{3 \cdot 2^t \varepsilon}{\log(S/\varepsilon)} && \text{by def. of } p_t \\
 745 \quad &\leq \frac{3\varepsilon}{\log(S/\varepsilon)} \sum_j \sum_{t=1}^{\lfloor \log q_j \rfloor} 2^t \\
 746 \quad &\leq \frac{3\varepsilon}{\log(S/\varepsilon)} \sum_j 2q_j \\
 747 \quad &= \frac{3\varepsilon}{\log(S/\varepsilon)} \cdot \frac{2}{72} \cdot \frac{\log(1/\varepsilon)}{\varepsilon} \\
 748 \quad &\leq 1/12
 \end{aligned}$$

749 \blacktriangleleft

750 We are now ready to prove Theorem 4.2.

751 **Proof of Theorem 4.2.** We want to show that any non-adaptive tester with perfect com-
 752 pleteness for $L(\mathcal{A})$ requires at least $\frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$ queries, by showing that any tester with
 753 fewer queries errs with probability at least $1/3$. We use Yao's minmax principle (Fact 4.19),
 754 and show that any **deterministic** non-adaptive algorithm T with perfect completeness that
 755 makes less than $\frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$ queries errs on u when $(0 : u) \sim \mathcal{D}$ with probability at least
 756 $1/3$.

757 Consider such an algorithm T . The probability that T makes an error on u is lower-
 758 bounded by the probability that u is ε -far from $L(\mathcal{A})$ and T accepts, which in turn is larger
 759 than the probability of $\mathcal{M} \cap \mathcal{F}$. By Corollary 4.25, we have $\mathbb{P}(\mathcal{F}) \geq 5/12$, and by Lemma 4.27,
 760 $\mathbb{P}(\mathcal{M})$ is at least $11/12$. Therefore, we have

$$761 \quad \mathbb{P}(T \text{ errs}) \geq \mathbb{P}(\mathcal{M} \cap \mathcal{F}) \geq 1 - 7/12 - 1/12 = 1/3.$$

762 This concludes the proof of Theorem 4.2. ◀

763 **5 The Case of General NFAs**

764 In this section we extend the previous characterisation to all finite automata, proving our
 765 main theorem, stated as Theorem 2.7 in the overview section. To do so, we generalise the
 766 notion of blocking factor: we introduce *blocking sequences*, which are sequences of factors that
 767 witness the fact that we cannot take any path through the strongly connected components of
 768 the automaton. We define a suitable partial order on blocking sequences, which extends the
 769 factor relation on words to those sequences. The classification between trivial, easy and hard
 770 of a language can be characterised by the set of *minimal blocking sequences* of an automaton
 771 recognising it. This is expressed by the following theorem, where $MBS(\mathcal{A})$ stands for the set
 772 of *minimal blocking sequences* of \mathcal{A} .

773 The statement we will prove is the following:

774 **► Theorem 5.1.** *Let L be an infinite language recognised by the trim NFA \mathcal{A} . The complexity*
 775 *of testing L is characterized by $MBS(\mathcal{A})$ as follows:*

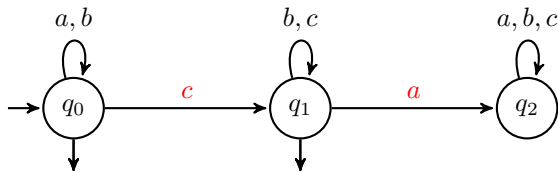
- 776 1. *L is hard to test if and only if $MBS(\mathcal{A})$ is infinite.*
- 777 2. *L is easy to test if and only if $MBS(\mathcal{A})$ is finite and nonempty.*
- 778 3. *L is trivial if and only if $MBS(\mathcal{A})$ is empty.*

779 Recall that we only consider infinite languages in this classification.

This section uses the `knowledge` package, to help the reader keep track of the various notions. Some *important terms* are coloured in red when we define them. Occurrences of those *important terms* are coloured in blue. The reader can click on those (or just hover over them on some PDF readers) to see the definition.

780

781 The rest of this section is dedicated to the proof of Theorem 5.1. Before we get into the
 782 proof, let us go through some examples, which illustrate some of the main difficulties. In all
 783 that follows we will abbreviate “strongly connected component” as SCC. We call an SCC
 784 trivial if it is just a single state with no self-loop.



785 **■ Figure 3** An automaton recognising the language $(a + b)^*(b + c)^*$.

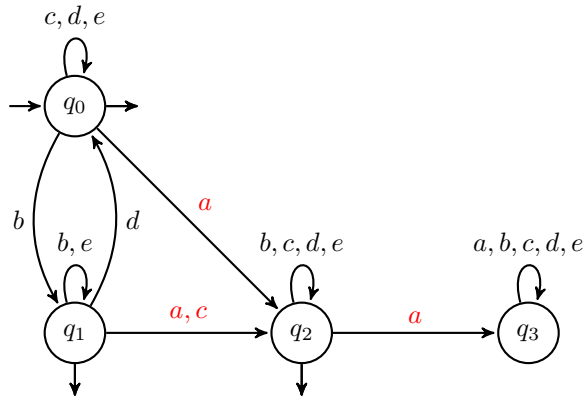
786 **► Example 5.2.** Observe the automaton in Figure 3. It has two SCCs, plus a sink state. The set of *minimal blocking factors* of its language is infinite: it is the set cb^*a . Yet, it is

787 easy: Given a word w , it suffices to sample $O(1/\varepsilon)$ positions at random and reject if we
 788 see a c appearing before an a . Clearly if the word is in the language, every c must be after
 789 every a , thus we accept. On the other hand, suppose the word is ε -far from the language.
 790 Let u be the maximal prefix of w containing less than $\varepsilon|w|/2$ occurrences of c . If $u = w$
 791 then we can turn every c in w into an a to make it accepted, and thus $d(w, L) < \varepsilon|w|/2$, a
 792 contradiction. Hence we can write w as ucv . If v contains less than $\varepsilon|w|/2$ occurrences of a
 793 then $d(w, L) < \varepsilon|w|$, again a contradiction.

794 Otherwise, u contains $\varepsilon|w|/2$ occurrences of c and v contains $\varepsilon|w|/2$ occurrences of a .
 795 Then the probability that when picking $\varepsilon|w|$ letters at random we sample one of the c in u
 796 and one of the a in v is lower-bounded by a positive constant. In conclusion, we reject with
 797 constant probability when the word is ε -far from the language.

798 The crucial point in the following proof is the use of blocking *sequences* instead of blocking
 799 *factors*. A blocking sequence is a list of factors that are blocking for SCCs of the automaton,
 800 so that seeing this sequence as disjoint factors of a word guarantees that it is rejected.
 801 Blocking sequences come with a natural notion of minimality, which lets us characterise
 802 languages that are *easy* as those that admit finitely many *minimal blocking sequences*.

803 In the example above, a (unique) minimal blocking sequence is (c, a) .



■ **Figure 4** An automaton recognising the language $[\varepsilon + ((c + d + e)^*b(b + e)^*d)^*a](b + c + d + e)^*$.

804 ► **Example 5.3.** In Figure 4 we display an automaton with two SCCs and a sink state. The
 805 first SCC has blocking factors $be^*c + a$ and the second one just a . This automaton is *easy*:
 806 intuitively, a word that is ε -far from this language has to contain many a , as otherwise we
 807 can make it accepted by deleting all a , thanks to the second SCC. As a is also a blocking
 808 factor of the first SCC, we only need to look for two a s in the word.

809 The family of unbounded blocking factors of the first SCC is made irrelevant by the fact
 810 that a word far from the language must contain many a anyway.

811 We fix an NFA $\mathcal{A} = (Q, \Sigma, \delta, q_{init}, q_f)$. Once again note that it has a single final state q_f .
 812 Let \mathcal{S} be its set of SCCs. We define the partial order relation $\leq_{\mathcal{A}}$ on \mathcal{S} as: $S \leq_{\mathcal{A}} T$ if and
 813 only if T is reachable from S . We write $<_{\mathcal{A}}$ for its strict part $\leq_{\mathcal{A}} \setminus \geq_{\mathcal{A}}$.

814 We define p as the least common multiple of the lengths of all simple cycles of \mathcal{A} . Given
 815 a number $k \in \{0, \dots, p - 1\}$, we say that a state t is k -reachable from a state s if there is
 816 a path from s to t of length k modulo p . In what follows, we use “positional words” for
 817 p -positional words with this value of p .

818 ▶ **Remark 5.4.** In the rest of this section we will not try to optimise the constants in the
 819 formulas. They will, in fact, become quite large in some of the proofs. We make this choice
 820 to make the proofs more readable, although some of them are already technical.

821 For instance, the choice of p as the lcm of the lengths of simple cycles is not optimal: we
 822 could use, for instance, the lcm of the periodicities of the SCCs.

823 ▶ **Definition 5.5.** A **portal** is a 4-tuple $s, x \rightsquigarrow t, y \in (Q \times \{0, \dots, p-1\})^2$, such that s and t
 824 are in the same SCC. It describes the first and last states visited by a path in an SCC, and
 825 the times at which it first and lasts visits that SCC (modulo p).

826 The *positional language* of a portal is the set

$$827 \quad \mathcal{PL}(s, x \rightsquigarrow t, y) = \{(x : w) \mid t \in \delta(s, w) \wedge x + |w| = y \pmod{p}\}.$$

828 Portals were already defined in [5], in a slightly different way. Our definition will allow us to
 829 express blocking sequences more naturally.

830 ▶ **Definition 5.6.** A *positional word* $(n : u)$ is **blocking** for a portal $s, x \rightsquigarrow t, y$ if it is not a
 831 factor of any word of $\mathcal{PL}(s, x \rightsquigarrow t, y)$. In other words, there is no path that starts in s and
 832 ends in t , of length $y - x$ modulo p , which reads u after $n - x$ steps modulo p .

833 ▶ **Remark 5.7.** There is an NFA with $\leq p|\mathcal{A}|$ states recognising $\mathcal{PL}(s, x \rightsquigarrow t, y)$: it simply
 834 simulates the SCC of s while keeping track of the number of letters read, plus x , modulo p .
 835 Its set of states is thus a subset of $\{0, \dots, p-1\} \times Q$.

836 It is strongly connected: say we read a word u from (s, x) and reach (s', x') . There is
 837 a path from s' to s in \mathcal{A} , labelled by a word v . Hence we can reach (s, x) from (s', x') by
 838 reading $v(uv)^{p-1}$.

839 Its periodicity is p . Hence we can use all results we obtained on strongly connected NFAs
 840 on **portals**, with $p|\mathcal{A}|$ as the number of states and p as the periodicity.

841 ▶ **Definition 5.8.** An **SCC-path** π of \mathcal{A} is a sequence of portals linked by transitions

$$842 \quad \pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} s_1, x_1 \rightsquigarrow t_1, y_1 \cdots \xrightarrow{a_k} s_k, x_k \rightsquigarrow t_k, y_k,$$

843 such that for all $i \in \{1, \dots, k\}$, $x_i = y_{i-1} + 1 \pmod{p}$, $s_i \in \delta(t_{i-1}, a_i)$, and $t_{i-1} <_{\mathcal{A}} s_i$.

844 It is a description of the states and times at which a path through the automaton enters
 845 and leaves the SCCs.

846 The language $\mathcal{L}(\pi)$ of an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \cdots s_k, x_k \rightsquigarrow t_k, y_k$ is the set

$$847 \quad \mathcal{L}(\pi) = \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0) a_1 \mathcal{L}(s_1, x_1 \rightsquigarrow t_1, y_1) a_2 \cdots \mathcal{L}(s_k, x_k \rightsquigarrow t_k, y_k)$$

848 We say that π is **accepting** if $x_0 = 0$, $s_0 = q_{init}$, $t_k = q_f$ and $\mathcal{L}(\pi)$ is non-empty.

▶ **Fact 5.9.**

$$849 \quad \mathcal{L}(\mathcal{A}) = \bigcup_{\pi \text{ accepting}} \mathcal{L}(\pi).$$

850 **Proof.** Let $w = b_1 \cdots b_\ell \in \mathcal{L}(\mathcal{A})$. There exists $\rho = q_0 \xrightarrow{b_1} q_1 \cdots \xrightarrow{b_\ell} q_\ell$ an accepting run in \mathcal{A} .

851 Let $i_1 < \dots < i_k$ be the sequence of indices such that $\{i_1, \dots, i_k\} = \{i \in \{1, \dots, m\} \mid$
 852 $q_{i-1} <_{\mathcal{A}} q_i\}$. We also define $i_0 = 0$ and $i_{k+1} = \ell + 1$. In other words, those are the indices at
 853 which ρ enters a new SCC.

854 We define the SCC-path $\pi(\rho)$ as follows:

$$855 \quad \pi(\rho) = q_0, 0 \rightsquigarrow q_{i_1-1}, y_0 \xrightarrow{a_{i_1}} q_{i_1}, x_1 \rightsquigarrow q_{i_2-1}, y_1 \cdots \xrightarrow{a_{i_k}} q_{i_k}, x_k \rightsquigarrow q_\ell, y_k$$

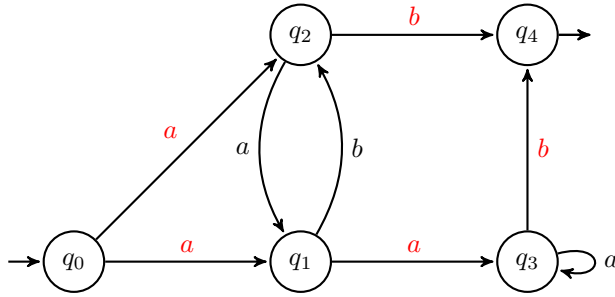
856 with $x_j = m + i_j \pmod{p}$ and $y_j = m + i_{j+1} - 1 \pmod{p}$ for all $j \in \{0, \dots, k\}$. Clearly
 857 $w \in \mathcal{L}(\pi(\rho))$ and $\pi(\rho)$ is an accepting SCC-path.

858 Conversely, let $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \cdots s_k, x_k \rightsquigarrow t_k, y_k$ be an accepting SCC-path in \mathcal{A}
 859 and let $w \in \mathcal{L}(\pi)$.

860 For all $j \in \{0, \dots, k\}$, there is a word w_j labelling a path from s_j to t_j , such that
 861 $w = w_0 a_1 \cdots w_k$. By gluing those paths and the transitions $t_{j-1} \xrightarrow{a_j} s_j$, we obtain an
 862 accepting run for w in \mathcal{A} . \blacktriangleleft

863 Decomposing \mathcal{A} as a union of SCC-paths allows us to use them as an intermediate step.
 864 We define blocking sequences for SCC-paths before defining them on automata.

865 **► Definition 5.10.** We say that a sequence $((n_1 : u_1), \dots, (n_\ell : u_\ell))$ of positional factors is
 866 **blocking** for an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \cdots s_k, x_k \rightsquigarrow t_k, y_k$ if there is a sequence of
 867 indices $i_0 \leq i_1 \leq \cdots \leq i_k$ such that $(n_{i_j} : u_{i_j})$ is blocking for $s_j, x_j \rightsquigarrow t_j, y_j$, for all j .



868 **Figure 5** Automaton used for Example 5.11.

869 **► Example 5.11.** Take a look at the automaton displayed in Figure 5. It has four SCCs,
 870 including two trivial ones $\{q_0\}$ and $\{q_4\}$. The lcm of the lengths of its simple cycles is $p = 2$.

871 It has six accepting SCC-paths:

- 872 $\blacksquare q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_1, 1 \rightsquigarrow q_1, 1 \xrightarrow{a} q_3, 0 \rightsquigarrow q_3, 0 \xrightarrow{b} q_4, 1 \rightsquigarrow q_4, 1$
- 873 $\blacksquare q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_1, 1 \rightsquigarrow q_1, 1 \xrightarrow{a} q_3, 0 \rightsquigarrow q_3, 1 \xrightarrow{b} q_4, 0 \rightsquigarrow q_4, 0$
- 874 $\blacksquare q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_2, 1 \rightsquigarrow q_1, 0 \xrightarrow{a} q_3, 1 \rightsquigarrow q_3, 0 \xrightarrow{b} q_4, 1 \rightsquigarrow q_4, 1$
- 875 $\blacksquare q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_2, 1 \rightsquigarrow q_1, 0 \xrightarrow{a} q_3, 1 \rightsquigarrow q_3, 1 \xrightarrow{b} q_4, 0 \rightsquigarrow q_4, 0$
- 876 $\blacksquare q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_1, 1 \rightsquigarrow q_2, 0 \xrightarrow{b} q_4, 1 \rightsquigarrow q_4, 1$
- 877 $\blacksquare q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_2, 1 \rightsquigarrow q_2, 1 \xrightarrow{b} q_4, 0 \rightsquigarrow q_4, 0$

878 The language of the first SCC-path is $a(ba)^*a(a^2)^*b$. A blocking sequence for this SCC-
 879 path is $(0 : aa), (0 : b)$, which is in fact blocking for all those SCC-paths. Another one is
 880 $(1 : ab)$.

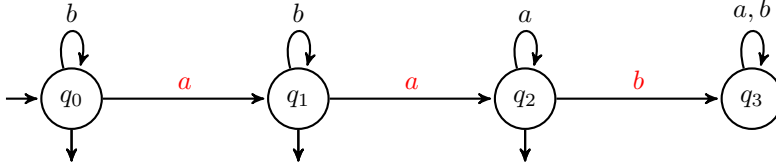
881 On the other hand, $(0 : ab)$ is not blocking for this path, as $(0 : ab)$ is not a blocking
 882 factor for the portal $q_1, 1 \rightsquigarrow q_1, 1$. It is, however, a blocking sequence for the third, fourth
 883 and last SCC-paths.

884 This is because if we enter the SCC $\{q_1, q_2\}$ through q_1 , a factor ab can only appear
 885 after an even number of steps, while if we enter through q_2 , it can only appear after an odd
 886 number of steps.

886 ► **Example 5.12.** The automaton \mathcal{A} displayed in Figure 6 only has cycles of length 1, hence
 887 $p = 1$. They are totally ordered by $\leq_{\mathcal{A}}$.

888 Observe that the sequence $((0 : a), (0 : b))$ is a blocking sequence for the SCC-path
 889 $\pi = q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_1, 0 \rightsquigarrow q_1, 0 \xrightarrow{a} q_2, 0 \rightsquigarrow q_2, 0$. Indeed, a is blocking for the first two
 890 portals, and b for the third. We can verify Lemma 5.15 here: If a word contains $|Q| = 4$
 891 disjoint sequences $((0 : a), (0 : b))$, then in particular it must contain factors a , a and b in
 892 that order.

893 Even two blocking sequences would be enough here, but note that containing one blocking
 894 sequence is not enough: the word aba contains $((a : 0), (b : 0))$, yet it is in the language of π .



■ **Figure 6** An automaton recognising the language $b^* + b^*ab^*a^*$.

895 In order to smoothen the proofs of the following results, let us start with two technical
 896 lemmas expressing two basic properties of the Hamming distance with respect to \mathcal{A} . The
 897 first one states that, for all SCC-path π and ℓ large enough, whether $\mathcal{L}(\pi)$ contains a word
 898 of length ℓ only depends on the value $\ell \pmod p$.

899 ► **Lemma 5.13.** *Let $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$ be an SCC-path and
 900 $r \in \{0, \dots, p-1\}$. If there exists a word $w \in \mathcal{L}(\pi)$ with $|w| = r \pmod p$ and $|w| \geq |\mathcal{A}|$
 901 then for all $\ell \in \mathbb{N}$ such that $\ell = r \pmod p$ and $\ell \geq p|\mathcal{A}| + 3|\mathcal{A}|^3$ there exists $w' \in \mathcal{L}(\pi)$ with
 902 $|w'| = \ell$.*

903 **Proof.** Suppose there exists $w \in \mathcal{L}(\pi)$ with $|w| = r \pmod p$ and $|w| \geq |\mathcal{A}|$. Then we can
 904 decompose it as $w = w_0 a_1 \dots w_k$ with $w_i \in \mathcal{L}(s_i, x_i \rightsquigarrow t_i, y_i)$ for all i . For each $i \in \{0, \dots, k\}$,
 905 let S_i be the SCC of s_i , and p_i the periodicity S_i . For all i such that the S_i is not trivial, by
 906 Fact 3.3, there is a word v'_i labelling a path from s_i to t_i of length ℓ_i with $\ell_i = m_{i+1} - m_i$
 907 $\pmod{p_i}$ and $\ell_i \leq 3|\mathcal{A}|^2$. If $i = k$, we set $m_{k+1} = r$. Since the SCC of s_i is not trivial, there is
 908 a simple cycle from s_i to itself. Let c_i be its length and u_i the word it reads. Since p_i divides
 909 p , we know that $m_{i+1} - m_i - \ell_i = r_i p_i \pmod p$ for some $r_i \in \{0, \dots, p/p_i - 1\}$. The word
 910 $w'_i = u_i^{r_i} v'_i$ labels a path from s_i to t_i , of length $\ell_i + r_i p_i = m_{i+1} - m_i \pmod p$. Furthermore
 911 we have $|w'_i| \leq p + 3|\mathcal{A}|^2$. If s_i is in a trivial SCC, then w_i is the empty word γ . In that case
 912 we set $w'_i = \gamma$. We set $w' = w'_1 a_1 w'_2 \dots a_{k-1} w'_k$. We have $w' \in \mathcal{L}(\pi)$, $|w'| \leq p|\mathcal{A}| + 3|\mathcal{A}|^3$ and
 913 $|w'| = r \pmod p$.

914 Since $w \in \mathcal{L}(\pi)$ and $|w| \geq |\mathcal{A}|$, the run reading w has to go through a cycle, hence there
 915 must be an i such that S_i is non-trivial. Let u be a word labelling a simple cycle from s_i
 916 to itself. Since $|u|$ divides p , for any $\ell \in \mathbb{N}$ such that $\ell = r \pmod p$ and $\ell \geq p|\mathcal{A}| + 3|\mathcal{A}|^3$
 917 we can find a word of length ℓ in $\mathcal{L}(\pi)$ by adding this cycle enough times in the run of w'
 918 constructed before. ◀

919 Our second technical lemma expresses that adding p letters to a word can only increase
 920 the distance by p .

921 ► **Lemma 5.14.** *Let $s, x \rightsquigarrow t, y$ be a portal such that the SCC of s and t is non-trivial, and
 922 w a word such that $d(w, \mathcal{L}(s, x \rightsquigarrow t, y)) < +\infty$. Let $u \in \Sigma^p$. Then we have $d(wu, \mathcal{L}(s, x \rightsquigarrow$
 923 $t, y)) \leq d(w, \mathcal{L}(s, x \rightsquigarrow t, y)) + p$.*

924 **Proof.** As $d(w, \mathcal{L}(s, x \rightsquigarrow t, y)) < +\infty$, there exists $w' \in \mathcal{L}(s, x \rightsquigarrow t, y)$ such that $d(w, w') =$
 925 $d(w, \mathcal{L}(s, x \rightsquigarrow t, y))$. Thus there is a path of length $y - x \pmod{p}$ from s to t reading w' . As
 926 the SCC of t is non-trivial, there is a cycle from t to itself. Let v be a word labelling a simple
 927 cycle from t to itself. By definition of p , $|v|$ divides p , thus there exists k such that $k|v| = p$. In
 928 consequence, the word $w'v^k$ is in $\mathcal{L}(s, x \rightsquigarrow t, y)$. Furthermore, since $d(w, w') = d(w, \mathcal{L}(s, x \rightsquigarrow$
 929 $t, y))$, we have $d(wu, \mathcal{L}(s, x \rightsquigarrow t, y)) \leq d(wu, w'v^k) \leq d(w, \mathcal{L}(s, x \rightsquigarrow t, y)) + p$. \blacktriangleleft

930 We say that blocking sequences of a word are disjoint if they appear on disjoint sets of
 931 positions.

932 \blacktriangleright **Lemma 5.15.** *If $(0 : w)$ contains $|Q|$ disjoint blocking sequences for an SCC-path $\pi =$
 933 $s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$, then $w \notin \mathcal{L}(\pi)$.*

934 **Proof.** We prove a slightly stronger statement by induction on k :

935 If $(m : w)$ contains k disjoint blocking sequences for an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1}$
 936 $\dots s_k, x_k \rightsquigarrow t_k, y_k$, with $x_0 = m$, then no word of $\mathcal{L}(\pi)$ has w as a suffix.

937 The base case is trivial as the empty SCC-path has an empty language.

938 Now let $k > 0$ and suppose this proposition holds for $k - 1$. Consider an SCC-path
 939 $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$ and disjoint blocking sequences $\sigma_1 \dots, \sigma_k$. Let
 940 $(m : w) = (m : w_-)(m_v : v)(m_+ : w_+)$ with $(m_v : v)$ the first factor from one of the blocking
 941 sequences that is blocking for (m_1, s_1, t_1) . Let σ_i be the blocking sequence in which it
 942 appears.

943 Since the blocking sequences are disjoint, for every blocking sequence other than σ ,
 944 its part appearing in w_+ must be a blocking sequence for π , and thus also for $s_1, x_1 \rightsquigarrow$
 945 $t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$. Hence w_+ contains $k - 1$ disjoint blocking sequences for
 946 $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$. By induction hypothesis, no word of $\mathcal{L}(s_1, x_1 \rightsquigarrow$
 947 $t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k)$ has w_+ as a suffix. Let u be a word having w as a suffix.
 948 Suppose by contradiction that $u \in \mathcal{L}(\pi)$. Then $u = u_- a_1 u_+$ with $u_- \in \mathcal{L}((m_1, s_1, t_1))$ and
 949 $u_+ \in \mathcal{L}((m_2, s_2, t_2), \dots, (m_k, s_k, t_k))$. Further, since w_+ is a suffix of w which is a suffix of
 950 u , we have $u = u_p w_+$ for some prefix u_p . Since w_+ cannot be a suffix of u_+ , u_p must be a
 951 prefix of u_- , meaning that $(m_v : v)$ must appear as a factor of $(m : u_-)$. As $(m_v : v)$ is a
 952 blocking factor for (m_1, s_1, t_1) , this contradicts the fact that u_- should be read entirely in
 953 the SCC of s_1 . As a result, $u \notin \mathcal{L}(\pi)$.

954 This concludes our induction. \blacktriangleleft

955 The following lemma expresses a sort of converse implication: if a word is far from the
 956 language then it contains many blocking sequences. Let $B = p|\mathcal{A}| + 3|\mathcal{A}|^2$.

957 In the following results we will often use terms like “ $(x : w)$ contains at least N_0 blocking
 958 factors for $s_0, x_0 \rightsquigarrow t_0, y_0, \dots, N_k$ blocking factors for $s_k, x_k \rightsquigarrow t_k, y_k$, in that order, all disjoint”.
 959 This means that we can cut the word $(x : w)$ in k parts $(x : w) = (x_0 : w_0) \dots (x_k : w_k)$,
 960 where for all i we have N_i disjoint blocking factors for $s_i, x_i \rightsquigarrow t_i, y_i$.

961 \blacktriangleright **Lemma 5.16.** *Let $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$ be an SCC-path. If $|w| \geq$
 962 $\max(\frac{6p^2|\mathcal{A}|^2}{\varepsilon}, (k+2)(B+p), \frac{(2k+4)p}{\varepsilon})$ and $+\infty > d(w, \mathcal{L}(\pi)) \geq \varepsilon|w|$ then $(x_0 : w)$ contains
 963 at least $\frac{\varepsilon|w|}{12p^2|\mathcal{A}|^2(k+2)}$ blocking factors for $s_0, x_0 \rightsquigarrow t_0, y_0, \dots, \frac{\varepsilon|w|}{12p^2|\mathcal{A}|^2(k+2)}$ blocking factors for
 964 $s_k, x_k \rightsquigarrow t_k, y_k$, in that order, all disjoint.*

965 **Proof.** We prove this by induction on k using Lemma 4.8. For $k = 0$ we can directly apply
 966 Lemma 4.8, in light of Remark 5.7.

967 Let $k > 0$, suppose the lemma holds for $k - 1$. Since $+\infty > d(w, \mathcal{L}(\pi))$, there is a word
 968 of length $|w|$ in $\mathcal{L}(\pi)$, hence we must have $|w| = y_k - x_0 \pmod{p}$.

969 Our goal is now to cut w in two parts with an intermediate letter, $w = w_- a w_+$, so that:

970 ■ $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \geq \frac{\varepsilon|w|}{2k+4}$, and we can apply Lemma 4.8
 971 ■ $d(w_+, \mathcal{L}(s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \cdots s_k, x_k \rightsquigarrow t_k, y_k)) \geq \frac{(k+1)\varepsilon|w|}{k+2}$ and we can apply the
 972 induction hypothesis

973 To do so, we use an intermediate value argument: We show that w has a short prefix
 974 which is very close to $\mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)$, and a large prefix which is far from it.

975 Then, we use Lemma 5.14, which says that extending a prefix with p letters can only
 976 change the distance to $\mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)$ by p . We then argue that there is an intermediate
 977 prefix w_- which is (roughly) at distance $\frac{\varepsilon|w|}{2k+4}$ from $\mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)$. We split w into
 978 w_-aw_+ , with a single letter. As $d(w, \mathcal{L}(\pi)) \geq \varepsilon|w|$, we infer that w_+ must be at distance at
 979 least $\varepsilon|w| - \frac{\varepsilon|w|}{2k+4}$ from $\mathcal{L}(s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \cdots s_k, x_k \rightsquigarrow t_k, y_k)$, which suffices to conclude.

980 Let us now detail the proof. We define $\pi_+ = s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \cdots s_k, x_k \rightsquigarrow t_k, y_k$.

981 \triangleright Claim 5.17. There is a prefix w' of w such that $|w'| \leq B + p$, $|w'| = y_0 - x_0 \pmod{p}$ and
 982 $d(w', \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \leq B + p$.

983 Proof. Let w' be the prefix of w such that $|w'| = y_0 - x_0 \pmod{p}$ and $p|\mathcal{A}| + 3|\mathcal{A}|^2 \leq |w'| < B +$
 984 p . It exists as $|w| \geq (k+2)(B+p) \geq B+p$. By Lemma 5.13, $d(w_p, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) < +\infty$
 985 and therefore $d(w', \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \leq |w_p| \leq B + p$. \triangleleft

986 \triangleright Claim 5.18. There is a prefix w'' of w such that $|w''| > B + p$, $|w''| = y_0 - x_0 \pmod{p}$
 987 and $d(w'', \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \geq \varepsilon|w| - B - p - 1$.

988 Proof. Let $w = w''aw_s$ such that $p|\mathcal{A}| + 3|\mathcal{A}|^2 \leq |w_s| \leq B + p$ and $|w_s| = y_k - x_1 \pmod{p}$.
 989 This decomposition exists as $|w| \geq (k+2)(B+p) \geq B+p$. We have $|w''| = y_0 - x_0 \pmod{p}$.
 990 Furthermore, as $B \leq |w_+|$, by Lemma 5.13, $d(w_+, \mathcal{L}(\pi_+)) < +\infty$. As a consequence,
 991 $d(w_+, \mathcal{L}(\pi_+)) \leq |w_+| \leq B + p$.

992 As $d(w_+, \mathcal{L}(\pi_+)) \leq |w_+| \leq B+p$ and $+\infty > d(w, \mathcal{L}(\pi)) \geq \varepsilon n$, we must have $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow$
 993 $t_0, y_0)) \geq \varepsilon n - B - p - 1$. \triangleleft

994 \triangleright Claim 5.19. There exist words w_-, w_+ and a letter a such that $w = w_-aw_+$ and
 995 $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \geq \frac{\varepsilon|w|}{2k+4}$ and $d(w_+, \mathcal{L}(\pi_+)) \geq \frac{(k+1)\varepsilon|w|}{k+2}$.

996 Proof. By the two previous claim, w has a prefix w' of length $\geq B$ at distance $\leq B + p$ from
 997 $\mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)$, and a longer prefix w'' at distance $\geq \varepsilon|w| - B - p - 1$ from $\mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)$.
 998 Furthermore, as $|w| \geq 2\frac{B+p+1}{\varepsilon}$, we have $\varepsilon|w| - B - p - 1 \geq \frac{\varepsilon|w|}{k+2}$.

999 In consequence, there must exist w_- a prefix of w and u a word of length p such that
 1000 $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) < \frac{\varepsilon|w|}{k+2} \leq d(w_-u, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0))$.

1001 By Lemma 5.14, we have $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \geq \frac{\varepsilon|w|}{k+2} - p \geq \frac{\varepsilon|w|}{2k+4}$. As $|w| \geq \frac{(2k+4)p}{\varepsilon}$,
 1002 we have $\frac{\varepsilon|w|}{k+2} - p \geq \frac{\varepsilon|w|}{2k+4}$ and thus $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) \geq \frac{\varepsilon|w|}{2k+4}$.

1003 On the other hand, as $d(w_-, \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)) < \frac{\varepsilon|w|}{k+2}$ and $+\infty > d(w, \mathcal{L}(\pi)) \geq \varepsilon|w|$, we
 1004 must have $d(w_+, \mathcal{L}(\pi_+)) \geq \frac{(k+1)\varepsilon|w|}{k+2}$. \triangleleft

1005 By the claim above and Lemma 4.8 we have that w_- contains at least $\frac{\varepsilon|w|}{12p^2|\mathcal{A}|^2(k+2)}$
 1006 blocking factors for $s_0, x_0 \rightsquigarrow t_0, y_0$. On the other hand, by induction hypothesis, w_+ contains
 1007 at least $\frac{\varepsilon|w|}{12p^2|\mathcal{A}|^2(k+2)}$ blocking factors for $s_1, x_1 \rightsquigarrow t_1, y_1, \dots, \frac{\varepsilon|w|}{12p^2|\mathcal{A}|^2(k+2)}$ blocking factors for
 1008 $s_k, x_k \rightsquigarrow t_k, y_k$, in that order, all disjoint.

1009 By combining the two, we obtain the lemma. \blacktriangleleft

1010 A *blocking sequence* for \mathcal{A} is a sequence $((n_1 : u_1), \dots, (n_\ell : u_\ell))$ that is blocking for
 1011 all SCC-paths of \mathcal{A} . As an example, observe that the sequences $(0 : ab), (1 : ab)$ and
 1012 $(0 : aa), (0 : b)$ are both blocking for the automaton displayed in Figure 5 (see Example 5.11).

1013 The goal of the next two lemmas is to show that we can reduce property testing of $\mathcal{L}(\mathcal{A})$
 1014 to a search for blocking sequences in the word:

- 1015 ■ If we find a few blocking sequences in a word then we can answer no as it is not in the
 1016 language (Lemma 5.20).
- 1017 ■ A word that is far from the language contains many blocking sequences (Lemma 5.21).
 1018 Hence if we do not find blocking sequences in the word then it is unlikely to be far from
 1019 the language.

1020 ► **Lemma 5.20.** *If w contains $|\mathcal{A}|$ disjoint blocking sequences for \mathcal{A} then $w \notin \mathcal{L}(\mathcal{A})$.*

1021 **Proof.** Let π be an accepting SCC-path through \mathcal{A} . By definition a blocking sequence for \mathcal{A}
 1022 is a blocking sequence for π . As w contains $|\mathcal{A}|$ disjoint blocking sequences for \mathcal{A} , it contains
 1023 $|\mathcal{A}|$ disjoint blocking sequences for π , hence $w \notin \mathcal{L}(\pi)$ by Lemma 5.15.

1024 As a result, w is not in the language of any accepting SCC-path of \mathcal{A} , and thus not in
 1025 $\mathcal{L}(\mathcal{A})$. ◀

1026 Before going into the next proof, we start by observing that an SCC-path has at most
 1027 $|\mathcal{A}|$ terms, and thus there are at most $(|\mathcal{A}|^2 p^2 |\Sigma| + 1)^{|\mathcal{A}|}$ SCC-paths in \mathcal{A} .

1028 Let $C = (|\mathcal{A}|^2 p^2 |\Sigma| + 1)^{|\mathcal{A}|}$.

1029 ► **Lemma 5.21.** *If $+\infty > d(w, \mathcal{L}(\mathcal{A})) \geq \varepsilon|w|$ and $|w| \geq \max(\frac{6p^2}{\varepsilon}, (k+2)(B+p), \frac{(2k+4)p}{\varepsilon})$
 1030 then w contains $\frac{\varepsilon|w|}{12C|\mathcal{A}|^3 p^2}$ disjoint blocking sequences for \mathcal{A} .*

1031 **Proof.** For each accepting SCC-path π , as $\mathcal{L}(\pi) \subseteq \mathcal{L}(\mathcal{A})$, $d(w, \mathcal{L}(\pi)) \geq d(w, \mathcal{L}(\mathcal{A}))$. Thus, \mathcal{A}
 1032 must have $\frac{\varepsilon|w|}{12|\mathcal{A}|^2 p^2}$ disjoint blocking sequences for π , by Lemma 5.16. It remains to prove that
 1033 $\frac{\varepsilon|w|}{12|\mathcal{A}|^2 p^2}$ disjoint blocking sequences for each π implies $\frac{\varepsilon|w|}{12C|\mathcal{A}|^3 p^2}$ disjoint blocking sequences
 1034 for \mathcal{A} . Given a set of SCC-paths Π , we define $|\Pi|$ as the sum of the lengths of its elements.
 1035 We say that a sequence is blocking for Π if it is blocking for all its elements.

1036 We now prove the following statement by induction on $|\Pi|$: Let Π be a set of SCC-paths
 1037 through \mathcal{A} , and let w be a word with $\frac{\varepsilon|w|}{12|\mathcal{A}|^2 p^2}$ disjoint blocking sequences for each $\pi \in \Pi$.

1038 Then w contains $\frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ disjoint blocking sequences for Π .

1039 The base case is immediate as w contains arbitrarily many disjoint occurrences of the
 1040 empty word, which is a blocking sequence for \emptyset .

1041 Let $w = w_- w_+$ where w_- is the minimal prefix of w containing $\frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ disjoint
 1042 blocking factors for the first element of some $\pi \in \Pi$. That is, $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots$
 1043 $\dots s_k, x_k \rightsquigarrow t_k, y_k$ and w_- contains $\frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ disjoint blocking factors for $s_0, x_0 \rightsquigarrow t_0, y_0$.

1044 Then, by minimality of w_- , w_+ must have $\frac{(|\Pi|-1)\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ many disjoint blocking sequences
 1045 for $\pi' = s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$ and for each $\pi'' \neq \pi$. We can then apply the
 1046 induction hypothesis on w_+ , with $\varepsilon' = \frac{(|\Pi|-1)\varepsilon}{|\Pi|}$ and $\Pi' = \Pi \setminus \{\pi\} \cup \{\pi'\}$: it must contain

1047 $\frac{\varepsilon'|w|}{12|\Pi'|\mathcal{A}|^2 p^2} = \frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ disjoint blocking sequences for Π' .

1048 Appending a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$ in front of any of those blocking sequences
 1049 for Π' yields a blocking sequence for Π . In consequence, we can form $\frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ disjoint
 1050 blocking sequences for Π by matching the $\frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ blocking factors for $s_0, x_0 \rightsquigarrow t_0, y_0$ in

1051 w_- with the $\frac{\varepsilon|w|}{12|\Pi||\mathcal{A}|^2 p^2}$ blocking sequences for Π' in w_+ .

1052 This concludes our induction. To obtain the lemma, we simply apply this property with
 1053 Π the set of accepting SCC-paths of \mathcal{A} . ◀

1054 We define a partial order \trianglelefteq on sequences of positional factors. It is an extension of the
 1055 factor order on blocking factors. It will let us define minimal blocking sequences, with which
 1056 we characterise hard languages.

1057 ► **Definition 5.22.** We have $(n_1 : u_1), \dots, (n_k : u_k) \trianglelefteq (n'_1 : u'_1), \dots, (n'_\ell : u'_\ell)$ when there
 1058 exists a sequence of indices $i_1 \leq i_2 \leq \dots \leq i_k$ such that $(n_{i_j} : u_{i_j})$ is a factor of $(n'_j : u'_j)$ for
 1059 all j .

1060 A blocking sequence $(n_1 : u_1), \dots, (n_k : u_k)$ for \mathcal{A} is **minimal** if it is minimal for \trianglelefteq
 1061 among blocking sequences of \mathcal{A} .

1062 ► **Remark 5.23.** If $\sigma \trianglelefteq \sigma'$ and σ is a blocking sequence for an SCC-path π then σ' is also a
 1063 blocking sequence for π .

1064 The *left effect* of a sequence σ on an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$ is
 1065 the maximal index i such that the sequence is blocking for $s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_i, x_i \rightsquigarrow t_i, y_i$
 1066 (-1 if there is no such i). It is written $(\sigma \gg \pi)$. Similarly, the *right effect* of a sequence on π
 1067 is the minimal index i such that the sequence is blocking for $(m_i, s_i), \dots, (m_k, s_k)$ ($k + 1$ if
 1068 there is no such i). It is written $(\pi \ll \sigma)$.

1069 ► **Remark 5.24.** A sequence σ is blocking for an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow$
 1070 t_k, y_k if and only if $(\sigma \gg \pi) = k$ if and only if $(\pi \ll \sigma) = 0$.

1071 Also, given two sequences σ^l, σ^r , the sequence $\sigma^l \sigma^r$ is blocking for π if and only if
 1072 $(\sigma^l \gg \pi) \geq (\pi \ll \sigma^r) - 1$.

1073 We make the remark that minimal blocking sequences have a bounded number of terms.
 1074 This is because if we build the sequence from left to right by adding terms one by one, the
 1075 minimality implies that at each step the left effect on some SCC-path should increase. As the
 1076 number and lengths of SCC-paths are bounded, so is the number of terms in the sequence.

1077 ► **Lemma 5.25.** A minimal blocking sequence for \mathcal{A} has at most $|Q|(p|Q|)^{2|Q|}$ terms.

1078 **Proof.** The number of SCC-paths in \mathcal{A} is bounded by $(p|Q|)^{2|Q|}$, as each path has at most
 1079 $|Q|$ portals and there are at most $p^2|Q|^2$ portals. Let $\sigma = (n_1 : u_1), \dots, (n_\ell : u_\ell)$ be a minimal
 1080 blocking sequence for \mathcal{A} . For all $i \in \{1, \dots, \ell\}$ we write σ_i for $(n_1 : u_1), \dots, (n_i : u_i)$.

1081 For all $i \in \{1, \dots, \ell - 1\}$ and SCC-path π , we have $(\sigma_i \gg \pi) \leq (\sigma_{i+1} \gg \pi)$. Furthermore,
 1082 for all i there must exist π such that $(\sigma_i \gg \pi) < (\sigma_{i+1} \gg \pi)$: Otherwise we could remove
 1083 $(n_{i+1} : u_{i+1})$ and the sequence would still be blocking for all SCC-paths of \mathcal{A} , contradicting
 1084 the minimality of σ .

1085 As there are at most $(p|Q|)^{2|Q|}$ SCC-paths, each of length at most $|Q|$, ℓ must be at most
 1086 $|Q|(p|Q|)^{2|Q|}$. ◀

1087 We now have all the tools to present the proof that languages recognised by automata
 1088 with bounded minimal blocking sequences are exactly easy languages. Let us start with the
 1089 easier direction.

1090 ► **Lemma 5.26.** If \mathcal{A} has finitely many minimal blocking sequences, then it is easy.

1091 **Proof.** As the length of minimal blocking sequences of \mathcal{A} is bounded, so is the number of
 1092 minimal blocking sequences. Let K be the bound on the length and P the bound on the
 1093 number of minimal blocking sequences.

1094 Let us first sketch the proof before detailing the formulas. We infer from the fact that
 1095 there are boundedly many blocking sequences that if a word w is ε -far from the language of
 1096 \mathcal{A} then it must contain $O(\varepsilon|w|)$ many times the same minimal sequence σ .

1097 Since each positional word in this sequence has length at most K , by sampling $O(\frac{1}{\varepsilon})$
 1098 factors of length K uniformly at random, we can show a positive constant lower bound on
 1099 the probability to find σ . We can repeat this step to obtain a probability $> 1/2$ to find $|\mathcal{A}|$
 1100 times the sequence σ . This proves that $w \notin \mathcal{L}(\mathcal{A})$ by Lemma 5.20.

1101 We now develop the formal proof, starting with a proof that a word that is ε -far from
 1102 $\mathcal{L}(\mathcal{A})$ must contain many times some minimal blocking sequence σ . The next claim shows
 1103 that having many sequences \sqsubseteq -greater than a sequence σ implies having many occurrences
 1104 of σ .

1105 \triangleright Claim 5.27. Let $\sigma = (n_1 : u_1) \cdots (n_k : u_k)$ be a blocking sequence for \mathcal{A} and let $M \in \mathbb{N}$. If
 1106 a positional word $(m : w)$ contains M disjoint blocking sequences for \mathcal{A} that are all greater
 1107 or equal to σ for \sqsubseteq , then $(m : w)$ contains at least $\frac{M}{k}$ occurrences of $(n_1 : u_1)$, ..., $\frac{M}{k}$
 1108 occurrences of $(n_k : u_k)$, in that order, all disjoint.

1109 Proof. We proceed by induction on k . If $k = 0$ the claim is immediate.

1110 Let $k > 0$, suppose the claim holds for sequences of length $k-1$. Suppose $(m : w)$ contains
 1111 M disjoint blocking sequences for \mathcal{A} that are all greater or equal to σ for \sqsubseteq . We can assume
 1112 without loss of generality that all those sequences have $(n_1 : u_1)$ as a factor of their first
 1113 element: If a sequence $\sigma' = (n'_1 : u'_1) \cdots (n'_l : u'_l)$ is such that $\sigma \sqsubseteq \sigma'$ and $(n_1 : u_1)$ is not a
 1114 factor of $(n'_1 : u'_1)$, then we must have $\sigma \sqsubseteq (n'_2 : u'_2) \cdots (n'_l : u'_l)$. Hence we can shorten the
 1115 sequences of timed words until they all have $(n_1 : u_1)$ as a factor of their first term.

1116 Let $(m : w_1)$ be the smallest prefix of w containing $\frac{M}{k}$ occurrences of $(n_1 : u_1)$. Let
 1117 $(m' : w')$ be such that $(m : w) = (m : w_1)(m' : w')$. As $(m : w)$ contains M disjoint blocking
 1118 sequences for \mathcal{A} which all have $(n_1 : u_1)$ as a factor of their first term, we can find at least
 1119 $M - \frac{M}{k}$ of them in $(m' : w')$. As they are all greater or equal to σ , they are also greater
 1120 or equal to $(n_2 : u_2), \dots, (n_k : u_k)$. By induction hypothesis, $(m' : w')$ contains at least
 1121 $\frac{1}{k-1}(M - \frac{M}{k}) = \frac{M}{k}$ disjoint occurrences of $(n_2 : u_2), \dots, (n_k : u_k)$, in that order, all disjoint.

1122 As a result, $(m : w)$ contains at least $\frac{M}{k}$ occurrences of σ . \triangleleft

1123 We can move on to the next step, which is to show that a word that is ε -far from $\mathcal{L}(\mathcal{A})$
 1124 contains many occurrences of some minimal blocking sequence σ .

1125 Let $D = 12C|\mathcal{A}|^4(p|\mathcal{A}|)^{2|\mathcal{A}|}p^2P$.

1126 \triangleright Claim 5.28. If $+\infty > d(w, \mathcal{L}(\mathcal{A})) \geq \varepsilon|w|$ and $|w| \geq \max(\frac{6p^2|\mathcal{A}|^2}{\varepsilon}, (k+2)(B+p), \frac{(2k+4)p}{\varepsilon})$
 1127 then there exists a minimal blocking sequence $\sigma = (n_1 : u_1) \cdots (n_k : u_k)$ for \mathcal{A} such that w
 1128 contains $\frac{\varepsilon|w|}{D}$ occurrences of $(n_1 : u_1)$, ..., $\frac{\varepsilon|w|}{D}$ occurrences of $(n_k : u_k)$, in that order, all
 1129 disjoint.

1130 Proof. We start by applying Lemma 5.21. We obtain that w contains $\frac{\varepsilon|w|}{12C|\mathcal{A}|^3p^2}$ disjoint
 1131 blocking sequences for \mathcal{A} .

1132 Each one of those sequences is greater or equal to a minimal blocking sequence of \mathcal{A} for
 1133 \sqsubseteq . As a result, there exist σ a minimal blocking sequence and $\frac{\varepsilon|w|}{12C|\mathcal{A}|^3p^2P}$ disjoint blocking
 1134 sequences in w that are all greater or equal to σ for \sqsubseteq . Furthermore, by Lemma 5.25, σ
 1135 has at most $|\mathcal{A}|(p|\mathcal{A}|)^{2|\mathcal{A}|}$ terms.

1136 We can then apply Claim 5.27 and obtain that w contains $\frac{\varepsilon|w|}{D}$ disjoint occurrences of
 1137 each term of $(n_1 : u_1)$, ..., $(n_k : u_k)$, in that order, all disjoint. \triangleleft

1138 Given a word of length n , we start by checking that \mathcal{A} accepts a word of length n . If not,
 1139 we reject.

1140 If $|w| \leq \max(\frac{6p^2|\mathcal{A}|^2}{\varepsilon}, (k+2)(B+p), \frac{(2k+4)p}{\varepsilon})$ then we read w entirely, and accept iff it
 1141 is in $\mathcal{L}(\mathcal{A})$.

1142 Otherwise, for each minimal blocking sequence σ , we sample uniformly at random $\frac{D}{\varepsilon}$
 1143 intervals of length K in w . We reject if we find $|\mathcal{A}|$ disjoint occurrences of σ . If we have gone
 1144 through every minimal blocking sequence without rejecting, we accept.

1145 If the word is in $\mathcal{L}(\mathcal{A})$, then by Lemma 5.15 it cannot contain $|Q|$ disjoint blocking
 1146 sequences, hence the algorithm will accept.

1147 If the word is ε -far from $\mathcal{L}(\mathcal{A})$ (but within a finite distance), then by Claim 5.28 there
 1148 exists a minimal blocking sequence $\sigma = (n_1 : u_1) \cdots (n_k : u_k)$ for \mathcal{A} such that w contains
 1149 $\frac{\varepsilon|w|}{D}$ occurrences of $(n_1 : u_1)$, ..., $\frac{\varepsilon|w|}{D}$ occurrences of $(n_k : u_k)$, in that order, all disjoint.
 1150 Recall that by Lemma 5.25, σ has at most $T = |\mathcal{A}|(p|\mathcal{A}|)^{2|\mathcal{A}|}$ terms, hence $k \leq T$. By
 1151 sampling $O(\frac{1}{\varepsilon})$ factors of length K at random, we have a constant positive lower bound
 1152 on the probability of finding $|Q|$ of those occurrences of $(n_i : u_i)$, for any i . From this we
 1153 infer that by sampling $O(\frac{1}{\varepsilon})$ factors of length K at random, we have a constant positive
 1154 lower bound on the probability of finding $|\mathcal{A}|$ occurrences of $(n_i : u_i)$ for each i , and thus $|\mathcal{A}|$
 1155 occurrences of σ .

1156 We can iterate this procedure a constant number of times to obtain a procedure using
 1157 $O(\frac{1}{\varepsilon})$ samples that accepts every word in the language and rejects with probability $> 1/2$
 1158 words that are ε -far from the language. ◀

1159 In order to prove a lower bound on the number of samples necessary to test a language
 1160 with infinitely many minimal blocking sequences, we proceed as follows. We exhibit a portal
 1161 with infinitely many minimal blocking factors $s, x \rightsquigarrow t, y$ and “isolate it” by constructing two
 1162 sequences of timed factors σ^l and σ^r such that for all $(n' : u')$, $\sigma^l(n' : u')\sigma^r$ is blocking for \mathcal{A}
 1163 if and only if $(n' : u')$ is blocking for $s, x \rightsquigarrow t, y$. Then we reduce the problem of testing the
 1164 language of this portal to the problem of testing the language of \mathcal{A} .

1165 For the next proof we define a partial order on portals: $s, x \rightsquigarrow t, y \preceq s', x' \rightsquigarrow t', y'$ if all
 1166 blocking factors of $s', x' \rightsquigarrow t', y'$ are also blocking factors of $s, x \rightsquigarrow t, y$. We write \succeq for
 1167 the reverse relation, \simeq for the equivalence relation $\preceq \cap \succeq$ and $\not\simeq$ for the complement
 1168 relation of \simeq .

1169 Additionally, given an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \cdots s_k, x_k \rightsquigarrow t_k, y_k$ and two
 1170 sequences of positional words σ^l, σ^r , we say that the portal $s_i, x_i \rightsquigarrow t_i, y_i$ *survives* (σ^l, σ^r) if
 1171 $(\sigma^l \gg \pi) < i < (\pi \ll \sigma^r)$.

1172 ▶ **Definition 5.29.** Let $s, x \rightsquigarrow t, y$ be a portal and σ^l and σ^r sequences of positional words.

1173 We define four properties that those objects may have:

1174 **P1** $\sigma^l \sigma^r$ is not blocking for \mathcal{A}

1175 **P2** $s, x \rightsquigarrow t, y$ has infinitely many minimal blocking factors

1176 **P3** for all accepting SCC-path π in \mathcal{A} , every portal in π which survives (σ^l, σ^r) is \simeq -
 1177 equivalent to $s, x \rightsquigarrow t, y$.

1178 ▶ **Lemma 5.30.** If \mathcal{A} has infinitely many minimal blocking sequences, then there exist a
 1179 portal $s, x \rightsquigarrow t, y$ and sequences σ^l and σ^r satisfying properties P1, P2 and P3.

1180 **Proof.** If \mathcal{A} has infinitely many minimal blocking sequences, let $(\sigma_j)_{j \in \mathbb{N}}$ be a family of
 1181 minimal blocking sequences such that the sum of the lengths of the elements of σ_j is at least
 1182 j for all j .

1183 By Lemma 5.25, a minimal blocking sequence has a bounded number of elements. We
 1184 can thus extract from this sequence another one $(\sigma'_j)_{j \in \mathbb{N}}$ such that each σ'_j contains a factor
 1185 of length at least j .

1186 For each j let i_j be the index in σ'_j of a factor of length at least j , and l_j and r_j respectively
 1187 the left effect of the $i_j - 1$ first factors and the right effect of the $k_j - i_j$ last ones, with

1188 k_j the length of σ'_j . As those objects are taken from bounded sets, we can obtain a third
 1189 sequence $(\bar{\sigma}_j)_{j \in \mathbb{N}}$ and α and K such that the i th element of each $\bar{\sigma}_j$ has length at least j
 1190 and the set of components for which it is blocking is K .

1191 For all j let (n_j, u_j) be the i th element of $\bar{\sigma}_j$. Define $\sigma^l = (n_1^l : u_1^l), \dots, (n_k^l : u_k^l)$ and
 1192 $\sigma^r = (n_1^r : u_1^r), \dots, (n_\ell^r : u_\ell^r)$ so that $\bar{\sigma}_1 = \sigma^l(n_1, u_1)\sigma^r$. For all j , $\sigma^l(n_j : u_j)\sigma^r$ is a minimal
 1193 blocking sequence.

1194 We call *surviving portals* the portals that survive (σ^l, σ^r) in at least one SCC-path.

1195 \triangleright Claim 5.31. There exists a surviving portal with infinitely many minimal blocking factors
 1196 that is minimal for \preceq among surviving portals.

1197 Proof. Suppose by contradiction that all \preceq -minimal surviving portals have finitely many
 1198 minimal blocking factors.

1199 For all j , $(n_j : u_j)$ must be blocking for all surviving portals (otherwise $\bar{\sigma}_j$ would not be
 1200 blocking for \mathcal{A}). Hence $(n_j : u_j)$ contains a blocking factor for each \preceq -minimal surviving
 1201 portal. As those factors are bounded while $(n_j : u_j)$ can get arbitrarily large, there exists
 1202 j such that $(n_j : u_j)$ can be split into two non-empty parts $(n_j : u_j^-)(n_j^+ : u_j^+)$ so that
 1203 each \preceq -minimal surviving portal has a minimal blocking factor in either $(n_j : u_j^-)$ or
 1204 $(n_j^+ : u_j^+)$. As a consequence, every surviving portal has a blocking factor in either $(n_j : u_j^-)$
 1205 or $(n_j^+ : u_j^+)$.

1206 Let P be the number of portals of \mathcal{A} . We obtain that $\sigma^l[(n_j : u_j^-)(n_j^+ : u_j^+)]^P\sigma^r$ is a
 1207 blocking sequence for \mathcal{A} , contradicting the minimality of $\sigma^l(n_j : u_j)\sigma^r$ for \preceq . In conclusion,
 1208 there is a \preceq -minimal surviving portal with infinitely many minimal blocking factors. \triangleleft

1209 Let $s, x \rightsquigarrow t, y$ be a \preceq -minimal surviving portal with infinitely many minimal blocking
 1210 factors: It satisfies P2.

1211 The following claim shows that there is a pair of sequences (σ^l, σ^r) such that properties
 1212 P1 and P3 are satisfied.

1213 \triangleright Claim 5.32. There exist σ^l, σ^r such that $\sigma^l\sigma^r$ is not a blocking sequence for \mathcal{A} , and for
 1214 all accepting SCC-path π in \mathcal{A} , every surviving portal in π is \simeq -equivalent to $s, x \rightsquigarrow t, y$.

1215 Proof. We start from the sequences σ^l, σ^r defined before and extend them so that they have
 1216 the desired property.

1217 For each $s', x' \rightsquigarrow t', y' \not\rightsquigarrow s, x \rightsquigarrow t, y$, since $s, x \rightsquigarrow t, y$ is \preceq -minimal we can pick a
 1218 positional word $(n : u)_{s', x' \rightsquigarrow t', y'}$ that is blocking for $s', x' \rightsquigarrow t', y'$ but not for $s, x \rightsquigarrow t, y$.

1219 We extend σ^l and σ^r as follows. While there is a surviving portal $s', x' \rightsquigarrow t', y'$ that is
 1220 not \simeq -equivalent to $s, x \rightsquigarrow t, y$:

- 1221 ■ We pick an SCC-path π such that $s', x' \rightsquigarrow t', y'$ survives in π .
- 1222 ■ Let $i_\ell = (\sigma^l \gg \pi)$ and $i_r = (\pi \ll \sigma^r)$
- 1223 ■ If for all $i \in \{i_\ell + 1, \dots, i_r - 1\}$, $s_i, x_i \rightsquigarrow t_i, y_i \not\rightsquigarrow s, x \rightsquigarrow t, y$ then we append at the
 1224 end of σ^l the sequence $(n : u)_{s_{i_\ell+1}, x_{i_\ell+1} \rightsquigarrow t_{i_\ell+1}, y_{i_\ell+1}}, \dots, (n : u)_{s_{i_r-1}, x_{i_r-1} \rightsquigarrow t_{i_r-1}, y_{i_r-1}}$. The
 1225 sequence $\sigma^l\sigma^r$ is now blocking for π . On the other hand, since we did not add any
 1226 blocking factor for $s, x \rightsquigarrow t, y$, there must still be a surviving portal that is \simeq -equivalent
 1227 to it.
- 1228 ■ If there is an $i \in \{i_\ell + 1, \dots, i_r - 1\}$ such that $s_i, x_i \rightsquigarrow t_i, y_i \simeq s, x \rightsquigarrow t, y$ then let c be the
 1229 maximal index in $\{i_\ell + 1, \dots, i\}$ such that (m_c, s_c, t_c) is not equivalent to $s, x \rightsquigarrow t, y$ for
 1230 \simeq , or i_ℓ if there is no such index. Symmetrically, let d the minimal index in $\{i, \dots, i_r - 1\}$
 1231 such that $(m_d, s_d, t_d) \not\rightsquigarrow s, x \rightsquigarrow t, y$, or i_r if there is no such index. We append at the
 1232 end of σ^l the sequence $(n : u)_{s_{i_\ell+1}, x_{i_\ell+1} \rightsquigarrow t_{i_\ell+1}, y_{i_\ell+1}}, \dots, (n : u)_{m_c, s_c, t_c}$. We append at the

1233 beginning of σ^r the sequence $(n : u)_{s_d, x_d \rightsquigarrow t_d, y_d}, \dots, (n : u)_{s_{i_r-1}, x_{i_r-1} \rightsquigarrow t_{i_r-1}, y_{i_r-1}}$. Now all
 1234 surviving portals in π are \simeq -equivalent to $s, x \rightsquigarrow t, y$, and $s_i, x_i \rightsquigarrow t_i, y_i$ still survives.

1235 We iterate this step until all surviving portals are \simeq -equivalent to $s, x \rightsquigarrow t, y$. We made
 1236 sure that at least one portal was still surviving after each step, hence in the end the sequence
 1237 $\sigma^l \sigma^r$ is not blocking for \mathcal{A} . \triangleleft

1238 \blacktriangleleft

1239 **► Lemma 5.33.** *Let $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell$ be an accepting SCC-path, and*
 1240 *let $i \in \{0, \dots, \ell\}$. Let $\sigma^l = (n_1^l : u_1^l), \dots, (n_k^l : u_k^l)$ a sequence such that $(\sigma^l \gg \pi) < i$ and*
 1241 *$N \in \mathbb{N}$.*

1242 *Then there is a word w^l of length at most $(3|\mathcal{A}|^3 + |\mathcal{A}|)(k+1) + N(2p^2 + p)k|\mathcal{A}| +$*
 1243 *$pN \sum_{i=1}^k |u_i^l|$ such that $|w^l| = x_i - x_0 \pmod{p}$, there is a run reading w^l from s_0 to s_i in \mathcal{A} ,*
 1244 *and $(x_0 : w)$ contains N times $(n_1^l : u_1^l)$, ..., N times $(n_k^l : u_k^l)$ as disjoint factors, in that*
 1245 *order.*

1246 **Proof.** We define w^l by induction on k . As π is accepting, by definition $\mathcal{L}(\pi) \neq \emptyset$, and thus
 1247 for all $j \in \{0, \dots, \ell\}$ there exists a word of length $y_j - x_j \pmod{p}$ labelling a path from
 1248 s_j to t_j . By Fact 3.3, there is such a word v_j of length at most $3|\mathcal{A}|^2$. As a result, for all
 1249 $z \in \{0, \dots, \ell\}$ we can form a word $w_z = v_0 a_1 v_1 \dots a_z$, of length at most $3|\mathcal{A}|^3 + |\mathcal{A}|$, labelling
 1250 a path of length $x_z \pmod{p}$ from q_{init} to s_z in \mathcal{A} . If $k = 0$, we can simply set $w^l = w_i$.

1251 Let $k > 0$, suppose the lemma holds for $k-1$. Let $j = ((n_1 : u_1^l) \gg \pi)$. As $((n_1 : u_1^l) \gg \pi) \leq$
 1252 $(\sigma^l \gg \pi) < i$, we have $j < i$. By definition, $(n_1 : u_1^l)$ is not blocking for $s_{j+1}, x_{j+1} \rightsquigarrow t_{j+1}, y_{j+1}$.
 1253 As a consequence, there is a word v_j labelling a path from s_j to t_j such that $(x_j : v_j)$ has
 1254 $(n_1 : u_1^l)$ as a factor. We can remove cycles of length 0 \pmod{p} in that path, before and
 1255 after reading $(x_j : v_j)$, so we can assume that $|v_j| \leq |u_1^l| + 2p|\mathcal{A}|$. As s_j and t_j are in the
 1256 same SCC, we can extend v_j into a word v'_j of length $\leq |v_j| + |\mathcal{A}| \leq |u_1^l| + (2p+1)|\mathcal{A}|$ that
 1257 labels a cycle from s_j to itself.

1258 Let $\sigma' = (n_2^l : u_2^l), \dots, (n_k^l : u_k^l)$ and $\pi' = s_{j+1}, x_{j+1} \rightsquigarrow t_{j+1}, y_{j+1} \xrightarrow{a_{j+2}} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell$.
 1259 By definition, we have $(\sigma' \gg \pi') = (\sigma^l \gg \pi) < i$. By induction hypothesis, there is a word w'
 1260 of length $\leq (3|\mathcal{A}|^3 + |\mathcal{A}|)k + N(2p^2 + p)(k-1)|\mathcal{A}| + pN \sum_{i=1}^{k-1} |u_i^l|$ such that $|w'| = x_i - x_j$
 1261 \pmod{p} , there is a run reading w' from s_j to s_i in \mathcal{A} , and $(x_j : w)$ contains N times $(n_2^l : u_2^l)$,
 1262 ..., N times $(n_k^l : u_k^l)$ as disjoint factors, in that order.

1263 We set $w^l = w_j(v'_j)^{pN}w'$. This word has length $x_i \pmod{p}$, and at most $|w_j| + pN|v'_j| +$
 1264 $|w'| \leq 3|\mathcal{A}|^3 + |\mathcal{A}| + pN(|u_1^l| + (2p+1)|\mathcal{A}|) + |w'| \leq (3|\mathcal{A}|^3 + |\mathcal{A}|)(k+1) + N(2p^2 + p)k|\mathcal{A}| +$
 1265 $pN \sum_{i=1}^k |u_i^l|$. It labels a path from s_0 to s_i , and contains N times $(n_1^l : u_1^l)$, ..., N times
 1266 $(n_k^l : u_k^l)$ as disjoint factors, in that order. \blacktriangleleft

1267 **► Lemma 5.34.** *Let $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell$ be an accepting SCC-path, and*
 1268 *let $i \in \{0, \dots, \ell\}$. Let $\sigma^r = (n_1^r : u_1^r), \dots, (n_k^r : u_k^r)$ a sequence such that $(\pi \ll \sigma^r) > i$ and*
 1269 *$N \in \mathbb{N}$.*

1270 *Then there is a word w^r of length at most $(3|\mathcal{A}|^3 + |\mathcal{A}|)(k+1) + N(2p^2 + p)k|\mathcal{A}| +$*
 1271 *$pN \sum_{i=1}^k |u_i^r|$ such that $|w^r| = y_\ell - y_i \pmod{p}$, there is a run reading w^r from t_i to t_ℓ in \mathcal{A} ,*
 1272 *and $(y_i : w^r)$ contains N times $(n_1^r : u_1^r)$, ..., N times $(n_k^r : u_k^r)$ as disjoint factors, in that*
 1273 *order.*

1274 **Proof.** By a symmetric proof to the one of the previous lemma. \blacktriangleleft

1275 Given a sequence σ , define $\|\sigma\|$ as the sum of the lengths of the terms of σ .

1276 ► **Lemma 5.35.** *If there exist $s, x \rightsquigarrow t, y$ and σ^l, σ^r satisfying properties P1, P2 and P3*
 1277 *then $\mathcal{L}(\mathcal{A})$ is hard.*

1278 **Proof.** A direct consequence of properties P1 and P3 is that for all $(n' : u')$, $\sigma^l(n' : u')\sigma^r$ is
 1279 blocking for \mathcal{A} if, and only if $(n' : u')$ is blocking for $s, x \rightsquigarrow t, y$.

1280 The proof goes as follows: we show that we can turn an algorithm testing $\mathcal{L}(\mathcal{A})$ with $f(\varepsilon)$
 1281 samples into an algorithm testing $\mathcal{L}(s, x \rightsquigarrow t, y)$ with $f(\varepsilon/X)$ samples with X a constant.
 1282 We then apply Theorem 4.18 to obtain the lower bound.

1283 Consider an algorithm testing $\mathcal{L}(\mathcal{A})$ with $f(\varepsilon)$ samples for some function f . We describe
 1284 an algorithm for testing $\mathcal{L}(s, x \rightsquigarrow t, y)$. Say we are given a threshold ε and a word v of
 1285 length n . First of all we can apply Lemmas 5.33 and 5.34 to compute two words w^l and
 1286 w^r of length at most $E + \varepsilon n F$ for some constants E and F such that we can read w^l from
 1287 q_{init} to s and w^r from t to q_f and w_l contains each element of σ^l at least εn times and w_r
 1288 contains each element of σ^r at least εn times. Let $w = w^l v w^r$. Suppose $|v| \geq \frac{6p^2|\mathcal{A}|^2}{\varepsilon}$ and
 1289 $d(v, \mathcal{L}(s, x \rightsquigarrow t, y)) < +\infty$.

- 1290 ■ If $v \in \mathcal{L}(\mathcal{A})$ then clearly $w \in \mathcal{L}(\mathcal{A})$.
- 1291 ■ If $d(v, \mathcal{L}(s, x \rightsquigarrow t, y)) \geq \varepsilon n$ then by Lemma 4.8 (in light of Remark 5.7), $(x : v)$ contains at
 1292 least $\frac{\varepsilon n}{6p^2|\mathcal{A}|^2}$ blocking factors for $s, x \rightsquigarrow t, y$. Then we have that w contains at least $\frac{\varepsilon n}{6p^2|\mathcal{A}|^2}$
 1293 disjoint blocking sequences for \mathcal{A} . As a result, $d(w, \mathcal{L}(\mathcal{A})) \geq \frac{\varepsilon n}{6p^2|\mathcal{A}|^2}$. We divide this by
 1294 the length of w , which is at most $2E + 2F\varepsilon n + n$. We obtain that $d(w, \mathcal{L}(\mathcal{A})) \geq \frac{\varepsilon}{X}|w|$ for
 1295 some constant X .

1296 Let us now describe the algorithm for testing $\mathcal{L}(s, x \rightsquigarrow t, y)$.

- 1297 ■ If $\mathcal{L}(s, x \rightsquigarrow t, y) \cap \Sigma^n = \emptyset$ then we reject.
- 1298 ■ If $|v| < \frac{6p^2|\mathcal{A}|^2}{\varepsilon}$ then we read v entirely and check that it is in $\mathcal{L}(s, x \rightsquigarrow t, y)$.
- 1299 ■ If $v \in \mathcal{L}(s, x \rightsquigarrow t, y)$ then we apply our algorithm for testing $\mathcal{L}(\mathcal{A})$ on $w = w^l v w^r$ with
 1300 parameter $\varepsilon' = \frac{\varepsilon}{X}$.

1301 The number of samples used on v is at most the number of samples needed on w , hence
 1302 $f(\varepsilon/X)$. We obtain a procedure to test $\mathcal{L}(s, x \rightsquigarrow t, y)$ using $f(\varepsilon/X)$ samples.

1303 By Theorem 4.18, $f(\varepsilon/X) = \Omega(\log(\varepsilon^{-1})/\varepsilon)$, hence $f(\varepsilon) = \Omega(\log(\varepsilon^{-1})/\varepsilon)$. This concludes
 1304 our proof. ◀

1305 ► **Proposition 5.36.** *If \mathcal{A} has infinitely many minimal blocking sequences, then $\mathcal{L}(\mathcal{A})$ is hard.*

1306 **Proof.** We combine Lemmas 5.30 and 5.35. ◀

1307 5.1 Trivial languages

1308 We now characterise trivial languages, as defined in [5]. The definition given there is that
 1309 a language is trivial if for all threshold $\varepsilon > 0$, above a certain length N , every word is at
 1310 distance $\leq \varepsilon|w|$ or $+\infty$ from the language. Hence, on words of length more than N , we do
 1311 not need to sample any letter: we just check if the language contains a word of length $|w|$. If
 1312 not, we answer no. If yes, then we know that w is ε -close to the language and we can answer
 1313 yes.

1314 We present here some other characterisations of this set of languages. They are exactly
 1315 the languages such that there is a bound B such that every word is at distance either $\leq B$
 1316 or $+\infty$ from the language.

1317 They are also the languages that are either finite or described by an automaton with a
 1318 blocking sequence.

1319 ► **Example 5.37.** A representative example of trivial language is $L_1 = a^*ba^*$, the set of
 1320 words containing a b over $\{a, b\}$.

1321 Given any word w , it is at distance at most 1 from L_1 : it suffices to make the first letter
 1322 a b to obtain a word of L_1 .

1323 In consequence, all words of length at least $\frac{1}{\varepsilon}$ are ε -close to the language, which allows us
 1324 to simply answer yes without sampling anything. For words of length $< \frac{1}{\varepsilon}$, we simply read
 1325 the word in full and check if it is in the language.

1326 Now consider the language $L_2 = L_1 \cap (\{a, b\}^2)^*$. It is still trivial, but now we have to take
 1327 into account the parity of the length of the input word: If $|w|$ is odd then $d(w, L_2) = +\infty$
 1328 and we can answer no. If $|w|$ is even then $d(w, L_2) \leq 1$ and we can answer yes as soon as
 1329 $|w| \geq \frac{1}{\varepsilon}$.

1330 ► **Lemma 5.38.** *Let \mathcal{A} a trim NFA. The following are equivalent:*

- 1331 1. *There exists $\varepsilon_0 > 0$, such that for infinitely many n there exist words in $\mathcal{L}(\mathcal{A}) \cap \Sigma^n$ and*
 1332 *there exists $w \in \Sigma^n$ such that $d(w, \mathcal{L}(\mathcal{A})) \geq \varepsilon_0 n$*
- 1333 2. *There exists a family of words $(w_i)_{i \in \mathbb{N}}$ such that for all i , $i \leq d(w_i, \mathcal{L}(\mathcal{A})) < +\infty$*
- 1334 3. *$\mathcal{L}(\mathcal{A})$ is infinite and \mathcal{A} admits a blocking sequence.*
- 1335 4. *$\mathcal{L}(\mathcal{A})$ is infinite and every portal appearing in an accepting SCC-path in \mathcal{A} has a blocking*
 1336 *factor.*

1337 **Proof.** ■ $1 \Rightarrow 2$ is immediate.

1338 ■ $2 \Rightarrow 3$: For all i , $i \leq d(w_i, \mathcal{L}(\mathcal{A})) < +\infty$ implies that $|w_i| \geq i$ and that there exists a
 1339 word $u_i \in \mathcal{L}(\mathcal{A})$ of length $|w_i|$.

1340 It remains to prove that \mathcal{A} has a blocking sequence. We use Lemma 5.21. Fix an arbitrary
 1341 ε , for instance $\varepsilon = 1/2$. Let i be such that $i \geq \max(\frac{6p^2}{\varepsilon}, (k+2)(B+p), \frac{(2k+4)p}{\varepsilon})$ and
 1342 $i > \frac{12C|\mathcal{A}|p^2}{\varepsilon}$.

1343 Then as $i \leq d(w_i, \mathcal{L}(\mathcal{A})) < +\infty$, we can apply Lemma 5.21 and obtain that w_i contains
 1344 $\frac{\varepsilon|w_i|}{12C|\mathcal{A}|p^2} > 1$ blocking sequences for \mathcal{A} . In particular, \mathcal{A} has a blocking sequence.

1345 ■ $3 \Rightarrow 1$: Let $\sigma = (n_1 : u_1), \dots, (n_k : u_k)$ be a blocking sequence for \mathcal{A} . As \mathcal{A} is infinite,
 1346 there exists an SCC-path π in \mathcal{A} and $w \in \mathcal{L}(\pi)$ with $|w| \geq |\mathcal{A}|$. By Lemma 5.13, for all
 1347 $\ell \geq p|\mathcal{A}| + 3|\mathcal{A}|^3$ such that $\ell = |w| \pmod{p}$ there exists $w' \in \mathcal{L}(\pi)$ with $|w'| = \ell$.

1348 For all $i \in \{1, \dots, k\}$ we define v_i as a word of length $\leq u_i + 2p$ such that $(0 : v_i)$ has
 1349 $(n_i : u_i)$ as a factor. For all $N \in \mathbb{N}$, we can then define the word $w_N = v_1^N \dots v_k^N a^{|w|}$ with
 1350 a an arbitrary letter. As it is of length $|w| \pmod{p}$, there is a word of the same length
 1351 in $\mathcal{L}(\mathcal{A})$. On the other hand, it contains N disjoint occurrences of σ , which is a blocking
 1352 sequence for \mathcal{A} . Let $\varepsilon_0 = \frac{1}{|u_1| + |u_2| + \dots + |u_k| + 2kp + |w|}$. We have $\varepsilon_0 |w_N| \leq N \leq d(w_N, \mathcal{L}(\mathcal{A}))$.

1353 ■ $3 \Rightarrow 4$: If \mathcal{A} has a blocking sequence, then every portal in \mathcal{A} appearing in an accepting
 1354 SCC-path has to have a blocking factor in that sequence.

1355 ■ $4 \Rightarrow 3$: If $\mathcal{L}(\mathcal{A})$ is infinite and every portal appearing in an accepting SCC-path in \mathcal{A} has
 1356 a blocking factor, then we can construct a blocking sequence for \mathcal{A} as follows. Let P be
 1357 the number of those portals in \mathcal{A} . Let σ be a sequence containing a blocking factor for
 1358 each of those portals. The sequence σ^P is blocking for \mathcal{A} .

1359 ◀

1360 This concludes the proof of Theorem 5.1.

1361 6 Hardness of classifying

1362 In the previous sections, we have shown that testing some regular languages (*easy* ones)
 1363 that requires fewer queries than testing others (*hard* ones). Therefore, given the task of

1364 testing a word for membership in $\mathcal{L}(\mathcal{A})$, it is natural to first try to determine if the language
 1365 of \mathcal{A} is easy, and if this is the case, run the appropriate ε -tester, that uses fewer queries.
 1366 In this section, we investigate the computational complexity of checking which class of the
 1367 trichotomy the language of a given automaton belongs to. We formalize this question as the
 1368 following decision problems:

1369 We show that, unfortunately, our combinatorial characterization based on minimal
 1370 blocking sequences does not lead to efficient algorithms: both problems are PSPACE-complete.

1371 ► **Theorem 6.1.** *The triviality and easiness problems are both PSPACE-complete, even for*
 1372 *strongly connected NFAs.*

1373 In the following section we show the PSPACE upper bounds on both problems (Proposi-
 1374 tions 6.8 and 6.9).

1375 6.1 A PSPACE upper-bound on classifying automata

1376 Let us first provide another characterisation of hard automata.

1377 ► **Lemma 6.2.** *Let $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell$ be an SCC-path, i an index,*
 1378 *Π a set of SCC-paths and $(\sigma_{\pi'})_{\pi' \in \Pi}$ a family of sequences of positional words such that*
 1379 *$(\sigma_{\pi'} \gg \pi') < i$ for all π' .*

1380 *There exists a sequence of positional words σ such that:*

- 1381 ■ $(\sigma \gg \pi) < i$
- 1382 ■ $(\sigma_{\pi'} \gg \pi') \leq (\sigma \gg \pi')$ for all $\pi' \in \Pi$.

1383 **Proof.** We prove this by induction on the sum of the lengths of the elements of Π . If Π is
 1384 empty then we can set σ as the empty sequence.

1385 If not, let π_{min} be such that the first term of $\sigma_{\pi_{min}}$ has the least left effect on π . Let
 1386 $\sigma_{\pi_{min}} = (n_1 : u_1), \dots, (n_k : u_k)$ and $\pi_{min} = s'_0, x'_0 \rightsquigarrow t'_0, y'_0 \xrightarrow{a_1} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$. Let
 1387 $j = ((n_1 : u_1) \gg \pi_{min})$ and $r = ((n_1 : u_1) \gg \pi)$.

1388 Let $\pi' = s'_{j+1}, x'_{j+1} \rightsquigarrow t'_{j+1}, y'_{j+1} \xrightarrow{a_1} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$. Define $\Pi' = \Pi \setminus \{\pi_{min}\} \cup \{\pi'\}$ if
 1389 $j < \ell$ and $\Pi' = \Pi \setminus \{\pi_{min}\}$ otherwise. In the first case the sequence associated with π' is
 1390 $\sigma_{\pi'} = (n_2 : u_2), \dots, (n_k : u_k)$.

1391 ▷ **Claim 6.3.** For all $\bar{\pi} \in \Pi \setminus \{\pi_{min}\}$, we have $(\sigma_{\bar{\pi}} \gg \pi) = r + (\sigma_{\bar{\pi}} \gg s_{r+1}, x_{r+1} \rightsquigarrow$
 1392 $t_{r+1}, y_{r+1} \xrightarrow{a_{r+2}} \dots s_k, x_k \rightsquigarrow t_k, y_k)$

1393 **Proof.** Since the first term of $\sigma_{\pi'}$ was the one with the least left effect on π , the first term of
 1394 every other sequence has a left effect at least r on it.

1395 Let $\bar{\pi} \in \Pi \setminus \{\pi_{min}\}$, let $\sigma_{\bar{\pi}} = (\bar{n}_1 : \bar{u}_1), \dots, (\bar{n}_m : \bar{u}_m)$. Let $z = ((\bar{n}_1 : \bar{u}_1) \gg \pi)$. This
 1396 means $(\bar{n}_1 : \bar{u}_1)$ is not a blocking factor for $s_{z+1}, x_{z+1} \rightsquigarrow t_{z+1}, y_{z+1}$.

1397 We have $(\sigma_{\bar{\pi}} \gg \pi) = z + ((\bar{n}_2 : \bar{u}_2), \dots, (\bar{n}_m : \bar{u}_m) \gg s_{z+1}, x_{z+1} \rightsquigarrow t_{z+1}, y_{z+1})$ and
 1398 $(\sigma_{\bar{\pi}} \gg s_{r+1}, x_{r+1} \rightsquigarrow t_{r+1}, y_{r+1} \xrightarrow{a_{r+2}} \dots) = z - r + ((\bar{n}_2 : \bar{u}_2), \dots, (\bar{n}_m : \bar{u}_m) \gg s_{z+1}, x_{z+1} \rightsquigarrow$
 1399 $t_{z+1}, y_{z+1}) = (\sigma_{\bar{\pi}} \gg \pi) - r$.

1400 ◁

1401 As a consequence of this claim, we have that $(\sigma_{\bar{\pi}} \gg s_{r+1}, x_{r+1} \rightsquigarrow t_{r+1}, y_{r+1} \xrightarrow{a_{r+2}}$
 1402 $\dots s_k, x_k \rightsquigarrow t_k, y_k) < i - r$ for all $\bar{\pi} \in \Pi \setminus \{\pi'\}$.

1403 By induction hypothesis, we obtain a sequence σ' such that

- 1404 ■ $(\sigma' \gg s_{r+1}, x_{r+1} \rightsquigarrow t_{r+1}, y_{r+1} \xrightarrow{a_1} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell) < i - r$
- 1405 ■ $(\sigma_{\pi'} \gg \pi') \leq (\sigma' \gg \pi')$ for all $\pi' \in \Pi'$.

1406 The sequence $(n_1 : u_1), \sigma'$ satisfies both conditions of the lemma. \blacktriangleleft

1407 **► Lemma 6.4.** *An automaton \mathcal{A} is hard if and only if there exists an accepting SCC-path π*
 1408 *containing a portal $s, x \rightsquigarrow t, y$ such that:*

- 1409 \blacksquare *$s, x \rightsquigarrow t, y$ has infinitely many minimal blocking factors.*
- 1410 \blacksquare *For all accepting SCC-path π' there exist sequences σ^l, σ^r such that:*
 - 1411 \blacksquare *$s, x \rightsquigarrow t, y$ survives (σ^l, σ^r) in π*
 - 1412 \blacksquare *All portals surviving (σ^l, σ^r) in π' are \simeq -equivalent to $s, x \rightsquigarrow t, y$*

1413 **Proof.** Let us start with the left-to-right direction. If \mathcal{A} is hard then by Lemma 5.26 it
 1414 has infinitely many minimal blocking sequences. Then by Lemma 5.30 we have a portal
 1415 $s, x \rightsquigarrow t, y$ and sequences σ^l, σ^r satisfying properties P1, P2 and P3.

1416 By P1, $\sigma^l \sigma^r$ is not blocking for \mathcal{A} , thus there exists an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1}$
 1417 $\dots s_k, x_k \rightsquigarrow t_k, y_k$ and an index i such that $(\sigma^l \gg \pi) < i < (\pi \ll \sigma^r)$.

1418 As a consequence, we have $s_i, x_i \rightsquigarrow t_i, y_i \simeq s, x \rightsquigarrow t, y$, by P3. We can assume without
 1419 loss of generality that $s_i, x_i \rightsquigarrow t_i, y_i = s, x \rightsquigarrow t, y$. As a result, for all accepting SCC-path
 1420 π' we have that $s, x \rightsquigarrow t, y$ survives (σ^l, σ^r) in π and all portals surviving (σ^l, σ^r) in π' are
 1421 \simeq -equivalent to $s, x \rightsquigarrow t, y$ (we use the same pair (σ^l, σ^r) for all π').

1422 Let us now prove the other direction. Suppose we have π and $s, x \rightsquigarrow t, y$ satisfying the
 1423 conditions of the lemma. We only need to construct two sequences σ^l, σ^r such that properties
 1424 P1 and P3 are satisfied. The result follows by Lemma 5.35.

1425 let Π be the set of accepting SCC-paths in \mathcal{A} . Consider families of sequences $(\sigma_{\pi'}^l)_{\pi' \in \Pi}$
 1426 and $(\sigma_{\pi'}^r)_{\pi' \in \Pi}$ such that for all $\pi' \in \Pi$:

- 1427 \blacksquare *$s, x \rightsquigarrow t, y$ survives $(\sigma_{\pi'}^l, \sigma_{\pi'}^r)$ in π*
- 1428 \blacksquare *All portals surviving $(\sigma_{\pi'}^l, \sigma_{\pi'}^r)$ in π' are \simeq -equivalent to $s, x \rightsquigarrow t, y$*

1429 Let i be the index of $s, x \rightsquigarrow t, y$ in π . By Lemma 6.2 we can build a sequence σ^l such that

- 1430 \blacksquare $(\sigma^l \gg \pi) < i$
- 1431 \blacksquare $(\sigma_{\pi'}^l \gg \pi') \leq (\sigma^l \gg \pi')$ for all $\pi' \in \Pi$.

1432 Using a symmetric argument, we build a sequence σ^r such that

- 1433 \blacksquare $i < (\pi \ll \sigma^r)$
- 1434 \blacksquare $(\pi' \ll \sigma_{\pi'}^r) \geq (\pi' \ll \sigma^r)$ for all $\pi' \in \Pi$.

1435 As a consequence, for all accepting SCC-path $\pi' \in \Pi$, all portals surviving (σ^l, σ^r) in π'
 1436 are \simeq -equivalent to $s, x \rightsquigarrow t, y$. Furthermore, $s, x \rightsquigarrow t, y$ survives (σ^l, σ^r) in π .

1437 We have shown that $s, x \rightsquigarrow t, y$ and (σ^l, σ^r) satisfy properties P1 and P3. P2 is immediate
 1438 by assumption. We simply apply Lemma 5.35 to obtain the result. \blacktriangleleft

1439 Next, we establish that the items listed in the previous lemma can all be checked in
 1440 polynomial space in $|\mathcal{A}|$.

1441 **► Lemma 6.5.** *Given a portal $s, x \rightsquigarrow t, y$, we can check whether it has infinitely many*
 1442 *minimal blocking factors in polynomial space in $|\mathcal{A}|$.*

1443 **Proof.** We start by defining a deterministic automaton \mathcal{B} recognising the set of positional
 1444 words that are factors of $\mathcal{PL}(s, x \rightsquigarrow t, y)$.

1445 For each $i \in \{0, \dots, p-1\}$ let Q_i be the set of states in the SCC of s that can be reached
 1446 in $i-x$ steps from s . It is easily computable using the partition of the states given by
 1447 Fact 3.3.

1448 Let \mathcal{A}_i be \mathcal{A} where the initial states are Q_i and every state in the SCC of s is final. It
 1449 recognises words that can be read from Q_i in \mathcal{A} without leaving the SCC.

1450 Then, we define \mathcal{B}_i as the automaton obtained by determinising \mathcal{A}_i . It has size at most
 1451 $2^{|\mathcal{A}|}$. From \mathcal{B}_i we easily obtain an automaton \mathcal{B}'_i of size $p|\mathcal{B}_i|$ recognising the set of positional
 1452 words $\{(i : w) \mid w \in \mathcal{L}(\mathcal{B}_i)\}$: we simply keep track in the states of the number of letters read,
 1453 plus i , modulo p .

1454 Lastly, we define \mathcal{B} as follows: We take all automata \mathcal{B}_i and merge their initial states into
 1455 one. Observe that \mathcal{B} is deterministic as all \mathcal{B}_i are, and for all letter (j, a) there is at most one
 1456 transition from the initial state reading (j, a) , which goes to a state of \mathcal{B}_j . This automaton
 1457 is of at most exponential size in $|\mathcal{A}|$. It recognises the set of positional words that are factors
 1458 of $\mathcal{PL}(s, x \rightsquigarrow t, y)$.

1459 We can complement it to obtain an automaton $\overline{\mathcal{B}}$ recognising the complement language,
 1460 i.e., the set of positional words that are not factors of $\mathcal{PL}(s, x \rightsquigarrow t, y)$. We have $|\overline{\mathcal{B}}| \leq |\mathcal{B}| + 1$.

1461 A positional word $(n : w)$ is a *minimal blocking factor* of $s, x \rightsquigarrow t, y$ if and only if it is
 1462 not a factor of $\mathcal{PL}(s, x \rightsquigarrow t, y)$ while removing its first or its last letter makes it a factor of
 1463 $\mathcal{PL}(s, x \rightsquigarrow t, y)$.

1464 The set of blocking factors can thus be recognised by an automaton of size $|\overline{\mathcal{B}}|^3$, which
 1465 runs $\overline{\mathcal{B}}$ on the input word, while running \mathcal{B} from the second to the last letter and from the
 1466 first to the second to last letter. The automaton accepts if all three runs are accepting. It is
 1467 of exponential size in $|\mathcal{A}|$.

1468 We simply need to check if this automaton has an infinite language, which is the case if and
 1469 only if it has a cycle reachable from the initial state and from which a final state is reachable.
 1470 This can be checked by exploring the state space of the automaton, in non-deterministic
 1471 polynomial space (in $|\mathcal{A}|$), and applying Savitch's theorem. ◀

1472 ▶ **Lemma 6.6.** *Given two SCC-paths π and π' , one can check in PSPACE whether there is a*
 1473 *sequence σ that is blocking for π and not π' .*

1474 **Proof.**

1475 ▷ **Claim 6.7.** There is a sequence σ that is blocking for $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow$
 1476 t_k, y_k and not $\pi' = s'_0, x'_0 \rightsquigarrow t'_0, y'_0 \xrightarrow{a'_1} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$ if and only if either:

- 1477 ■ there is a positional word $(n : w)$ that is a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$ and
 1478 not $s'_0, x'_0 \rightsquigarrow t'_0, y'_0$ and there is a sequence σ' that is blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2}$
 1479 $\dots s_k, x_k \rightsquigarrow t_k, y_k$ and not π' ,
- 1480 ■ or there is a positional word $(n : w)$ that is a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$ and $s'_0, x'_0 \rightsquigarrow$
 1481 t'_0, y'_0 and there is a sequence σ' that is blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$
 1482 and not $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$.

1483 **Proof.** The right-to-left direction is clear (just take $\sigma = (n : w), \sigma'$ in both cases).

1484 For the left-to-right direction, consider a sequence σ that is blocking for π and not π' , of
 1485 minimal length. Let σ_+ and $(n : w)$ be such that $\sigma = (n : w)\sigma_+$.

- 1486 ■ If $(n : w)$ is not blocking for $s_0, x_0 \rightsquigarrow t_0, y_0$ then σ_+ is blocking for π and not π' ,
 1487 contradicting the minimality of σ .
- 1488 ■ If $(n : w)$ is blocking for $s_0, x_0 \rightsquigarrow t_0, y_0$ and not $s'_0, x'_0 \rightsquigarrow t'_0, y'_0$ then we set $\sigma' = \sigma$. We
 1489 know that σ is not blocking for π' . On the other hand, as σ is blocking for π , it is also
 1490 blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$.
- 1491 ■ If $(n : w)$ is blocking for both $s_0, x_0 \rightsquigarrow t_0, y_0$ and $s'_0, x'_0 \rightsquigarrow t'_0, y'_0$ then we set $\sigma' = \sigma$.
 1492 As σ is blocking for π , it is also blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$. On
 1493 the other hand, if σ was blocking for $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$, then it would

1494 also be blocking for π' , a contradiction. Hence σ is not blocking for $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2}$
 1495 $\cdots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$
 1496 ◁

1497 The claim above lets us define a recursive algorithm.

- 1498 ■ First check if there is a positional word $(n : w)$ that is blocking for $s_0, x_0 \rightsquigarrow t_0, y_0$ and
 1499 not $s'_0, x'_0 \rightsquigarrow t'_0, y'_0$. If it is the case, make a recursive call to check if there is a sequence
 1500 σ' that is blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \cdots s_k, x_k \rightsquigarrow t_k, y_k$ and not π' . If it is the case,
 1501 answer yes.
- 1502 ■ Then check if there is a positional word $(n : w)$ that is a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$
 1503 and $s'_0, x'_0 \rightsquigarrow t'_0, y'_0$. If so, make a recursive call to check if there is a sequence σ' that is
 1504 blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \cdots s_k, x_k \rightsquigarrow t_k, y_k$ and not $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2} \cdots s'_\ell, x'_\ell \rightsquigarrow$
 1505 t'_ℓ, y'_ℓ . If it is the case, answer yes.

1506 If both items fail, answer no.

1507 The existence of those positional words can be checked in polynomial space using the
 1508 automaton \mathcal{B} constructed in the proof of Lemma 6.5. The depth of the recursive calls is at
 1509 most the sum of the lengths of π and π' , which is bounded by $2|\mathcal{A}|$. In consequence, this
 1510 algorithm runs in polynomial space.

1511 ◀

1512 ► **Proposition 6.8.** *The following problem is in PSPACE: Given an automaton \mathcal{A} , is it hard?*

1513 **Proof.** We use Lemma 6.4. We guess an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \cdots s_k, x_k \rightsquigarrow t_k, y_k$
 1514 and an index i .

1515 We check that $s_i, x_i \rightsquigarrow t_i, y_i$ has infinitely many minimal blocking factors, using Lemma 6.5.

1516 We then enumerate all SCC-paths in \mathcal{A} . For each one $\pi' = s'_0, x'_0 \rightsquigarrow t'_0, y'_0 \xrightarrow{a'_1} \cdots s'_\ell, x'_\ell \rightsquigarrow$
 1517 t'_ℓ, y'_ℓ we guess indices j^l and j^r . We check that every portal $s'_j, x'_j \rightsquigarrow t'_j, y'_j$ with $j^l < j < j^r$
 1518 is \simeq -equivalent to $s, x \rightsquigarrow t, y$.

1519 Then, we use Lemma 6.6 to check that there is a sequence σ^l that is blocking for
 1520 $s'_0, x'_0 \rightsquigarrow t'_0, y'_0 \xrightarrow{a'_1} \cdots s'_{j^l}, x'_{j^l} \rightsquigarrow t'_{j^l}, y'_{j^l}$ and not $s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \cdots s_i, x_i \rightsquigarrow t_i, y_i$.

1521 Symmetrically, we check that there is a sequence σ^r that is blocking for $s'_{j^r}, x'_{j^r} \rightsquigarrow$
 1522 $t'_{j^r}, y'_{j^r} \xrightarrow{a'_1} \cdots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$ and not $s_i, x_i \rightsquigarrow t_i, y_i \xrightarrow{a_{i+1}} \cdots s_k, x_k \rightsquigarrow t_k, y_k$.

1523 If all those tests succeed, we answer yes, otherwise we answer no. This algorithm is
 1524 correct and complete by Lemma 6.4. ◀

1525 Our last result is the PSPACE upper bound on the complexity of checking if a language
 1526 is trivial. It is based on the characterisation of trivial languages given by Lemma 5.38.

1527 ► **Proposition 6.9.** *One can check if an automaton has a trivial language in PSPACE.*

1528 **Proof.** By Lemma 5.38, it suffices to enumerate all accepting SCC-paths in the automaton,
 1529 and then check that all portals appearing in them have a blocking factor. This is feasible in
 1530 PSPACE, using the automaton $\bar{\mathcal{B}}$ from the proof of Lemma 6.5. ◀

1531 **6.2 Hardness of classifying automata**

1532 We prove hardness of the triviality problem and easiness problems, concluding on their
 1533 PSPACE-completeness. We reduce from the universality problem for NFAs, which is well-
 1534 known to be PSPACE-complete (see e.g. [1, Theorem 10.14]).

1535 ► **Lemma 6.10.** *The triviality problem is PSPACE-hard.*

1536 **Proof.** Consider an NFA $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ on an alphabet Σ . Without loss of generality,
 1537 we assume that \mathcal{A} is trim (up to removing unreachable or non-co-reachable states) and that
 1538 it accepts all words of length less than 2: this can be checked in polynomial time. Let $\#$ and
 1539 $!$ be two letters that are not in Σ . We apply the following transformations to \mathcal{A} :

- 1540 ■ add a transition labeled by $!$ from every final state to the initial state q_0
- 1541 ■ add a self-loop labeled by $\#$ to each state.

1542 We call the resulting automaton $\mathcal{B} = (Q, \Sigma \cup \{!, \#\}, \delta', q_0, F)$. Note that \mathcal{B} is strongly
 1543 connected: consider any two states $q, q' \in Q$, we show that q' is reachable from q . As \mathcal{A} is
 1544 trim, there exists $q_f \in F$ that is reachable from q , and q' is reachable from the initial state
 1545 q_0 . Furthermore, we have put a $!$ transition from q_f to q_0 , hence q' is reachable from q .

1546 Recall that the language of a strongly connected automaton is trivial if and only if it
 1547 has no minimal blocking factor. Hence, to complete this reduction, we need to show that
 1548 $\text{MBF}(\mathcal{B})$ is empty if and only if \mathcal{A} is universal.

1549 First, let us describe the language recognized by \mathcal{B} . It is given by

$$1550 \mathcal{L}(\mathcal{B}) = \{u_1!u_2!\cdots!u_n \mid \forall i, u_i \in (\Sigma \cup \{\#\})^* \wedge \pi_\Sigma(u_i) \in \mathcal{L}(\mathcal{A})\},$$

1551 where $\pi_\Sigma(u)$ is the word in Σ^* obtained by removing all letters not in Σ from u .

1552 ▷ **Claim 6.11.** If \mathcal{A} is universal, then \mathcal{B} is also universal.

1553 **Proof.** Indeed, any word u in $\mathcal{L}(\mathcal{B})$ can be uniquely decomposed into $u = u_1!u_2!\cdots!u_n$ where
 1554 each u_i does not contain the letter “!”. As $\#$ is idempotent on \mathcal{B} , $\delta'(q_0, u_i)$ is equal to
 1555 $\delta(q_0, \pi_\Sigma(u_i))$ for every i . Since \mathcal{A} is universal, each of the $\delta'(q_0, u_i)$ contains a final state,
 1556 hence $\delta'(q_0, u_i) = \{q_0\}$. Therefore, the set $\delta'(q_0, u)$ is equal to $\delta'(q_0, u_n)$, which contains a
 1557 final state, and u is in $\mathcal{L}(\mathcal{B})$, which shows that \mathcal{B} is universal. ◀

1558 This shows that if \mathcal{A} is universal, then $\text{MBF}(\mathcal{B})$ is empty.

1559 For the converse, we show that a word $w \in \Sigma^*$ not in $\mathcal{L}(\mathcal{A})$ induces minimum blocking
 1560 factors for \mathcal{B} . Consider such a w of minimal size. As we assumed that \mathcal{A} accepts all words of
 1561 size less than 2, $|w| \geq 2$. Let u, v be words of length at least 1 such that $w = uv$. For all
 1562 $n \in \mathbb{N}$, at least one of $u\#^n v, !u\#^n v, u\#^n v!, !u\#^n v!$ is a minimal blocking factor (depending
 1563 respectively on whether w is not a factor of any word of $\mathcal{L}(\mathcal{A})$ or is a prefix/suffix of a word
 1564 of $\mathcal{L}(\mathcal{A})$ or not). As a consequence, \mathcal{B} has infinitely many blocking factors, and is thus hard
 1565 to test by Theorem 4.2.

1566 In summary, \mathcal{A} is universal if and only if \mathcal{B} is trivial to test. This shows the PSPACE-
 1567 hardness of the triviality problem. ◀

1568 The above proof can be extended to show the PSPACE-hardness of the easiness problem.

1569 ► **Corollary 6.12.** *The easiness problem is PSPACE-hard.*

1570 **Proof.** We proceed as in the proof of Lemma 6.10: given an automaton \mathcal{A} over an alphabet
 1571 Σ , we build an automaton \mathcal{B} over the alphabet $\Sigma \cup \{!, \#\}$ such that if \mathcal{A} is universal, $\text{MBF}(\mathcal{B})$
 1572 is empty, and if \mathcal{A} is not universal, then $\text{MBF}(\mathcal{B})$ is infinite.

1573 To show the hardness of the easiness problem, let b denote a new letter not in $\Sigma \cup \{!, \#\}$
 1574 and consider the automaton \mathcal{B}' equal to \mathcal{B} but taken over the alphabet $\Sigma \cup \{!, \#, b\}$. As
 1575 there are no transitions labeled by b in \mathcal{B}' , the word b is always a minimum blocking factor
 1576 of \mathcal{B}' . As a result, we have $\text{MBF}(\mathcal{B}') = \text{MBF}(\mathcal{B}) \cup \{b\}$, hence \mathcal{A} is universal if and only if
 1577 $\text{MBF}(\mathcal{B}')$ is finite but non-empty: by Theorem 4.2, this is equivalent to $\mathcal{L}(\mathcal{B}')$ is easy to test.
 1578 Therefore, the easiness problem is also PSPACE-hard. ◀

1579 This concludes the proof of Theorem 6.1

1580 — References —

- 1581 1 Alfred V. Aho and John E. Hopcroft. *The Design and Analysis of Computer Algorithms*.
 1582 Addison-Wesley Longman Publishing Co., Inc., 1974.
- 1583 2 Maryam Aliakbarpour, Ilias Diakonikolas, and Ronitt Rubinfeld. Differentially private identity
 1584 and equivalence testing of discrete distributions. In *International Conference on Machine
 1585 Learning, ICML, 2018*. URL: <https://proceedings.mlr.press/v80/aliakbarpour18a.html>.
- 1586 3 Noga Alon, Richard A. Duke, Hanno Lefmann, Vojtech Rödl, and Raphael Yuster. The
 1587 algorithmic aspects of the regularity lemma. *Journal of Algorithms*, 16(1):80–109, 1994.
 1588 doi:10.1006/JAGM.1994.1005.
- 1589 4 Noga Alon, Eldar Fischer, Michael Krivelevich, and Mario Szegedy. Efficient testing of large
 1590 graphs. *Combinatorica*, 20(4):451–476, 2000. doi:10.1007/s004930070001.
- 1591 5 Noga Alon, Michael Krivelevich, Ilan Newman, and Mario Szegedy. Regular languages are
 1592 testable with a constant number of queries. *SIAM Journal on Computing*, 30(6):1842–1862,
 1593 2001. doi:10.1109/SFFCS.1999.814641.
- 1594 6 Noga Alon and Asaf Shapira. Every monotone graph property is testable. *SIAM Journal of
 1595 Computing*, 38(2):505–522, 2008. doi:10.1137/050633445.
- 1596 7 Gabriel Bathie and Tatiana Starikovskaya. Property testing of regular languages with ap-
 1597 plications to streaming property testing of visibly pushdown languages. In *International
 1598 Colloquium on Automata, Languages, and Programming, ICALP, 2021*. doi:10.4230/LIPIcs.
 1599 ICALP.2021.119.
- 1600 8 Manuel Blum, Michael Luby, and Ronitt Rubinfeld. Self-testing/correcting with applications
 1601 to numerical problems. *Journal of Computer and System Sciences*, 47(3):549–595, 1993.
 1602 doi:10.1016/0022-0000(93)90044-W.
- 1603 9 Ilias Diakonikolas and Daniel M Kane. A new approach for testing properties of discrete
 1604 distributions. In *IEEE Symposium on Foundations of Computer Science, FOCS*, pages 685–694.
 1605 IEEE, 2016. doi:10.1109/FOCS.2016.78.
- 1606 10 Yevgeniy Dodis, Oded Goldreich, Eric Lehman, Sofya Raskhodnikova, Dana Ron, and Alex
 1607 Samorodnitsky. Improved testing algorithms for monotonicity. In *International Workshop
 1608 on Randomization and Approximation Techniques in Computer Science, RANDOM*, pages
 1609 97–108. Springer, 1999. doi:10.1007/978-3-540-48413-4_10.
- 1610 11 Funda Ergün, Sampath Kannan, S.Ravi Kumar, Ronitt Rubinfeld, and Mahesh Viswanathan.
 1611 Spot-checkers. *Journal of Computer and System Sciences*, 60(3):717–751, 2000. doi:10.1006/
 1612 jcsc.1999.1692.
- 1613 12 Nathanaël François, Frédéric Magniez, Michel de Rougemont, and Olivier Serre. Streaming
 1614 Property Testing of Visibly Pushdown Languages. In *European Symposium on Algorithms, ESA*,
 1615 volume 57, pages 43:1–43:17, Dagstuhl, Germany, 2016. Schloss Dagstuhl – Leibniz-Zentrum
 1616 für Informatik. doi:10.4230/LIPIcs.ESA.2016.43.
- 1617 13 Oded Goldreich. *Introduction to property testing*. Cambridge University Press, 2017. doi:
 1618 10.1017/9781108135252.

- 1619 **14** Oded Goldreich, Shari Goldwasser, and Dana Ron. Property testing and its connection to
1620 learning and approximation. *Journal of the ACM*, 45(4):653–750, jul 1998. doi:10.1145/
1621 285055.285060.
- 1622 **15** Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The*
1623 *collected works of Wassily Hoeffding*, pages 409–426, 1994.
- 1624 **16** Stefan Kiefer and Corto Mascle. On finite monoids over nonnegative integer matrices and
1625 short killing words. *Symposium on Theoretical Aspects of Computer Science, STACS*, 126,
1626 2019. doi:10.4230/LIPIcs.STACS.2019.43.
- 1627 **17** Stefan Kiefer and Corto N Mascle. On nonnegative integer matrices and short killing words.
1628 *SIAM Journal on Discrete Mathematics*, 35(2):1252–1267, 2021. doi:10.1137/19M1250893.
- 1629 **18** Frédéric Magniez and Michel de Rougemont. Property testing of regular tree languages.
1630 *Algorithmica*, 49(2):127–146, 2007. doi:10.1007/s00453-007-9028-3.
- 1631 **19** Liam Paninski. A coincidence-based test for uniformity given very sparsely sampled discrete
1632 data. *IEEE Transactions on Information Theory*, 54(10):4750–4755, 2008. doi:10.1109/TIT.
1633 2008.928987.
- 1634 **20** Ronitt Rubinfeld and Madhu Sudan. Robust characterizations of polynomials with applications
1635 to program testing. *SIAM Journal on Computing*, 25(2):252–271, 1996. doi:10.1137/
1636 S0097539793255151.
- 1637 **21** Andrew Chi-Chin Yao. Probabilistic computations: Toward a unified measure of complexity.
1638 In *Symposium on Foundations of Computer Science, SFCS*, pages 222–227. IEEE, 1977.
1639 doi:10.1109/SFCS.1977.24.

A Properties of minimal blocking factors

1640

1641 In this section, we discuss properties of the set of minimal blocking factors of an NFA. First,
1642 we show the set of minimal blocking factors of an automaton is a regular language.

1643 ► **Lemma A.1.** *Let $\mathcal{A} = (Q, \Sigma, \delta, I, F)$ be a strongly connected NFA with m states and let
1644 $\lambda = \lambda(\mathcal{A})$. For every $i \in \mathbb{Z}/\lambda\mathbb{Z}$, the set of minimal blocking factors of \mathcal{A} of the form $(i : u)$ is
1645 a regular language recognized by a NFA of size $2^{\mathcal{O}(m)}$.*

1646 **Proof.** We call blocking factors of \mathcal{A} of the form $(i : u)$ its *i*-blocking factors.

1647 We first show that the set of *i*-blocking factors of \mathcal{A} , but not necessarily minimal ones, is
1648 a regular language recognized by an NFA \mathcal{A}_i with $m + 1$ states. The result follows by using
1649 a standard construction for the automaton recognizing words in a regular language L that
1650 have no proper factor in a regular language L' , which gives an automaton of size $2^{\mathcal{O}(m)}$.

1651 Consider the NFA \mathcal{A}_i obtained by adding a new sink state \perp to \mathcal{A} , making it the only
1652 accepting state, with initial states Q_i . Formally, \mathcal{A}_i is defined as $\mathcal{A}_i = (Q \cup \{\perp\}, \Sigma, \delta', Q_i, \{\perp\})$,
1653 where δ' is defined as follows:

$$1654 \quad \forall p \in Q, \forall a \in \Sigma : \delta'(p, a) = \begin{cases} \{\perp\} & \text{if } \delta(p, a) = \emptyset, \\ \delta(p, a) & \text{otherwise.} \end{cases}$$

1655 This automaton³ recognizes the set of *i*-blocking factors of \mathcal{A} and has size $\mathcal{O}(m)$. Applying
1656 the aforementioned construction to $L = L' = \mathcal{L}(\mathcal{A}_i)$ yields the desired automaton, of size
1657 $2^{\mathcal{O}(m)}$. ◀

1658 It follows that the set of minimal blocking factors of \mathcal{A} is also a regular language.

1659 ► **Corollary A.2.** *Let \mathcal{A} be an NFA with m states. The set of minimal blocking factors of \mathcal{A}
1660 is a regular language recognized by an NFA of size $2^{\mathcal{O}(m)}$.*

1661 Therefore, if $\text{MBF}(\mathcal{A})$ is infinite, we can use Kleene's lemma to find an infinite family of
1662 minimal blocking factors with a shared structure $\{\phi\nu^r\chi, r \in \mathbb{N}\}$.

1663 ► **Lemma 4.20.** *If $\text{MBF}(\mathcal{A})$ is infinite, then there exist positional words ϕ, ν_+, ν_-, χ such
1664 that:*

- 1665 1. *the words ν_+ and ν_- have the same length,*
- 1666 2. *there exists a constant $S = 2^{\mathcal{O}(m)}$ such that $|\phi|, |\nu_+|, |\nu_-|, |\chi| \leq S$,*
- 1667 3. *there exists an index $i_* \in \mathbb{Z}/\lambda\mathbb{Z}$ and a state $q_* \in Q_{i_*}$ such that for every integer $r \geq 1$,*
1668 $\tau_{-,r} = \phi(\nu_-)^r z$ *is blocking for \mathcal{A} , and for every $s < r$, we have*

$$1669 \quad q_* \xrightarrow{\tau_{+,r,s}} q_* \text{ where } \tau_{+,r,s} = \phi(\nu_-)^j \nu_+ (\nu_-)^{r-1-s} \chi.$$

1670 *In particular, $\tau_{+,r,s}$ is not blocking for \mathcal{A} .*

1671 Note that here, the state q_* is the same for every integers r, s .

1672 **Proof.** As $\text{MBF}(\mathcal{A})$ is infinite, there must exist an i_* such that \mathcal{A} has infinitely many minimal
1673 i_* -blocking factors; we fix such an i_* in what follows.

1674 As the set of minimal i_* -blocking factors is an infinite regular language recognized by
1675 an NFA of size $S = 2^{\mathcal{O}(m)}$, by Kleene's Lemma, there exist positional words τ, μ, η , each of

³ Our definition of NFAs does not allow for multiple initial states. As there is no constraint of strong connectivity for \mathcal{A}_i , this can be solved using a simple construction that adds a new initial state.

length at most S with $|\mu| \geq 1$, such that for any non-negative integer k , $\tau\mu^k\eta$ is a minimal i_* -blocking factor. We can assume w.l.o.g. that neither τ nor η is empty, otherwise we set their value to μ : after this modification, $\tau\mu^k\eta$ is still a minimal i_* -blocking factor for every $k \geq 0$.

Notice that the word $\tau\mu^m$ is not a blocking factor, as a proper factor of the minimal blocking factor $\tau\mu^m\eta$. Therefore, by the pigeonhole principle, there exist integers $k_0, k_1 \geq 1$ with $k_0 + k_1 = m$ and states p, p_1 such that we have

$$p \xrightarrow{\tau\mu^{k_0}} p_1 \xrightarrow{\mu^{k_1}} p_1.$$

Note that, by Fact 3.3, $p_1 \xrightarrow{\mu^{k_1}} p_1$ implies that $k_1 \cdot |\mu| = 0 \pmod{\lambda}$.

Similarly, the word $\mu^m\eta$ is not a blocking factor, since it is a proper factor of the minimal i_* -blocking factor $\tau\mu^m\eta$. Again, there exist integers $k_2 \geq 1, k_3$ summing to m and states p_2 and q such that

$$p_2 \xrightarrow{\mu^{k_2}} p_2 \xrightarrow{\mu^{k_3}\eta} q.$$

Now, define $\phi = \tau\mu^{k_0}, \chi = \mu^{k_3}\eta$ and $\nu_- = \mu^K$, where $K = \rho \cdot k_1 \cdot k_2$. As there are transitions starting from p_1 and p_2 labeled by μ, p_1 and p_2 belong to the same periodicity class. Therefore, by Fact 3.3, as $K \geq \rho$ and $K \cdot |\mu| = 0 \pmod{\lambda}$, there exists a word ν_+ of length $K \cdot |\mu|$ such that $p_1 \xrightarrow{\nu_+} p_2$. This choice of ϕ, ν_+, ν_- and χ satisfies all the conditions of the lemma. \blacktriangleleft

B Hoeffding's inequality

► **Lemma B.1** ([15, Theorem 2]). *Let X_1, \dots, X_k be independent random variables such that for every $i = 1, \dots, k$, we have $a_i \leq X_i \leq b_i$, and let $S = \sum_{i=1}^k X_i$. Then, for any $t > 0$, we have*

$$\mathbb{P}(\mathbb{E}[S] - S \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^k (b_i - a_i)^2}\right).$$