

Asymptotic Approximation by Regular Languages

Ryoma Sin'ya 

Akita University, Akita, Japan

RIKEN AIP, Japan

ryoma@math.akita-u.ac.jp

Abstract

This paper investigates a new property of formal languages called REG-measurability where REG is the class of regular languages. Intuitively, a language L is REG-measurable if there exists an infinite sequence of regular languages that “converges” to L . A language without REG-measurability has a complex shape in some sense so that it can not be (asymptotically) approximated by regular languages. We show that several context-free languages are REG-measurable (including languages with transcendental generating function and transcendental density, in particular), while a certain simple deterministic context-free language and the set of primitive words are REG-immeasurable in a strong sense.

2012 ACM Subject Classification Theory of computation → Formal languages and automata theory

Keywords and phrases Automata, context-free languages, density, primitive words

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

Funding Ryoma Sin'ya: JSPS KAKENHI Grant Number JP19K14582

1 Introduction

Approximating a complex object by more simple objects is a major concept in both computer science and mathematics. In the theory of formal languages, various types of approximations have been investigated (*e.g.*, [15, 16, 10, 7, 5, 8]). For example, Kappes and Kintala [15] introduced *convergent-reliability* and *slender-reliability* which measure how a given deterministic automaton \mathcal{A} nicely approximates a given language L over an alphabet A . Formally \mathcal{A} is said to accept L convergent-reliability if the ratio of the number of *incorrectly* accepted/rejected words of length n

$$\#((L(\mathcal{A})\Delta L) \cap A^n) / \#(A^n)$$

tends to 0 if n tends to infinity, and is said to accept L slender-reliability if the number of incorrectly accepted/rejected words of length n is always bounded above by some constant c : *i.e.*, $\#((L(\mathcal{A})\Delta L) \cap A^n) \leq c$ for any n . Here $L(\mathcal{A})$ denotes the language accepted by \mathcal{A} , $\#(S)$ denotes the cardinality of the set S , \bar{L} denotes the complement of L and Δ denotes the symmetric difference. A slightly modified version of approximation is *bounded- ϵ -approximation* which was introduced by Eisman and Ravikumar. They say that two languages L_1 and L_2 provide a bounded- ϵ -approximation of language L if $L_1 \subseteq L \subseteq L_2$ holds and the ratio of their length- n difference satisfies

$$\#((L_2 \setminus L_1) \cap A^n) / \#(A^n) \leq \epsilon$$

for every sufficiently large $n \in \mathbb{N}$. Perhaps surprisingly, they showed that no pair of regular languages can provide a bounded- ϵ -approximation of the language $\{w \in \{a, b\}^* \mid w \text{ has more } a\text{'s than } b\text{'s}\}$ for any $0 \leq \epsilon < 1$ [10]. This result is a very strong *inapproximable* (by regular languages) example of certain non-regular languages. Also, there is a different framework of approximation so-called *minimal-cover* [8, 5], and a notion represents some *inapproximability* by regular languages so-called *REG-immunity* [12].



© Ryoma Sin'ya;
licensed under Creative Commons License CC-BY
42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:16



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

43 A model of approximation introduced in this paper is rather close to the work of Eisman
 44 and Ravikumar [10]. Instead of approximating by a *single* regular language, we consider an
 45 approximation of some non-regular language L by an *infinite sequence* of regular languages
 46 that “converges” to L . Intuitively, we say that L is REG-measurable if there exists an infinite
 47 sequence of pairs of regular languages $(K_n, M_n)_{n \in \mathbb{N}}$ such that $K_n \subseteq L \subseteq M_n$ holds for all
 48 n and the “size” of the difference $M_n \setminus K_n$ tends to 0 if n tends to infinity. The formal
 49 definition of “size” is formally described in the next section: we use a notion called *density*
 50 (of languages) for measuring the “size” of a language.

51 Although we used the term “approximation” in the title and there are various research on
 52 this topic in formal language theory, our work is strongly influenced by the work of Buck [4]
 53 which investigates, as the title said, *the measure theoretic approach to density*. In [4] the
 54 concept of *measure density* μ of subsets of natural numbers \mathbb{N} was introduced. Roughly
 55 speaking, Buck considered an arithmetic progression $X = \{cn + d \mid n \in \mathbb{N}\}$ (where $c, d \in \mathbb{N}$,
 56 c can be zero) as a “basic set” whose *natural density* as $\delta(X) = 1/c$ if $c \neq 0$ and $\delta(X) = 0$
 57 otherwise, then defined the *outer measure density* $\mu^*(S)$ of any subset $S \subseteq \mathbb{N}$ as

$$58 \quad \mu^*(S) = \inf \left\{ \sum_i \delta(X_i) \mid S \subseteq X \text{ and } X \text{ is a finite union of} \right.$$

$$59 \quad \left. \text{disjoint arithmetic progressions } X_1, \dots, X_k \right\}.$$

61 Then the *measure density* $\mu(S) = \mu^*(S)$ was introduced for the sets satisfying the condition

$$62 \quad \mu^*(S) + \mu^*(\bar{S}) = 1 \tag{1}$$

64 where $\bar{S} = \mathbb{N} \setminus S$. Technically speaking, the class \mathcal{D}_μ of all subsets of natural numbers
 65 satisfying Condition (1) is the *Carathéodory extension* of the class

$$66 \quad \mathcal{D}_0 \stackrel{\text{def}}{=} \{X \subseteq \mathbb{N} \mid X \text{ is a finite union of arithmetic progressions } \},$$

67 see Section 2 of [4] for more details. Notice that here we regard a singleton $\{d\}$ as an
 68 arithmetic progression (the case $c = 0$ for $\{cn + d \mid n \in \mathbb{N}\}$), any finite set belongs to \mathcal{D}_0 .
 69 Buck investigated several properties of μ and \mathcal{D}_μ , and showed that \mathcal{D}_μ *properly* contains \mathcal{D}_0 .

70 In the setting of formal languages, it is very natural to consider the class REG of regular
 71 languages as “basic sets” since it has various types of representation, good closure properties
 72 and rich decidable properties. Moreover, if we consider regular languages REG_A over a unary
 73 alphabet $A = \{a\}$, then REG_A is isomorphic to the class \mathcal{D}_0 ; it is well known that the Parikh
 74 image $\{|w| \mid w \in L\} \subseteq \mathbb{N}$ (where $|w|$ denotes the length of w) of every regular language L in
 75 REG_A is semilinear and hence it is just a finite union of arithmetic progressions. From this
 76 observation, investigating the densities of regular languages and its measure densities (*i.e.*,
 77 REG-measurability) for non-regular languages can be naturally considered as an adaptation
 78 of Buck’s study [4] for formal language theory.

79 Our contribution

80 In this paper we investigate REG-measurability (\simeq asymptotic approximability by regu-
 81 lar languages) of non-regular, mainly context-free languages. The main results consist of
 82 three kinds. We show that: (1) several context-free languages (including languages with
 83 *transcendental generating function* and *transcendental density*) are REG-measurable [The-
 84 orem 23–30]. (2) there are “very large/very small” (deterministic) context-free languages
 85 that are REG-immeasurable in a strong sense [Theorem 36]. (3) the set of *primitive words*

86 is “very large” and REG-immeasurable in a strong sense [Theorem 37–38]. Open problems
87 and some possibility of an application of the notion of measurability to classifying formal
88 languages will be stated in Section 6.

89 The paper is organised as follows. Section 2 provides mathematical background of
90 densities of formal languages. The formal definition of REG-approximability and REG-
91 measurability are introduced in Section 3. The scenario of Section 3 mostly follows one
92 of the measure density introduced by Buck [4] which was described above. In Section 4,
93 we will give several examples of REG-inapproximable but REG-measurable context-free
94 languages. These examples include, perhaps somewhat surprisingly, a language with a
95 *transcendental density* which have been considered as a very complex context-free language
96 from a combinatorial viewpoint. In Section 5, we consider the set of so-called *primitive*
97 *words* and its REG-measurability. Section 6 ends this paper with concluding remarks, some
98 future work and open problems. We assume that the reader has a basic knowledge of formal
99 language theory.

100 2 Densities of Formal Languages

101 For a set S , we write $\#(S)$ for the cardinality of S . The set of natural numbers including
102 0 is denoted by \mathbb{N} . For an alphabet A , we denote the set of all words (resp. all non-empty
103 words) over A by A^* (resp. A^+). We write ε for the empty word and write A^n (resp. $A^{<n}$)
104 for the set of all words of length n (resp. less than n). For a language L , we write $\text{Alph}(L)$
105 for the set of all letters appeared in L . For word $w \in A^*$ and a letter $a \in A$, $|w|_a$ denotes the
106 number of occurrences of a in w . A word v is said to be a *factor* of a word w if $w = xvy$ for
107 some $x, y \in A^*$, further said to be a *prefix* of w if $x = \varepsilon$. For a language $L \subseteq A^*$, we denote
108 by $\bar{L} = A^* \setminus L$ the complement of L .

109 A *language class* \mathcal{C} is a family of languages $\{\mathcal{C}_A\}_{A: \text{finite alphabet}}$ where $\mathcal{C}_A \subseteq 2^{A^*}$ for each
110 A and $\mathcal{C}_A \subseteq \mathcal{C}_B$ for each $A \subseteq B$. We simply write $L \in \mathcal{C}$ if $L \in \mathcal{C}_A$ for some alphabet A .
111 We denote by REG, DetCFL, UnCFL and CFL the class of regular languages, deterministic
112 context-free languages, unambiguous context-free languages and context-free languages,
113 respectively. A language L is said to be \mathcal{C} -*immune* if L is infinite and no infinite subset of L
114 belongs to \mathcal{C} .

115 ► **Definition 1.** Let $L \subseteq A^*$ be a language. The *natural density* $\delta_A(L)$ of L is defined as

$$116 \quad \delta_A(L) \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \frac{\#(L \cap A^n)}{\#(A^n)}$$

117 if the limit exists, otherwise we write $\delta_A(L) = \perp$ and say that L does not have a natural
118 density. The *density* $\delta_A^*(L)$ of L is defined as

$$119 \quad \delta_A^*(L) \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} \frac{\#(L \cap A^k)}{\#(A^k)}$$

120 if its exists, otherwise we write $\delta_A^*(L) = \perp$ and say that L does not have a density. A
121 language $L \subseteq A^*$ is called *null* if $\delta_A^*(L) = 0$, and conversely L is called *co-null* if $\delta_A^*(L) = 1$.

122 ► **Remark 2.** Notice that if L has a natural density (*i.e.*, $\delta_A(L) \neq \perp$), then it also has a
123 density and $\delta_A^*(L) = \delta_A(L)$ holds. But the converse is not true in general, *e.g.*, the case
124 $L = (AA)^*$ (see Example 4 below).

125 The following observation is basic.

23:4 Asymptotic Approximation by Regular Languages

126 ▷ **Claim 3.** Let $K, L \subseteq A^*$ with $\delta_A^*(K) = \alpha, \delta_A^*(L) = \beta$. Then we have:

- 127 1. $\alpha \leq \beta$ if $K \subseteq L$.
- 128 2. $\delta_A^*(L \setminus K) = \beta - \alpha$ if $K \subseteq L$.
- 129 3. $\delta_A^*(\overline{K}) = 1 - \alpha$.
- 130 4. $\delta_A^*(K \cup L) \leq \alpha + \beta$ if $\delta_A^*(K \cup L) \neq \perp$.
- 131 5. $\delta_A^*(K \cup L) = \alpha + \beta$ if $K \cap L = \emptyset$.

132 For more properties of δ_A^* , see Chapter 13 of [3].

133 ► **Example 4.** Here we enumerate a few examples of densities of languages.

- 134 ■ The set of all words A^* clearly satisfies $\delta_A(A^*) = 1$, and its complement \emptyset satisfies
- 135 $\delta_A(\emptyset) = 0$. It is also clear that every finite language is null.
- 136 ■ For the set $\{a\}A^*$ of all words starting with $a \in A$, we have $\#\(\{a\}A^* \cap A^n) / \#(A^n) =$
- 137 $\#(aA^{n-1}) / \#(A^n) = 1 / \#(A)$. Hence $\delta_A(\{a\}A^*) = 1 / \#(A)$.
- 138 ■ Consider $(AA)^*$ the set of all words with even length. Because

$$139 \quad \frac{\#\((AA)^* \cap A^n)}{\#(A^n)} = \begin{cases} 1 & \text{if } n \text{ is even,} \\ 0 & \text{if } n \text{ is odd.} \end{cases}$$

140 holds, its limit does not exist and thus $(AA)^*$ does not have a natural density $\delta_A((AA)^*) =$

141 \perp . However, it has a density $\delta_A^*((AA)^*) = 1/2$.

- 142 ■ The semi-Dyck language

$$143 \quad D \stackrel{\text{def}}{=} \{w \in \{a, b\}^* \mid |w|_a = |w|_b \text{ and } |u|_a \geq |u|_b \text{ for every prefix } u \text{ of } w\}$$

144 is non-regular but context-free. It is well known that the number of words in D of length

145 $2n$ is equal to the n -th Catalan number whose asymptotic approximation is $\Theta(4^n/n^{3/2})$.

146 Thus

$$147 \quad \frac{\#(D \cap A^n)}{\#(A^n)} = \begin{cases} \Theta(1/(n/2)^{3/2}) & \text{if } n \text{ is even,} \\ 0 & \text{if } n \text{ is odd.} \end{cases}$$

148 and we have $\delta_A(D) = 0$, *i.e.*, D is null.

149 Example 4 shows us that, for some regular language L , its natural density is either zero or

150 one, for some, like $L = \{a\}A^*$ (for $\#(A) \geq 2$), $\delta_A(L)$ could be a real number strictly between

151 zero and one, and for some, like $L = (AA)^*$, a natural density may not even exist. However,

152 the following theorem tells us that all regular languages *do* have densities.

153 ► **Theorem 5** (*cf.* Theorem III.6.1 of [21]). *Let $L \subseteq A^*$ be a regular language. Then there is*

154 *a positive integer c such that for all natural numbers $d < c$, the following limit exists*

$$155 \quad \lim_{n \rightarrow \infty} \frac{\#(L \cap A^{cn+d})}{\#(A^{cn+d})}$$

156 *and it is always rational, i.e., the sequence $(\#(L \cap A^n) / \#(A^n))_{n \in \mathbb{N}}$ has only finitely many*

157 *accumulation points and these are rational and periodic.*

158 ► **Corollary 6.** *Every regular language has a density and it is rational.*

159 ► **Corollary 7.** *For any regular language $L \subseteq A^*$, $\delta_A(L) = 0$ if and only if $\delta_A^*(L) = 0$.*

160 Furthermore, for *unambiguous* context-free languages, the following holds.

161 ► **Theorem 8** (Berstel [2]). *For any unambiguous context-free language L over A , its density*

162 *$\delta_A^*(L)$, if it exists (*i.e.*, $\delta_A^*(L) \neq \perp$), is always algebraic.*

163 In the next section we will introduce a language with a transcendental density, which should
 164 be inherently ambiguous due to Theorem 8.

165 We conclude the section by introducing the notion called *dense*: a property about some
 166 topological “largeness” of a language (*cf.* Chapter 2.5 of [3]).

167 ► **Definition 9.** A language $L \subseteq A^*$ is said to be *dense* if the set of all factors of L is equal
 168 to A^* . We say that a word $w \in A^*$ is a *forbidden word* (resp. *forbidden prefix*) of L if
 169 $L \cap A^*wA^* = \emptyset$ (resp. $L \cap wA^* = \emptyset$).

170 Observe that $L \subseteq A^*$ is dense if and only if no word is a forbidden word of L . The next
 171 theorem ties two different notions of “largeness” of languages in the regular case.

172 ► **Theorem 10** (S. [23]). *A regular language is non-null if and only if it is dense.*

173 The “only if”-part of Theorem 10 is nothing but the well-known so-called *infinite monkey*
 174 *theorem* (which states that L is not dense implies L is null), and this part is true for any
 175 (non-regular) languages. But we stress that “if”-part is *not true* beyond regular languages; for
 176 example the semi-Dyck language D is null *but dense* (which will be described in Proposition 12).
 177 We denote by REG^+ the family of non-null regular languages, which is equivalent to the
 178 family of regular languages with positive densities thanks to Corollary 6.

179 **3** Approximability and Measurability

180 Although we will mainly consider REG-measurability of non-regular languages in this paper,
 181 here we define two notions approximability and measurability in general setting, with few
 182 concrete examples.

183 ► **Definition 11.** Let \mathcal{C}, \mathcal{D} be classes of languages. A language L is said to be (\mathcal{C}, ϵ) -*lower-*
 184 *approximable* if there exists $K \in \mathcal{C}$ such that $K \subseteq L$ and $\delta_{\text{Alph}(L)}^*(L \setminus K) \leq \epsilon$. A language
 185 L is said to be (\mathcal{C}, ϵ) -*upper-approximable* if there exists $M \in \mathcal{C}$ such that $L \subseteq M$ and
 186 $\delta_{\text{Alph}(M)}^*(M \setminus L) \leq \epsilon$. A language L is said to be \mathcal{C} -*approximable* if L is both $(\mathcal{C}, 0)$ -lower
 187 and $(\mathcal{C}, 0)$ -upper-approximable. \mathcal{D} is said to be \mathcal{C} -approximable if every language in \mathcal{D} is
 188 \mathcal{C} -approximable.

189 The following proposition gives a simple REG-inapproximable example.

190 ► **Proposition 12.** *The semi-Dyck language D is REG-inapproximable.*

191 **Proof.** We already mentioned that D is null in Example 4, and thus D is $(\text{REG}, 0)$ -lower-
 192 approx by $\emptyset \subseteq D$. One can easily observe that D has no forbidden word: since for any
 193 $w \in A^*$ there exists a pair of natural numbers $(n, m) \in \mathbb{N}^2$ such that $a^n w b^m \in D$. Hence if a
 194 regular language L satisfies $D \subseteq L$, L has no forbidden word, too, and thus L is non-null by
 195 Theorem 10. Thus by Claim 3, $\delta_A^*(L \setminus D) = \delta_A^*(L) - \delta_A^*(D) = \delta_A^*(L) > 0$, which means that
 196 D can not be $(\text{REG}, 0)$ -upper-approximable. ◀

197 The proof of Proposition 12 only depends on the non-existence of forbidden words, hence we
 198 can apply the same proof to the next theorem.

199 ► **Theorem 13.** *Any null language having no forbidden word is $(\text{REG}, 0)$ -upper-inapproximable.*

200 Because D is deterministic context-free, in our term we have:

201 ► **Corollary 14.** *DetCFL is REG-inapproximable.*

23:6 Asymptotic Approximation by Regular Languages

202 Furthermore, by the combination of Theorem 8 and the next theorem, we will know that
 203 there exists a context-free language which can not be approximated by any unambiguous
 204 context-free language.

205 ► **Theorem 15** (Kemp [17]). *Let $A = \{a, b, c\}$. Define*

$$206 \quad S_1 \stackrel{\text{def}}{=} \{a\}\{b^i a^i \mid i \geq 1\}^* \quad S_2 \stackrel{\text{def}}{=} \{a^i b^{2i} \mid i \geq 1\}^* \{a\}^+,$$

207 *and*

$$208 \quad L_1 \stackrel{\text{def}}{=} S_1 \{c\} A^* \quad L_2 \stackrel{\text{def}}{=} S_2 \{c\} A^*.$$

209 *Then $K \stackrel{\text{def}}{=} L_1 \cup L_2$ is a context-free language with a transcendental natural density $\delta_A(K)$.*

210 ► **Corollary 16.** *CFL is UnCFL-inapproximable.*

211 We then introduce the notion of \mathcal{C} -measurability which is a formal language theoretic
 212 analogue of Buck's measure density [4].

213 ► **Definition 17.** Let \mathcal{C}, \mathcal{D} be classes of languages. For a language L , we define its \mathcal{C} -lower-
 214 density as

$$215 \quad \underline{\mu}_{\mathcal{C}}(L) \stackrel{\text{def}}{=} \sup\{\delta_A^*(K) \mid A = \text{Alph}(L), K \subseteq L, K \in \mathcal{C}_A, \delta_A^*(K) \neq \perp\}$$

216 and its \mathcal{C} -upper-density as

$$217 \quad \overline{\mu}_{\mathcal{C}}(L) \stackrel{\text{def}}{=} \inf\{\delta_A^*(K) \mid A = \text{Alph}(L), L \subseteq K, K \in \mathcal{C}_A, \delta_A^*(K) \neq \perp\}.$$

218 A language L is said to be \mathcal{C} -measurable if $\overline{\mu}_{\mathcal{C}}(L) = \underline{\mu}_{\mathcal{C}}(L)$ holds, and we simply write $\overline{\mu}_{\mathcal{C}}(L)$
 219 as $\mu_{\mathcal{C}}(L)$. \mathcal{D} is said to be \mathcal{C} -measurable if every language in \mathcal{D} is \mathcal{C} -measurable.

220 ► **Definition 18.** We call $\overline{\mu}_{\mathcal{C}}(L) - \underline{\mu}_{\mathcal{C}}(L)$ the \mathcal{C} -gap of a language L . We say that a language
 221 L has full \mathcal{C} -gap if its \mathcal{C} -gap equals to 1, i.e., $\overline{\mu}_{\mathcal{C}}(L) - \underline{\mu}_{\mathcal{C}}(L) = 1$.

222 In the next section, we describe several examples of both REG-measurable and REG-
 223 immeasurable languages. The REG-gap could be a good measure how much a given language
 224 has a complex shape from the viewpoint of regular languages.

225 The following lemmata are basic.

226 ► **Lemma 19.** *Let K, L be two languages.*

- 227 1. $\overline{\mu}_{\mathcal{C}}(K) \leq \overline{\mu}_{\mathcal{C}}(L)$ if $K \subseteq L$.
- 228 2. $\overline{\mu}_{\mathcal{C}}(K \cup L) \leq \overline{\mu}_{\mathcal{C}}(K) + \overline{\mu}_{\mathcal{C}}(L)$ if \mathcal{C} is closed under union.
- 229 3. $\overline{\mu}_{\mathcal{C}}(K) = \delta_A^*(K)$ if $K \in \mathcal{C}$ and $\delta_A^*(K) \neq \perp$.

230 ► **Lemma 20.** *Let \mathcal{C} be a language class such that \mathcal{C} is closed under complement and every
 231 language in \mathcal{C} has a density. A language $L \subseteq A^*$ is \mathcal{C} -measurable if and only if*

$$232 \quad \overline{\mu}_{\mathcal{C}}(L) + \overline{\mu}_{\mathcal{C}}(\overline{L}) = 1. \tag{2}$$

234 **Proof.** Let L be a language and $A = \text{Alph}(L)$. By definition, L satisfies Condition (2) if and
 235 only if

$$236 \quad \inf\{\delta_A^*(K) \mid L \subseteq K, K \in \mathcal{C}\} = 1 - \inf\{\delta_A^*(K) \mid \overline{L} \subseteq K, K \in \mathcal{C}\} \tag{3}$$

238 holds. On the other hand, L is measurable if and only if

$$239 \quad \inf\{\delta_A^*(K) \mid L \subseteq K, K \in \mathcal{C}\} = \sup\{\delta_A^*(K) \mid K \subseteq L, K \in \mathcal{C}\}. \tag{4}$$

241 For any language $K \in \mathcal{C}_A$ such that $K \subseteq L$ and $\delta_A^*(K) \neq \perp$, its complement \overline{K} satisfies
 242 $\overline{L} \subseteq \overline{K}$ and $\delta_A^*(\overline{K}) = 1 - \delta_A^*(K)$. This means that if \mathcal{C}_A is closed under complement then
 243 $\sup\{\delta_A^*(K) \mid K \subseteq L, K \in \mathcal{C}_A\} = 1 - \inf\{\delta_A^*(K) \mid \overline{L} \subseteq K, K \in \mathcal{C}_A\}$, holds, which immediately
 244 implies the equivalence of Condition (3) and Condition (4). ◀

4 REG-measurability on Context-free Languages

In this section we examine REG-measurability of several types of context-free languages. The first type of languages (Section 4.1) is null context-free languages. Although some null language can have a full REG-gap as stated in the next theorem, we will show that typical null context-free languages are REG-measurable.

► **Theorem 21.** *There is a recursive language L which is null but $\bar{\mu}_{\text{REG}}(L) = 1$.*

Proof. Let A be an alphabet with $\#(A) \geq 2$ and let $(\mathcal{A}_i)_{i \in \mathbb{N}}$ be an enumeration of automata over A such that $\text{REG}_A = \{L(\mathcal{A}_i) \mid i \in \mathbb{N}\}$; we can take such enumeration by enumerating some binary representation of automata via shortlex order $<_{\text{lex}}$. We will construct a null language L such that $\bar{\mu}_{\text{REG}}(L) = 1$, in particular, L is not a subset of every regular co-infinite language.

Consider the following program P which takes an input word w :

Step 1 set $i = 0$ and $\ell = 0$.

Step 2 check $L(\mathcal{A}_i)$ is co-infinite (*i.e.*, the complement $\overline{L(\mathcal{A}_i)}$ is infinite) or not.

Step 3 if $L(\mathcal{A}_i)$ is co-finite, then set $i = i + 1$ and go back to Step 2.

Step 4 otherwise, pick u such that u is the smallest (with respect to $<_{\text{lex}}$) word satisfying $|u| > \ell$ and $u \notin L(\mathcal{A}_i)$ (such u surely exists since $L(\mathcal{A}_i)$ is co-infinite).

Step 5 if $w = u$ then P accepts w and halts.

Step 6 if $w <_{\text{lex}} u$ then P rejects w and halts.

Step 7 if $u <_{\text{lex}} w$ then set $\ell = |u|$, $i = i + 1$ and go back to Step 2.

One can easily observe that all Steps are effective and P ultimately halts for any input word w because the length of the word u in Step 4 is strictly increasing until $u = w$ or $w <_{\text{lex}} u$. Thus the language $L \stackrel{\text{def}}{=} \{w \in A^* \mid P \text{ accepts } w\}$ is recursive. Moreover, L satisfies the following properties: (1) $L \not\subseteq R$ for any regular co-infinite language because by Step (4–5) P accepts some word $w \notin R$, and (2) $\delta_A(L) = 0$; by Step (5–6) and the length of u is strictly increasing, P rejects every word in A^n except for one single word u , for each n . Clearly, (2) implies $\delta_A(L) = 0$, and (1) implies $\bar{\mu}_{\text{REG}}(L) = 1$ since every language R with $\delta_A^*(R) < 1$ is co-infinite. ◀

The second type of languages (Section 4.2) is inherently ambiguous languages and the third type of languages (Section 4.3) includes Kemp's language K whose density is transcendental. The last type of languages (Section 4.4) is languages with full REG-gap, *i.e.*, strongly REG-immeasurable languages.

4.1 Null Context-free Languages

First we consider the following language with constraints on the number of occurrences of letters, which is a very typical example of a non-regular but context-free language.

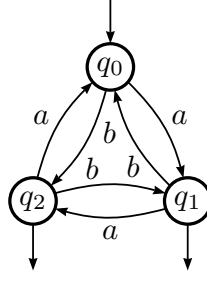
► **Definition 22.** For an alphabet A and letters $a, b \in A$ such that $a \neq b$, we define

$$L_A(a, b) \stackrel{\text{def}}{=} \{w \in A^* \mid |w|_a = |w|_b\}.$$

► **Theorem 23.** $L_A(a, b)$ is REG-measurable where $A = \{a, b\}$.

Proof. It is enough to show that the complement $L = \overline{L_A(a, b)}$ satisfies $\underline{\mu}_{\text{REG}}(L) = 1$. For each $k \geq 1$, we define

$$L_k \stackrel{\text{def}}{=} \{w \in A^* \mid |w|_a \neq |w|_b \pmod k\}.$$



■ **Figure 1** The deterministic automaton \mathcal{A}_3 in the Proof of Theorem 23. Here, the state q_0 having unlabelled incoming arrow is initial and the states q_1, q_2 having unlabelled outgoing arrow are final.

Clearly, $L_k \subseteq L$ holds. Each L_k is recognised by a k -states deterministic automaton

$$\mathcal{A}_k = (Q_k = \{q_0, \dots, q_{k-1}\}, \Delta_k : Q_k \times A \rightarrow Q_k, q_0, Q_k \setminus \{q_0\})$$

where

$$\Delta_k(q_i, a) = q_{i+1 \bmod k} \quad \Delta_k(q_i, b) = q_{i-1 \bmod k} \quad (\text{for each } i \in \{0, \dots, k-1\}),$$

q_0 is the initial state, and any other state $q \in Q_k \setminus \{q_0\}$ is a final state (the case $k = 3$ is depicted in Fig 1). The adjacency matrix of \mathcal{A}_k is

$$M_k = \begin{bmatrix} 0 & 1 & 0 & \cdots & \cdots & 1 \\ 1 & 0 & 1 & \ddots & & \vdots \\ 0 & 1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 1 & 0 \\ \vdots & & \ddots & 1 & 0 & 1 \\ 1 & \cdots & \cdots & 0 & 1 & 0 \end{bmatrix} = E_k + E_k^{k-1} \quad \text{where } E_k = \begin{bmatrix} 0 & 0 & 0 & \cdots & \cdots & 1 \\ 1 & 0 & 0 & \ddots & & \vdots \\ 0 & 1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 & 0 \\ \vdots & & \ddots & 1 & 0 & 0 \\ 0 & \cdots & \cdots & 0 & 1 & 0 \end{bmatrix}.$$

M_k is a special case of *circulant matrices*. A k -dimensional circulant matrix C_k is a matrix that can be represented by a polynomial of E_k :

$$C_k = p(E_k) = \sum_{n=0}^{k-1} c_n E_k^n$$

and it is well known that C_k can be diagonalised as, for a k -th root of unity $\xi_k = e^{-\frac{2\pi i}{k}}$ (where i is the imaginary unit),

$$\frac{1}{\sqrt{k}} F_k^H \cdot C_k \cdot \frac{1}{\sqrt{k}} F_k = \text{diag}(p(1), p(\xi_k^{-1}), p(\xi_k^{-2}), \dots, p(\xi_k^{-(k-1)}))$$

where $F_k = (f_{n,m})$ with $f_{n,m} = \xi_k^{(n-1)(m-1)}$ (for $1 \leq n, m \leq k$) is the k -dimensional *Fourier matrix*, F_k^H is its Hermitian transpose and $\text{diag}(\lambda_1, \dots, \lambda_k)$ is the diagonal matrix whose n -th diagonal element is λ_n (for $1 \leq n \leq k$) (cf. Section 5.2.1 of [18]). Hence, in the case of $M_k = p_{\mathcal{A}_k}(E_k) = E_k + E_k^{k-1}$, we have

$$\frac{1}{\sqrt{k}} F_k^H \cdot M_k \cdot \frac{1}{\sqrt{k}} F_k = \text{diag}(2, \xi_k^{-1} + \xi_k, \xi_k^{-2} + \xi_k^2, \dots, \xi_k^{-(k-1)} + \xi_k^{k-1}) \quad (5)$$

306 because, for any $n \geq 0$, $p_{\mathcal{A}_k}(\xi_k^{-n}) = \xi_k^{-n} + \xi_k^{-n(k-1)} = \xi_k^{-n} + \xi_k^n$ holds.

307 Let $\Lambda_k = \text{diag}(2, \xi_k^{-1} + \xi_k, \xi_k^{-2} + \xi_k^2, \dots, \xi_k^{-(k-1)} + \xi_k^{k-1})$. Because \mathcal{A}_k is deterministic and
 308 the final states are all but q_0 , the number of words of length n in L_k is exactly the number
 309 of paths from q_0 to any other state in \mathcal{A}_k . For the k -dimensional vectors $\mathbf{e} = (1, 0, 0, \dots, 0)$
 310 and $\mathbf{1} = (1, 1, 1, \dots, 1)$, from Equation (5) we have

$$\begin{aligned}
 311 \quad \#(L_k \cap A^n) &= \mathbf{e} \cdot M_k^n \cdot (\mathbf{1} - \mathbf{e})^T \\
 312 \quad &= \frac{1}{k} \mathbf{e} \cdot F_k \cdot \Lambda_k^n \cdot F_k^H (\mathbf{1} - \mathbf{e})^T \\
 313 \quad &= \frac{1}{k} \mathbf{1} \cdot \Lambda_k^n \cdot \left(k-1, \sum_{j=1}^{k-1} \xi_k^{-j}, \sum_{j=1}^{k-1} \xi_k^{-2j}, \dots, \sum_{j=1}^{-(k-1)} \xi_k^{-(k-1)j} \right)^T \\
 314 \quad &= \frac{1}{k} \left(2^n(k-1) + (\xi_k^{-1} + \xi_k)^n \sum_{j=1}^{k-1} \xi_k^{-j} + \dots + (\xi_k^{-(k-1)} + \xi_k^{k-1})^n \sum_{j=1}^{k-1} \xi_k^{-(k-1)j} \right). \quad (6) \\
 315
 \end{aligned}$$

316 If k is odd $k = 2m + 1$, then for any $1 \leq j \leq k-1$, $\xi_k^{-j} + \xi_k^j$ is a real number whose
 317 absolute value is strictly smaller than 2; because ξ_k^{-j} is the complex conjugate of ξ_k^j and
 318 hence $|\xi_k^{-j} + \xi_k^j| = |2\text{Re}(\xi_k^j)| < 2$ for odd k . Hence from Equation (6) we can deduce that

$$319 \quad \#(L_k \cap A^n) = \frac{k-1}{k} 2^n + o(2^n)$$

320 where $o(2^n)$ means some function such that $\lim_{n \rightarrow \infty} o(2^n)/2^n = 0$. Thus we have $\delta_A(L_k) =$
 321 $\frac{k-1}{k}$ for odd $k = 2m + 1$, which tends to 1 if k tends to infinity, *i.e.*, $\mu_{\text{REG}}(L) = 1$. This
 322 completes the proof. ◀

323 By Theorem 23, it is also true that any subset of $L_{\{a,b\}}(a, b)$ is REG-measurable. In
 324 particular, we have:

325 ▶ **Corollary 24.** *The semi-Dyck language $D \subseteq L_{\{a,b\}}(a, b)$ is REG-measurable.*

326 The next example is the set of all palindromes.

327 ▶ **Theorem 25.** $P_A \stackrel{\text{def}}{=} \{w \in A^* \mid w = \text{rev}(w)\}$ is REG-measurable.

328 **Proof.** Because the case $\#(A) = 1$ is trivial ($P_A = A^*$), we assume that $\#(A) \geq 2$. It is
 329 enough to show that the complement $\overline{P_A}$ is REG-measurable.

330 For each $k \geq 1$, we define

$$331 \quad L_k \stackrel{\text{def}}{=} \{w_1 A^* w_2 \mid w_1, w_2 \in A^k, w_1 \neq \text{rev}(w_2)\}.$$

332 One can easily observe that $L_k \subseteq \overline{P_A}$ for each $k \geq 1$. Moreover, for any $n > 2k$, the number
 333 of words in L_k of length n is

$$334 \quad \#(L_k \cap A^n) = \#(A)^k \cdot \#(A)^{n-2k} \cdot (\#(A)^k - 1) = \#(A)^n - \#(A)^{n-k}.$$

335 From this we can conclude that $\delta_A(L_k) = 1 - \#(A)^{-k}$ and it tends to 1 if k tends to infinity.
 336 Thus we have $\mu_{\text{REG}}(\overline{P_A}) = 1$. ◀

337 **4.2 Some Inherently Ambiguous Languages**

338 There are REG-measurable inherently ambiguous context-free languages. Since every *bounded*
 339 *language* $L \subseteq w_1^* \cdots w_k^*$ is trivially REG-measurable ($\mu_{\text{REG}}(L) = 0$), a typical example of an
 340 inherently ambiguous context-free language $\{a^i b^j c^k \mid i = j \text{ or } i = k\}$ is REG-measurable.

341 Some more complex examples of inherently ambiguous languages are the following
 342 languages with constraints on the number of occurrences of letters investigated by Flajolet [13]:

343
$$\mathbf{O}_3 \stackrel{\text{def}}{=} \{w \in \{a, b, c\}^* \mid |w|_a = |w|_b \text{ or } |w|_a = |w|_c\},$$

 344
$$\mathbf{O}_4 \stackrel{\text{def}}{=} \{w \in \{x, \bar{x}, y, \bar{y}\}^* \mid |w|_x = |w|_{\bar{x}} \text{ or } |w|_y = |w|_{\bar{y}}\}.$$

346 **► Theorem 26.** \mathbf{O}_3 and \mathbf{O}_4 are REG-measurable.

347 **Proof.** Let $A = \{a, b, c\}$. For the case \mathbf{O}_3 , in a very similar way to Theorem 23, we
 348 can construct a sequence of automata $(\mathcal{A}_k^{ab})_{k \in \mathbb{N}}$ such that each automaton \mathcal{A}_k^{ab} satisfies
 349 $L(\mathcal{A}_k^{ab}) \subseteq \overline{L_A(a, b)}$ and its adjacency matrix is of the form

350
$$M_k^{ab} = M_k + I_k = \begin{bmatrix} 1 & 1 & 0 & \cdots & \cdots & 1 \\ 1 & 1 & 1 & \ddots & & \vdots \\ 0 & 1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 1 & 0 \\ \vdots & & \ddots & 1 & 1 & 1 \\ 1 & \cdots & \cdots & 0 & 1 & 1 \end{bmatrix}$$

351

352 where M_k is the adjacency matrix stated in Theorem 23 and I_k is the k -dimensional identity
 353 matrix. The automaton \mathcal{A}_k^{ab} is obtained by just adding self-loop labeled by c for each state
 354 $q \in Q_k$ of \mathcal{A}_k in Theorem 23. This sequence of automata ensures that the language $L_A(a, b)$
 355 is REG-measurable ($\bar{\mu}_{\text{REG}}(L_A(a, b)) = 0$, in particular). The same argument is applicable to
 356 the language $L_A(a, c)$, thus these union $\mathbf{O}_3 = L_A(a, b) \cup L_A(a, c)$ is also REG-measurable by
 357 Lemma 19. The case \mathbf{O}_4 can be achieved in the same manner. ◀

358 Next we consider the so-called *Goldstine language*

359
$$\mathbf{G} \stackrel{\text{def}}{=} \{a^{n_1} b a^{n_2} b \cdots a^{n_p} b \mid p \geq 1, n_i \neq i \text{ for some } i\}.$$

360 While \mathbf{G} can be accepted by a non-deterministic pushdown automaton, its generating function
 361 is not algebraic [14] and thus it is an inherently ambiguous context-free language due to the
 362 well-known Chomsky–Schützenberger theorem stating that the generating function of every
 363 unambiguous context-free language is algebraic [6].

364 **► Theorem 27.** \mathbf{G} is REG-measurable.

365 **Proof.** Let $A = \{a, b\}$. Observe that $\mathbf{G} \subseteq A^* b$ and $\bar{\mu}_{\text{REG}}(\mathbf{G}) \leq \delta_A(A^* b) = 1/2$. Let

366
$$L_{\mathbf{G}} = \{u \in A^* \mid uA^*\{b\} \cap \overline{\mathbf{G}} = \emptyset\}$$

367 be the set of all forbidden prefixes of the complement $\overline{\mathbf{G}}$. For each $k \geq 1$, we define

368
$$L_k \stackrel{\text{def}}{=} \{uA^*\{b\} \mid u \in L_{\mathbf{G}} \cap A^k\}.$$

369 If a word u is in $L_{\mathbf{G}}$, then by definition of $L_{\mathbf{G}}$, uvb is always in \mathbf{G} for any word v , thus
 370 $L_k \subseteq \mathbf{G}$ holds for each k . Any word in $\overline{L_{\mathbf{G}}} = A^* \setminus L_{\mathbf{G}}$ is a prefix of the infinite word

371 $a^{n_1}ba^{n_2}ba^{n_3}b \dots$ ($n_i = i$ for each $i \in \mathbb{N}$) thus $\#(L_G \cap A^n) = \#(A^n) - 1$ holds for each $n \geq 1$.
 372 Hence we have

$$373 \quad \delta_A(L_k) = \lim_{n \rightarrow \infty} \frac{\#(L_k \cap A^n)}{\#(A^n)} = \lim_{n \rightarrow \infty} \frac{(\#(A^k) - 1) \cdot \#(A^{n-k-1})}{\#(A^n)}$$

$$374 \quad = (\#(A)^k - 1) \cdot \#(A)^{-k-1} = 2^{-1} - 2^{-k-1}.$$

376 This implies that $\delta_A(L_k)$ tends to $1/2$. Thus $\mu_{\text{REG}}(\mathbf{G}) = 1/2$. \blacktriangleleft

377 In general, for an infinite word $w \in A^\omega$, the set

$$378 \quad \text{Copref}(w) \stackrel{\text{def}}{=} A^* \setminus \{u \in A^* \mid u \text{ is a prefix of } w\}$$

379 is called the *coprefix language of w* . The proof of Theorem 27 uses a key property that \mathbf{G} can
 380 be characterised by using the coprefix language of the infinite word $w = a^{n_1}ba^{n_2}ba^{n_3}b \dots$ as
 381 $\mathbf{G} = \text{Copref}(w) \cap \{a, b\}^* \{b\}$ which was pointed out in [1]. Thus by the same argument, we
 382 can say that any coprefix language L is REG-measurable ($\mu_{\text{REG}}(L) = 1$, in particular).

383 For coprefix languages, the following nice ‘‘gap theorem’’ holds.

384 **► Theorem 28** (Autebert–Flajolet–Gaborro [1]). *Let $w \in A^\omega$ be an infinite word generated by*
 385 *an iterated morphism, i.e., $w = h(w) = h^\omega(a)$ for some monoid morphism $h : A^* \rightarrow A^*$ and*
 386 *letter $a \in A$. Then for the coprefix language $L = \text{Copref}(w)$ there are only two possibilities:*

- 387 1. L is a regular language.
- 388 2. L is an inherently ambiguous context-free language.

389 This means that we can construct, by finding some suitable morphism h , many examples of
 390 inherently ambiguous context-free languages.

391 4.3 K: A Language with Transcendental Density

392 We now show the fact that the language \mathbf{K} defined by Kemp [17] (recall that the definition of
 393 \mathbf{K} appeared in Thorem 15) is REG-measurable. We will actually show a more general result
 394 regarding the following type of languages.

395 **► Definition 29.** Let $L \subseteq A^*$ be a language and $c \notin A$ be a letter. We call the language
 396 $L\{c\}(A \cup \{c\})^*$ over $A \cup \{c\}$ *suffix extension of L by c* .

397 **► Theorem 30.** *The suffix extension $L' \subseteq (A \cup \{c\})^*$ of any language $L \subseteq A^*$ by $c \notin A$ is*
 398 *REG-measurable.*

399 **Proof.** Let $B = A \cup \{c\}$ and $k = \#(B)$. We first show that L' has a natural density. For
 400 any words $u, v \in L$ with $u \neq v$, two languages $u\{c\}B^*$ and $v\{c\}B^*$ are disjoint, and clearly

$$401 \quad \#(u\{c\}B^* \cap B^n) / \#(B^n) = \#(u\{c\}B^{n-|u|-1}) / \#(B^n) = k^{n-|u|-1} / k^n = k^{-(|u|+1)}$$

402 holds for $n > |u|$ thus $\delta_B(u\{c\}B^*) = k^{-(|u|+1)}$. The natural density of L' is

$$403 \quad \delta_B(L') = \lim_{n \rightarrow \infty} \frac{\#(L' \cap B^n)}{\#(B^n)} = \lim_{n \rightarrow \infty} \frac{\#(\bigcup_{w \in L} (w\{c\}B^* \cap B^n))}{\#(B^n)}$$

$$404 \quad = \lim_{n \rightarrow \infty} \frac{\sum_{w \in L} \#(w\{c\}B^* \cap B^n)}{\#(B^n)} = \lim_{n \rightarrow \infty} \sum_{w \in (L \cap A^{<n})} k^{-(|w|+1)}. \quad (7)$$

406 Because the sequence $(\sum_{w \in (L \cap A^{<n})} k^{-(|w|+1)})_{n \in \mathbb{N}}$ is non-decreasing and bounded above by
 407 1, the limit (7) exists, say $\delta_B(L') = \alpha$.

23:12 Asymptotic Approximation by Regular Languages

408 For each $n \in \mathbb{N}$, the language $L_n \stackrel{\text{def}}{=} \bigcup_{w \in L \cap A^{<n}} w\{c\}B^*$ is regular (since $L \cap A^{<n}$ is
 409 finite), $L_n \subseteq L'$ and $\delta_B(L_n) = \sum_{w \in (L \cap A^{<n})} k^{-(|w|+1)}$. Hence $\mu_{\text{REG}}(L') = \alpha$. By similar
 410 argument, for each $n \in \mathbb{N}$, we can claim that the language $K_n \stackrel{\text{def}}{=} B^* \setminus \bigcup_{w \in \bar{L} \cap A^{<n}} w\{c\}B^*$
 411 satisfies $K_n \supseteq L'$ and $\delta_B(K_n)$ tends to α if n tends to infinity. Thus $\mu_{\text{REG}}(L') = \alpha$. ◀

412 Since K is the suffix extensions of the union $S_1 \cup S_2$ in Theorem 15, we have:

413 ▶ **Corollary 31.** K is REG-measurable.

414 ▶ **Remark 32.** Theorem 30 indicates that REG-measurability is a quite relaxed property
 415 in some sense: even for a non-recursively-enumerable language, its suffix extension is still
 416 non-recursively-enumerable but REG-measurable. Moreover, because the class of recursively
 417 enumerable languages is just a countable set, there exist *uncountably many* REG-measurable
 418 non-recursively-enumerable languages.

419 The same proof method works for the *prefix extension* and the *infix extension* (see the
 420 full version [22] for details).

421 The same proof method works for the *prefix extension* and the *infix extension*.

422 ▶ **Theorem 33.** Let $c \notin A$ and $A' = A \cup \{c\}$. The prefix extension $L' = A'^*\{c\}L$ of any
 423 language $L \subseteq A^*$ is REG-measurable. Also, the infix extension $L'' = A'^*\{c\}L\{c\}A'^*$ of any
 424 language $L \subseteq A^*$ is REG-measurable, $\mu_{\text{REG}}(L'') = 0$ if $L = \emptyset$, $\mu_{\text{REG}}(L'') = 1$ otherwise, in
 425 particular.

426 **Proof.** The prefix extension of L is just the reverse of the suffix extension of L , the same
 427 proof method trivially works. For the infix extension $L'' = A'^*\{c\}L\{c\}A'^*$, if $L = \emptyset$ then L''
 428 is also empty and thus $\mu_{\text{REG}}(L'') = 0$. Further, if $L \neq \emptyset$ then there is a word $w \in L$ and
 429 thus $A'^*cwcA'^* \subseteq L''$ holds, which means that $\delta_{A'}(A'^*cwcA'^*) = 1$ by the infinite monkey
 430 theorem and we have $\mu_{\text{REG}}(L'') = 1$. ◀

4.4 Languages with Full REG-Gap

432 In Section 4.1, we showed that the language $L_{\{a,b\}}(a, b)$ is REG-measurable. On the other
 433 hand, by the result of Eisman–Ravikumar [10], we will know that the closely related language

$$434 \quad M \stackrel{\text{def}}{=} \{w \in \{a, b\}^* \mid |w|_a > |w|_b\},$$

435 sometimes called the *majority language*, is not REG-measurable. This contrast is interesting.

436 ▶ **Theorem 34** (Eisman–Ravikumar [10, 11]). Let $A = \{a, b\}$ and $L \subseteq A^*$ be a regular
 437 language. Then $M \subseteq L$ implies

$$438 \quad \limsup_{n \rightarrow \infty} \{\#(\bar{L} \cap A^n) / \#(A^n)\} = 0.$$

439 One can easily observe that $\limsup_{n \rightarrow \infty} \{\#(\bar{L} \cap A^n) / \#(A^n)\} = 0$ if and only if $\delta_A(\bar{L}) = 0$,
 440 which means that any regular superset of M is co-null. Thus the above theorem implies that
 441 both M and \bar{M} are REG^+ -immune, hence we have:

442 ▶ **Corollary 35.** M has full REG-gap.

443 By using the infinite monkey theorem and some probabilistic arguments, we can generalise
 444 the previous theorem as follows.

445 ► **Theorem 36.** For any $m \geq 1$, the following language over $A = \{a, b\}$

$$446 \quad M_m \stackrel{\text{def}}{=} \{w \in A^* \mid |w|_a > m \cdot |w|_b\}$$

447 has full REG-gap, and $\delta_A(M_m) = 1/2$ if $m = 1$ otherwise $\delta_A(M_m) = 0$.

448 **Proof.** First we prove that any non-null regular language L can not be a subset of M_m . Let $\eta : A^* \rightarrow M$ be the syntactic morphism η and monoid M of L , and let $c = \max_{m \in M} \min_{w \in \eta^{-1}(m)} |w|$ (this is well-defined natural number since M is finite). By the infinite monkey theorem, 450 L is not null implies that L has no forbidden word, and thus for the word b^{2c} there exist 451 two words x and y such that $xb^{2c}y$ is in L . We can assume that $|x|, |y| \leq c$ without loss of 452 generality by the definition of c , which implies $|xb^{2c}y|_a \leq |x| + |y| = 2c \leq |xb^{2c}y|_b$ hence 453 $xb^{2c}y \notin M_m$. Thus $L \not\subseteq M_m$ and $\underline{\mu}_{\text{REG}}(M_m) = 0$. By using same argument, we can prove 454 that $\overline{\mu}_{\text{REG}}(M_m) = 1$ and hence M_m has full REG-gap. 455

456 In the case $m = 1$, $\delta_A(M_1) = \delta_A(M) = 1/2$ is obvious. It is enough to show that 457 $\delta_A(M_2) = 0$ holds (since $M_m \subseteq M_2$ for any $m \leq 2$). Indeed, we have

$$458 \quad \delta_A(M_2) = \lim_{n \rightarrow \infty} \frac{\#\{w \in A^n \mid |w|_a > 2|w|_b\}}{2^n} = \lim_{n \rightarrow \infty} \frac{\#\{w \in A^n \mid |w|_a > 2n/3\}}{2^n}$$

$$459 \quad = \lim_{n \rightarrow \infty} \Pr(|\overline{X}_n - n/2| > n/6) = 0$$

461 where $\Pr(|\overline{X}_n - n/2| > n/6)$ means the probability that the absolute value of the difference 462 of the number \overline{X}_n of the occurrences of a 's in a randomly chosen word of length n and its 463 mean value $n/2$ is larger than $n/6$; it tends to zero by the weak law of large numbers. ◀

464 5 REG-Immeasurability of Primitive Words

465 A non-empty word $w \in A^+$ is said to be primitive if $u^n = w$ implies $u = w$ for any $u \in A^+$ 466 and $n \in \mathbb{N}$. The set of all primitive words over A is denoted by Q_A . Because the case 467 $\#(A) = 1$ is meaningless ($Q_A = A$ in this case), hereafter we always assume $\#(A) \geq 2$. 468 Whether Q_A is context-free or not is a well-known long-standing open problem posed by 469 Dömösi, Horváth and Ito [9]. Reis and Shyr [20] proved $Q_A^2 = A^+ \setminus \{a^n \mid a \in A, n \neq 2\}$, 470 which intuitively means that every non-empty word w not a power of a letter is a product of 471 two primitive words. From this result one may think that Q_A is “very large” in some sense. 472 Actually, Q_A is somewhat “large” (it is dense in the sense of Definition 9), but we can show 473 more stronger property as follows.

474 ► **Theorem 37.** $\delta_A(Q_A) = 1$.

475 **Proof.** It is enough to show that $\delta_A(\overline{Q_A}) = 0$ holds. One can easily observe that any natural 476 number $n \in \mathbb{N}$ has at most $2\sqrt{n}$ divisors. In addition, for any non-primitive word $w = v^m$ of 477 length n is uniquely determined by v (since $m = n/|v|$) and $|v| \leq n/2$. Hence the number of 478 non-primitive words of length n satisfies

$$479 \quad \#\overline{(Q_A \cap A^n)} \leq 2\sqrt{n} \sum_{i=0}^{\lfloor n/2 \rfloor} \#(A^i) \leq 2\sqrt{n} \cdot \#(A)^{\lfloor n/2 \rfloor + 1}.$$

480 By using the above estimation, we can deduce that

$$481 \quad \frac{\#\overline{(Q_A \cap A^n)}}{\#(A^n)} \leq \frac{2\sqrt{n} \cdot \#(A)^{\lfloor n/2 \rfloor + 1}}{\#(A)^n} \leq \frac{2\sqrt{n}}{\#(A)^{n/2 - 1}}$$

482 and it tends to 0 if n tends to infinity (since we assume $\#(A) \geq 2$). Thus $\delta_A(\overline{Q_A}) = 0$. ◀

23:14 Asymptotic Approximation by Regular Languages

483 While Q_A is “very large” (co-null) as stated above, we can also prove that Q_A is REG^+ -
 484 immune. The proof relies on an analysis of the structure of the syntactic monoid of a non-null
 485 regular language. We assume that the reader has a basic knowledge of semigroup theory
 486 (*cf.* [19]): Green’s relations $\mathcal{J}, \mathcal{R}, \mathcal{L}, \mathcal{H}$ and a direct consequence of Green’s theorem (an
 487 \mathcal{H} -class H in a semigroup S is a subgroup of S if and only if H contains an idempotent), in
 488 particular.

489 ► **Theorem 38.** *Any non-null regular language contains infinitely many non-primitive words,*
 490 *and hence $\mu_{\text{REG}}(Q_A) = 0$.*

491 **Proof.** Let L be a regular language over A with a positive density $\delta_A(L) > 0$. We consider
 492 $\eta : A^* \rightarrow M$ the syntactic morphism η and the syntactic monoid M of L , and let S be a
 493 subset of M satisfying $\eta^{-1}(S) = L$. L is regular means that M is finite, and hence M has at
 494 least one $\leq_{\mathcal{J}}$ -minimal element.

495 We first show that S contains a $\leq_{\mathcal{J}}$ -minimal element t . This is rather clear because,
 496 for any non- $\leq_{\mathcal{J}}$ -minimal element s , its language $\eta^{-1}(s) \subseteq A^*$ is null: s is non- $\leq_{\mathcal{J}}$ -minimal
 497 means that there is an other element t such that $t <_{\mathcal{J}} s$ (*i.e.*, $MtM \subsetneq MsM$), whence
 498 $s \notin MtM$ which implies that any word $w \in \eta^{-1}(t)$ is a forbidden word of $\eta^{-1}(s)$. Thus by
 499 the infinite monkey theorem $\eta^{-1}(s)$ is null.

500 Clearly, we have $t^n \leq_{\mathcal{J}} t$ and thus $t \mathcal{J} t^n$ holds for any $n > 1$ by the $\leq_{\mathcal{J}}$ -minimality of t .
 501 $t \mathcal{J} t^n$ implies that there is a pair of words x, y such that $xt^n y = t$. Since M is finite, x^m is
 502 an idempotent for some $m > 0$ (*i.e.*, $x^{2m} = x^m$). Thus we obtain $t = xt^n y = x(t)t^{n-1}y =$
 503 $x^2(t)(t^{n-1}y)^2 = \dots = x^m t (t^{n-1}y)^m = x^m x^m t (t^{n-1}y)^m = x^m t$ whence $t = t^n (y(t^{n-1}y)^{m-1})$.
 504 It follows that $t \mathcal{R} t^n$. Dually, we also obtain $t \mathcal{L} t^n$ and hence we can deduce that $t \mathcal{H} t^n$ holds.
 505 By the finiteness of M , there exists some $n > 0$ such that t^n is an idempotent. Thanks to
 506 Green’s theorem, the \mathcal{H} -equivalent class H_t of t is a subgroup of M with the identity element
 507 t^n . Because η is surjective, we can take a word w' from $\eta^{-1}(t)$. Let $t' = \eta(w'a) = t\eta(a)$ for
 508 some letter $a \in A$, then by the $\leq_{\mathcal{J}}$ -minimality of t , we can take some words $x, y \in A^*$ so that
 509 $\eta(xw'ay) = \eta(x)t'\eta(y) = t$. Hence we can deduce that $\eta^{-1}(t)$ contains a non-empty word
 510 $w = xw'ay$. Then for any $\varepsilon \neq w \in \eta^{-1}(t)$ and $m \geq 1$, we have

$$511 \quad \eta(w^{mn+1}) = t^{mn+1} = (t^n)^m \cdot t = t \in S$$

512 which means that $L \supseteq \eta^{-1}(t)$ contains infinitely many non-primitive words w^{mn+1} . ◀

513 ► **Corollary 39** (of Theorem 37 and 38). Q_A has full REG-gap.

514 ► **Remark 40.** We emphasise that the assumption “ L is non-null” in Theorem 38 is quite tight,
 515 since a slightly weaker assumption “ L is of exponential growth” (*i.e.*, $\#(L \cap A^n)$ is exponential
 516 for n) does not imply that L contains non-primitive words. A trivial counterexample is
 517 $L_0 = \{a, b\}^* c$ over $A = \{a, b, c\}$: $\#(L_0 \cap A^n) = 2^{n-1}$ ($n \geq 1$) is exponential but L_0 only
 518 consists of primitive words. L_0 has a cc as a forbidden word, hence it is null by the infinite
 519 monkey theorem. Thus L_0 is not a counterexample of Theorem 38.

520 6 Conclusion and Open Problems

521 In this paper we proposed REG-measurability and showed that several context-free languages
 522 are REG-measurable, excluding M_m . Interestingly, it is shown that, like G and K , languages
 523 that have been considered as complex from a combinatorial viewpoint are, actually, easy
 524 to asymptotically approximate by regular languages. It is also interesting that a modified
 525 majority language M_2 is just a deterministic context-free but it is complex from a measure

526 theoretic viewpoint. Its complement $\overline{M_2}$ is also deterministic context-free, and actually it is
 527 co-null but REG^+ -immune (*i.e.*, has full REG -gap). This means that $\overline{M_2}$ is as complex as
 528 Q_A from a viewpoint of REG -measurability.

529 The following fundamental problems are still open and we consider these to be future
 530 work.

531 ► **Problem 41.** *Can we give an alternative characterisation of the null (resp. co-null)*
 532 *context-free languages (like Theorem 10)?*

533 ► **Problem 42.** *Can we give an alternative characterisation of the REG -measurable context-*
 534 *free languages?*

535 ► **Problem 43.** *Can we find a language class that can “separate” Q_A and CFL? *i.e.*, is there*
 536 *\mathcal{C} such that Q_A has full \mathcal{C} -gap but no co-null context-free language has full \mathcal{C} -gap, or Q_A is*
 537 *\mathcal{C} -immeasurable but any co-null context-free language is \mathcal{C} -measurable?*

538 The our results (Theorem 36, 37 and 38) tell us that the class REG of regular languages can
 539 not separate Q_A and CFL. However, it is still open whether the situation is the same or not
 540 when $\mathcal{C} = \text{DetCFL}, \text{UnCFL}, \text{CFL}$ or other extension of regular languages. Notice that *if* the
 541 answer of Problem 43 is “yes”, then Q_A is not context-free.

542

543 **Acknowledgement:** The author would like to thank Takanori Maehara (RIKEN AIP) and
 544 Fazekas Szilárd (Akita University) whose helpful discussion were an enormous help to me.
 545 The author also thank to anonymous reviewers for many valuable comments. This work was
 546 supported by JSPS KAKENHI Grant Number JP19K14582.

547 ——— References ———

- 548 1 Jean-Michel Autebert, Philippe Flajolet, and Joaquim Gabarro. Prefixes of infinite words and
 549 ambiguous context-free languages. *Information Processing Letters*, 25(4):211–216, 1987.
- 550 2 Jean Berstel. Sur la densité asymptotique de langages formels. In *International Colloquium*
 551 *on Automata, Languages and Programming*, pages 345–358, France, 1973. North-Holland.
- 552 3 Jean Berstel, Dominique Perrin, and Christophe Reutenauer. *Codes and Automata*. Encyclo-
 553 *pedia of Mathematics and its Applications*. Cambridge University Press, 2009.
- 554 4 Robert C. Buck. The measure theoretic approach to density. *American Journal of Mathematics*,
 555 68(4):560–580, 1946.
- 556 5 Cezar Câmpeanu, Nicolae Sântean, and Sheng Yu. Minimal cover-automata for finite languages.
 557 *Theoretical Computer Science*, 267(1):3–16, 2001.
- 558 6 N. Chomsky and M.P. Schützenberger. The algebraic theory of context-free languages*. In
 559 *Computer Programming and Formal Systems*, volume 35, pages 118–161. Elsevier, 1963.
- 560 7 Brendan Cordy and Kai Salomaa. On the existence of regular approximations. *Theoretical*
 561 *Computer Science*, 387(2):125–135, 2007.
- 562 8 Michael Domaratzki. Minimal covers of formal languages. Master’s thesis, University of
 563 Waterloo, 2001.
- 564 9 Pál Dömösi, Sándor Horváth, and Masami Ito. On the connection between formal languages
 565 and primitive words. pages 59–67, 1991.
- 566 10 Gerry Eisman and Bala Ravikumar. Approximate recognition of non-regular languages by
 567 finite automata. In *Twenty-Eighth Australasian Computer Science Conference (ACSC2005)*,
 568 volume 38 of *CRPIT*, pages 219–228, Newcastle, Australia, 2005. ACS.
- 569 11 Gerry Eisman and Bala Ravikumar. On approximating non-regular languages by regular
 570 languages. *Fundamenta Informaticae*, 110:125–142, 2011.
- 571 12 P. Flajolet and J. M. Steyaert. On sets having only hard subsets. In *International Colloquium*
 572 *on Automata, Languages and Programming*, pages 446–457. North-Holland, 1974.

23:16 Asymptotic Approximation by Regular Languages

- 573 13 Philippe Flajolet. Ambiguity and transcendence. In *Automata, Languages and Programming*,
574 pages 179–188, Berlin, Heidelberg, 1985. Springer Berlin Heidelberg.
- 575 14 Philippe Flajolet. Analytic models and ambiguity of context-free languages. *Theoretical*
576 *Computer Science*, 49(2):283–309, 1987.
- 577 15 Martin Kappes and Chandra M. R. Kintala. Tradeoffs between reliability and conciseness
578 of deterministic finite automata. *Journal of Automata, Languages and Combinatorics*, 9(2–
579 3):281–292, 2004.
- 580 16 Martin Kappes and Frank Nießner. Succinct representations of languages by dfa with different
581 levels of reliability. *Theoretical Computer Science*, 330(2):299–310, 2005.
- 582 17 Rainer Kemp. A note on the density of inherently ambiguous context-free languages. *Acta*
583 *Informatica*, 14(3):295–298, 1980.
- 584 18 Piet van Mieghem. *Graph Spectra for Complex Networks*. Cambridge University Press, 2010.
- 585 19 Jean-Éric Pin. *Mathematical foundations of automata theory*, 2012.
- 586 20 C.M. Reis and H.J. Shyr. Some properties of disjunctive languages on a free monoid. *Inform-*
587 *ation and Control*, 37(3):334–344, 1978.
- 588 21 Arto Salomaa and Matti Soittola. *Automata Theoretic Aspects of Formal Power Series*.
589 Springer-Verlag New York, Inc., Secaucus, NJ, USA, 1978.
- 590 22 Ryoma Sin’ya. Asymptotic approximation by regular languages (full version). URL: <http://www.math.akita-u.ac.jp/~ryoma/misc/measure.pdf>.
591
- 592 23 Ryoma Sin’ya. An automata theoretic approach to the zero-one law for regular languages:
593 Algorithmic and logical aspects. In *Proceedings Sixth International Symposium on Games,*
594 *Automata, Logics and Formal Verification, GandALF 2015*, pages 172–185, 2015.