

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.*, [14, 15, 10, 7, 5, 8]). For example, Kappes and Kintala [14] 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 $0 < \epsilon < 1$ [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* [5, 8].



© 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



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

23:2 Asymptotic Approximation by Regular Languages

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

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

$$57 \quad \mu^*(S) = \inf \left\{ \sum_i \delta(X_i) \mid S \subseteq X \text{ and } X \text{ is a finite union of} \right.$$

$$58 \quad \left. \text{disjoint arithmetic progressions } X_1, \dots, X_k \right\}.$$

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

$$61 \quad \mu^*(S) + \mu^*(\bar{S}) = 1 \tag{1}$$

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

$$65 \quad \mathcal{D}_0 \stackrel{\text{def}}{=} \{X \subseteq \mathbb{N} \mid X \text{ is a finite union of arithmetic progressions } \},$$

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

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

78 Our contribution

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

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

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

99 2 Densities of Formal Languages

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

108 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
 109 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 .
 110 We denote by REG, DetCFL, UnCFL and CFL the class of regular languages, deterministic
 111 context-free languages, unambiguous context-free languages and context-free languages,
 112 respectively. A language L is said to be \mathcal{C} -*immune* if L is infinite and no infinite subset of L
 113 belongs to \mathcal{C} .

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

$$115 \quad \delta_A(L) \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \frac{\#(L \cap A^n)}{\#(A^n)}$$

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

$$118 \quad \delta_A^*(L) \stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^{n-1} \frac{\#(L \cap A^k)}{\#(A^k)}$$

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

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

124 The following observation is basic.

23:4 Asymptotic Approximation by Regular Languages

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

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

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

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

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

$$138 \quad \frac{\#((AA)^* \cap A^n)}{\#(A^n)} = \begin{cases} 1 & \text{if } n \text{ is even,} \\ 0 & \text{if } n \text{ is odd.} \end{cases}$$

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

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

- 141 ■ The semi-Dyck language

$$142 \quad D \stackrel{\text{def}}{=} \{w \in \{a, b\}^* \mid |w|_a = |w|_b \text{ and } |u|_a \geq |u|_b \text{ for every prefix } u \text{ of } w\}$$

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

144 $2n$ is equal to the n -th Catalan number whose asymptotic approximation is $\Theta(4^n/n^{3/2})$.

145 Thus

$$146 \quad \frac{\#(D \cap A^n)}{\#(A^n)} = \begin{cases} \Theta(1/(n/2)^{3/2}) & \text{if } n \text{ is even,} \\ 0 & \text{if } n \text{ is odd.} \end{cases}$$

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

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

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

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

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

152 ► **Theorem 5** (*cf.* Theorem III.6.1 of [20]). *Let $L \subseteq A^*$ be a regular language. Then there is*

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

$$154 \quad \lim_{n \rightarrow \infty} \frac{\#(L \cap A^{cn+d})}{\#(A^{cn+d})}$$

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

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

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

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

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

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

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

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

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

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

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

171 ► **Theorem 10** (S. [21]). *A regular language is non-null if and only if it is dense.*

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

178 3 Approximability and Measurability

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

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

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

189 ► **Proposition 12.** *The semi-Dyck language D is REG-inapprox.*

190 **Proof.** We already mentioned that D is null in Example 4, and thus D is $(\text{REG}, 0)$ -lower-
 191 approx by $\emptyset \subseteq D$. One can easily observe that D has no forbidden word: since for any
 192 $w \in A^*$ 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
 193 regular language L satisfies $D \subseteq L$, L has no forbidden word, too, and thus L is non-null by
 194 Theorem 10. Thus by Claim 3, $\delta_A^*(L \setminus D) = \delta_A^*(L) - \delta_A^*(D) = \delta_A^*(L) > 0$, which means that
 195 D can not $(\text{REG}, 0)$ -upper-approx. ◀

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

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

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

200 ► **Corollary 14.** *DetCFL is REG-inapprox.*

23:6 Asymptotic Approximation by Regular Languages

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

204 ► **Theorem 15** (Kemp [16]). *Let $A = \{a, b, c\}$. Define*

$$205 \quad S_1 \stackrel{\text{def}}{=} \{a\}\{b^i a^i \mid i \geq 1\}^* \quad S_2 \stackrel{\text{def}}{=} \{a^i b^{2i} \mid i \geq 1\}^* \{a\}^+,$$

206 *and*

$$207 \quad L_1 \stackrel{\text{def}}{=} S_1 \{c\} A^* \quad L_2 \stackrel{\text{def}}{=} S_2 \{c\} A^*.$$

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

209 ► **Corollary 16.** *CFL is UnCFL-inapprox.*

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

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

$$214 \quad \underline{\mu}_{\mathcal{C}}(L) \stackrel{\text{def}}{=} \sup\{\delta_A^*(K) \mid A = \text{Alph}(L), K \subseteq L, K \in \mathcal{C}_A, \delta_A^*(K) \neq \perp\}$$

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

$$216 \quad \overline{\mu}_{\mathcal{C}}(L) \stackrel{\text{def}}{=} \inf\{\delta_A^*(K) \mid A = \text{Alph}(L), L \subseteq K, K \in \mathcal{C}_A, \delta_A^*(K) \neq \perp\}.$$

217 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 $\overline{\mu}_{\mathcal{C}}(L)$
 218 as $\mu_{\mathcal{C}}(L)$. \mathcal{D} is said to be \mathcal{C} -measurable if every language in \mathcal{D} is \mathcal{C} -measurable.

219 ► **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
 220 L has full \mathcal{C} -gap if its \mathcal{C} -gap equals to 1, i.e., $\overline{\mu}_{\mathcal{C}}(L) - \underline{\mu}_{\mathcal{C}}(L) = 1$.

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

224 The following lemmata are basic.

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

- 226 1. $\overline{\mu}_{\mathcal{C}}(K) \leq \overline{\mu}_{\mathcal{C}}(L)$ if $K \subseteq L$.
- 227 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.
- 228 3. $\overline{\mu}_{\mathcal{C}}(K) = \delta_A^*(K)$ if $K \in \mathcal{C}$ and $\delta_A^*(K) \neq \perp$.

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

$$231 \quad \overline{\mu}_{\mathcal{C}}(L) + \overline{\mu}_{\mathcal{C}}(\overline{L}) = 1. \tag{2}$$

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

$$235 \quad \inf\{\delta_A^*(K) \mid L \subseteq K, K \in \mathcal{C}\} = 1 - \inf\{\delta_A^*(K) \mid \overline{L} \subseteq K, K \in \mathcal{C}\} \tag{3}$$

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

$$238 \quad \inf\{\delta_A^*(K) \mid L \subseteq K, K \in \mathcal{C}\} = \sup\{\delta_A^*(K) \mid K \subseteq L, K \in \mathcal{C}\}. \tag{4}$$

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

244 4 REG-measurability on Context-free Languages

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

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

250 **Proof.** Let A be an alphabet with $\#(A) \geq 2$ and let $(\mathcal{A}_i)_{i \in \mathbb{N}}$ be an enumeration of automata
 251 over A such that $\text{REG}_A = \{L(\mathcal{A}_i) \mid i \in \mathbb{N}\}$; we can take such enumeration by enumerating
 252 some binary representation of automata via shortlex order $<_{\text{lex}}$. We will construct a null
 253 language L such that $\bar{\mu}_{\text{REG}}(L) = 1$, in particular, L intersects with every regular infinite
 254 language.

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

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

257 **Step 2** check $L(\mathcal{A}_i)$ is infinite or not.

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

259 **Step 4** otherwise, pick u such that u is the smallest (with respect to $<_{\text{lex}}$) word satisfying
 260 $|u| > \ell$ and $u \in L$ (such u surely exists since $L(\mathcal{A}_i)$ is infinite).

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

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

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

264 One can easily observe that all Steps are effective and P ultimately halts for any input
 265 word w because the length of the word u in Step 4 is strictly increasing until $u = w$ or
 266 $w <_{\text{lex}} u$. Thus the language $L \stackrel{\text{def}}{=} \{w \in A^* \mid P \text{ accepts } w\}$ is recursive, (1) $L \cap R \neq \emptyset$ for
 267 any regular infinite language because by Step (4–5) P accepts some word $w \in R$, and (2)
 268 $\delta_A(L) = 0$; by Step (5–6) and the length of u is strictly increasing, P rejects every word in
 269 A^n except for one single word u , for each n . Thus $\delta_A(L) = 0$ and $\bar{\mu}_{\text{REG}}(L) = 1$. ◀

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

274 4.1 Null Context-free Languages

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

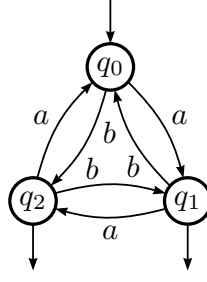
277 ► **Definition 22.** For an alphabet A and letters $a, b \in A$ such that $a \neq b$, we define

$$278 \quad L_A(a, b) \stackrel{\text{def}}{=} \{w \in A^* \mid |w|_a = |w|_b\}.$$

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

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

$$282 \quad L_k \stackrel{\text{def}}{=} \{w \in A^* \mid |w|_a \neq |w|_b \pmod k\}.$$



■ **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 [17]). Hence, in the case of $M_k = p_{\mathcal{A}_k}(E_k) = E_k + E_k^{k-1}$, we have

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

303 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.

304 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
 305 the final states are all but q_0 , the number of words of length n in L_k is exactly the number
 306 of paths from q_0 to any other state in \mathcal{A}_k . For the k -dimensional vectors $\mathbf{e} = (1, 0, 0, \dots, 0)$
 307 and $\mathbf{1} = (1, 1, 1, \dots, 1)$, from Equation (5) we have

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

313 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
 314 absolute value is strictly smaller than 2; because ξ_k^{-j} is the complex conjugate of ξ_k^j and
 315 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

$$316 \quad \#(L_k \cap A^n) = \frac{k-1}{k} 2^n + o(2^n)$$

317 where $o(2^n)$ means some function such that $\lim_{n \rightarrow \infty} o(2^n)/2^n = 0$. Thus we have $\delta_A(L_k) =$
 318 $\frac{k-1}{k}$ for odd $k = 2m + 1$, which tends to 1 if k tends infinity, *i.e.*, $\mu_{\text{REG}}(L) = 1$. This
 319 completes the proof. ◀

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

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

323 The next example is the set of all palindromes.

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

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

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

$$328 \quad L_k \stackrel{\text{def}}{=} \{w_1 A^* w_2 \mid w_1, w_2 \in A^k, w_1 \neq \text{rev}(w_2)\}.$$

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

$$331 \quad \#(L_k \cap A^n) = \#(A)^k \cdot \#(A)^{n-2k} \cdot (\#(A)^k - 1) = \#(A)^n - \#(A)^{n-k}.$$

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

334 **4.2 Some Inherently Ambiguous Languages**

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

338 Some more complex examples of inherently ambiguous languages are the following
 339 languages with constraints on the number of occurrences of letters investigated by Flajolet [12]:

$$340 \quad \mathcal{O}_3 \stackrel{\text{def}}{=} \{w \in \{a, b, c\}^* \mid |w|_a = |w|_b \text{ or } |w|_a = |w|_c\},$$

$$341 \quad \mathcal{O}_4 \stackrel{\text{def}}{=} \{w \in \{x, \bar{x}, y, \bar{y}\}^* \mid |w|_x = |w|_{\bar{x}} \text{ or } |w|_y = |w|_{\bar{y}}\}.$$

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

344 **Proof.** Let $A = \{a, b, c\}$. For the case \mathcal{O}_3 , in a very similar way to Theorem 23, we
 345 can construct a sequence of automata $(\mathcal{A}_k^{ab})_{k \in \mathbb{N}}$ such that each automaton \mathcal{A}_k^{ab} satisfies
 346 $L(\mathcal{A}_k^{ab}) \subseteq \overline{L_A(a, b)}$ and its adjacency matrix is of the form

$$347 \quad M_k^{ab} = M_k + I_k = \begin{bmatrix} 1 & 1 & 0 & \cdots & \cdots & 1 \\ 1 & 1 & 1 & \ddots & & \vdots \\ 0 & 1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 1 & 0 \\ \vdots & & \ddots & 1 & 1 & 1 \\ 1 & \cdots & \cdots & 0 & 1 & 1 \end{bmatrix}$$

348 where M_k is the adjacency matrix stated in Theorem 23 and I_k is the k -dimensional identity
 349 matrix. The automaton \mathcal{A}_k^{ab} is obtained by just adding self-loop labeled by c for each state
 350 $q \in Q_k$ of \mathcal{A}_k in Theorem 23. This sequence of automata ensures that the language $L_A(a, b)$
 351 is REG-measurable ($\bar{\mu}_{\text{REG}}(L_A(a, b)) = 0$, in particular). The same argument is applicable to
 352 the language $L_A(a, c)$, thus these union $\mathcal{O}_3 = L_A(a, b) \cup L_A(a, c)$ is also REG-measurable by
 353 Lemma 19. The case \mathcal{O}_4 can be archived in the same manner. ◀

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

$$356 \quad \mathbf{G} \stackrel{\text{def}}{=} \{a^{n_1} b a^{n_2} b \cdots a^{n_p} b \mid p \geq 1, n_i \neq i \text{ for some } i\}.$$

357 While \mathbf{G} can be accepted by a non-deterministic pushdown automaton, its generating function
 358 is not algebraic [13] and thus it is an inherently ambiguous context-free language due to the
 359 well-known Chomsky–Schützenberger theorem stating that the generating function of every
 360 unambiguous context-free language is algebraic [6].

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

362 **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

$$363 \quad L_{\mathbf{G}} = \{u \in A^* \mid uA^*\{b\} \cap \overline{\mathbf{G}} = \emptyset\}$$

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

$$365 \quad L_k \stackrel{\text{def}}{=} \{uA^*\{b\} \mid u \in L_{\mathbf{G}} \cap A^k\}.$$

366 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
 367 $L_k \subseteq \mathbf{G}$ holds for each k . Any word in $\overline{L_{\mathbf{G}}} = A^* \setminus L_{\mathbf{G}}$ is a prefix of the infinite word

368 $a^{n_1}ba^{n_2}ba^{n_3}b \dots$ ($n_i = i$ for each $i \in \mathbb{N}$) thus $\#(L_G \cap A^n) = \#(A^n) - 1$ holds for each $n \geq 1$.
 369 Hence we have

$$\begin{aligned} 370 \quad \delta_A(L_k) &= \lim_{n \rightarrow \infty} \frac{\#(L_k \cap A^n)}{\#(A^n)} = \lim_{n \rightarrow \infty} \frac{(\#(A^k) - 1) \cdot \#(A^{n-k-1})}{\#(A^n)} \\ 371 \quad &= (\#(A)^k - 1) \cdot \#(A)^{-k-1} = 2^{-1} - 2^{-k-1}. \end{aligned}$$

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

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

$$375 \quad \text{Copref}(w) \stackrel{\text{def}}{=} A^* \setminus \{u \in A^* \mid u \text{ is a prefix of } w\}$$

376 is called the *coprefix language of w* . The proof of Theorem 27 uses a key property that \mathbf{G} can
 377 be characterised by using the coprefix language of the infinite word $w = a^{n_1}ba^{n_2}ba^{n_3}b \dots$ as
 378 $\mathbf{G} = \text{Copref}(w) \cap \{a, b\}^* \{b\}$ which was pointed out in [1]. Thus by the same argument, we
 379 can say that any coprefix language L is REG-measurable ($\mu_{\text{REG}}(L) = 1$, in particular).

380 For coprefix languages, the following nice “gap theorem” holds.

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

- 384 1. L is a regular language.
- 385 2. L is an inherently ambiguous context-free language.

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

388 4.3 K: A Language with Transcendental Density

389 We now show the fact that the language \mathbf{K} defined by Kemp [16] (recall that the definition of
 390 \mathbf{K} appeared in Theorem 15) is REG-measurable. We will actually show a more general result
 391 regarding the following type of languages.

392 **► Definition 29.** Let $L \subseteq A^*$ be a language and $c \notin A$ be a letter. We call the language
 393 $L\{c\}(A \cup \{c\})^*$ over $A \cup \{c\}$ *suffix extension of L by c* .

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

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

$$398 \quad \#(u\{c\}B^* \cap B^n) / \#(B^n) = \#(u\{c\}B^{n-|u|-1}) / \#(B^n) = k^{n-|u|-1} / k^n = k^{-(|u|+1)}$$

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

$$\begin{aligned} 400 \quad \delta_B(L') &= \lim_{n \rightarrow \infty} \frac{\#(L' \cap B^n)}{\#(B^n)} = \lim_{n \rightarrow \infty} \frac{\#(\bigcup_{w \in L} (w\{c\}B^* \cap B^n))}{\#(B^n)} \\ 401 \quad &= \lim_{n \rightarrow \infty} \frac{\sum_{w \in L} \#(w\{c\}B^* \cap B^n)}{\#(B^n)} = \lim_{n \rightarrow \infty} \sum_{w \in (L \cap A^{<n})} k^{-(|w|+1)}. \end{aligned} \quad (7)$$

403 Because the sequence $(\sum_{w \in (L \cap A^{<n})} k^{-(|w|+1)})_{n \in \mathbb{N}}$ is non-decreasing and bounded above by
 404 1, the limit (7) exists, say $\delta_B(L') = \alpha$.

23:12 Asymptotic Approximation by Regular Languages

405 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
 406 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
 407 argument, for each $n \in \mathbb{N}$, we can claim that the language $K_n \stackrel{\text{def}}{=} A^* \setminus \bigcup_{w \in \bar{L} \cap A^{<n}} w\{c\}B^*$
 408 satisfies $K_n \supseteq L'$ and $\delta_B(K_n)$ tends to α if n tends to infinity. Thus $\mu_{\text{REG}}(L') = \alpha$. ◀

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

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

411 ▶ **Remark 32.** Theorem 30 indicates that REG-measurability is a quite relaxed property
 412 in some sense: even for a non-recursively-enumerable language, its suffix extension is still
 413 non-recursively-enumerable but REG-measurable.

414 The same proof method works for the *prefix extension*, and the *infix extension* is also
 415 REG-measurable.

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

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

4.4 Languages with Full REG-Gap

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

$$428 \quad \mathbb{M} \stackrel{\text{def}}{=} \{w \in \{a, b\}^* \mid |w|_a > |w|_b\},$$

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

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

$$432 \quad \limsup_{n \rightarrow \infty} \{\#(\bar{L} \cap A^n) / \#(A^n)\} = 0.$$

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

436 ▶ **Corollary 35.** \mathbb{M} has full REG-gap.

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

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

$$440 \quad \mathbb{M}_m \stackrel{\text{def}}{=} \{w \in A^* \mid |w|_a > m \cdot |w|_b\}$$

441 has full REG-gap, and $\delta_A(\mathbb{M}_m) = 1/2$ if $m = 1$ otherwise $\delta_A(\mathbb{M}_m) = 0$.

442 **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, 443
444 (this is well-defined natural number since M is finite). By the infinite monkey theorem, L is not null implies that L has no forbidden word, and thus for the word b^{2c} there exist 445
446 two words x and y such that $xb^{2c}y$ is in L . We can assume that $|x|, |y| \leq c$ without loss of 447
448 generality by the definition of c , which implies $|xb^{2c}y|_a \leq |x| + |y| = 2c \leq |xb^{2c}y|_b$ hence 449
450 $xb^{2c}y \notin M_m$. Thus $L \not\subseteq M_m$ and $\mu_{\text{REG}}(M_m) = 0$. By using same argument, we can prove 451
452 that $\bar{\mu}_{\text{REG}}(M_m) = 0$ and hence M_m has full REG-gap.

450 In the case $m = 1$, $\delta_A(M_1) = \delta_A(M) = 1/2$ is obvious. It is enough to show that 451
452 $\delta_A(M_2) = 0$ holds (since $M_m \subseteq M_2$ for any $m \leq 2$). Indeed, we have

$$452 \quad \delta_A(M_2) = \lim_{n \rightarrow \infty} \frac{\#\{w \in A^n \mid |w|_a > 2|w|_b\}}{2^n} = \lim_{n \rightarrow \infty} \frac{\#\{w \in A^n \mid |w|_a > 2n/3\}}{2^n}$$

$$453 \quad = \lim_{n \rightarrow \infty} \Pr(|\bar{X}_n - n/2| > n/6) = 0$$

455 where $\Pr(|\bar{X}_n - n/2| > n/6)$ means the probability that the absolute value of the difference 456
457 of the number \bar{X}_n of the occurrences of a 's in a randomly chosen word of length n and its 458
459 mean value $n/2$ is larger than $n/6$; its tends to zero by the weak law of large numbers. ◀

458 5 REG-Immesurability of Primitive Words

459 A non-empty word $w \in A^+$ is said to be primitive if $u^n = w$ implies $u = w$ for any $u \in A^+$ 460
461 and $n \in \mathbb{N}$. The set of all primitive words over A is denoted by Q_A . Because the case 462
463 $\#(A) = 1$ is meaningless ($Q_A = A$ in this case), hereafter we always assume $\#(A) \geq 2$. 464
465 Whether Q_A is context-free or not is a well-known long-standing open problem posed by 466
467 Dömösi, Horváth and Ito [9]. Reis and Shyr [19] proved $Q_A^2 = A^+ \setminus \{a^n \mid a \in A, n \neq 2\}$, 468
469 which intuitively means that every non-empty word w not a power of a letter is a product of 470
471 two primitive words. From this result one may think that Q_A is “very large” in some sense. 472
473 Actually, Q_A is somewhat “large” (it is dense in the sense of Definition 9), but we can show 474
475 more stronger property as follows.

468 ▶ **Theorem 37.** $\delta_A(Q_A) = 1$.

469 **Proof.** It is enough to show that $\delta_A(\overline{Q_A}) = 0$ holds. One can easily observe that any natural 470
471 number $n \in \mathbb{N}$ has at most $2\sqrt{n}$ divisors. In addition, for any non-primitive word $w = v^m$ of 472
473 length n is uniquely determined by v (since $m = n/|v|$) and $|v| \leq n/2$. Hence the number of 474
475 non-primitive words of length n satisfies

$$473 \quad \#\overline{(Q_A \cap A^n)} \leq 2\sqrt{n} \sum_{i=0}^{\lfloor n/2 \rfloor} \#(A^i) \leq 2\sqrt{n} \cdot \#(A)^{\lfloor n/2 \rfloor + 1}.$$

474 By using the above estimation, we can deduce that

$$475 \quad \frac{\#\overline{(Q_A \cap A^n)}}{\#(A^n)} \leq \frac{2\sqrt{n} \cdot \#(A)^{\lfloor n/2 \rfloor + 1}}{\#(A)^n} \leq \frac{2\sqrt{n}}{\#(A)^{n/2 - 1}}$$

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

477 While Q_A is “very large” (co-null) as stated above, we can also prove that Q_A is REG⁺- 478
479 immune. The proof relies on an analysis of the structure of the syntactic monoid of a non-null 480
481 regular language. We assume that the reader has a basic knowledge of semigroup theory

23:14 Asymptotic Approximation by Regular Languages

(cf. [18]): Green's relations $\mathcal{J}, \mathcal{R}, \mathcal{L}, \mathcal{H}$ and a direct consequence of Green's theorem (an \mathcal{H} -class H in a semigroup S is a subgroup of S if and only if H contains an idempotent), in particular.

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

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

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

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 . $t \mathcal{J} t^n$ implies that there is a pair of words x, y such that $xt^n y = t$. Since M is finite, x^m is an idempotent for some $m > 0$ (i.e., $x^{2m} = x^m$). Thus we obtain $t = xt^n y = x(t)t^{n-1}y = x^2(t)(t^{n-1}y)^2 = \dots = x^m t (t^{n-1}y)^m = x^m x^m t (t^{n-1}y)^m = x^m t$ whence $t = t^n (y(t^{n-1}y)^{m-1})$. 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. By the finiteness of M , there exists some $n > 0$ such that t^n is an idempotent. Thanks to Green's theorem, the \mathcal{H} -equivalent class H_t of t is a subgroup of M with the identity element 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 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 $\eta(xw'ay) = \eta(x)t'\eta(y) = t$. Hence we can deduce that $\eta^{-1}(t)$ contains a non-empty word $w = xw'ay$. Then for any $\varepsilon \neq w \in \eta^{-1}(t)$ and $m \geq 1$, we have

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

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

► **Corollary 39** (of Theorem 37 and 38). \mathbb{Q}_A has full REG-gap.

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

6 Conclusion and Open Problems

In this paper we proposed REG-measurability and showed that several context-free languages are REG-measurable, excluding M_m . Interestingly, it is shown that, like G and K , languages that have been considered as complex from a combinatorial viewpoint are, actually, easy to asymptotically approximate by regular languages. It is also interesting that a modified majority language M_2 is just a deterministic context-free but it is complex from a measure theoretic viewpoint. Its complement $\overline{M_2}$ is also deterministic context-free, and actually it is co-null but REG^+ -immune (i.e., has full REG-gap). This means that $\overline{M_2}$ is as complex as \mathbb{Q}_A from a viewpoint of REG-measurability.

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

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

527 ▶ **Problem 42.** *Can we give an alternative characterisation of the REG-measurable context-*
528 *free languages?*

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

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

536 **Acknowledgement:** The author would like to thank Takanori Maehara (RIKEN AIP)
537 whose helpful discussion were an enormous help to me. The author also thank to anonymous
538 reviewers for many valuable comments. This work was supported by JSPS KAKENHI Grant
539 Number JP19K14582.

540 ——— References ———

- 541 1 Jean-Michel Autebert, Philippe Flajolet, and Joaquim Gabarro. Prefixes of infinite words and
542 ambiguous context-free languages. *Information Processing Letters*, 25(4):211 – 216, 1987.
- 543 2 Jean Berstel. Sur la densité asymptotique de langages formels. In *International Colloquium*
544 *on Automata, Languages and Programming (ICALP, 1972)*, pages 345–358, France, 1973.
545 North-Holland.
- 546 3 Jean Berstel, Dominique Perrin, and Christophe Reutenauer. *Codes and Automata*. Encyclo-
547 *pedia of Mathematics and its Applications*. Cambridge University Press, 2009.
- 548 4 Robert C. Buck. The measure theoretic approach to density. *American Journal of Mathematics*,
549 68(4):560–580, 1946.
- 550 5 Cezar Câmpeanu, Nicolae Sântean, and Sheng Yu. Minimal cover-automata for finite languages.
551 *Theoretical Computer Science*, 267(1):3 – 16, 2001.
- 552 6 N. Chomsky and M.P. Schützenberger. The algebraic theory of context-free languages*. In
553 *Computer Programming and Formal Systems*, volume 35, pages 118 – 161. Elsevier, 1963.
- 554 7 Brendan Cordy and Kai Salomaa. On the existence of regular approximations. *Theoretical*
555 *Computer Science*, 387(2):125 – 135, 2007.
- 556 8 Michael Domaratzki. Minimal covers of formal languages. Master’s thesis, University of
557 Waterloo, 2001.
- 558 9 Pál Dömösi, Sándor Horváth, and Masami Ito. On the connection between formal languages
559 and primitive words. pages 59–67, 1991.
- 560 10 Gerry Eisman and Bala Ravikumar. Approximate recognition of non-regular languages by
561 finite automata. In *Twenty-Eighth Australasian Computer Science Conference (ACSC2005)*,
562 volume 38 of *CRPIT*, Newcastle, Australia, 2005. ACS.
- 563 11 Gerry Eisman and Bala Ravikumar. On approximating non-regular languages by regular
564 languages. *Fundamenta Informaticae*, 110:125–142, 2011.
- 565 12 Philippe Flajolet. Ambiguity and transcendence. In *Automata, Languages and Programming*,
566 pages 179–188, Berlin, Heidelberg, 1985. Springer Berlin Heidelberg.
- 567 13 Philippe Flajolet. Analytic models and ambiguity of context-free languages. *Theoretical*
568 *Computer Science*, 49(2):283 – 309, 1987.

23:16 Asymptotic Approximation by Regular Languages

- 569 **14** Martin Kappes and Chandra M. R. Kintala. Tradeoffs between reliability and conciseness
570 of deterministic finite automata. *Journal of Automata, Languages and Combinatorics*, 9(2–
571 3):281–292, 2004.
- 572 **15** Martin Kappes and Frank Nießner. Succinct representations of languages by dfa with different
573 levels of reliability. *Theoretical Computer Science*, 330(2):299 – 310, 2005.
- 574 **16** Rainer Kemp. A note on the density of inherently ambiguous context-free languages. *Acta*
575 *Informatica*, 14(3):295–298, 1980.
- 576 **17** Piet van Mieghem. *Graph Spectra for Complex Networks*. Cambridge University Press, 2010.
- 577 **18** Jean-Éric Pin. *Mathematical foundations of automata theory*, 2012.
- 578 **19** C.M. Reis and H.J. Shyr. Some properties of disjunctive languages on a free monoid. *Inform-*
579 *ation and Control*, 37(3):334 – 344, 1978.
- 580 **20** Arto Salomaa and Matti Soittola. *Automata Theoretic Aspects of Formal Power Series*.
581 Springer-Verlag New York, Inc., Secaucus, NJ, USA, 1978.
- 582 **21** Ryoma Sin'ya. An automata theoretic approach to the zero-one law for regular languages:
583 Algorithmic and logical aspects. In *Proceedings Sixth International Symposium on Games,*
584 *Automata, Logics and Formal Verification, GandALF 2015*, pages 172–185, 2015.