

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

MPI-SWS, Kaiserslautern, Germany

Abstract

Property testing is concerned with the design of algorithms making a sub-linear 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 an optimal 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 introduced by Goldreich, Goldwasser and Ron [19] in 1998: it is the study of randomized approximate decision procedures that must distinguish objects that have a given property from those that are *far* from having it. Because of this relaxation on the specification, the field focuses on very efficient decision procedures, typically with sublinear (or even constant) running time – in particular, the algorithm does not even have the time to read the whole input.

In a seminal paper, Alon et al. [5] introduced property testing of formal languages: given a language L of finite words, the goal is to determine whether an input word u belongs to the language or is ε -far¹ from it, where ε is the precision parameter. The model assumes random access to the input word: a *query* specifies a position in the word and asks for the letter at that position, and the *query complexity* of the algorithm is the worst-case number of queries it makes to the input. Alon et al. [5] showed a surprising result: under the Hamming distance, all regular languages 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. In that paper, they also identified the class of *trivial* regular languages, those for which the answer is always *yes* or always *no* for sufficiently large n , e.g. finite languages or the set of words starting with an a , and showed that testing membership in a *non-trivial* regular language requires $\Omega(1/\varepsilon)$ queries.

¹ Informally, u is ε -far from L means that even by changing an ε -fraction of the letters of u , we cannot obtain a word in L . See Section 2 for a formal definition.

2 The Trichotomy of Regular Property Testing

42 The results of Alon et al. [5] leave a multiplicative gap of $\mathcal{O}(\log^3(1/\varepsilon))$ between the best
43 upper and lower bounds. We set out to improve our understanding of property testing of
44 regular languages by closing this gap. Bathie and Starikovskaya obtained in 2021 [9] the first
45 improvement over the result of Alon et al. [5] in more than 20 years:

46 ► **Fact 1.1** (From [9, Theorem 5]). *Under the edit distance, every regular language can be*
47 *tested with $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.*

48 Testers under the edit distance are weaker than testers under the Hamming distance, hence
49 this result does not exactly improve the result of Alon et al. [5]. We overcome this shortcoming
50 later in this article: Theorem 4.12 extends the above result to the case of the Hamming
51 distance.

52 Bathie and Starikovskaya also showed that this upper bound is tight, in the sense that
53 there is a regular language L_0 for which this complexity cannot be further improved, thereby
54 closing the query complexity gap.

55 ► **Fact 1.2** (From [9, Theorem 15]). *There is a regular language L_0 with query complexity*
56 *$\Omega(\log(\varepsilon^{-1})/\varepsilon)$ under the edit distance², for all small enough $\varepsilon > 0$.*

57 Furthermore, it is easy to find specific non-trivial regular languages for which there is an
58 algorithm using only $\mathcal{O}(1/\varepsilon)$ queries. Some examples of languages that can require $\mathcal{O}(1/\varepsilon)$
59 queries are a^* over the alphabet $\{a, b\}$, $(ab)^*$ or $(aa + bb)^*$.

60 Hence, these results combined with those of Alon et al. [5] show that there exist trivial
61 languages (that require 0 queries for large enough n), *easy* languages (with query complexity
62 $\Theta(1/\varepsilon)$) and *hard* languages (with query complexity $\Theta(\log(\varepsilon^{-1})/\varepsilon)$). This raises the question
63 of whether there exist languages with a different query complexity (e.g. $\Theta(\log \log(\varepsilon^{-1})/\varepsilon)$),
64 or if every regular is either trivial, easy or hard. This further asks the question of giving a
65 characterization of the languages that belong to each class, inspired by the recent success of
66 exact characterizations of the complexity of sliding window [16] recognition and dynamic
67 membership [7] of regular languages.

68 In this article, we answer both questions: we show a trichotomy of the complexity of
69 testing regular languages under the Hamming distance³, showing that there are only the
70 three aforementioned complexity classes (trivial, easy and hard), we give a characterization
71 of all three classes using a combinatorial object called *blocking sequences*, and show that this
72 characterization can be decided in polynomial space (and that it is complete for PSPACE).
73 This trichotomy theorem closes a line of work on improving query complexity for property
74 testers and identifying easier subclasses of regular languages.

75 1.1 Related work

76 A very active branch of property testing focuses on graph properties, for instance one can
77 test whether a given graph appears as a subgraph [3] or as an induced subgraph [4], and
78 more generally every monotone graph property can be tested with one-sided error [6]. Other
79 families of objects heavily studied under this algorithmic paradigm include probabilistic
80 distributions [25, 11] combined with privacy constraints [2], numerical functions [10, 28], and
81 programs [13, 12]. We refer to the book of Goldreich [18] for an overview of the field of
82 property testing.

² Note that, as opposed to testers, lower bounds for the edit distance are stronger than lower bounds of the Hamming distance.

³ We consider one-sided error testers, also called testing with perfect completeness, see definitions below.

83 **Testing formal languages.** Building upon the seminal work of Alon et al. [5], Magniez
 84 et al. [23] gave a tester using $\mathcal{O}(\log^2(\varepsilon^{-1})/\varepsilon)$ queries for regular languages under the edit
 85 distance with moves, and François et al. [15] gave a tester using $\mathcal{O}(1/\varepsilon^2)$ queries for the case
 86 of the weighted edit distance. Alon et al. [5] also show that context-free languages cannot be
 87 tested with a constant number of queries, and subsequent work has considered testing specific
 88 context-free languages such as the DYCK languages [26, 14] or regular tree languages [23].
 89 Property testing of formal languages has been investigated in other settings: Ganardi et
 90 al. [17] studied the question of testing regular languages in the so-called “sliding window
 91 model”, while others considered property testing for subclasses of context-free languages in
 92 the streaming model: Visibly Pushdown languages [15], DYCK languages [21, 22, 24] or DLIN
 93 and $\text{LL}(k)$ [8]. A recent application of property testing of regular languages was to detect
 94 race conditions in execution traces of distributed systems [29].

95 1.2 Main result and overview of the paper

96 ► **Definition 1.3** (Hard, easy and trivial languages). *Let L be a regular language. We say that:*

- 97 ■ L is hard if the optimal query complexity for a property tester for L is $\Theta(\log(\varepsilon^{-1})/\varepsilon)$.
- 98 ■ L is easy if the optimal query complexity for a property tester for L is $\Theta(1/\varepsilon)$.
- 99 ■ L is trivial if there exists $\varepsilon_0 > 0$ such that for all positive $\varepsilon < \varepsilon_0$, there is a property
 100 tester and some $n \in \mathbb{N}$ such that the tester makes 0 queries on words of length $\geq n$.

101 Our characterisation of those three classes uses the notion of *blocking sequence* of a
 102 language L . Intuitively, they are sequences of words such that finding those words as factors
 103 of a word w proves that w is not in L . We also define an order on them, which gives us a
 104 notion of *minimal* blocking sequence.

105 ► **Theorem 1.4.** *Let L be an infinite regular language recognized by an NFA \mathcal{A} and let
 106 $MBS(\mathcal{A})$ denote the set of minimal blocking sequences of \mathcal{A} . The complexity of testing L is
 107 characterized by $MBS(\mathcal{A})$ as follows:*

- 108 1. L is trivial if and only if $MBS(\mathcal{A})$ is empty;
- 109 2. L is easy if and only if $MBS(\mathcal{A})$ is finite and nonempty;
- 110 3. L is hard if and only if $MBS(\mathcal{A})$ is infinite.

111 In the case where L is recognised by a strongly connected automaton, blocking sequences
 112 can be replaced by *blocking factors*. A blocking factor is a single word that is not a factor of
 113 any word in L .

114 Section 2 defines the necessary terms and notations. The rest of the paper is structured
 115 as follows. In Sections 3 and 4, we delimit the set of hard languages, that is, the ones that
 116 require $\Theta(\log(\varepsilon^{-1})/\varepsilon)$ queries. More precisely, Section 3 focuses on the subcase of languages
 117 defined by *strongly connected automata*.

- 118 ■ First, we show that the language of every strongly connected automaton is testable with
 119 $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.
- 120 ■ Then, we show that if the language of a strongly connected automaton has infinitely
 121 many blocking factors then it requires $\Omega(\log(\varepsilon^{-1})/\varepsilon)$ queries.

122 Then, Section 4 extends those results to all automata. The interplay with the previous
 123 section is different for the upper and the lower bound

- 124 ■ We show the upper bound of $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries to all automata by a natural extension
 125 of the proof in the strongly connected case.

4 The Trichotomy of Regular Property Testing

126 ■ We obtain a lower bound of $\Omega(\log(\varepsilon^{-1})/\varepsilon)$ queries for languages with infinitely many
127 minimal blocking sequences by a reduction to the strongly connected case. We show that
128 we can use a tester for the general automaton to test one of its hard strongly connected
129 components.

130 Section 5 completes the trichotomy, by characterising the easy and trivial languages.

131 ■ We show that languages of automata with finitely many blocking sequences can be tested
132 with $\mathcal{O}(1/\varepsilon)$ queries.

133 ■ Finally, we prove that if an automaton has at least one blocking sequence, then it requires
134 $\Omega(1/\varepsilon)$ queries to be tested. By contrast, automata with no blocking sequence recognise
135 trivial languages

136 Once we have the trichotomy, it is natural to ask whether it is effective: given an
137 automaton \mathcal{A} , can we determine if its language is trivial, easy or hard? The answer is yes,
138 and we show that all three decision problems are PSPACE-complete in Section 6.

2 Preliminaries

140 **Words and automata.** We write Σ^* (resp. Σ^+) for the set of finite words (resp. non-empty
141 words) over the alphabet Σ . The length of a word u is denoted $|u|$, and its i th letter is
142 denoted $u[i]$. The empty word is denoted γ . Given $u \in \Sigma^*$ and $0 \leq i, j \leq |u| - 1$, define
143 $u[i..j]$ as the word $u[i]u[i+1]\dots u[j]$ if $i \leq j$ and γ otherwise. Further, $u[i..j)$ denotes the
144 word $u[i..j-1]$. A word w is a *factor* (resp. *prefix*, *suffix*) of u if there exist indices i, j such
145 that $w = u[i..j]$ (resp. with $i = 0, j = |u| - 1$). We use $w \preceq u$ to denote “ w is a factor of u ”.
146 Furthermore, if w is a factor of u and $w \neq u$, we say that w is a *proper factor* of u .

147 A *nondeterministic finite automaton* (NFA) \mathcal{A} is a transition system defined by a tuple
148 $(Q, \Sigma, \delta, q_0, F)$, with Q a finite set of states, Σ a finite alphabet, $\delta : Q \times \Sigma \rightarrow 2^Q$ the transition
149 function, $q_0 \in Q$ the initial state and $F \subseteq Q$ the set of final states. The semantics is as
150 usual [27]. When there is a path from a state p to a state q in \mathcal{A} , we say that q is reachable
151 from p and that p is co-reachable from q . In this article, we assume w.l.o.g. that all NFA \mathcal{A}
152 are *trim*, i.e., every state is reachable from the initial state and co-reachable from some final
153 state.

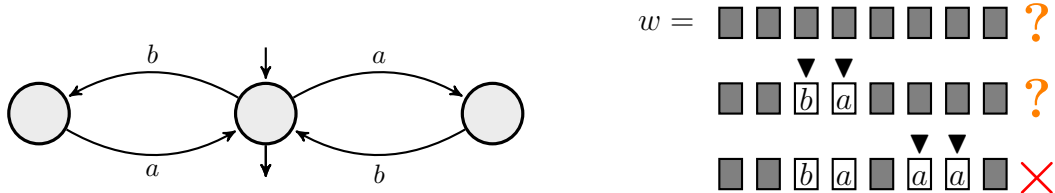
154 **Property testing.** Let us start with the notion of a property tester for a language L : the
155 goal is to determine whether an input word u belongs to the language L , or whether it is
156 ε -far from a distance d .

157 ► **Definition 2.1.** Let L be a language, let u be a word of length n , let $\varepsilon > 0$ be a precision
158 parameter and let $d : \Sigma^* \times \Sigma^* \rightarrow \mathbb{N} \cup \{+\infty\}$ be a metric. We say that the word u is ε -far
159 from L w.r.t. d if $d(u, L) \geq \varepsilon n$, where

$$160 \quad d(u, L) := \inf_{v \in L} d(u, v).$$

161 We assume random access to the input word: a query specifies a position in the word and
162 asks for the letter in this position.

163 Throughout this work and unless explicitly stated otherwise, we will consider the case
164 where d is the Hamming distance, defined for two words u and v as the number of positions
165 at which they differ if they have the same length, and as $+\infty$ otherwise. Here, $d(u, L) \geq \varepsilon n$
166 means that one cannot change an ε -fraction of the letters in u to obtain a word in L .



■ **Figure 1** An automaton \mathcal{A} recognizing the language $L = (ab)^* + (ba)^*$. A witness that a word is not in this language is a factor aa or bb . We query factors of size 2, and reject if we find one of those factors. We can show that a word ε -far from L must have many such factors, and infer that we can get a constant probability of rejecting it by sampling only $\mathcal{O}(1/\varepsilon)$ factors.

167 A ε -property tester (or simply a *tester*) T for a language L is a randomized algorithm
 168 that, given an input word u , always answers “yes” if $u \in L$ and answers “no” with probability
 169 bounded away from 0 when u is ε -far from L .

170 ► **Definition 2.2.** A property tester for the language L with precision $\varepsilon > 0$ is a randomized
 171 algorithm T that, for any input u of length n , given random access to u , satisfies the following
 172 properties:

173
$$\text{if } u \in L, \text{ then } T(u) = 1, \tag{1}$$

174
$$\text{if } u \text{ is } \varepsilon\text{-far from } L, \text{ then } \mathbb{P}(T(u) = 0) \geq 2/3. \tag{2}$$

175 The query complexity of T is a function of n and ε that counts the maximum number of
 176 queries that T makes over all inputs of length n and over all possible random choices.

177 We measure the complexity of a tester by its *query complexity*. Let us emphasize that
 178 throughout this article we focus on so-called “testers with perfect completeness”, or “one-sided
 179 error”: if a word is in the language, the tester answers positively (with probability 1). In
 180 particular our characterization applies to this class. Because they are based on the notion of
 181 blocking factors that we will discuss below, all known testers for regular languages [5, 23, 15, 9]
 182 have perfect completeness.

183 ► **Remark 2.3.** If L is finite, then it is trivial: since there is a bound B on the lengths of its
 184 words, it suffices to read entirely words of length $\leq B$ and answer No for words of greater
 185 length. For that reason, in the rest of the paper we only consider infinite languages.

186 **Graphs and periodicity.** We now recall tools introduced by Alon et al. [5] to deal with
 187 periodicity in finite automata.

188 Let $G = (V, E)$ with $E \subseteq V^2$ be a directed graph. A *strongly connected component* (or
 189 SCC) of G is a maximal set of vertices that are all reachable from each other. It is *trivial* if
 190 it contains a single state with no self-loop on it. We say that G is *strongly connected* if its
 191 entire set of vertices is an SCC.

192 The period $\lambda = \lambda(G)$ of a non-trivial strongly connected graph G is the greatest common
 193 divisor of the length of the cycles in G . Following the work of Alon et al. [5], we will use the
 194 following property of directed graphs.

195 ► **Fact 2.4** (From [5, Lemma 2.3]). Let $G = (V, E)$ be a non-empty, non-trivial, strongly
 196 connected graph with finite period $\lambda = \lambda(G)$. Then there exists a partition $V = Q_0 \sqcup \dots \sqcup Q_{\lambda-1}$
 197 and a reachability constant $\rho = \rho(G)$ that does not exceed $3|V|^2$ such that:

6 The Trichotomy of Regular Property Testing

- 198 1. For every $0 \leq i, j \leq \lambda - 1$ and for every $s \in Q_i, t \in Q_j$, the length of any directed path
 199 from s to t in G is equal to $(j - i) \bmod \lambda$.
- 200 2. For every $0 \leq i, j \leq \lambda - 1$, for every $s \in Q_i, t \in Q_j$ and for every integer $r \geq \rho$, if
 201 $r = (j - i) \bmod \lambda$, then there exists a directed path from s to t in G of length r .

202 The sets $(Q_i : i = 0, \dots, \lambda - 1)$ are the *periodicity classes* of G . In what follows, we will
 203 slightly abuse notation and use Q_i even when $i \geq \lambda$ to mean $Q_{i \bmod \lambda}$.

204 An automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ defines an underlying graph $G = (Q, E)$ where
 205 $E = \{(p, q) \in Q^2 \mid \exists a \in \Sigma : q \in \delta(p, a)\}$. In what follows, we naturally extend the notions
 206 defined above to finite automata through this graph G . Note that the numbering of the
 207 periodicity classes is defined up to a shift mod λ : we can thus always assume that Q_0 is the
 208 class that contains the initial state q_0 . The period of \mathcal{A} is written $\lambda(\mathcal{A})$.

209 **Positional words and positional languages.** Take a look at the language $L = (ab)^*$. The
 210 word $v = ab$ can appear as a factor of a word $u \in L$ if v occurs at an even position in u .
 211 However, if v occurs at an *odd* position in u , then u cannot be in L . Therefore, v can be
 212 used to witness that u is not in L , but only if we find it at an odd position. This example
 213 motivates the introduction of *p-positional words*, which additionally encode information
 214 about the index of each letter modulo an integer p .

215 More generally, we will associate to each regular language a period p , and working with
 216 p -positional words will allow us to define blocking factors in a position-dependent way without
 217 explicitly considering the index at which the factor occurs.

218 **► Definition 2.5 (Positional words).** Let p be a positive integer. A p -positional word is a
 219 word over the alphabet $\mathbb{Z}/p\mathbb{Z} \times \Sigma$ of the form $(n \bmod p, a_0)((n + 1) \bmod p, a_1) \cdots ((n + \ell)$
 220 $\bmod p, a_\ell)$. If $u = a_0 \cdots a_\ell$, we call this word $(n : u)$ (when p is clear from context).

221 Any strongly connected finite automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$ can naturally be extended
 222 into an automaton $\hat{\mathcal{A}}$ over $\lambda(\mathcal{A})$ -positional words with $\lambda(\mathcal{A})|Q|$ states, recognising $\{(0 : u) \mid$
 223 $u \in \mathcal{L}(\mathcal{A})\}$. It suffices to keep track in the states of the current state of \mathcal{A} and the number of
 224 letters read modulo $\lambda(\mathcal{A})$. We call the language recognized by $\hat{\mathcal{A}}$ the *positional language* of
 225 \mathcal{A} , and denote it $\mathcal{TL}(\mathcal{A})$. This definition is motivated by the following property:

226 **► Property 2.6.** For any word $u \in \Sigma^*$, we have $u \in \mathcal{L}(\mathcal{A})$ if and only if $(0 : u) \in \mathcal{TL}(\mathcal{A})$.

227 Positional words make it easier to manipulate factors with positional information, hence
 228 we phrase our property testing results in terms of positional languages. Note that a tester for
 229 $\mathcal{TL}(\mathcal{A})$ immediately gives a property tester for $\mathcal{L}(\mathcal{A})$, as one can simulate queries to $(0 : u)$
 230 with queries to u by simply pairing the index of the query modulo $\lambda(\mathcal{A})$ with its result.

3 Hard Languages for Strongly Connected NFAs

232 Before considering the case of arbitrary NFAs, we first study the case of strongly connected
 233 NFAs. We define the set of *minimal blocking factors* of \mathcal{A} , which are factor-minimal $\lambda(\mathcal{A})$ -
 234 positional words that witness the fact that a word does not belong to $\mathcal{TL}(\mathcal{A})$. We show that all
 235 strongly connected NFA have a query complexity in $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$, and that the ones which
 236 have infinitely many minimal blocking factors have a query complexity of $\Theta(\log(\varepsilon^{-1})/\varepsilon)$.

237 In this section, we consider a fixed NFA \mathcal{A} and simply use “positional words” to refer to
 238 λ -positional words, where $\lambda = \lambda(\mathcal{A})$ is the period of \mathcal{A} .

239 ► **Definition 3.1** (Blocking factors). *Let \mathcal{A} be a strongly connected NFA. A positional word τ*
 240 *is a blocking factor of \mathcal{A} if for any other positional word μ we have $\tau \preceq \mu \Rightarrow \mu \notin \mathcal{TL}(\mathcal{A})$.*

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

243 Intuitively and in terms of automata, $(i : u)$ is blocking for \mathcal{A} if u does not label any path in
 244 \mathcal{A} from a state of Q_i . (This property is formally established later in Lemma A.1.) The main
 245 result of this section is the following:

246 ► **Theorem 3.2.** *Let L be an infinite language recognised by a strongly connected NFA \mathcal{A} . If*
 247 *$MBF(\mathcal{A})$ is infinite, then L is hard, i.e., it has query complexity $\Theta(\log(\varepsilon^{-1})/\varepsilon)$*

248 This result gives both an upper bound of $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ and a lower bound of $\Omega(\log(\varepsilon^{-1})/\varepsilon)$
 249 on the query complexity of a tester for L : we prove the upper bound in Section 3.1 and the
 250 lower bound in Section 3.2.

251 3.1 An efficient property tester for strongly connected NFAs.

252 In this section, we show the following result.

253 ► **Theorem 3.3.** *Let \mathcal{A} be a strongly connected NFA. For any $\varepsilon > 0$, there exists an ε -property*
 254 *tester for $\mathcal{L}(\mathcal{A})$ that uses $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.*

255 As mentioned in the overview, this result supersedes the one that was obtained in [9]:
 256 while both testers use the same number of queries, the tester in [9] works under the edit
 257 distance, while the one of Theorem 3.3 is designed for the Hamming distance. As the edit
 258 distance never exceeds the Hamming distance, the set of words that are ε -far with respect to
 259 the former is contained in the set of words ε -far for the latter. Therefore, an ε -tester for the
 260 Hamming distance is also an ε -tester for the edit distance, and our result supersedes and
 261 generalizes theirs. Our proof is similar to that of [9], with one notable improvement: we use
 262 a new method for sampling factors in u , which greatly simplifies the correctness analysis.

263 The algorithm for Theorem 3.3 is Algorithm 1. The procedure is fairly simple: the
 264 algorithm samples at random factors of various lengths in u , and rejects if and only if at
 265 least one of these factors is blocking. On the other hand, the correctness of the tester is
 266 far from trivial. The lengths and the number of factors of each lengths are chosen so that
 267 the number of queries is $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ and the probability of finding a blocking factor is
 268 maximized, regardless of their repartition in u .

269 We now turn to formally establishing these properties.

270 ▷ **Claim 3.4.** The tester given in Algorithm 1 makes $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries to u .

271 Next, we show that if u is ε -far from $\mathcal{L}(\mathcal{A})$, then $(0 : u)$ contains $\Omega(\varepsilon n)$ blocking factors,
 272 each of length $\mathcal{O}(1/\varepsilon)$ (see Appendix A.2 for a proof).

273 ► **Lemma 3.5.** *Let $\varepsilon > 0$, let u be a word of length $n \geq 6m^2/\varepsilon$ and assume that $\mathcal{L}(\mathcal{A})$*
 274 *contains at least one word of length n . If u is ε -far from $\mathcal{L}(\mathcal{A})$, then the positional word*
 275 *$(0 : u)$ contains at least $\varepsilon n/(12m^2)$ disjoint blocking factors of length at most $12m^2/\varepsilon$.*

276 For the correctness analysis, we assume that u is ε -far from $\mathcal{L}(\mathcal{A})$, and show that
 277 Algorithm 1 finds at least one of the blocking factors given by Lemma 3.5 with probability
 278 at least $2/3$.

279 ► **Lemma 3.6.** *In the last Else block, if u is ε -far from $\mathcal{L}(\mathcal{A})$, then Algorithm 1 rejects with*
 280 *probability at least $2/3$.*

■ **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:   return ( $l : u[l..r]$ )
5: function TESTER( $u, \varepsilon$ )
6:    $\beta \leftarrow 12m^2/\varepsilon$ 
7:   if  $\mathcal{L}(\mathcal{A}) \cap \Sigma^n = \emptyset$  then
8:     Reject
9:   else if  $n < \beta$  then
10:    Query all of  $u$  and run  $\mathcal{A}$  on it
11:    Accept if and only if  $\mathcal{A}$  accepts
12:   else
13:      $T \leftarrow \lceil \log(\beta) \rceil$ 
14:     for  $t = 0$  to  $T$  do
15:        $\ell_t \leftarrow 2^t, r_t \leftarrow \lceil 2 \ln(3)\beta/\ell_t \rceil$ 
16:       Call SAMPLE( $u, \ell_t$ )  $r_t$  times
17:     Reject if and only if any call to SAMPLE returned a blocking factor for  $\mathcal{A}$ .

```

281 **Proof sketch.** Assume that u is ε -far from $\mathcal{L}(\mathcal{A})$, that $\mathcal{L}(\mathcal{A}) \cap \Sigma^n$ is not empty (i.e. $\mathcal{L}(\mathcal{A})$
282 contains a word of length n) and that $n \geq \beta$. For $t = 0, \dots, T$, let B_t denote the subset of \mathcal{B}
283 of blocking factors of length at most $\ell_t = 2^t$.

284 For each t , if in a call to SAMPLE(u, ℓ_t) the random position i is inside a blocking factor
285 of B_t , then the returned factor is blocking. Therefore, as the factors given by Lemma 3.5
286 are disjoint, each one of the $r_t = 2 \ln(3)\beta/\ell_t$ calls to SAMPLE(u, ℓ_t) has a probability
287 $p_t \geq \frac{1}{n} \sum_{\tau \in B_t} |\tau|$ to return a blocking factor.

288 Careful calculations let us show that the probability that no blocking factor is found is at
289 most $1/3$. They are detailed in Appendix A.2. ◀

290 3.2 Lower bound when there are infinitely many minimal blocking words

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

293 We start from the parity language P consisting of words that contain an even number
294 of b 's, over the alphabet $\{a, b\}$. Distinguishing $u \in P$ from $u \notin P$ requires $\Omega(n)$ queries, as
295 changing the letter at a single position can change membership in P . However, P is trivial
296 to test, as any word is at distance at most 1 from P , for the same reason. Now, consider
297 language L consisting of words over $\{a, b, c, d\}$ such that between a c and the next d , there
298 is a word in P . Intuitively, this language encodes multiple instances of P , which are hard
299 to identify for property testers. In [9], the authors proved a lower bound of $\Omega(\log(\varepsilon^{-1})/\varepsilon)$
300 on the query complexity of any property tester for L , matching the upper bound in the
301 same paper. The upper bound of [9] is for the edit distance, we show that the result can be
302 extended to the Hamming distance.

303 The minimal blocking factors of L include all words for the form cvd where $v \notin P$: there
304 are infinitely many such words. This is no coincidence: we show in this article than the
305 lower bound from [9] can be lifted to any languages with infinitely many minimal blocking
306 sequences (or minimal blocking factors in the case of strongly connected automata).

307 **► Theorem 3.7.** *Let \mathcal{A} be a strongly connected NFA. If $MBF(\mathcal{A})$ is infinite, then there exists*
 308 *$\varepsilon_0 > 0$ such that for all $\varepsilon < \varepsilon_0$, every ε -property tester for $L = \mathcal{L}(\mathcal{A})$ uses $\Omega(\log(\varepsilon^{-1})/\varepsilon)$*
 309 *queries.*

310 The proof of this result is a full generalization of the lower bound against the “repeated
 311 parity” language that was published in [9, Theorem 15].

312 Our proof is based on (a consequence of) Yao’s Minmax Principle [30]: if there is a
 313 distribution \mathcal{D} over inputs such that any *deterministic* algorithm that makes at most q
 314 queries errs on $u \sim \mathcal{D}$ with probability at least p , then any *randomized* with q queries errs
 315 with probability at least p on some input u .

316 To prove Theorem 3.7, we first exhibit such a distribution \mathcal{D} for $q = \Theta(\log(\varepsilon^{-1})/\varepsilon)$ (see
 317 Appendix A.3.2 for its construction). We take the following steps:

- 318 1. we show that with high probability, an input u sampled w.r.t. \mathcal{D} is either in or ε -far from
 319 L (Lemma 3.8),
- 320 2. we show that with high probability, any deterministic tester that makes fewer than
 321 $c \cdot \log(\varepsilon^{-1})/\varepsilon$ queries (for a suitable constant c) cannot distinguish whether the instance
 322 u is positive or ε -far, hence it errs with large probability.
- 323 3. we combine the above two results to prove Theorem 3.7 via Yao’s Minmax principle.

324 Constructing a Hard Distribution \mathcal{D}

325 To sample an input w.r.t. \mathcal{D} , we first sample a uniformly random bit π : if $\pi = 1$, we construct
 326 a word u that belongs to L , and if $\pi = 0$, u will be far from L with high probability. The
 327 input word is then conceptually divided into εn disjoint intervals. For the j -th interval,
 328 we sample a random variable $\kappa_j \in \{0, 1, \dots, \log(1/\varepsilon)\}$: this random variable describes the
 329 content of the interval. If $\pi = 1$, we fill the interval with *non-blocking* factors of length
 330 $\mathcal{O}(2^{\kappa_j})$, and if $\pi = 0$, we instead fill the interval with *blocking* factors of the same length,
 331 chosen to be very similar to the non-blocking factors, making it hard to distinguish between
 332 the two cases with few queries. By carefully choosing the distribution for κ_j , we can ensure
 333 that when $\pi = 0$, the resulting instance u is ε -far from L with high probability. (Proofs for
 334 this section can be found in Appendix A.3.)

335 **► Lemma 3.8.** *Conditioned on $\pi = 0$, the probability of the event $\mathcal{F} = \{u \text{ is } \varepsilon\text{-far from}$*
 336 *$\mathcal{TL}(\mathcal{A})\}$ goes to 1 as n goes to infinity.*

337 **► Corollary 3.9.** *For large enough n , we have $\mathbb{P}(\mathcal{F}) \geq 5/12$.*

338 Intuitively, our distribution is hard to test because positive and negative instance are very
 339 similar. Therefore, a tester with few queries will likely not be able to tell them apart: the
 340 perfect completeness constraint forces the tester to accept in that case. Below, we establish
 341 this result formally.

342 **► Lemma 3.10.** *Let T be a deterministic tester with perfect completeness (i.e. one sided*
 343 *error, always accepts $\tau \in \mathcal{TL}(\mathcal{A})$) and let q_j denote the number of queries that it makes in*
 344 *the j -th interval. Conditioned on the event $\mathcal{M} = \{\forall j \text{ s.t. } \kappa_j > 0, q_j < 2^{\kappa_j}\}$, the probability*
 345 *that T accepts u is 1.*

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

347 **► Lemma 3.11.** *Let T be a deterministic tester, let q_j denote the number of queries that*
 348 *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.*
 349 *$\sum_j q_j \leq \frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$. The probability of the event $\mathcal{M} = \{\forall j \text{ s.t. } \kappa_j > 0, q_j < 2^{\kappa_j}\}$ is at*
 350 *least $11/12$.*

351 We are now ready to prove Theorem 3.7.

352 **Proof of Theorem 3.7.** We want to show that any tester with perfect completeness for
 353 $\mathcal{L}(\mathcal{A})$ requires at least $\frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$ queries, by showing that any tester with fewer queries
 354 errs with probability at least $1/3$. We show that any **deterministic** algorithm T with
 355 perfect completeness that makes less than $\frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$ queries errs on u when $u \sim \mathcal{D}$ with
 356 probability at least $1/3$, and conclude using Yao’s Minmax principle.

357 Consider such an algorithm T . The probability that T makes an error on u is lower-
 358 bounded by the probability that u is ε -far from $\mathcal{L}(\mathcal{A})$ and T accepts, which in turn is larger
 359 than the probability of $\mathcal{M} \cap \mathcal{F}$. By Corollary 3.9, we have $\mathbb{P}(\mathcal{F}) \geq 5/12$, and by Lemma 3.11,
 360 $\mathbb{P}(\mathcal{M})$ is at least $11/12$. Therefore, we have

$$361 \quad \mathbb{P}(T \text{ errs}) \geq \mathbb{P}(\mathcal{M} \cap \mathcal{F}) \geq 1 - 7/12 - 1/12 = 4/12 = 1/3.$$

362 This concludes the proof of Theorem 3.7, and consequently of Theorem 3.2. ◀

363 4 Characterisation of Hard Languages for All NFAs

364 In this section we extend the results of the previous section to all finite automata. This
 365 extension is based on a generalization blocking factors: we introduce *blocking sequences*,
 366 which are sequences of factors that witness the fact that we cannot take any path through
 367 the strongly connected components of the automaton. For the lower bound, we define a
 368 suitable partial order on blocking sequences, which extends the factor relation on words to
 369 those sequences, and allows us to define *minimal* blocking sequences.

370 4.1 Blocking sequences

371 We start with an example that highlights why we need to use blocking sequences instead of
 372 blocking factors.

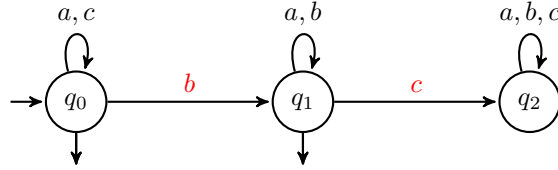
373 ▶ **Example 4.1.** Let us study the automaton \mathcal{A} depicted in Figure 2. The set of minimal
 374 blocking factors of $\mathcal{L}(\mathcal{A})$ is infinite: it is the language ba^*c . Yet, L is easy to test: We sample
 375 $\mathcal{O}(1/\varepsilon)$ letters at random, answer “no” if the sample contains a c occurring after a b , and
 376 “yes” otherwise. To prove that this yields a property tester, we rely on the following property:
 377

378 ▶ **Property 4.2.** *If u is ε -far from L , then it can be decomposed $u = u_1u_2$ where u_1 contains*
 379 *$\Omega(\varepsilon n)$ letters b and u_2 contains $\Omega(\varepsilon n)$ letters c .*

380 The pair of factors (c, a) is an example of blocking sequence: a word that contains an
 381 occurrence of the first followed by an occurrence of the second cannot be in L . We can
 382 also show that a word ε -far from L must contain many disjoint blocking sequences – this
 383 property (Lemma 4.11) underpins the algorithm for general regular languages.

384 What this example shows is that blocking factors are not enough: we need to consider
 385 sequences of factors, yielding the notion of *blocking sequences*. Intuitively, a blocking sequence
 386 for L is a sequence $\sigma = (v_1, \dots, v_k)$ of (positional) words such that if each word of the sequence
 387 appears in u , in the same order as in σ , then u is not in L . While L has infinitely many
 388 minimal blocking factors, it has a single minimal blocking sequence $\sigma = (b, c)$.

⁴ This is not quite the definition, but it conveys the right intuition.



■ **Figure 2** An automaton \mathcal{A} recognizing the language $L = (a + c)^*(a + b)^*$.

4.1.1 Portals and SCC-paths

Intuitively, blocking sequences are sequences of blocking factors of successive strongly connected components. To formalize this intuition, we use *portals*, which describe how a run in the automaton traverses a strongly connected component, and *SCC-paths*, that describe a succession of portals. Appendix B presents several examples motivating and illustrating the definitions below.

We fix an NFA $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$. Let \mathcal{S} be the set of SCCs of \mathcal{A} . We define p as the least common multiple of the lengths of all simple cycles of \mathcal{A} . Given a number $k \in \mathbb{Z}/p\mathbb{Z}$, we say that a state t is k -reachable from a state s if there is a path from s to t of length k modulo p . From now on, we use “positional words” for p -positional words.

► **Definition 4.3** (Portal). *A portal is a 4-tuple $P = s, x \rightsquigarrow t, y \in (Q \times \mathbb{Z}/p\mathbb{Z})^2$, such that s and t are in the same SCC. It describes the first and last states visited by a path in an SCC, and the positions x, y (modulo p) at which it first and lasts visits that SCC.*

The positional language of a portal is the set

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

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

► **Definition 4.4.** *A positional word $(n : u)$ is blocking for a portal P if it is not a factor of any word of $\mathcal{L}(P)$. In other words, there is no path that starts in s and ends in t , of length $y - x$ modulo p , which reads u after $n - x$ steps modulo p .*

Portals describe the behavior of a run in a single strongly connected component.

► **Definition 4.5** (SCC-path). *An SCC-path π of \mathcal{A} is a sequence of portals linked by transitions $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} s_1, x_1 \rightsquigarrow t_1, y_1 \cdots \xrightarrow{a_k} s_k, x_k \rightsquigarrow t_k, y_k$, such that for all $i \in \{1, \dots, k\}$, $x_i = y_{i-1} + 1 \pmod{p}$ and $s_i \in \delta(t_{i-1}, a_i)$.*

Intuitively, an SCC-path is a description of the states and positions at which a path through the automaton enters and leaves each SCC. See Appendix B for some examples.

► **Definition 4.6.** *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 $\mathcal{L}(\pi) = \mathcal{L}(s_0, x_0 \rightsquigarrow t_0, y_0)a_1\mathcal{L}(s_1, x_1 \rightsquigarrow t_1, y_1)a_2 \cdots \mathcal{L}(s_k, x_k \rightsquigarrow t_k, y_k)$*

We say that π is accepting if $x_0 = 0$, $s_0 = q_0$, $t_k \in F$ and $\mathcal{L}(\pi)$ is non-empty.

► **Remark 4.7.** We have $\mathcal{L}(\mathcal{A}) = \bigcup_{\pi \text{ accepting}} \mathcal{L}(\pi)$.

Decomposing \mathcal{A} as a union of SCC-paths allows us to use them as an intermediate step to define blocking sequences. We earlier defined blocking factors for portals: we now generalize this definition to blocking sequences for SCC-paths, to finally define blocking sequences for automata.

423 ► **Definition 4.8.** We say that a sequence (μ_1, \dots, μ_ℓ) of positional factors is blocking for an
 424 SCC-path $\pi = P_0 \xrightarrow{a_1} \dots P_k$ if there is a sequence of indices $i_0 \leq i_1 \leq \dots \leq i_k$ such that for
 425 every j , μ_{i_j} is blocking for P_j .

426 Crucially, in this definition, consecutive indices i_j and i_{j+1} can be equal, i.e. a single
 427 factor of the sequence may be blocking for multiple consecutive SCCs in the SCC-path. This
 428 choice is motivated by Example B.1 in the appendix, where the language is easy because
 429 consecutive SCCs share blocking factors.

430 ► **Definition 4.9** (Blocking sequence for \mathcal{A}). A sequence of positional words $\sigma = (\mu_1, \dots, \mu_\ell)$
 431 is blocking for \mathcal{A} if it is blocking for all SCC-paths of \mathcal{A} .

432 See Example B.2) for an example illustrating these definitions.

433 The goal of the next two lemmas is to show that we can reduce property testing of $\mathcal{L}(\mathcal{A})$
 434 to a search for blocking sequences in the word. They are proven in Appendix C.

435 ■ If we find a few blocking sequences in a word then it is not in the language and we can
 436 answer no (Lemma 4.10).

437 ■ A word that is far from the language contains many blocking sequences (Lemma 4.11).
 438 Hence if we do not find blocking sequences in the word, it is unlikely to be far from
 439 the language (Lemma 4.11). In fact, the lemma is more precise: we can find blocking
 440 sequences organised in a particular way.

441 ► **Lemma 4.10.** If μ contains $|\mathcal{A}|$ disjoint blocking sequences for \mathcal{A} then $\mu \notin \mathcal{L}(\mathcal{A})$.

442 ► **Lemma 4.11.** Let μ be a word of length n . There exist constants E, K such that if
 443 $+\infty > d(\mu, \mathcal{L}(\mathcal{A})) \geq \varepsilon n$ and $n \geq E/\varepsilon$, then μ can be partitioned into $\mu = \mu_0 \mu_1 \dots \mu_K$
 444 such that for every $i = 0, \dots, K$, μ_i contains at least $\frac{\varepsilon n}{E}$ disjoint occurrences of words
 445 $\nu_{i,1}, \nu_{i,2}, \dots$, each of length $\mathcal{O}(1/\varepsilon)$, such that for any choice of j_0, j_1, \dots, j_K , the sequence
 446 $(\nu_{0,j_0}, \nu_{1,j_1}, \dots, \nu_{K,j_K})$ is a blocking sequence for \mathcal{A} .

447 4.2 An efficient property tester

448 In this section, we show that for any regular language L and any small enough $\varepsilon > 0$, there
 449 is an ε -property tester for L that uses $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.

450 ► **Theorem 4.12.** For any NFA \mathcal{A} and any small enough $\varepsilon > 0$, there exists an ε -property
 451 tester for $\mathcal{L}(\mathcal{A})$ that uses $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries.

452 The property tester behind this theorem uses the property tester for strongly connected NFAs
 453 as a subroutine, and its correctness is based on Lemma 4.11, an extension of Lemma 3.5 to
 454 blocking sequences.

455 Recall that by Lemma 4.10, if we can find $N = |\mathcal{A}|$ disjoint blocking sequences in μ , then
 456 we know $\mu \notin L$, and the tester can reject. By Lemma 4.11, if we find N disjoint $\nu_{i,j}$'s in
 457 each of the μ_i , then we have found N disjoint blocking sequences. Since all of the $\nu_{i,j}$ have
 458 length at most $\mathcal{O}(1/\varepsilon)$ and each μ_i contains $\Omega(\varepsilon n)$ of them, we can adapt the algorithm of
 459 Theorem 3.3 by tweaking the constants so that it has a constant probability of finding at
 460 least one $\nu_{i,j}$ in a fixed μ_i , using $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries. By repeating $\mathcal{O}(N)$ times for each
 461 of the $K + 1$ possible values of i , the algorithm finds the N $\nu_{i,j}$ for each i with probability at
 462 least $2/3$, and uses a total of $\mathcal{O}(\log(\varepsilon^{-1})/\varepsilon)$ queries, as N and K are constants.

4.3 Lower bound

In order to characterise hard languages for all automata, we define a partial order \preceq on sequences of positional factors. It is an extension of the factor order on blocking factors. It will let us define *minimal blocking sequences*, which we use to characterize the complexity of testing a language.

► **Definition 4.13** (Minimal blocking sequence). *Let $\sigma = (\mu_1, \dots, \mu_k)$ and $\sigma' = (\mu'_1, \dots, \mu'_t)$ be blocking sequences of \mathcal{A} . We have $\sigma \preceq \sigma'$ if there exists a sequence of indices $i_1 \leq i_2 \leq \dots \leq i_k$ such that μ_j is a factor of μ'_{i_j} for all $j = 1, \dots, k$.*

A blocking sequence σ of \mathcal{A} is minimal if it is a minimal element of \preceq among blocking sequences of \mathcal{A} . The set of minimal blocking sequences of \mathcal{A} is written $MBS(\mathcal{A})$.

► **Remark 4.14.** If $\sigma \preceq \sigma'$ and σ is a blocking sequence for an SCC-path π then σ' is also a blocking sequence for π .

To prove a lower bound on the number of queries necessary to test a language when $MBS(\mathcal{A})$ is infinite, we present a reduction to the strongly connected case. Under the assumption that \mathcal{A} has infinitely many minimal blocking sequences, we exhibit a portal P of \mathcal{A} with infinitely many minimal blocking factors and “isolate it” by constructing two sequences of timed factors σ_l and σ_r such that for all μ , $\sigma_l, (\mu), \sigma_r$ is blocking for \mathcal{A} if and only if μ is a blocking factor of P . Then we reduce the problem of testing the language of this portal to the problem of testing the language of P .

To define “isolating P ” formally, we define the left (and right) effect of a sequence on an SCC-path. Intuitively, the left effect of a sequence σ on an SCC-path π is the index of the first portal in π where a run can be after reading σ , because all previous portals have been blocked. The right effect represents the same in reverse, starting from the end of the run.

More formally, the *left effect* of a sequence σ on an SCC-path $\pi = P_0 \xrightarrow{a_1} \dots P_k$ is the largest index i such that the sequence is blocking for $P_0 \xrightarrow{a_1} \dots P_i$ (-1 if there is no such i). We denote it by $(\sigma \gg \pi)$. Similarly, the *right effect* of a sequence on π is the smallest index i such that the sequence is blocking for $P_i \xrightarrow{a_{i+1}} \dots P_k$ ($k+1$ if there is no such i); we denote it by $(\pi \ll \sigma)$.

► **Remark 4.15.** A sequence σ is blocking for an SCC-path $\pi = P_0 \xrightarrow{a_1} \dots P_k$ if and only if $(\sigma \gg \pi) = k$, if and only if $(\pi \ll \sigma) = 0$.

Also, given two sequences σ_l, σ_r , the sequence $\sigma_l \sigma_r$ is blocking for π if and only if $(\sigma_l \gg \pi) > (\pi \ll \sigma_r)$.

For the next lemma we define a partial order on portals: $P \preceq P'$ if all blocking factors of P' are also blocking factors of P . We write \succeq for the reverse relation, \simeq for the equivalence relation $\preceq \cap \succeq$ and $\not\simeq$ for the complement relation of \simeq .

Additionally, given an SCC-path $\pi = P_0 \xrightarrow{a_1} \dots P_k$ and two sequences of positional words σ_l, σ_r , we say that the portal P_i *survives* (σ_l, σ_r) if $(\sigma_l \gg \pi) < i < (\pi \ll \sigma_r)$.

► **Definition 4.16.** *Let P be a portal and σ_l and σ_r sequences of positional words.*

We define three properties that those objects may have:

P1 $\sigma_l \sigma_r$ is not blocking for \mathcal{A}

P2 P has infinitely many minimal blocking factors

P3 for all accepting SCC-path π in \mathcal{A} , every portal in π which survives (σ_l, σ_r) is \simeq -equivalent to P .

We use the existence of P , σ_l and σ_r with those properties as an intermediate step to show that if $MBS(\mathcal{A})$ is infinite then $\mathcal{L}(\mathcal{A})$ is hard.

508 ► **Lemma 4.17.** *If \mathcal{A} has infinitely many minimal blocking sequences, then there exist a*
 509 *portal P and sequences σ_l and σ_r satisfying properties P1, P2 and P3.*

510 ► **Lemma 4.18.** *If there exist $P = s, x \rightsquigarrow t, y$ and σ_l, σ_r satisfying properties P1, P2 and*
 511 *P3 then $\mathcal{L}(\mathcal{A})$ is hard.*

512 ► **Proposition 4.19.** *If \mathcal{A} has infinitely many minimal blocking sequences, then $\mathcal{L}(\mathcal{A})$ is hard.*

513 **Proof.** We combine Lemmas 4.17 and 4.18. ◀

514 **5 Trivial and Easy languages**

515 **5.1 Upper bound for easy languages**

516 We establish that if $\text{MBS}(\mathcal{A})$ is finite then \mathcal{A} is testable with $\mathcal{O}(1/\varepsilon)$ queries.

517 ► **Lemma 5.1.** *If \mathcal{A} has finitely many minimal blocking sequences, then there is a tester for*
 518 *$\mathcal{L}(\mathcal{A})$ using $\mathcal{O}(1/\varepsilon)$ queries.*

519 **Proof sketch.** The full proof is in Appendix C. As the number of minimal blocking sequences
 520 of \mathcal{A} is bounded, so is the maximum length of a minimal blocking sequence. Let K be the
 521 upper bound on the length and P the bound on the number of minimal blocking sequences.

522 We infer from the fact that there are finitely many blocking sequences that if a word μ is
 523 ε -far from the language of \mathcal{A} then it must contain $\mathcal{O}(\varepsilon|\mu|)$ disjoint occurrence of the same
 524 minimal sequence σ .

525 Since each positional word in this sequence has length at most K , by sampling $\mathcal{O}(\frac{1}{\varepsilon})$
 526 factors of length K uniformly at random, we can show a positive constant lower bound on
 527 the probability to find σ . We can repeat this step to obtain a probability $> 1/2$ of finding
 528 $|\mathcal{A}|$ occurrences the sequence σ . This witnesses that $\mu \notin \mathcal{L}(\mathcal{A})$, by Lemma 4.10. ◀

529 This already gives us a clear dichotomy: all languages either require $\Theta(\log(\varepsilon^{-1})/\varepsilon)$ queries
 530 to be tested, or can be tested with $\mathcal{O}(1/\varepsilon)$ queries.

531 **5.2 Separation between trivial and easy languages**

532 At this point we can revisit the class of *trivial* regular languages identified by Alon et al. [5].

533 An example of a trivial language is L_2 consisting of words containing at least one a over
 534 the alphabet $\{a, b\}$. For any word u , replacing any letter by a yields a word in L_2 , hence
 535 $d(u, L_2) \leq 1$. Therefore, for $n > 1/\varepsilon$, no word of length n is ε -far from L_2 , and the trivial
 536 property tester that answers “yes” without sampling any letter is correct.

537 We give a combinatorial characterization of trivial languages (previously identified by
 538 Alon et al. [5]) based on their set of blocking factors. Our notion of *trivial* languages coincides
 539 with the notion of trivial languages of Alon et al. [5], recalled below.

540 ► **Definition 5.2** ([5, Definition 3.1]). *A language L is non-trivial if there exists a constant*
 541 *$\varepsilon_0 > 0$, so that for infinitely many values of n the set $L \cap \Sigma^n$ is non-empty, and there exists*
 542 *a word $w \in \Sigma^n$ so that $d(w, L) \geq \varepsilon_0 n$.*

543 We show that triviality is equivalent to having no blocking sequence. Recall that we focus
 544 on infinite languages, since we know that all finite ones are trivial (Remark 2.3).

545 ► **Lemma 5.3.** *$\text{MBS}(\mathcal{A})$ is empty if and only if $L = \mathcal{L}(\mathcal{A})$ is trivial.*

546 **Proof sketch.** The full proof is in Appendix C. Suppose $\text{MBS}(\mathcal{A})$ is empty. By Lemma 4.11,
 547 for all $\varepsilon > 0$ long enough words cannot be ε -far from L , hence L is trivial.

548 Suppose $\text{MBS}(\mathcal{A})$ contains a blocking sequence $\sigma = (\mu_1, \dots, \mu_k)$. For all $N \in \mathbb{N}$ we
 549 construct a word w_N containing $k + 1 + N$ factors μ_1, \dots, μ_k , plus a
 550 few letters for padding. This word has length $\mathcal{O}(N)$. Furthermore, by changing N letters we
 551 always obtain a word with $k + 1$ disjoint occurrences of σ , and thus not in L by Lemma C.4.
 552 As a result, $d(w_N, L) > N$ and thus w_N is ε -far from L for some ε independent of N . ◀

553 It is easy to see that if a language is trivial in the above sense, then for large enough
 554 input length n , membership in L only depends n , and the algorithm does not need to query
 555 the input. Alon et al. [5] show that if a language is non-trivial, then testing it requires $\Omega(1/\varepsilon)$
 556 queries for small enough $\varepsilon > 0$. As a corollary of that lower bound, we obtain that if $\text{MBS}(\mathcal{A})$
 557 is non-empty, then testing $\mathcal{L}(\mathcal{A})$ requires $\Omega(1/\varepsilon)$ queries.

558 6 Hardness of classifying

559 In the previous sections, we have shown that regular languages could be classified in three
 560 classes, which characterise their query complexity. In this section, we investigate the
 561 computational complexity of checking which class of the trichotomy the language of a given
 562 automaton belongs to. We formalize this question as the following decision problems:

563 ▶ **Problem 6.1** (Triviality problem). Given an finite automaton \mathcal{A} , is $\mathcal{L}(\mathcal{A})$ trivial?

564 ▶ **Problem 6.2** (Easiness problem). Given an finite automaton \mathcal{A} , is $\mathcal{L}(\mathcal{A})$ easy?

565 ▶ **Problem 6.3** (Hardness problem). Given an finite automaton \mathcal{A} , is $\mathcal{L}(\mathcal{A})$ hard?

566 We show that, our combinatorial characterization based on minimal blocking sequences
 567 is effective, in the sense that all three problems are decidable. However, it does not lead to
 568 efficient algorithms, as all three problems are PSPACE-complete.

569 ▶ **Theorem 6.4.** *The triviality, easiness and hardness problems are all PSPACE-complete,*
 570 *even for strongly connected NFAs.*

571 The proof is presented in Appendix F. The upper bound requires us to give another
 572 characterisation of hard languages, as it is not clear that the one we provided before can
 573 be turned into a PSPACE algorithm. The hardness is proven by a reduction from NFA
 574 universality. We transform a given automaton so that it remains universal if it was at the
 575 beginning, but any word out of the language yields an infinite family of minimal blocking
 576 factors in the new automaton.

577 7 Conclusion

578 We presented an effective classification of regular languages in three classes, each associated
 579 with an optimal query complexity for property testing. We thus close a line of research
 580 aiming to determine the optimal complexity of regular languages. All our results are with
 581 respect to the Hamming distance. We conjecture that they can be adapted to the edit
 582 distance. We use non-deterministic automata to represent regular languages. A natural open
 583 question is the complexity of classifying languages represented by *deterministic* automata.

584 ——— **References** ———

- 585 1 Alfred V. Aho and John E. Hopcroft. *The Design and Analysis of Computer Algorithms*.
586 Addison-Wesley Longman Publishing Co., Inc., 1974.
- 587 2 Maryam Aliakbarpour, Ilias Diakonikolas, and Ronitt Rubinfeld. Differentially private identity
588 and equivalence testing of discrete distributions. In *International Conference on Machine*
589 *Learning, ICML*, 2018. URL: <https://proceedings.mlr.press/v80/aliakbarpour18a.html>.
- 590 3 Noga Alon, Richard A. Duke, Hanno Lefmann, Vojtech Rödl, and Raphael Yuster. The
591 algorithmic aspects of the regularity lemma. *Journal of Algorithms*, 16(1):80–109, 1994.
592 doi:10.1006/JAGM.1994.1005.
- 593 4 Noga Alon, Eldar Fischer, Michael Krivelevich, and Mario Szegedy. Efficient testing of large
594 graphs. *Combinatorica*, 20(4):451–476, 2000. doi:10.1007/s004930070001.
- 595 5 Noga Alon, Michael Krivelevich, Ilan Newman, and Mario Szegedy. Regular languages are
596 testable with a constant number of queries. *SIAM Journal on Computing*, 30(6):1842–1862,
597 2001. doi:10.1109/SFFCS.1999.814641.
- 598 6 Noga Alon and Asaf Shapira. Every monotone graph property is testable. *SIAM Journal of*
599 *Computing*, 38(2):505–522, 2008. doi:10.1137/050633445.
- 600 7 Antoine Amarilli, Louis Jachiet, and Charles Paperman. Dynamic membership for regular
601 languages. In Nikhil Bansal, Emanuela Merelli, and James Worrell, editors, *48th International*
602 *Colloquium on Automata, Languages, and Programming, ICALP 2021, July 12–16, 2021,*
603 *Glasgow, Scotland (Virtual Conference)*, volume 198 of *LIPICs*, pages 116:1–116:17. Schloss
604 Dagstuhl - Leibniz-Zentrum für Informatik, 2021. URL: [https://doi.org/10.4230/LIPICs.](https://doi.org/10.4230/LIPICs.ICALP.2021.116)
605 *ICALP.2021.116*, doi:10.4230/LIPICs.ICALP.2021.116.
- 606 8 Ajesh Babu, Nutan Limaye, Jaikumar Radhakrishnan, and Girish Varma. Streaming algorithms
607 for language recognition problems. *Theoretical Computer Science*, 494:13–23, 2013.
- 608 9 Gabriel Bathie and Tatiana Starikovskaya. Property testing of regular languages with ap-
609 plications to streaming property testing of visibly pushdown languages. In *International*
610 *Colloquium on Automata, Languages, and Programming, ICALP*, 2021. doi:10.4230/LIPICs.
611 *ICALP.2021.119*.
- 612 10 Manuel Blum, Michael Luby, and Ronitt Rubinfeld. Self-testing/correcting with applications
613 to numerical problems. *Journal of Computer and System Sciences*, 47(3):549–595, 1993.
614 doi:10.1016/0022-0000(93)90044-w.
- 615 11 Ilias Diakonikolas and Daniel M Kane. A new approach for testing properties of discrete
616 distributions. In *IEEE Symposium on Foundations of Computer Science, FOCS*, pages 685–694.
617 IEEE, 2016. doi:10.1109/FOCS.2016.78.
- 618 12 Yevgeniy Dodis, Oded Goldreich, Eric Lehman, Sofya Raskhodnikova, Dana Ron, and Alex
619 Samorodnitsky. Improved testing algorithms for monotonicity. In *International Workshop*
620 *on Randomization and Approximation Techniques in Computer Science, RANDOM*, pages
621 97–108. Springer, 1999. doi:10.1007/978-3-540-48413-4_10.
- 622 13 Funda Ergün, Sampath Kannan, S.Ravi Kumar, Ronitt Rubinfeld, and Mahesh Viswanathan.
623 Spot-checkers. *Journal of Computer and System Sciences*, 60(3):717–751, 2000. doi:10.1006/
624 *jcss.1999.1692*.
- 625 14 Eldar Fischer, Frédéric Magniez, and Tatiana Starikovskaya. Improved bounds for testing
626 Dyck languages. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on*
627 *Discrete Algorithms*, pages 1529–1544. SIAM, 2018.
- 628 15 Nathanaël François, Frédéric Magniez, Michel de Rougemont, and Olivier Serre. Streaming
629 property testing of visibly pushdown languages. In *Proceedings of the 24th Annual European*
630 *Symposium on Algorithms*, volume 57 of *LIPICs*, pages 43:1–43:17, 2016. doi:10.4230/LIPICs.
631 *ESA.2016.43*.
- 632 16 Moses Ganardi, Danny Hucce, and Markus Lohrey. Randomized sliding window algorithms
633 for regular languages. In *Proceedings of the 45th International Colloquium on Automata,*
634 *Languages, and Programming*, volume 107 of *LIPICs*, pages 127:1–127:13, 2018. doi:10.4230/
635 *LIPICs.ICALP.2018.127*.

- 636 **17** Moses Ganardi, Danny Hucke, Markus Lohrey, and Tatiana Starikovskaya. Sliding Window
637 Property Testing for Regular Languages. In Pinyan Lu and Guochuan Zhang, editors, *30th*
638 *International Symposium on Algorithms and Computation (ISAAC 2019)*, volume 149 of
639 *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 6:1–6:13, Dagstuhl, Germany,
640 2019. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. URL: [https://drops.dagstuhl.de/
641 entities/document/10.4230/LIPIcs.ISAAC.2019.6](https://drops.dagstuhl.de/entities/document/10.4230/LIPIcs.ISAAC.2019.6), doi:10.4230/LIPIcs.ISAAC.2019.6.
- 642 **18** Oded Goldreich. *Introduction to property testing*. Cambridge University Press, 2017. doi:
643 10.1017/9781108135252.
- 644 **19** Oded Goldreich, Shari Goldwasser, and Dana Ron. Property testing and its connection to
645 learning and approximation. *Journal of the ACM*, 45(4):653–750, jul 1998. doi:10.1145/
646 285055.285060.
- 647 **20** Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *The*
648 *collected works of Wassily Hoeffding*, pages 409–426, 1994.
- 649 **21** Rahul Jain and Ashwin Nayak. The space complexity of recognizing well-parenthesized
650 expressions in the streaming model: The index function revisited. *IEEE Transactions on*
651 *Information Theory*, 60(10):6646–6668, Oct 2014. doi:10.1109/TIT.2014.2339859.
- 652 **22** Andreas Krebs, Nutan Limaye, and Srikanth Srinivasan. Streaming algorithms for recognizing
653 nearly well-parenthesized expressions. In *Proc. of MFCS 2011*, volume 6907 of *LNCS*, pages
654 412–423. Springer, 2011. doi:10.1007/978-3-642-22993-0_38.
- 655 **23** Frédéric Magniez and Michel de Rougemont. Property testing of regular tree languages.
656 *Algorithmica*, 49(2):127–146, 2007. doi:10.1007/s00453-007-9028-3.
- 657 **24** Frédéric Magniez, Claire Mathieu, and Ashwin Nayak. Recognizing well-parenthesized
658 expressions in the streaming model. *SIAM J. Comput.*, 43(6):1880–1905, 2014. doi:
659 10.1137/130926122.
- 660 **25** Liam Paninski. A coincidence-based test for uniformity given very sparsely sampled discrete
661 data. *IEEE Transactions on Information Theory*, 54(10):4750–4755, 2008. doi:10.1109/TIT.
662 2008.928987.
- 663 **26** Michal Parnas, Dana Ron, and Ronitt Rubinfeld. Testing membership in parenthesis languages.
664 *Random Struct. Algorithms*, 22(1):98–138, 2003. doi:10.1002/rsa.10067.
- 665 **27** Jean-Éric Pin, editor. *Handbook of Automata Theory*. European Mathematical Society
666 Publishing House, Zürich, Switzerland, 2021. URL: <https://doi.org/10.4171/Automata>,
667 doi:10.4171/AUTOMATA.
- 668 **28** Ronitt Rubinfeld and Madhu Sudan. Robust characterizations of polynomials with applications
669 to program testing. *SIAM Journal on Computing*, 25(2):252–271, 1996. doi:10.1137/
670 S0097539793255151.
- 671 **29** Mosaad Al Thokair, Minjian Zhang, Umang Mathur, and Mahesh Viswanathan. Dynamic
672 race detection with $O(1)$ samples. *Proc. ACM Program. Lang.*, 7(POPL):1308–1337, 2023.
673 doi:10.1145/3571238.
- 674 **30** Andrew Chi-Chin Yao. Probabilistic computations: Toward a unified measure of complexity.
675 In *Symposium on Foundations of Computer Science, SFCS*, pages 222–227. IEEE, 1977.
676 doi:10.1109/SFCS.1977.24.

677 **A** The Case of Strongly Connected NFAs

678 **A.1** Positional words, blocking factors and strongly connected NFAs

679 We establish some properties of positional words that will help us ensure that we are creating
680 well-formed positional words. We fix a strongly connected automaton $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$
681 and sets of states $Q_0, \dots, Q_{\lambda(\mathcal{A})-1}$ as in Fact 2.4, with the initial state q_0 in Q_0 .

682 **► Lemma A.1.** *A positional word $\tau = (i : u)$ is a blocking factor for \mathcal{A} iff for every states*
683 *$p \in Q_i, q \in Q$, we have $p \not\stackrel{u}{\rightarrow} q$.*

684 **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
685 not blocking, i.e., there exists $\mu \in \mathcal{TL}(\mathcal{A})$ such that $\tau \preceq \mu$. As \mathcal{A} is strongly connected,
686 there exist words w, w' such that $q_0 \stackrel{w}{\rightarrow} p$ and $q \stackrel{w'}{\rightarrow} q_f$ for some $q_f \in F$. By Fact 2.4, the
687 positional word $(0 : wuw')$ contains τ as a factor. Since wuw' labels a path from q_0 to q_f ,
688 we have $(0 : wuw') \in \mathcal{TL}(\mathcal{A})$, and τ is not blocking.

689 For the converse, assume that τ is not blocking: we show that there exists two states
690 $p \in Q_i, q \in Q$ such that $p \stackrel{u}{\rightarrow} q$. As τ is non-blocking, there exists a positional word $\mu = (0 : w)$
691 such that $\tau \preceq \mu$ and there exists a final state r such that $q_0 \stackrel{\mu}{\rightarrow} r$, and equivalently, $q_0 \stackrel{w}{\rightarrow} r$.
692 Since $\tau \preceq (0 : w)$, there exists words v, v' such that $w = vv'$ and the length of v is equal to
693 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
694 particular, we have $p \stackrel{u}{\rightarrow} q$. It only remains to show that p is in Q_i : this follows by Fact 2.4
695 since $|v| = i \pmod{\lambda}$ and $q_0 \in Q_0$. ◀

696 Next, we show that the Hamming distance between u and $\mathcal{L}(\mathcal{A})$ is the same as the
697 (Hamming) distance between $(0 : u)$ and $\mathcal{TL}(\mathcal{A})$.

698 **▷ Claim A.2.** For any word $u \in \Sigma^*$, we have $d(u, \mathcal{L}(\mathcal{A})) = d((0 : u), \mathcal{TL}(\mathcal{A}))$.

699 **Proof.** It suffices to observe that for all $v \in \mathcal{L}(\mathcal{A})$, the i th letter differs in u and v if and only
700 if it differs in $(0 : u)$ and $(0 : v)$. Hence $d(u, v) = d((0 : u), (0 : v))$ for all $v \in \mathcal{L}(\mathcal{A})$, and
701 therefore $d(u, \mathcal{L}(\mathcal{A})) = d((0 : u), \mathcal{TL}(\mathcal{A}))$. ◀

702 The above claim allows us to interchangeably use the statements “ u is ε -far from $\mathcal{L}(\mathcal{A})$ ” and
703 “ $(0 : u)$ is ε -far from $\mathcal{TL}(\mathcal{A})$ ”.

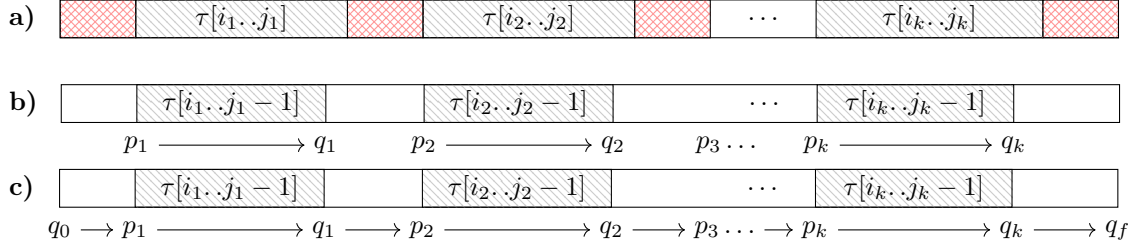
704 **A.2** Upper Bound for Strongly Connected NFAs

705 Alon et al. [5, Lemma 2.6] first noticed that if a word u is ε -far from $\mathcal{L}(\mathcal{A})$, then it contains
706 $\Omega(\varepsilon n)$ short factors that witness the fact that u is not in $\mathcal{L}(\mathcal{A})$. We start by translating
707 the lemma of Alon et al. on “short witnesses” to the framework of blocking factors. More
708 precisely, we show that if u is ε -far from $\mathcal{L}(\mathcal{A})$, then $(0 : u)$ contains many disjoint blocking
709 factors.

710 **► Lemma A.3.** *Let $\varepsilon > 0$, let u be a word of length $n \geq 6m^2/\varepsilon$ and assume that $\mathcal{L}(\mathcal{A})$*
711 *contains at least one word of length n . If $\tau = (0 : u)$ is ε -far from $\mathcal{TL}(\mathcal{A})$, then τ contains at*
712 *least $\varepsilon n/(6m^2)$ disjoint blocking factors.*

713 **Proof.** We build a set \mathcal{P} of disjoint blocking factors of τ as follows: we process u from left to
714 right, starting at index $i_1 = \rho$. Next, at iteration t , set j_t to be the smallest integer greater
715 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
716 no such integer, we stop the process. Otherwise, we add $\tau[i_t..j_t + \rho - 1]$ to the set \mathcal{P} , and
717 iterate starting from the index $i_{t+1} = j_t + \rho$.

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



■ **Figure 3 a)** The decomposition process returns k factors $\tau[i_1, j_1], \dots, \tau[i_k, j_k]$ (represented as diagonally hatched in gray regions), separated together and with the start of the word by padding regions of $\rho - 1$ letters (red cross-hatched regions). **b)** After removing the last letter, each previously blocking factor now labels a path 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.

722 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
 723 not blocking, and therefore $\tau[i_t..j_t - 1]$ labels a run from some state $p_t \in Q_k$ to some state
 724 $q_t \in Q_{k+j_t-i_t}$. Therefore, using the strong connectivity of \mathcal{A} and Fact 2.4, we can substitute
 725 the letters in $\tau[j_t..j_t + \rho - 1]$ to obtain a factor that labels a path from q_t to p_{t+1} . After this
 726 transformation, the word $\tau[i_t..j_t + \rho - 1]$ labels a path from p_t to p_{t+1} . Using the ρ letters
 727 at the start and the end of the word, we add transitions from an initial state to p_1 and from
 728 q_k to a final state: the assumption that $\mathcal{L}(\mathcal{A})$ contains a word of length n ensures that Q_n
 729 contains a final state, hence this is always possible. The resulting word is in $\mathcal{TL}(\mathcal{A})$ and was
 730 obtained from τ using $(k+1)\rho \leq 3(k+1)m^2$ substitutions. As τ is ε -far from $\mathcal{TL}(\mathcal{A})$, we
 731 obtain the following bound on k :

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

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

736 ► **Lemma 3.5.** *Let $\varepsilon > 0$, let u be a word of length $n \geq 6m^2/\varepsilon$ and assume that $\mathcal{L}(\mathcal{A})$
 737 contains at least one word of length n . If u is ε -far from $\mathcal{L}(\mathcal{A})$, then the positional word
 738 $(0 : u)$ contains at least $\varepsilon n/(12m^2)$ disjoint blocking factors of length at most $12m^2/\varepsilon$.*

739 **Proof.** Let u, \mathcal{A} be a word and an automaton satisfying the above hypotheses. By Lemma A.3,
 740 $(0 : u)$ contains at least $\varepsilon n/(6m^2)$ disjoint blocking factors. As these factors are disjoint, at
 741 most half of them (that is, $\varepsilon n/(12m^2)$ of them) can have length greater than $12m^2/\varepsilon$, as the
 742 sum of their lengths cannot exceed n . \blacktriangleleft

743 ► **Lemma 3.6.** *In the last Else block, if u is ε -far from $\mathcal{L}(\mathcal{A})$, then Algorithm 1 rejects with
 744 probability at least $2/3$.*

745 **Proof.** Assume that u is ε -far from $\mathcal{L}(\mathcal{A})$. As we are in the last **Else** block of Algorithm 1,
 746 $\mathcal{L}(\mathcal{A}) \cap \Sigma^n$ is not empty (i.e., $\mathcal{L}(\mathcal{A})$ contains a word of length n) and $n \geq \beta$, therefore the
 747 conditions of Lemma 3.5 are satisfied. Let \mathcal{B} denote the set of minimal blocking factors in
 748 $(0 : u)$ given by Lemma 3.5: we have $|\mathcal{B}| \geq n/\beta$. We conceptually divide the blocking factors
 749 in \mathcal{B} into different categories depending on their length: for $t = 0, \dots, T$, let B_t denote the
 750 subset of \mathcal{B} of blocking factors of length at most $\ell_t = 2^t$. We then carefully analyse the
 751 probability that randomly sampled factors of length $2\ell_t$ contains a blocking factor from B_t ,
 752 and show that over all t , at least one blocking factor is found with probability at least $2/3$.

753 \triangleright **Claim A.4.** If in a call to **SAMPLE**, the value i is such that there exists indices $l, r, l \leq i \leq r$,
 754 such that $(0 : u)[l, r]$ is a blocking factor of \mathcal{A} of length at most ℓ , then the factor η returned
 755 by the function is blocking for \mathcal{A} .

756 As the factors given by Lemma 3.5 are disjoint, the probability p_t that the factor returned
 757 by **SAMPLE** is blocking is lower bounded by

$$758 \quad p_t \geq \frac{1}{n} \sum_{\tau \in B_t} |\tau|$$

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

$$\begin{aligned}
 761 \quad (1-p) &= \prod_{t=0}^T (1-p_t)^{r_t} \leq \exp\left(-\sum_{t=0}^T p_t r_t\right) \\
 762 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{t=0}^T \frac{1}{\ell_t} \sum_{\tau \in B_t} |\tau|\right) \\
 763 \quad &= \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{\tau \in \mathcal{B}} |\tau| \sum_{t=\lceil \log |\tau| \rceil}^T 2^{-t}\right) \\
 764 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{\tau \in \mathcal{B}} |\tau| \cdot 2^{-\lceil \log |\tau| \rceil}\right) \\
 765 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \sum_{\tau \in \mathcal{B}} |\tau| \frac{1}{2^{|\tau|}}\right) \\
 766 \quad &= \exp\left(-\frac{2 \ln(3)\beta}{n} \cdot \frac{|\mathcal{B}|}{2}\right) \\
 767 \quad &\leq \exp\left(-\frac{2 \ln(3)\beta}{n} \cdot \frac{n}{2\beta}\right) \\
 768 \quad &\leq \exp(-\ln(3)) = 1/3
 \end{aligned}$$

769 It follows that $p \geq 2/3$, and Algorithm 1 satisfies Definition 2.2. \blacktriangleleft

770 **A.3 Lower Bound for strongly connected automata with infinitely many** 771 **blocking factors**

772 Let \mathcal{A} be a strongly connected NFA. Let $\widehat{\mathcal{A}}$ be an automaton recognising $\mathcal{TL}(\mathcal{A})$, as defined
 773 just after Definition 2.5. As \mathcal{A} is strongly connected, so is $\widehat{\mathcal{A}}$: from every reachable state we
 774 can go back to the initial state by following a path to the initial state of \mathcal{A} and adding cycles
 775 to adjust the modulo. In this whole section we will thus assume that \mathcal{A} recognises $\mathcal{TL}(\mathcal{A})$
 776 instead of just $\mathcal{L}(\mathcal{A})$.

777 A.3.1 The structure of $\text{MBF}(\mathcal{A})$

778 We prove that the set of minimal blocking factors of an automaton is a regular language,
779 albeit recognized by an automaton that is possibly exponentially larger than \mathcal{A} . We first
780 prove the result for blocking factors of the form $(i : u)$ for a fixed $i \in \mathbb{Z}/\lambda\mathbb{Z}$.

781 ► **Lemma A.5.** *Let $\mathcal{A} = (Q, \Sigma, \delta, I, F)$ be a strongly connected NFA with m states and let
782 $\lambda = \lambda(\mathcal{A})$. For every $i \in \mathbb{Z}/\lambda\mathbb{Z}$, the set of minimal blocking factors of \mathcal{A} of the form $(i : u)$ is
783 a regular language recognized by a NFA of size $2^{\mathcal{O}(m)}$.*

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

785 We first show that the set of *i-blocking factors* of \mathcal{A} , but not necessarily minimal ones, is
786 a regular language recognized by an NFA \mathcal{A}_i with $m + 1$ states.

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

$$790 \quad \forall p \in Q, \forall a \in \Sigma : \delta'(p, a) = \begin{cases} \{\perp\} & \text{if } \delta(p, a) = \emptyset, \\ \delta(p, a) & \text{otherwise.} \end{cases}$$

791 This automaton⁵ recognizes the set of *i-blocking factors* of \mathcal{A} and has size $\mathcal{O}(m)$. We can
792 then construct \mathcal{B}_i a deterministic automaton for the same language.

793 The result follows by using a standard construction: build an automaton that runs \mathcal{B}_i
794 on the input word, and \mathcal{B}_{i+1} on the word obtained by ignoring the first letter. It accepts if
795 the word is accepted by \mathcal{B}_i , and the words obtained respectively by ignoring the first and
796 last letters are rejected by respectively \mathcal{B}_{i+1} and \mathcal{B}_i . A word is accepted if and only if it is a
797 minimal blocking factor.

798 This yields the desired automaton, of size $2^{\mathcal{O}(m)}$.
799 ◀

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

801 ► **Corollary A.6.** *Let \mathcal{A} be an NFA with m states. The set of minimal blocking factors of \mathcal{A}
802 is a regular language recognized by an NFA of size $2^{\mathcal{O}(m)}$.*

803 Therefore, if $\text{MBF}(\mathcal{A})$ is infinite, we can use Kleene's lemma to find an infinite family of
804 minimal blocking factors with a shared structure $\{\phi\nu^r\chi, r \in \mathbb{N}\}$. We will use this property
805 later, when proving a lower bound against the language of automata with infinitely many
806 blocking factors.

807 ► **Lemma A.7.** *If $\text{MBF}(\mathcal{A})$ is infinite, then there exist positional words $\phi, \nu_+, \nu_-, z, \chi$ such
808 that:*

- 809 1. *the words ν_+ and ν_- have the same length,*
- 810 2. *there exists an index $i_* \in \mathbb{Z}/\lambda\mathbb{Z}$ and a state $q_* \in Q_{i_*}$ such that for every integer $r \geq 1$,*
811 *the positional word $\tau_{-,r} = \phi(\nu_-)^r z$ is blocking for \mathcal{A} , and for every $s < r$, we have*

$$812 \quad q_* \xrightarrow{\tau_{+,r,s}} q_* \text{ where } \tau_{+,r,s} = \phi(\nu_-)^s \nu_+ (\nu_-)^{r-1-s} \chi.$$

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

⁵ 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.

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

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

817 As the set of minimal i_* -blocking factors is an infinite regular language recognized by
818 an NFA of size $S = 2^{\mathcal{O}(m)}$, by Kleene's Lemma, there exist positional words τ, μ, η , each of
819 length at most S with $|\mu| \geq 1$, such that for any non-negative integer k , $\tau\mu^k\eta$ is a minimal
820 i_* -blocking factor. We can assume w.l.o.g. that neither τ nor η is empty, otherwise we set
821 their value to μ : after this modification, $\tau\mu^k\eta$ is still a minimal i_* -blocking factor for every
822 $k \geq 0$.

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

$$826 \quad p \xrightarrow{\tau\mu^{k_0}} p_1 \xrightarrow{\mu^{k_1}} p_1.$$

827 Note that, by Fact 2.4, $p_1 \xrightarrow{\mu^{k_1}} p_1$ implies that $k_1 \cdot |\mu| = 0 \pmod{\lambda}$.

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

$$831 \quad p_2 \xrightarrow{\mu^{k_2}} p_2 \xrightarrow{\mu^{k_3}\eta} q.$$

832 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
833 transitions starting from p_1 and p_2 labelled by μ , p_1 and p_2 belong to the same periodicity
834 class. Therefore, by Fact 2.4, as $K \geq \rho$ and $K \cdot |\mu| = 0 \pmod{\lambda}$, there exists a word ν_+ of
835 length $K \cdot |\mu|$ such that $p_1 \xrightarrow{\nu_+} p_2$. This choice of ϕ, ν_+, ν_- and χ satisfies all the conditions
836 of the lemma. \blacktriangleleft

837 A.3.2 Constructing a hard distribution

838 Let $\varepsilon > 0$ be sufficiently small and let n be a large enough integer. In what follows, m
839 denotes the number of states of \mathcal{A} . To construct the hard distribution \mathcal{D} , we will use an
840 infinite family of blocking factors that share a common structure, given by Lemma A.7.

841 Let S be the length of ν_- and ν_+ . The crucial property here is that $\tau_{-,r}$ and $\tau_{+,r,s}$
842 are very similar: they have the same length, differ in at most S letters, yet one of them is
843 blocking and the other is not.

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

846 Let π be a uniformly random bit. If $\pi = 1$, we will construct a positive instance
847 $\mu \in \mathcal{TL}(\mathcal{A})$, and otherwise the instance will be ε -far from $\mathcal{TL}(\mathcal{A})$ with high probability. We
848 divide the interval $[0..n-1]$ into $k = \varepsilon n$ intervals of length $\ell = 1/\varepsilon$, plus small initial and
849 final segments μ_i and μ_f of length $\mathcal{O}(\rho)$ to be specified later. For the sake of simplicity, we
850 assume that k and ℓ are integers and that λ divides ℓ . For $j = 1, \dots, k$, let a_j, b_j denote
851 the endpoints of the j -th interval. For each interval, we sample independently at random a
852 variable κ_j with the following distribution:

$$853 \quad \kappa_j = \begin{cases} t, & \text{w.p. } p_t = 3 \cdot 2^t S \varepsilon / \log((S\varepsilon)^{-1}) \text{ for } t = 1, 2, \dots, \log((S\varepsilon)^{-1}), \\ 0, & \text{w.p. } p_0 = 1 - \sum_{t=1}^{\log((S\varepsilon)^{-1})} p_t. \end{cases} \quad (3)$$

854 The event $\kappa_j > 0$ means that the j -th interval is filled with with $N \approx 2^{-\kappa_j}/\varepsilon$ “special”
 855 factors. When $\pi = 0$, these “special” factors will be minimal blocking factors $\tau_{-,r}$ for $r = 2^{\kappa_j}$,
 856 whereas when $\pi = 1$, they will instead be similar non-blocking factors $\tau_{+,r,s}$ for a uniformly
 857 random s : they will be hard to distinguish with few queries. On the other hand, the event
 858 $\kappa_j = 0$ means that the j -th interval contains no specific information. More precisely, we
 859 choose a positional word η_* of length ℓ such that $q_* \xrightarrow{w_*} q_*$: by Fact 2.4, this is possible as
 860 $\ell = 0 \pmod{\lambda}$. Then, if $\kappa_j = 0$, we set $\mu[a_j..b_j] = \eta_*$, regardless of the value of π .

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

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

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

868 A.3.3 Properties of the distribution \mathcal{D}

869 We first state the important properties of the distribution \mathcal{D}

870 ► **Observation A.8.** *If ε is small enough, \mathcal{D} is well-defined, i.e. for every t between 0 and*
 871 *$\log((S\varepsilon)^{-1})$, we have $0 \leq p_t \leq 1$.*

872 ► **Observation A.9.** *If $\pi = 1$, then $\mu \in \mathcal{TL}(\mathcal{A})$.*

873 We recall Hoeffding’s inequality, which will be used in the next proof.

874 ► **Lemma A.10** ([20, Theorem 2]). *Let X_1, \dots, X_k be independent random variables such*
 875 *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 $t > 0$,*
 876 *we have*

$$877 \quad \mathbb{P}(\mathbb{E}[S] - S \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^k (b_i - a_i)^2}\right).$$

878 ► **Lemma 3.8.** *Conditioned on $\pi = 0$, the probability of the event $\mathcal{F} = \{u \text{ is } \varepsilon\text{-far from}$*
 879 *$\mathcal{TL}(\mathcal{A})\}$ goes to 1 as n goes to infinity.*

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

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

889 ► **Claim A.11.** Let $B = \sum_{j=1}^k B_j$, and let $E = \mathbb{E}[B]$. We have $E \geq 2\varepsilon n$.

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

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

896 ◀

897 We will now show that $\mathbb{P}(B < \varepsilon n)$ goes to 0 as n goes to infinity. By Claim A.11, we have
 898 $B < \varepsilon n \Rightarrow E - B \geq \varepsilon n$, and therefore $\mathbb{P}(B < \varepsilon n) \leq \mathbb{P}(E - B \geq \varepsilon n)$. The random variable
 899 B is the sum of k independent random variables, each taking values between 0 and $1/(S\varepsilon)$.
 900 Therefore, by Hoeffding's Inequality (Lemma A.10), we have

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

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

905 **► Lemma 3.10.** *Let T be a deterministic tester with perfect completeness (i.e. one sided*
 906 *error, always accepts $\tau \in \mathcal{TL}(\mathcal{A})$) and let q_j denote the number of queries that it makes in*
 907 *the j -th interval. Conditioned on the event $\mathcal{M} = \{\forall j \text{ s.t. } \kappa_j > 0, q_j < 2^{\kappa_j}\}$, the probability*
 908 *that T accepts u is 1.*

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

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

921 **► Lemma 3.11.** *Let T be a deterministic tester, let q_j denote the number of queries that*
 922 *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.*
 923 $\sum_j q_j \leq \frac{1}{72} \cdot \log(\varepsilon^{-1})/\varepsilon$. *The probability of the event $\mathcal{M} = \{\forall j \text{ s.t. } \kappa_j > 0, q_j < 2^{\kappa_j}\}$ is at*
 924 *least 11/12.*

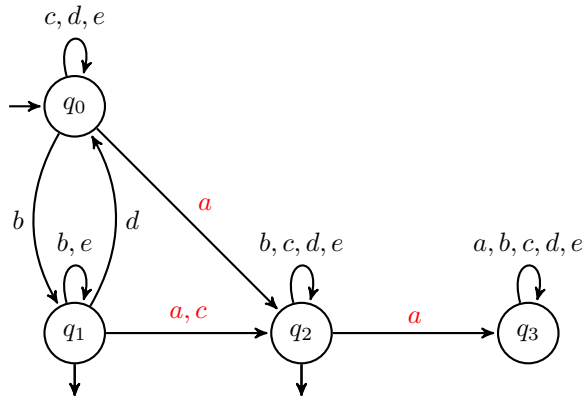
925 **Proof.** We show that the probability of $\overline{\mathcal{M}}$, the complement of \mathcal{M} , is at most $1/12$. We
 926 have:

$$\begin{aligned}
 927 \quad \mathbb{P}(\overline{\mathcal{M}}) &= \mathbb{P}(\exists j : \kappa_j > 0 \wedge q_j \geq 2^{\kappa_j}) \\
 928 \quad &\leq \sum_j \mathbb{P}(\kappa_j > 0 \wedge q_j \geq 2^{\kappa_j}) && \text{by union bound} \\
 929 \quad &\leq \sum_j \sum_{t=1}^{\lfloor \log q_j \rfloor} p_t \\
 930 \quad &= \sum_j \sum_{t=1}^{\lfloor \log q_j \rfloor} \frac{3 \cdot 2^t \varepsilon}{\log(S/\varepsilon)} && \text{by def. of } p_t \\
 931 \quad &\leq \frac{3\varepsilon}{\log(S/\varepsilon)} \sum_j \sum_{t=1}^{\lfloor \log q_j \rfloor} 2^t \\
 932 \quad &\leq \frac{3\varepsilon}{\log(S/\varepsilon)} \sum_j 2q_j \\
 933 \quad &= \frac{3\varepsilon}{\log(S/\varepsilon)} \cdot \frac{2}{72} \cdot \frac{\log(1/\varepsilon)}{\varepsilon} \\
 934 \quad &\leq 1/12
 \end{aligned}$$

935

936 **B** Useful examples for blocking sequences

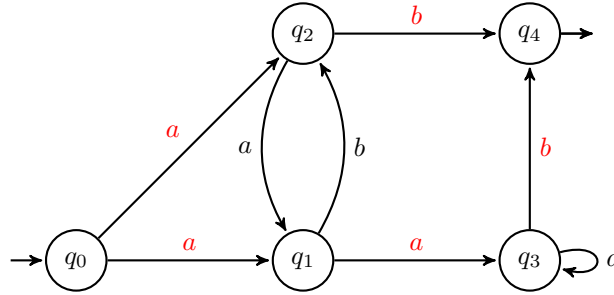
937 The goal of this appendix is to provide a few examples to guide the intuition of the reader
 938 through the definitions of Section 4.1.



939 **Figure 4** An automaton \mathcal{A}_2 recognizing the language $L_2 = [(c+d+e)^*b(b+e)^*d]^*a(b+c+d+e)^*$.

939 **► Example B.1.** Consider the automaton in Figure 4: it has two SCCs and a sink state.
 940 The minimal blocking factors of the first SCC are given by $B_1 = be^*c + a$, and $B_2 = \{a\}$
 941 for the second SCC. This automaton is easy to test: intuitively, a word that is ε -far from
 942 this language has to contain many a , as otherwise we can make it accepted by deleting all a ,

943 thanks to the second SCC. However, a is also a blocking factor of the first SCC, therefore, as
 944 soon as we find a 's in the word, we know that it is not in L_2 . The crucial facts here are that
 945 the set B_2 of minimal blocking factors of the second SCC is finite and it is a subset of B_1 :
 946 the infinite nature of B_1 is made irrelevant because any word far from the language contains
 947 many a 's.



■ **Figure 5** Automaton used for Example B.2.

948 ► **Example B.2.** Consider the automaton displayed in Figure 5. The lcm of the lengths of
 949 its simple cycles is $p = 2$. This automaton has six accepting SCC-paths, including

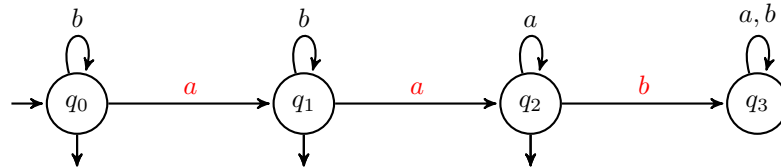
$$950 \quad \pi_1 = q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_1, 1 \rightsquigarrow q_1, 1 \xrightarrow{a} q_3, 0 \rightsquigarrow q_3, 0 \xrightarrow{b} q_4, 1 \rightsquigarrow q_4, 1$$

$$951 \quad \pi_2 = q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_2, 1 \rightsquigarrow q_1, 0 \xrightarrow{a} q_3, 1 \rightsquigarrow q_3, 0 \xrightarrow{b} q_4, 1 \rightsquigarrow q_4, 1$$

952 The language of the π_1 is $a(ba)^*a(a^2)^*b$. A blocking sequence for this SCC-path is
 953 $((0 : aa), (0 : b))$, which is in fact blocking for all of the SCC-paths.

954 On the other hand, $((0 : ab))$ is not blocking for π_1 , as $(0 : ab)$ is not a blocking factor
 955 for the portal $q_1, 1 \rightsquigarrow q_1, 1$. It is, however, a blocking sequence for π_3 . This is because if
 956 we enter the SCC $\{q_1, q_2\}$ through q_1 , a factor ab can only appear after an even number of steps,
 957 while if we enter through q_2 , it can only appear after an odd number of steps.

958 ► **Example B.3.** Consider the automaton \mathcal{A} displayed in Figure 6: it only has cycles of
 959 length 1, hence $p = 1$. Its states are totally ordered by $\leq_{\mathcal{A}}$.



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

960 Observe that the sequence $\sigma = ((0 : a), (0 : b))$ is a blocking sequence for the SCC-path
 961 $\pi = q_0, 0 \rightsquigarrow q_0, 0 \xrightarrow{a} q_1, 0 \rightsquigarrow q_1, 0 \xrightarrow{a} q_2, 0 \rightsquigarrow q_2, 0$. Indeed, a is blocking for the first two
 962 portals, and b for the third.

963 We can verify Lemma C.4 on this example: if a word u contains $|Q| = 4$ disjoint
 964 occurrences of σ , then in particular it must contain factors a, a and b in that order, hence u

965 is not in $\mathcal{L}(\mathcal{A})$. Two occurrences of blocking sequences would be enough here, but note that
 966 one occurrence is not enough: the word $(0 : aba)$ contains $((0 : a), (0 : b))$, yet it is in the
 967 language of π .

968 **C** Properties of blocking sequences

969 **► Lemma C.1.** *Let \mathcal{A} be an automaton and P a portal in \mathcal{A} . There is a strongly connected*
 970 *NFA with at most $p|\mathcal{A}|$ states that recognizes $L' = \mathcal{L}(P)$.*

971 **Proof.** Let S denote the SCC of s and t in \mathcal{A} , and let λ denote its periodicity. By definition
 972 of p , λ divides p : let k be the integer such that $p = \lambda k$. The automaton \mathcal{A}' for L' simulates
 973 the behavior of \mathcal{A} restricted to S starting from the state s , while keeping track of the number
 974 of letters read modulo p , starting from x . More precisely, let $Q_0, \dots, Q_{\lambda-1}$ be the partition of
 975 the states of S given by Fact 2.4. The set of states of \mathcal{A}' is given by

$$976 \quad Q' = \{(s', i + j\lambda) \mid s' \in Q_i, i = 0, \dots, \lambda - 1, j = 0, \dots, k - 1\}.$$

977 It is a subset of $Q \times \mathbb{Z}/p\mathbb{Z}$, hence it has cardinality at most $p|\mathcal{A}|$. The transitions in \mathcal{A}' are
 978 of the form $(s_1, i + j\lambda) \xrightarrow{(i,a)} (s_2, i + j\lambda + 1 \pmod{p})$ for any s_1, s_2 such that $s_1 \xrightarrow{a} s_2$ in \mathcal{A} .

979 Furthermore, \mathcal{A}' is strongly connected. Let i_1, i_2 be indices of periodicity classes of S ,
 980 and let $s_1 \in Q_{i_1}, s_2 \in Q_{i_2}$ be states of S . We show that for any $j_1, j_2 < k$, there is a path
 981 from $\sigma_1 = (s_1, i_1 + j_1\lambda)$ to $\sigma_2 = (s_2, i_2 + j_2\lambda)$ in \mathcal{A}' . Let ℓ be a sufficiently large integer equal
 982 to $(i_2 - i_1) + (j_2 - j_1)\lambda \pmod{p}$. As λ divides p , ℓ is equal to $(i_2 - i_1) \pmod{\lambda}$. By taking
 983 ℓ larger than the reachability constant of S , Fact 2.4 gives us that there is a path of length ℓ
 984 from s_1 to s_2 in S , labeled by some word u . The positional word $(x : u)$ labels a transition
 985 from σ_1 to σ_2 in \mathcal{A}' , hence it is strongly connected.

986 Note that the periodicity of \mathcal{A}' is p , hence we can apply the results we obtained on
 987 strongly connected NFAs in Section 3 to portals, with $p|\mathcal{A}|$ as the number of states and p as
 988 the periodicity. ◀

989 In order to simplify the exposition of the proofs of later results, let us start with two
 990 technical lemmas expressing two basic properties of the Hamming distance with respect to
 991 \mathcal{A} . The first one extends Kleene's Lemma on languages of SCC-paths: for large enough ℓ ,
 992 whether $\mathcal{L}(\pi)$ contains a word of length ℓ only depends on the value of ℓ modulo p .

993 **► Lemma C.2.** *Let $\pi = P_0 \xrightarrow{a_1} \dots \xrightarrow{a_k} P_k$ be an SCC-path. There exists a constant B such that,*
 994 *for all $\ell \geq B$, if there is a word μ of length ℓ in $\mathcal{L}(\pi)$, then there exists a word μ' of length*
 995 *$\ell - p$ and a word μ'' of length $\ell + p$ in $\mathcal{L}(\pi)$.*

996 **Proof.** Recall the definition of $\mathcal{L}(\pi)$ (Definition 4.6):

$$997 \quad \mathcal{L}(\pi) = L_0 a_1 L_1 a_2 \dots L_k, \text{ where } L_i = \mathcal{L}(P_i) \text{ for } i = 0, \dots, k.$$

998 It follows that a word $\mu \in \mathcal{L}(\pi)$ can be written as $\mu = \mu_1 a_1 \mu_2 \dots \mu_k$ with $\mu_i \in L_i$. Each L_i is
 999 recognized by a strongly connected automaton \mathcal{A}_i with at most $p|\mathcal{A}|$ states. Let $B = 5(p|\mathcal{A}|)^2$.
 1000 If the length ℓ of μ exceeds B , then the run of μ in each of the \mathcal{A}_i 's contains simple loops
 1001 with sum of lengths greater than $p + 3(p|\mathcal{A}|)^2$. Let $\ell_0 + p$ denote the sum of the length of
 1002 these simple cycles: by construction ℓ_0 is greater than $3(p|\mathcal{A}|)^2$. We remove these simple
 1003 cycles from the run: the resulting run is still in $\mathcal{L}(\pi)$. Next, select any non-trivial SCC S_i in
 1004 π and let s be a state of S_i used by the run. As $\ell_0 \geq 3(p|\mathcal{A}|)^2$, by Fact 2.4, there is a path
 1005 of length ℓ_0 from s to itself in \mathcal{A}_i . Adding this path to the run yields an accepting run of
 1006 length $\ell - (\ell_0 + p) + \ell_0 = \ell - p$: the word labeling this run is the desired word μ' .

1007 To obtain μ'' , consider any simple cycle in the run of μ in \mathcal{A} , and let m denote the length
 1008 of this cycle. By definition of p , m divides p . Iterating this cycle p/m times yields a word μ''
 1009 of length $\ell + p$ that is in $\mathcal{L}(\pi)$. ◀

1010 ▶ **Corollary C.3.** *Let π be an SCC path. For large enough ℓ , whether there is an word of*
 1011 *length ℓ in $\mathcal{L}(\pi)$ only depends on the value of $\ell \pmod{p}$.*

1012 We say that two occurrences of blocking sequences in a word μ are disjoint if the
 1013 occurrences of their factors appear on disjoint sets of positions in μ .

1014 ▶ **Lemma C.4.** *Let $\pi = P_0 \xrightarrow{a_1} P_1 \cdots P_k$ be an SCC-path. If μ contains $k + 1$ disjoint*
 1015 *occurrences of blocking sequences for π , then $\mu \notin \mathcal{L}(\pi)$.*

1016 **Proof.** We prove the statement by induction on k .

1017 For $k = 0$, the language of π is exactly the language of P_0 , and blocking sequences of π
 1018 contain at least a blocking factor of P_0 . Therefore, a word containing a blocking sequence is
 1019 not in $\mathcal{L}(P_0) = \mathcal{L}(\pi)$.

1020 Now let $k > 0$ and suppose that the induction hypothesis holds for $k - 1$. Assume for the
 1021 sake of contradiction that μ is in $\mathcal{L}(\pi)$ and that it contains $k + 1$ occurrences of blocking
 1022 sequences. As it is in $\mathcal{L}(\pi)$, we can decompose μ into $\mu = \mu_0 a_0 \mu'$ with $\mu_0 \in \mathcal{L}(P_0)$ and
 1023 $\mu' \in \mathcal{L}(\pi')$, where $\pi' = P_1 \xrightarrow{a_2} \cdots P_k$ is the tail of π . Now, we can assume w.l.o.g. that the
 1024 first element ν of each of the $k + 1$ blocking sequence occurring in μ is blocking for P_0 , as
 1025 otherwise the removal of ν still leaves a blocking sequence. It follows that each ν does not
 1026 occur completely in μ_0 , as μ_0 is in $\mathcal{L}(P_0)$. Furthermore, as the sequences are disjoint, at
 1027 most one ν can overlap with $\mu_0 a_0$, hence at least k blocking sequences are fully contained in
 1028 μ' . Our induction hypothesis implies that μ' is not in $\mathcal{L}(\pi')$, a contradiction. This concludes
 1029 our induction. ◀

1030 The next lemma expresses a kind of converse implication, which generalizes Lemma A.3
 1031 from the strongly connected case: if a word is far from the language, then it contains many
 1032 blocking sequences.

1033 ▶ **Lemma C.5.** *Let $\pi = P_0 \xrightarrow{a_1} \cdots P_k$ be an SCC-path, let $L = \mathcal{L}(\pi)$, and let μ be a positional*
 1034 *word of length n such that $d(\mu, L)$ is finite. There is a constant C such that if $n \geq 2C/\varepsilon$*
 1035 *and μ is ε -far from L , then μ can be partitioned into $\mu = \mu_0 \mu_1 \cdots \mu_k$ such that for every*
 1036 *$i = 0, \dots, k$, μ_i contains at least $\frac{\varepsilon n}{C}$ blocking factors for P_i , for some constant C .*

1037 **Proof.** We proceed similarly to the proof of Lemma A.3, and only sketch this proof. Starting
 1038 from the left end of μ , we accumulate letters until we find a factor blocking for P_0 , and
 1039 iterate again starting from ρ positions later, where ρ is the reachability constant of a strongly
 1040 connected automaton recognizing $\mathcal{L}(P_0)$. When we have found at least $K = \frac{\varepsilon n}{C(k+2)}$ blocking
 1041 factors (C is to be determined later) for $\mathcal{L}(P_i)$, this position marks the end of μ_i , and we
 1042 iterate with the next portal in π .

1043 Let us assume that the process ends (i.e. we reach the right end of μ) before finding
 1044 enough blocking factors for all portals. We show that in this case, the distance between μ
 1045 and L is at most εn . Assume w.l.o.g. that we stop before finding enough blocking factors for
 1046 P_i . As in the proof of Lemma A.3, we replace the last letter of each blocking factor and use
 1047 the padding between them to make the run accepted by the automaton: this uses at most
 1048 $((i + 1) \cdot K + 2)\rho \leq \varepsilon n$ substitutions, if we set $C = 4(k + 1)\rho$. Therefore, if μ is ε -far from
 1049 $\mathcal{L}(\pi)$, then the decomposition process finds at least K blocking factors for P_i in μ_i for each
 1050 i . ◀

1051 We make the remark that minimal blocking sequences have a bounded number of terms.
 1052 This is because if we build the sequence from left to right by adding terms one by one, the
 1053 minimality implies that at each step we should block a previously unblocked portal.

1054 ► **Lemma C.6.** *A minimal blocking sequence for \mathcal{A} contains at most $p^2|Q|^2$ terms.*

1055 **Proof.** First, remark that there are at most $p^2|Q|^2$ portals in \mathcal{A} . Let $\sigma = (\mu_1 \dots, \mu_\ell)$ be a
 1056 minimal blocking sequence for \mathcal{A} . For all $i = 1, \dots, \ell$, we define $\sigma_i = (\mu_1 \dots, \mu_i)$, and σ_0 is
 1057 the empty sequence.

1058 Then, for each i , we consider the set \mathcal{S}_i of portals P such that for all accepting SCC-path
 1059 π of \mathcal{A} containing P , the prefix of π ending at P is blocked by σ_i . We have $\mathcal{S}_0 = \emptyset$, and \mathcal{S}_ℓ is
 1060 the set of all portals of \mathcal{A} .

1061 We claim that for every $i < \ell$, \mathcal{S}_i is a proper subset of \mathcal{S}_{i+1} . Otherwise, if $\mathcal{S}_i = \mathcal{S}_{i+1}$, then
 1062 removing μ_{i+1} from σ gives a blocking sequence σ' of \mathcal{A} , such that $\sigma' \preceq \sigma$, contradicting the
 1063 minimality of σ .

1064 Therefore, it follows that $\ell \leq p^2|Q|^2$. ◀

1065 ► **Lemma 4.10.** *If μ contains $|\mathcal{A}|$ disjoint blocking sequences for \mathcal{A} then $\mu \notin \mathcal{L}(\mathcal{A})$.*

1066 **Proof.** Let π be an accepting SCC-path through \mathcal{A} . Note that number ℓ of portals in π is
 1067 at most $|Q|$. By definition, a blocking sequence for \mathcal{A} is a blocking sequence for π . As μ
 1068 contains $|Q|$ disjoint blocking sequences for \mathcal{A} , it contains $\ell \leq |Q|$ disjoint blocking sequences
 1069 for π , hence $\mu \notin \mathcal{L}(\pi)$ by Lemma C.4. As a result, μ is not in the language of any accepting
 1070 SCC-path of \mathcal{A} , and thus not in $\mathcal{L}(\mathcal{A})$. ◀

1071 ► **Lemma 4.11.** *Let μ be a word of length n . There exist constants E, K such that if
 1072 $+\infty > d(\mu, \mathcal{L}(\mathcal{A})) \geq \varepsilon n$ and $n \geq E/\varepsilon$, then μ can be partitioned into $\mu = \mu_0 \mu_1 \dots \mu_K$
 1073 such that for every $i = 0, \dots, K$, μ_i contains at least $\frac{\varepsilon n}{E}$ disjoint occurrences of words
 1074 $\nu_{i,1}, \nu_{i,2}, \dots$, each of length $\mathcal{O}(1/\varepsilon)$, such that for any choice of j_0, j_1, \dots, j_K , the sequence
 1075 $(\nu_{0,j_0}, \nu_{1,j_1}, \dots, \nu_{K,j_K})$ is a blocking sequence for \mathcal{A} .*

1076 **Proof.** We start by observing that an SCC-path contains at most $|\mathcal{A}|$ portals, hence there are
 1077 at most $(|\mathcal{A}|^2 p^2 |\Sigma| + 1)^{|\mathcal{A}|}$ SCC-paths in \mathcal{A} . Let $K \leq (|\mathcal{A}|^2 p^2 |\Sigma| + 1)^{|\mathcal{A}|}$ denote the number
 1078 of accepting SCC-paths of \mathcal{A} , and let Π be the set of these SCC paths.

1079 For each accepting SCC-path π , as $\mathcal{L}(\pi) \subseteq \mathcal{L}(\mathcal{A})$, $d(\mu, \mathcal{L}(\pi)) \geq d(\mu, \mathcal{L}(\mathcal{A}))$, hence μ
 1080 contains $\frac{\varepsilon n}{C}$ disjoint blocking sequences for π , by Lemma C.5. We can apply this result as
 1081 long as $D \geq 2C$: we define $D = \max(2C, K \cdot C)$. As the blocking sequences are disjoint, the
 1082 sum of their lengths does not exceed n , thus up to doubling C , we have that $\frac{\varepsilon n}{C}$ of them have
 1083 a total length of $\mathcal{O}(1/\varepsilon)$.

1084 It remains to prove that $\frac{\varepsilon n}{C}$ disjoint blocking sequences for each π implies $\frac{\varepsilon n}{D}$ disjoint
 1085 blocking sequences for \mathcal{A} . We decompose μ into η_1, \dots, η_K and index the accepting SCC-paths
 1086 $\Pi = \{\pi_1, \dots, \pi_K\}$ such that each η_i contains a $1/K$ -fraction of the paths that are blocking
 1087 for π_i , the i -th accepting SCC path of \mathcal{A} . We proceed as follows: for all $i < K$ we define
 1088 η_i so that $\eta_1 \dots \eta_i$ is the shortest prefix of μ such that there exist $\pi_i \in \Pi \setminus \{\pi_1, \dots, \pi_{i-1}\}$
 1089 such that there are $\geq \frac{\varepsilon n i}{KC}$ disjoint blocking sequences for π_i in $\eta_1 \dots \eta_i$. We define η_K so
 1090 that $\eta_1 \dots \eta_K = \mu$. As all $\pi \in \Pi$ have at least $\frac{\varepsilon n}{C}$ blocking sequences in μ , this sequence is
 1091 well-defined. Furthermore, the π_i are distinct by definition.

1092 It follows that η_i contains $\frac{\varepsilon n}{KC} \geq \frac{\varepsilon n}{D}$ disjoint blocking sequences for π_i . Finally, concatenating
 1093 a blocking sequence from each of the η_i yields a blocking sequence for \mathcal{A} .

1094 Note that the resulting blocking sequences $\sigma = (\xi_0, \xi_1, \dots, \xi_t)$ contain at most $t \leq K \cdot |\mathcal{A}|$
 1095 elements, as the blocking sequences for the π_i contain at most $|\mathcal{A}|$ elements. The word μ_0 is

1096 the prefix of μ that contains the element ξ_0 from the first $\frac{\varepsilon n}{E}$ sequences σ , where $E = D \cdot |\mathcal{A}|$;
 1097 the selected ξ_0 are the $(\nu_{0,j})_j$. We then iterate on the remainder of the word, with μ_i
 1098 containing the elements ξ_i from $\frac{\varepsilon n}{E}$ subsequent sequences σ . ◀

1099 **D Languages with infinitely many blocking sequences are hard**

1100 **D.1 Proof of Lemma 4.17**

1101 ▶ **Lemma 4.17.** *If \mathcal{A} has infinitely many minimal blocking sequences, then there exist a*
 1102 *portal P and sequences σ_l and σ_r satisfying properties P1, P2 and P3.*

1103 **Proof.** As \mathcal{A} has infinitely many minimal blocking sequences, there exists an accepting
 1104 SCC-path π with infinitely many minimal blocking sequences. By Lemma C.6, the number
 1105 of terms in a minimal blocking sequence is bounded, therefore there is a portal in π with
 1106 infinitely many minimal blocking factors.

1107 If \mathcal{A} has infinitely many minimal blocking sequences, let $(\sigma_j)_{j \in \mathbb{N}}$ be a family of minimal
 1108 blocking sequences such that the sum of the lengths of the terms of σ_j is at least j for all j .

1109 By Lemma C.6, a minimal blocking sequence has a bounded number of elements. We can
 1110 thus extract from this sequence another one $(\sigma'_j)_{j \in \mathbb{N}}$ such that each σ'_j contains a factor of
 1111 length at least j .

1112 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
 1113 the left effect of the $i_j - 1$ first factors and the right effect of the $k_j - i_j$ last ones, with
 1114 k_j the length of σ'_j . As those objects are taken from bounded sets, we can obtain a third
 1115 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
 1116 and the set of components for which it is blocking is K .

1117 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
 1118 $\sigma_r = (n_1^r : u_1^r), \dots, (n_\ell^r : u_\ell^r)$ 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
 1119 blocking sequence.

1120 We call surviving portals the portals that survive (σ_l, σ_r) in at least one SCC-path.

1121 ▷ **Claim D.1.** *There exists a surviving portal with infinitely many minimal blocking factors*
 1122 *that is minimal for \preceq among surviving portals.*

1123 **Proof.** Suppose by contradiction that all \preceq -minimal surviving portals have finitely many
 1124 minimal blocking factors.

1125 For all j , $(n_j : u_j)$ must be blocking for all surviving portals (otherwise $\bar{\sigma}_j$ would not be
 1126 blocking for \mathcal{A}). Hence $(n_j : u_j)$ contains a blocking factor for each \preceq -minimal surviving
 1127 portal. As those factors are bounded while $(n_j : u_j)$ can get arbitrarily large, there exists j
 1128 such that $(n_j : u_j)$ can be split into two non-empty parts $(n_j : u_j^-)(n_j^+ : u_j^+)$ so that each
 1129 \preceq -minimal surviving portal has a minimal blocking factor in either $(n_j : u_j^-)$ or $(n_j^+ : u_j^+)$. As
 1130 a consequence, every surviving portal has a blocking factor in either $(n_j : u_j^-)$ or $(n_j^+ : u_j^+)$.

1131 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
 1132 blocking sequence for \mathcal{A} , contradicting the minimality of $\sigma_l(n_j : u_j)\sigma_r$ for \preceq . In conclusion,
 1133 there is a \preceq -minimal surviving portal with infinitely many minimal blocking factors. ◀

1134 Let $s, x \rightsquigarrow t, y$ be a \preceq -minimal surviving portal with infinitely many minimal blocking
 1135 factors: It satisfies P2.

1136 The following claim shows that there is a pair of sequences (σ_l, σ_r) such that properties
 1137 P1 and P3 are satisfied.

1138 \triangleright Claim D.2. There exist σ_l, σ_r such that $\sigma_l \sigma_r$ is not a blocking sequence for \mathcal{A} , and for all
 1139 accepting SCC-path π in \mathcal{A} , every surviving portal in π is \simeq -equivalent to $s, x \rightsquigarrow t, y$.

1140 Proof. We start from the sequences σ_l, σ_r defined before and extend them so that they have
 1141 the desired property.

1142 For each $s', x' \rightsquigarrow t', y' \not\sim s, x \rightsquigarrow t, y$, since $s, x \rightsquigarrow t, y$ is \preceq -minimal we can pick a positional
 1143 word $(n : u)_{s', x' \rightsquigarrow t', y'}$ that is blocking for $s', x' \rightsquigarrow t', y'$ but not for $s, x \rightsquigarrow t, y$.

1144 We extend σ_l and σ_r as follows. While there is a surviving portal $s', x' \rightsquigarrow t', y'$ that is
 1145 not \simeq -equivalent to $s, x \rightsquigarrow t, y$:

- 1146 ■ We pick an SCC-path π such that $s', x' \rightsquigarrow t', y'$ survives in π .
- 1147 ■ Let $i_\ell = (\sigma_l \gg \pi)$ and $i_r = (\pi \ll \sigma_r)$
- 1148 ■ If for all $i \in \{i_\ell + 1, \dots, i_r - 1\}$, $s_i, x_i \rightsquigarrow t_i, y_i \not\sim s, x \rightsquigarrow t, y$ then we append at the
 1149 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
 1150 sequence $\sigma_l \sigma_r$ is now blocking for π . On the other hand, since we did not add any blocking
 1151 factor for $s, x \rightsquigarrow t, y$, there must still be a surviving portal that is \simeq -equivalent to it.
- 1152 ■ 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
 1153 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
 1154 \simeq , or i_ℓ if there is no such index. Symmetrically, let d the minimal index in $\{i, \dots, i_r - 1\}$
 1155 such that $(m_d, s_d, t_d) \not\sim s, x \rightsquigarrow t, y$, or i_r if there is no such index. We append at the
 1156 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
 1157 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
 1158 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.

1159 We iterate this step until all surviving portals are \simeq -equivalent to $s, x \rightsquigarrow t, y$. We made
 1160 sure that at least one portal was still surviving after each step, hence in the end the sequence
 1161 $\sigma_l \sigma_r$ is not blocking for \mathcal{A} . \triangleleft

1162 \blacktriangleleft

1163 D.2 Proof of Lemma 4.18

1164 We show the following statement.

1165 \blacktriangleright **Lemma 4.18.** *If there exist $P = s, x \rightsquigarrow t, y$ and σ_l, σ_r satisfying properties P1, P2 and*
 1166 *P3 then $\mathcal{L}(\mathcal{A})$ is hard.*

1167 \blacktriangleright **Lemma D.3.** *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*
 1168 *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*
 1169 *$N \in \mathbb{N}$.*

1170 *Then there is a word w^l of length at most $(3|\mathcal{A}|^3 + |\mathcal{A}|)(k + 1) + N(2p^2 + p)k|\mathcal{A}| +$*
 1171 *$pN \sum_{i=1}^k |u_i^l|$ 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} ,*
 1172 *and $(x_0 : w)$ contains N times $(n_1^l : u_1^l), \dots, N$ times $(n_k^l : u_k^l)$ as disjoint factors, in that*
 1173 *order.*

1174 **Proof.** We define w^l by induction on k . As π is accepting, by definition $\mathcal{L}(\pi) \neq \emptyset$, and thus
 1175 for all $j \in \{0, \dots, \ell\}$ there exists a word of length $y_j - x_j \pmod{p}$ labelling a path from
 1176 s_j to t_j . By Fact 2.4, there is such a word v_j of length at most $3|\mathcal{A}|^2$. As a result, for all
 1177 $z \in \{0, \dots, \ell\}$ 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
 1178 a path of length $x_z \pmod{p}$ from q_0 to s_z in \mathcal{A} . If $k = 0$, we can simply set $w^l = w_i$.

1179 Let $k > 0$, suppose the lemma holds for $k - 1$. Let $j = ((n_1 : u_1^l) \gg \pi)$. As $((n_1 : u_1^l) \gg \pi) \leq$
 1180 $(\sigma_l \gg \pi) < i$, we have $j < i$. By definition, $(n_1 : u_1^l)$ is not blocking for

1181 $s_{j+1}, x_{j+1} \rightsquigarrow t_{j+1}, y_{j+1}$. As a consequence, there is a word v_j labelling a path from s_j to t_j
 1182 such that $(x_j : v_j)$ has $(n_1 : u_1^l)$ as a factor. We can remove cycles of length $0 \pmod{p}$ in that
 1183 path, before and 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
 1184 are in the same SCC, we can extend v_j into a word v_j' of length $\leq |v_j| + |\mathcal{A}| \leq |u_1^l| + (2p+1)|\mathcal{A}|$
 1185 that labels a cycle from s_j to itself.

1186 Let $\sigma' = (n_2^l : u_2^l), \dots, (n_k^l : u_k^l)$ and $\pi' = s_{j+1}, x_{j+1} \rightsquigarrow t_{j+1}, y_{j+1} \xrightarrow{a_{j+2}} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell$.
 1187 By definition, we have $(\sigma' \gg \pi') = (\sigma_l \gg \pi) < i$. By induction hypothesis, there is a word
 1188 w' of length $\leq (3|\mathcal{A}|^3 + |\mathcal{A}|)k + N(2p^2 + p)(k-1)|\mathcal{A}| + pN \sum_{i=1}^{k-1} |u_i|$ such that $|w'| = x_i - x_j$
 1189 \pmod{p} , there is a run reading w' from s_j to s_i in \mathcal{A} , and $(x_j : w')$ contains N times $(n_2^l : u_2^l)$,
 1190 \dots , N times $(n_k^l : u_k^l)$ as disjoint factors, in that order.

1191 We set $w^l = w_j(v_j')^{pN} w'$. This word has length $x_i \pmod{p}$, and at most $|w_j| + pN|v_j'| +$
 1192 $|w'| \leq 3|\mathcal{A}|^3 + |\mathcal{A}| + pN(|u_1^l| + (2p+1)|\mathcal{A}|) + |w'| \leq (3|\mathcal{A}|^3 + |\mathcal{A}|)(k+1) + N(2p^2 + p)k|\mathcal{A}| +$
 1193 $pN \sum_{i=1}^k |u_i^l|$. It labels a path from s_0 to s_i , and contains N times $(n_1^l : u_1^l)$, \dots , N times
 1194 $(n_k^l : u_k^l)$ as disjoint factors, in that order. \blacktriangleleft

1195 \blacktriangleright **Lemma D.4.** *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*
 1196 *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*
 1197 *$N \in \mathbb{N}$.*

1198 *Then there is a word w^r of length at most $(3|\mathcal{A}|^3 + |\mathcal{A}|)(k+1) + N(2p^2 + p)k|\mathcal{A}| +$*
 1199 *$pN \sum_{i=1}^k |u_i^r|$ 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} ,*
 1200 *and $(y_i : w^r)$ contains N times $(n_1^r : u_1^r)$, \dots , N times $(n_k^r : u_k^r)$ as disjoint factors, in that*
 1201 *order.*

1202 **Proof.** By a proof symmetric to the one of the previous lemma. \blacktriangleleft

1203 Given a sequence σ , define $\|\sigma\|$ as the sum of the lengths of the terms of σ .

1204 **Proof of Lemma 4.18.** A direct consequence of properties P1 and P3 is that for all $(n' : u')$,
 1205 $\sigma_l(n' : u')\sigma_r$ is blocking for \mathcal{A} if, and only if $(n' : u')$ is blocking for P .

1206 The proof goes as follows: we show that we can turn an algorithm testing $\mathcal{L}(\mathcal{A})$ with $f(\varepsilon)$
 1207 samples into an algorithm testing $\mathcal{L}(P)$ with $f(\varepsilon/X)$ samples with X a constant. We then
 1208 apply Theorem 3.7 to obtain the lower bound.

1209 Consider an algorithm testing $\mathcal{L}(\mathcal{A})$ with $f(\varepsilon)$ samples for some function f . We describe
 1210 an algorithm for testing $\mathcal{L}(P)$. Say we are given a threshold ε and a word v of length n . First
 1211 of all we can apply Lemmas D.3 and D.4 to compute two words w^l and w^r of length at most
 1212 $E + \varepsilon nF$ for some constants E and F such that we can read w^l from q_0 to s and w^r from t
 1213 to q_f and w_l contains each element of σ_l at least εn times and w_r contains each element of
 1214 σ_r at least εn times. Let $w = w^l v w^r$. Suppose $|v| \geq \frac{6p^2|\mathcal{A}|^2}{\varepsilon}$ and $d(v, \mathcal{L}(P)) < +\infty$.

- 1215 \blacksquare If $v \in \mathcal{L}(\mathcal{A})$ then clearly $w \in \mathcal{L}(\mathcal{A})$.
- 1216 \blacksquare If $d(v, \mathcal{L}(P)) \geq \varepsilon n$ then by Lemma A.3 (in light of Lemma C.1), $(x : v)$ contains at least
 1217 $\frac{\varepsilon n}{6p^2|\mathcal{A}|^2}$ blocking factors for P . Then we have that w contains at least $\frac{\varepsilon n}{6p^2|\mathcal{A}|^2}$ disjoint
 1218 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 the length
 1219 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 some
 1220 constant X .

1221 Let us now describe the algorithm for testing $\mathcal{L}(P)$.

- 1222 \blacksquare If $\mathcal{L}(P) \cap \Sigma^n = \emptyset$ then we reject.
- 1223 \blacksquare If $|v| < \frac{6p^2|\mathcal{A}|^2}{\varepsilon}$ then we read v entirely and check that it is in $\mathcal{L}(P)$.

1224 ■ If $v \in \mathcal{L}(P)$ then we apply our algorithm for testing $\mathcal{L}(\mathcal{A})$ on $w = w^l v w^r$ with parameter
 1225 $\varepsilon' = \frac{\varepsilon}{X}$.

1226 The number of samples used on v is at most the number of samples needed on w , hence
 1227 $f(\varepsilon/X)$. We obtain a procedure to test $\mathcal{L}(P)$ using $f(\varepsilon/X)$ samples. By Lemma C.1, $\mathcal{L}(P)$
 1228 is recognised by a strongly connected automaton. Furthermore by P2, P has infinitely many
 1229 blocking factors. By Theorem 3.7, $f(\varepsilon/X) = \Omega(\log(\varepsilon^{-1})/\varepsilon)$, hence $f(\varepsilon) = \Omega(\log(\varepsilon^{-1})/\varepsilon)$.
 1230 This concludes our proof. ◀

1231 **E** Trivial and easy languages

1232 **E.1** Proof of Lemma 5.1

1233 ► **Lemma 5.1.** *If \mathcal{A} has finitely many minimal blocking sequences, then there is a tester for*
 1234 *$\mathcal{L}(\mathcal{A})$ using $\mathcal{O}(1/\varepsilon)$ queries.*

1235 **Proof.** Let us start by proving that a word that is ε -far from $\mathcal{L}(\mathcal{A})$ must contain many times
 1236 some minimal blocking sequence σ . We first check whether \mathcal{A} contains a word of length n .
 1237 If not, we reject, and otherwise, the distance between μ and $\mathcal{L}(\mathcal{A})$ is finite. Therefore, by
 1238 Lemma 4.11, if μ is ε -far from $\mathcal{L}(\mathcal{A})$, then μ contains $N = \Omega(\varepsilon n)$ disjoint blocking sequences,
 1239 organised as stated in the lemma. We can extract from each group of $|\mathcal{A}|$ such blocking
 1240 sequence a minimal blocking sequence, hence μ contains at least N disjoint occurrences
 1241 of minimal blocking sequences. By pigeonhole principle, there exists a minimal blocking
 1242 sequence σ such that μ contains $N/P = \Omega(\varepsilon n)$ disjoint occurrences of σ .

1243 If $|\mu| \leq \frac{D}{\varepsilon}$ then we read μ entirely, and accept iff it is in $\mathcal{L}(\mathcal{A})$: this uses $\mathcal{O}(1/\varepsilon)$ queries.

1244 Otherwise, for each minimal blocking sequence σ , we sample uniformly at random $\mathcal{O}(1/\varepsilon)$
 1245 intervals of length K in μ . We reject if we find $|\mathcal{A}|$ disjoint occurrences of σ . If we have gone
 1246 through every minimal blocking sequence without rejecting, we accept.

1247 If $\mu \in \mathcal{L}(\mathcal{A})$, then by Lemma C.4 it cannot contain $|\mathcal{A}|$ disjoint blocking sequences, hence
 1248 the algorithm will accept.

1249 If the word is ε -far from $\mathcal{L}(\mathcal{A})$ (but within a finite distance), then there exists a minimal
 1250 blocking sequence σ such that μ contains $\Omega(\varepsilon n)$ disjoint occurrences of σ . By sampling
 1251 $\mathcal{O}(1/\varepsilon)$ factors of length K at random, we have a probability bounded away from zero of
 1252 finding $|\mathcal{A}|$ of those occurrences, and subsequently rejecting μ .

1253 We can iterate this procedure a constant number of times to obtain a procedure using
 1254 $\mathcal{O}(\frac{1}{\varepsilon})$ samples that accepts every word in the language and rejects with probability at least
 1255 $1/2$ words that are ε -far from the language. ◀

1256 **E.2** Proof of Lemma 5.3

1257 ► **Lemma 5.3.** *$MBS(\mathcal{A})$ is empty if and only if $L = \mathcal{L}(\mathcal{A})$ is trivial.*

1258 We prove the two directions separately.

1259 ► **Lemma E.1.** *If $MBS(\mathcal{A})$ is empty, then $L = \mathcal{L}(\mathcal{A})$ is trivial in the sense of Definition 5.2.*

1260 **Proof.** We showed in Lemma 4.11 that if μ is long enough and ε -far from L , then μ contains
 1261 $\Omega(\varepsilon n)$ disjoint blocking sequences for \mathcal{A} . As there are no blocking sequences for \mathcal{A} , long
 1262 enough words cannot be ε -far from L , hence it is trivial in the sense of Definition 5.2. ◀

1263 Alon et al. [5] show that if a language is non-trivial (in their sense), then testing it requires
 1264 $\Omega(1/\varepsilon)$ queries for small enough $\varepsilon > 0$. To finish our characterization of trivial languages, we
 1265 show that if $MBS(\mathcal{A})$ is not empty, then $L = \mathcal{L}(\mathcal{A})$ is non-trivial in the sense of Alon.

1266 ► **Lemma E.2.** *Let \mathcal{A} be a trim NFA such that $L = \mathcal{L}(\mathcal{A})$ is infinite. If \mathcal{A} admits a blocking*
 1267 *sequence, then there exists $\varepsilon_0 > 0$, such that for infinitely many n there exist words in*
 1268 *$\mathcal{L}(\mathcal{A}) \cap \Sigma^n$ and there exists $w \in \Sigma^n$ such that $d(w, \mathcal{L}(\mathcal{A})) \geq \varepsilon_0 n$*

1269 **Proof.** Let $\sigma = (\mu_1, \dots, \mu_k)$ be a blocking sequence for \mathcal{A} , and let C be the maximum length
 1270 of a μ_i 's. As L is infinite, there exists an SCC-path π in \mathcal{A} and $w \in \mathcal{L}(\pi)$ with $|w| \geq |\mathcal{A}|$. By
 1271 Corollary C.3, for all sufficiently large ℓ such that $\ell = |w| \pmod{p}$, there exists $w' \in \mathcal{L}(\pi)$
 1272 with $|w'| = \ell$.

1273 For all $i = 1, \dots, k$, let ν_i be a shortest word of the form $(0 : v_i)$ and of length ℓ_i equal to
 1274 0 modulo p , such that μ_i is a factor of ν_i . By minimality, ℓ_i is at most $C + 2p$. Then, for
 1275 any integer $N \in \mathbb{N}$, let $w_N = \nu_1^N \dots \nu_k^N (0 : a^{|w|})$, where a is an arbitrary letter.

1276 As w_N is of length $|w| \pmod{p}$, there is a word of the same length in $\mathcal{L}(\mathcal{A})$, i.e. $\mathcal{L}(\mathcal{A}) \cap \Sigma^n$
 1277 is nonempty. On the other hand, it contains N disjoint occurrences of σ , which is a blocking
 1278 sequence for \mathcal{A} . Therefore, the distance between w_N and \mathcal{A} is at least $N - |\mathcal{A}|$. Furthermore,
 1279 the length of w_N is less than $|w| + N(C + 2p)$. Therefore, if we let $\varepsilon_0 = \frac{1}{C+2p+w}$, then for
 1280 $N \geq 2|\mathcal{A}|$, we have $\varepsilon_0 |w_N| \leq N - |\mathcal{A}| \leq d(w_N, \mathcal{L}(\mathcal{A}))$, i.e. w_N is ε -far from L . ◀

1281 **F Complexity proofs**

1282 In this section we prove Theorem 6.4:

1283 ► **Theorem 6.4.** *The triviality, easiness and hardness problems are all PSPACE-complete,*
 1284 *even for strongly connected NFAs.*

1285 In Appendix F.1 we show the PSPACE upper bounds on the hardness and triviality
 1286 problems (Propositions F.7 and F.9). The upper bound on the easiness problem follows
 1287 immediately, as the three properties form a trichotomy.

1288 In Appendix F.2, we show that all three problems are PSPACE-hard (Lemma F.11 and
 1289 Corollary F.13).

1290 **F.1 A PSPACE upper-bound**

1291 **F.1.1 Testing hardness**

1292 A naive algorithm to check hardness of a language $\mathcal{L}(\mathcal{A})$ would be to construct an automaton
 1293 recognising blocking sequences of $\mathcal{L}(\mathcal{A})$. However, there are subtleties in the construction
 1294 due to the fact that the number of SCC-paths can be exponential in the $|\mathcal{A}|$. For instance,
 1295 an automaton computing the set of left effects of a given sequence on SCC-paths of \mathcal{A} would
 1296 be doubly exponential. This would a priori not give a PSPACE algorithm, since we obtain a
 1297 doubly-exponential state space.

1298 We choose to prove another characterisation of automata with hard languages, which
 1299 straightforwardly gives a recursive PSPACE algorithm to test it.

1300 ► **Lemma F.1.** *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,*
 1301 *Π a set of SCC-paths and $(\sigma_{\pi'})_{\pi' \in \Pi}$ a family of sequences of positional words such that*
 1302 *$(\sigma_{\pi'} \gg \pi) < i$ for all π' .*

1303 *There exists a sequence of positional words σ such that:*

- 1304 ■ $(\sigma \gg \pi) < i$
- 1305 ■ $(\sigma_{\pi'} \gg \pi') \leq (\sigma \gg \pi')$ for all $\pi' \in \Pi$.

1306 **Proof.** We prove this by induction on the sum of the lengths of the elements of Π . If Π is
 1307 empty then we can set σ as the empty sequence.

1308 If not, let π_{min} be such that the first term of $\sigma_{\pi_{min}}$ has the least left effect on π . Let
 1309 $\sigma_{\pi_{min}} = (n_1 : u_1), \dots, (n_k : u_k)$ 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
 1310 $j = ((n_1 : u_1) \gg \pi_{min})$ and $r = ((n_1 : u_1) \gg \pi)$.

1311 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
 1312 $j < \ell$ and $\Pi' = \Pi \setminus \{\pi_{min}\}$ otherwise. In the first case the sequence associated with π' is
 1313 $\sigma_{\pi'} = (n_2 : u_2), \dots, (n_k : u_k)$.

1314 \triangleright **Claim F.2.** For all $\bar{\pi} \in \Pi \setminus \{\pi_{min}\}$, we have $(\sigma_{\bar{\pi}} \gg \pi) = r + (\sigma_{\bar{\pi}} \gg s_{r+1}, x_{r+1} \rightsquigarrow$
 1315 $t_{r+1}, y_{r+1} \xrightarrow{a_{r+2}} \dots s_k, x_k \rightsquigarrow t_k, y_k)$

1316 **Proof.** Since the first term of $\sigma_{\pi'}$ was the one with the least left effect on π , the first term of
 1317 every other sequence has a left effect at least r on it.

1318 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
 1319 means $(\bar{n}_1 : \bar{u}_1)$ is not a blocking factor for $s_{z+1}, x_{z+1} \rightsquigarrow t_{z+1}, y_{z+1}$.

1320 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
 1321 $(\sigma_{\bar{\pi}} \gg s_{r+1}, x_{r+1} \rightsquigarrow t_{r+1}, y_{r+1} \xrightarrow{a_{r+2}} \dots) = z - r + ((\bar{n}_2 : \bar{u}_2), \dots, (\bar{n}_m : \bar{u}_m) \gg s_{z+1}, x_{z+1} \rightsquigarrow$
 1322 $t_{z+1}, y_{z+1}) = (\sigma_{\bar{\pi}} \gg \pi) - r$.

1323 \triangleleft

1324 As a consequence of this claim, we have that $(\sigma_{\bar{\pi}} \gg s_{r+1}, x_{r+1} \rightsquigarrow t_{r+1}, y_{r+1} \xrightarrow{a_{r+2}}$
 1325 $\dots s_k, x_k \rightsquigarrow t_k, y_k) < i - r$ for all $\bar{\pi} \in \Pi \setminus \{\pi'\}$.

1326 By induction hypothesis, we obtain a sequence σ' such that

- 1327 $\blacksquare (\bar{\sigma} \gg s_{r+1}, x_{r+1} \rightsquigarrow t_{r+1}, y_{r+1} \xrightarrow{a_1} \dots s_\ell, x_\ell \rightsquigarrow t_\ell, y_\ell) < i - r$
- 1328 $\blacksquare (\sigma_{\pi'} \gg \pi') \leq (\sigma' \gg \pi')$ for all $\pi' \in \Pi'$.

1329 The sequence $(n_1 : u_1), \sigma'$ satisfies both conditions of the lemma. \blacktriangleleft

1330 \blacktriangleright **Lemma F.3.** An automaton \mathcal{A} is hard if and only if there exists an accepting SCC-path π
 1331 containing a portal $s, x \rightsquigarrow t, y$ such that:

- 1332 $\blacksquare s, x \rightsquigarrow t, y$ has infinitely many minimal blocking factors.
- 1333 \blacksquare For all accepting SCC-path π' there exist sequences σ^l, σ^r such that:
 - 1334 $\blacksquare s, x \rightsquigarrow t, y$ survives (σ^l, σ^r) in π
 - 1335 \blacksquare All portals surviving (σ^l, σ^r) in π' are \simeq -equivalent to $s, x \rightsquigarrow t, y$

1336 **Proof.** Let us start with the left-to-right direction. If \mathcal{A} is hard then by Lemma 5.1 it
 1337 has infinitely many minimal blocking sequences. Then by Lemma 4.17 we have a portal
 1338 $s, x \rightsquigarrow t, y$ and sequences σ^l, σ^r satisfying properties P1, P2 and P3.

1339 By P1, $\sigma^l \sigma^r$ is not blocking for \mathcal{A} , thus there exists an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1}$
 1340 $\dots s_k, x_k \rightsquigarrow t_k, y_k$ and an index i such that $(\sigma^l \gg \pi) < i < (\pi \ll \sigma^r)$.

1341 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
 1342 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
 1343 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
 1344 \simeq -equivalent to $s, x \rightsquigarrow t, y$ (we use the same pair (σ^l, σ^r) for all π').

1345 Let us now prove the other direction. Suppose we have π and $s, x \rightsquigarrow t, y$ satisfying the
 1346 conditions of the lemma. We only need to construct two sequences σ^l, σ^r such that properties
 1347 P1 and P3 are satisfied. The result follows by Lemma 4.18.

1348 Let Π be the set of accepting SCC-paths in \mathcal{A} . Consider families of sequences $(\sigma^l_{\pi'})_{\pi' \in \Pi}$
 1349 and $(\sigma^r_{\pi'})_{\pi' \in \Pi}$ such that for all $\pi' \in \Pi$:

- 1350 ■ $s, x \rightsquigarrow t, y$ survives $(\sigma_{\pi'}^l, \sigma_{\pi'}^r)$ in π
 1351 ■ All portals surviving $(\sigma_{\pi'}^l, \sigma_{\pi'}^r)$ in π' are \simeq -equivalent to $s, x \rightsquigarrow t, y$

1352 Let i be the index of $s, x \rightsquigarrow t, y$ in π . By Lemma F.1 we can build a sequence σ^l such
 1353 that

- 1354 ■ $(\sigma^l \gg \pi) < i$
 1355 ■ $(\sigma_{\pi'}^l \gg \pi') \leq (\sigma^l \gg \pi')$ for all $\pi' \in \Pi$.

1356 Using a symmetric argument, we build a sequence σ^r such that

- 1357 ■ $i < (\pi \ll \sigma^r)$
 1358 ■ $(\pi' \ll \sigma_{\pi'}^r) \geq (\pi' \ll \sigma^r)$ for all $\pi' \in \Pi$.

1359 As a consequence, for all accepting SCC-path $\pi' \in \Pi$, all portals surviving (σ^l, σ^r) in π'
 1360 are \simeq -equivalent to $s, x \rightsquigarrow t, y$. Furthermore, $s, x \rightsquigarrow t, y$ survives (σ^l, σ^r) in π .

1361 We have shown that $s, x \rightsquigarrow t, y$ and (σ^l, σ^r) satisfy properties P1 and P3. P2 is immediate
 1362 by assumption. We simply apply Lemma 4.18 to obtain the result. ◀

1363 Next, we establish that the items listed in the previous lemma can all be checked in
 1364 polynomial space in $|\mathcal{A}|$.

1365 ▶ **Lemma F.4.** *Given a portal P , we can check whether it has infinitely many minimal*
 1366 *blocking factors in space polynomial in $|\mathcal{A}|$.*

1367 **Proof.** Recall that, by Lemma C.1, $L = \mathcal{L}(P)$ is recognized by a strongly connected auto-
 1368 maton \mathcal{A}' with at most $p|\mathcal{A}|$ states. While this number may be exponential in $|\mathcal{A}|$, the
 1369 transition function of \mathcal{A}' can be computed in polynomial space from the polynomial-sized
 1370 representation of a state. Furthermore, the same property holds for the construction used in
 1371 Lemma A.5, as in the determinization step, all states share the index modulo p .

1372 We then simply need to check if the resulting automaton has an infinite language, which
 1373 is the case if and only if it has a cycle reachable from the initial state and from which a final
 1374 state is reachable. This can be checked by exploring the state space of the automaton, in
 1375 non-deterministic polynomial space (in $|\mathcal{A}|$), and applying Savitch's theorem. ◀

1376 ▶ **Lemma F.5.** *Given two SCC-paths π and π' , one can check in PSPACE whether there is a*
 1377 *sequence σ that is blocking for π and not π' .*

1378 **Proof.** We prove an intermediate result that yields a recursive PSPACE algorithm.

1379 ▷ **Claim F.6.** There is a sequence σ that is blocking for $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow$
 1380 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:

- 1381 ■ there is a positional word $(n : w)$ that is a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$ and
 1382 not $s'_0, x'_0 \rightsquigarrow t'_0, y'_0$ and there is a sequence σ' that is blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2}$
 1383 $\dots s_k, x_k \rightsquigarrow t_k, y_k$ and not π' ,
 1384 ■ or there is a positional word $(n : w)$ that is a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$ and $s'_0, x'_0 \rightsquigarrow$
 1385 t'_0, y'_0 and there is a sequence σ' that is blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$
 1386 and not $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$.

1387 **Proof.** The right-to-left direction is clear (just take $\sigma = (n : w), \sigma'$ in both cases).

1388 For the left-to-right direction, consider a sequence σ that is blocking for π and not π' , of
 1389 minimal length. Let σ_+ and $(n : w)$ be such that $\sigma = (n : w)\sigma_+$.

- 1390 ■ 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' ,
 1391 contradicting the minimality of σ .
- 1392 ■ 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
 1393 know that σ is not blocking for π' . On the other hand, as σ is blocking for π , it is also
 1394 blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$.
- 1395 ■ 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 $\sigma' = \sigma$.
 1396 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
 1397 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
 1398 also be blocking for π' , a contradiction. Hence σ is not blocking for $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2}$
 1399 $\dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$

1400

◁

1401 The claim above lets us define a recursive algorithm.

- 1402 ■ First check if there is a positional word $(n : w)$ that is blocking for $s_0, x_0 \rightsquigarrow t_0, y_0$ and
 1403 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
 1404 σ' that is blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$ and not π' . If it is the case,
 1405 answer yes.
- 1406 ■ Then check if there is a positional word $(n : w)$ that is a blocking factor for $s_0, x_0 \rightsquigarrow t_0, y_0$
 1407 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
 1408 blocking for $s_1, x_1 \rightsquigarrow t_1, y_1 \xrightarrow{a_2} \dots s_k, x_k \rightsquigarrow t_k, y_k$ and not $s'_1, x'_1 \rightsquigarrow t'_1, y'_1 \xrightarrow{a'_2} \dots s'_\ell, x'_\ell \rightsquigarrow$
 1409 t'_ℓ, y'_ℓ . If it is the case, answer yes.

1410 If both items fail, answer no.

1411 The existence of those positional words can be checked in polynomial space using the
 1412 automaton \mathcal{B} constructed in the proof of Lemma F.4. The depth of the recursive calls is at
 1413 most the sum of the lengths of π and π' , which is bounded by $2|\mathcal{A}|$. In consequence, this
 1414 algorithm runs in polynomial space.

1415

◀

1416 ► **Proposition F.7.** *The hardness problem is in PSPACE.*

1417 **Proof.** We use Lemma F.3. We guess an SCC-path $\pi = s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_k, x_k \rightsquigarrow t_k, y_k$
 1418 and an index i .

1419 We check that $s_i, x_i \rightsquigarrow t_i, y_i$ has infinitely many minimal blocking factors, using Lemma F.4.

1420 We then enumerate all SCC-paths in \mathcal{A} . For each one $\pi' = s'_0, x'_0 \rightsquigarrow t'_0, y'_0 \xrightarrow{a'_1} \dots s'_\ell, x'_\ell \rightsquigarrow$
 1421 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 $j^l < j < j^r$
 1422 is \simeq -equivalent to $s, x \rightsquigarrow t, y$.

1423 Then, we use Lemma F.5 to check that there is a sequence σ^l that is blocking for
 1424 $s'_0, x'_0 \rightsquigarrow t'_0, y'_0 \xrightarrow{a'_1} \dots s'_{j^l}, x'_{j^l} \rightsquigarrow t'_{j^l}, y'_{j^l}$ and not $s_0, x_0 \rightsquigarrow t_0, y_0 \xrightarrow{a_1} \dots s_i, x_i \rightsquigarrow t_i, y_i$.

1425 Symmetrically, we check that there is a sequence σ^r that is blocking for $s'_{j^r}, x'_{j^r} \rightsquigarrow$
 1426 $t'_{j^r}, y'_{j^r} \xrightarrow{a'_{j^r}} \dots s'_\ell, x'_\ell \rightsquigarrow t'_\ell, y'_\ell$ and not $s_i, x_i \rightsquigarrow t_i, y_i \xrightarrow{a_{i+1}} \dots s_k, x_k \rightsquigarrow t_k, y_k$.

1427 If all those tests succeed, we answer yes, otherwise we answer no. This algorithm is
 1428 correct and complete by Lemma F.3. ◀

1429 **F.1.2 Testing triviality**

1430 We show the PSPACE upper bound on the complexity of checking if a language is trivial.
 1431 It is based on the characterisation of trivial languages given by Lemma 5.3, and uses the
 1432 following result.

1433 ► **Lemma F.8.** *Given a portal P , we can check whether it has a blocking factor in space*
 1434 *polynomial in $|\mathcal{A}|$.*

1435 **Proof.** We proceed as in the proof of Lemma F.4, except that we only need to check whether
 1436 some final state is reachable from the initial state. ◀

1437 ► **Proposition F.9.** *The triviality problem is in PSPACE.*

1438 **Proof.** Recall that $\mathcal{L}(\mathcal{A})$ is trivial if and only if \mathcal{A} has no blocking sequences.

1439 ▷ **Claim F.10.** \mathcal{A} has a blocking sequence if and only if for all accepting SCC-path π of \mathcal{A} ,
 1440 every portal in π has a blocking factor.

1441 **Proof.** If \mathcal{A} has a blocking sequence then by definition it is blocking for all accepting SCC-
 1442 paths, and thus contains blocking factors for all portals along those SCC-paths.

1443 Suppose that every portal π of every accepting SCC-path has a blocking factor. Then by
 1444 taking the sequence of those blocking factors we obtain a blocking sequence for π . Then,
 1445 since there are finitely many accepting SCC-paths, we can concatenate all those sequences to
 1446 obtain a sequence that is blocking for all accepting SCC-paths. By definition, that sequence
 1447 is blocking for \mathcal{A} . ◀

1448 Therefore, it suffices to enumerate all accepting SCC-paths π in the automaton, and then
 1449 check whether all portals in π have at least one blocking factor, using Lemma F.8. ◀

1450 **F.2 Hardness of classifying automata**

1451 We prove hardness of the triviality problem and easiness problems, concluding on their
 1452 PSPACE-completeness. We reduce from the universality problem for NFAs, which is well-
 1453 known to be PSPACE-complete (see e.g. [1, Theorem 10.14]).

1454 ► **Lemma F.11.** *The triviality and hardness problems are PSPACE-hard.*

1455 **Proof.** Consider an NFA $\mathcal{A} = (Q, \Sigma, \delta, q_0, F)$. Without loss of generality, we assume that
 1456 \mathcal{A} is trim (up to removing unreachable or non-co-reachable states) and that it accepts all
 1457 words of length less than 2: this can be checked in polynomial time and does not affect the
 1458 PSPACE-hardness of universality. Let $\#$ and $!$ be two letters that are not in Σ . We apply
 1459 the following transformations to \mathcal{A} :

- 1460 ■ add a transition labeled by $!$ from every final state to the initial state q_0
- 1461 ■ add a self-loop labeled by $\#$ to each state.

1462 We call the resulting automaton $\mathcal{B} = (Q, \Sigma \cup \{!, \#\}, \delta', q_0, F)$. Note that \mathcal{B} is strongly
 1463 connected: consider any two states $q, q' \in Q$, we show that q' is reachable from q . As \mathcal{A} is
 1464 trim, there exists $q_f \in F$ that is reachable from q , and q' is reachable from the initial state
 1465 q_0 . Furthermore, we have put a $!$ transition from q_f to q_0 , hence q' is reachable from q .

1466 As \mathcal{A} is strongly connected, it is easy to check that its set of minimal blocking sequences
 1467 is exactly its set of minimal blocking factors.

1468 Therefore the language of a strongly connected automaton is trivial if and only if it has
 1469 no minimal blocking factor, and hard if and only if it has infinitely many minimal blocking
 1470 factors. To complete the proof, we show that $\text{MBF}(\mathcal{B})$ is empty when \mathcal{A} is universal and
 1471 infinite otherwise.

1472 First, let us describe the language recognized by \mathcal{B} . It is given by

$$1473 \quad \mathcal{L}(\mathcal{B}) = \{u_1!u_2!\cdots!u_n \mid \forall i, u_i \in (\Sigma \cup \{\#\})^* \wedge \pi_\Sigma(u_i) \in \mathcal{L}(\mathcal{A})\},$$

1474 where $\pi_\Sigma(u)$ is the word in Σ^* obtained by removing all $\#$ from u .

1475 \triangleright Claim F.12. If \mathcal{A} is universal, then \mathcal{B} is also universal.

1476 Proof. Indeed, any word in u in $\Sigma \cup \{\#, !\}$ can be decomposed into $u = u_1!u_2!\cdots!u_n$ where
 1477 each u_i does not contain the letter “!”. As $\#$ is idempotent on \mathcal{B} , $\delta'(q_0, u_i)$ is equal to
 1478 $\delta(q_0, \pi_\Sigma(u_i))$ for every i . Since \mathcal{A} is universal, each of the $\delta(q_0, u_i)$ contains a final state, and
 1479 thus $u_i \in \mathcal{L}(\mathcal{A})$.

1480 By the description of $\mathcal{L}(\mathcal{B})$ given above, we have $u \in \mathcal{B}$. \triangleleft

1481 This shows that if \mathcal{A} is universal, then so is \mathcal{B} and thus $\text{MBF}(\mathcal{B})$ is empty.

1482 Now we show that a word $w \in \Sigma^*$ not in $\mathcal{L}(\mathcal{A})$ induces infinitely many minimal blocking
 1483 factors for \mathcal{B} . Consider such a w of minimal size. As we assumed that \mathcal{A} accepts all words of
 1484 size less than 2, $|w| \geq 2$. Let u, v be words of length at least 1 such that $w = uv$. For all
 1485 $n \in \mathbb{N}$, at least one of $u\#^n v, !u\#^n v, u\#^n v!, !u\#^n v!$ is a minimal blocking factor (depending
 1486 respectively on whether w is not a factor of any word of $\mathcal{L}(\mathcal{A})$ or is a prefix/suffix of a word
 1487 of $\mathcal{L}(\mathcal{A})$ or not). As a consequence, \mathcal{B} has infinitely many blocking factors, and is thus hard
 1488 to test by Theorem 3.2.

1489 In summary, \mathcal{A} is universal if and only if \mathcal{B} is trivial to test, and \mathcal{A} is *not* universal if and
 1490 only if \mathcal{B} is hard to test. This shows the PSPACE-hardness of both the triviality problem
 1491 and the hardness problem. \blacktriangleleft

1492 The above proof can be extended to show the PSPACE-hardness of the easiness problem.

1493 \blacktriangleright **Corollary F.13.** *The easiness problem is PSPACE-hard.*

1494 **Proof.** We proceed as in the proof of Lemma F.11: given an automaton \mathcal{A} over an alphabet
 1495 Σ , we build an automaton \mathcal{B} over the alphabet $\Sigma \cup \{!, \#\}$ such that if \mathcal{A} is universal, $\text{MBF}(\mathcal{B})$
 1496 is empty, and if \mathcal{A} is not universal, then $\text{MBF}(\mathcal{B})$ is infinite.

1497 To show the hardness of the easiness problem, let b denote a new letter not in $\Sigma \cup \{\#, !\}$
 1498 and consider the automaton \mathcal{B}' equal to \mathcal{B} but taken over the alphabet $\Sigma \cup \{\#, !, b\}$. As
 1499 there are no transitions labeled by b in \mathcal{B}' , the word b is always a minimum blocking factor
 1500 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
 1501 $\text{MBF}(\mathcal{B}')$ is finite but non-empty: by Theorem 3.2, this is equivalent to $\mathcal{L}(\mathcal{B}')$ is easy to test.
 1502 Therefore, the easiness problem is also PSPACE-hard. \blacktriangleleft

1503 This concludes the proof of Theorem 6.4