


1 Revisiting the growth of polyregular functions

2 Output languages, weighted automata and unary inputs

3 Sandra Kiefer ✉

4 University of Oxford, United Kingdom

5 Lê Thành Dũng (Tito) Nguyễn ✉ 🏠 

6 École normale supérieure de Lyon, France

7 Cécilia Pradic 🏠

8 Swansea University, United Kingdom

9 — Abstract —

10 Polyregular functions are the class of string-to-string functions definable by pebble transducers (an
11 extension of finite automata) or equivalently by MSO interpretations (a logical formalism). Their
12 output length is bounded by a polynomial in the input length: a function computed by a k -pebble
13 transducer or by a k -dimensional MSO interpretation has growth rate $O(n^k)$.

14 Bojańczyk has recently shown that the converse holds for MSO interpretations, but not for
15 pebble transducers. We give significantly simplified proofs of those two results, extending the former
16 to first-order interpretations by reduction to an elementary property of \mathbb{N} -weighted automata. For
17 any k , we also prove the stronger statement that there is some quadratic polyregular function whose
18 output language differs from that of any k -fold composition of macro tree transducers (and which
19 therefore cannot be computed by any k -pebble transducer).

20 In the special case of unary input alphabets, we show that k pebbles suffice to compute polyregular
21 functions of growth $O(n^k)$. This is obtained as a corollary of a basis of simple word sequences whose
22 ultimately periodic combinations generate all polyregular functions with unary input. Finally, we
23 study polyregular and polyblind functions between unary alphabets (i.e. integer sequences), as well
24 as their first-order subclasses.

25 **2012 ACM Subject Classification** Theory of computation → Transducers

26 **Keywords and phrases** integer-valued polynomials, k -iterative languages, macro tree transducers,
27 (comparison-free) polyregular functions

28 **Funding** Lê Thành Dũng (Tito) Nguyễn: Supported by the LABEX MILYON (ANR-10-LABX-0070)
29 of Université de Lyon, within the program “Investissements d’Avenir” (ANR-11-IDEX-0007) operated
30 by the French National Research Agency (ANR).

31 **Acknowledgements** We would like to thank Mikołaj Bojańczyk for giving us much of the initial
32 inspiration that led to this paper; Nathan Lhote for suggesting to us the approach of Section 3; and
33 Gaëtan Douéneau-Tabot for stimulating discussions.

34 Contents

35	1 Introduction	2
36	2 A simple proof that inner squaring requires 3 pebbles	4
37	2.1 Regular functions with origin information	5
38	2.2 Pebble transducers in an abstract style and polyblind functions	5
39	2.3 Proof of Theorem 2.1	7
40	3 Dimension minimization for first-order interpretations	9
41	3.1 First-order interpretations	9
42	3.2 Reducing dimension minimization to a lemma on FO queries	11

2 Revisiting the growth of polyregular functions

43	3.3	N-weighted automata	12
44	3.4	Proof of Lemma 3.10	13
45	4	Quadratic polyregular functions vs macro tree transducers	14
46	4.1	Compositions of macro tree transducers (MTTs)	14
47	4.2	Proof of Theorem 4.1	15
48	5	Polyregular word sequences	17
49	5.1	For-transducers in an abstract style	17
50	5.2	Minimizing the number of loops for polyregular sequences	19
51	5.3	Polyregular integer sequences (unary output too)	21

Warning: This paper is unfinished: among other things, Section 5 is still very messy.

1 Introduction

Many works about *transducers* – automata that can produce output string/trees, not just recognize languages – are concerned with:

- either well-known classes of *linearly growing* functions, i.e. $|f(x)| = O(|x|)$ – in the string-to-string case, those are the sequential, rational and regular functions, see e.g. [39];
- or devices with possibly (hyper)exponential growth rates: HDTOL systems [29, 30, 15], macro tree transducers (MTTs) [26], compositions of MTTs (see below)...

A middle ground is occupied by the *polyregular functions*, surveyed in [5], which derive their name from their *polynomial* growth rate. They are the functions computed by string-to-string *pebble transducers*; although they have been around for two decades (starting with [38], see also [25, 24, 18]), it is only in 2018–2019 that Bojańczyk et al. [3, 6] introduced several alternative characterizations. This had led to renewed interest in this robust function class (e.g. [9, 33]), including applications to linguistics [43].

► **Example 1.1** (from [5, §6.2]). The *inner squaring* function, defined as

$$\begin{aligned} \mathbf{innsq} : \{a, b, \#\}^* &\rightarrow \{a, b, \#\}^* \\ w_0\#\dots\#w_n &\mapsto (w_0)^n\#\dots\#(w_n)^n \quad (w_0, \dots, w_n \in \{a, b\}^*) \end{aligned}$$

(e.g. $\mathbf{innsq}(ba\#abba\#b) = baba\#abbaabba\#bb$), is polyregular and $|\mathbf{innsq}(w)| = O(|w|^2)$.

A straightforward property of k -pebble¹ transducers is that their growth rate is² $O(n^k)$, and this is the best possible bound. Hence the question of *pebble minimization*: if a polyregular function (defined by some ℓ -pebble transducer) has growth $O(n^k)$, is it always computable by some k -pebble transducer? Definitely not, as Bojańczyk recently showed [4]: no number of pebbles suffices to compute all polyregular functions with *quadratic* growth. Moreover, he uses the tools from [4] to show that \mathbf{innsq} requires 3 pebbles in [5, Theorem 6.3].³ (That said, it is known that pebble minimization holds in some restricted settings [14].)

The same paper [4] also contains a proof of a positive result: the polyregular functions with growth $O(n^k)$ are exactly those defined by *k-dimensional MSO interpretations* (a logical

¹ We follow the convention of newer papers for two-way string transducers are 1-pebble transducers, while older papers would have called them 0-pebble transducers; thus, $k \geq 1$ here.

² This applies to strings; for trees, (output *height*) = O ((input *size*) ^{k}) [25, Lemma 7].

³ The paper [4] proposes a different quadratic example that requires 3 pebbles, called “block squaring”.

79 formalism introduced in [6]; as usual, MSO stands for Monadic Second-Order logic). But
 80 this proof is rather technical, and the one for pebble non-minimization even more so.

81 **Contributions (1): simpler proofs** In this paper, we first propose faster ways to reprove
 82 the results of [4], sometimes obtaining new results along the way:

- 83 ■ In Section 2, we give a short proof that `innsq` requires 3 pebbles. It only depends on a
 84 few familiar properties of regular (not polyregular) functions, such as their closure under
 85 composition and an elementary pumping lemma for their output languages.
- 86 ■ We prove dimension minimization for *first-order* (FO) interpretations in Section 3, using
 87 a strategy that could directly apply *mutatis mutandis* to MSO interpretations (though
 88 we chose to treat the FO case because it has not been explicitly proved in the literature).
 89 In fact, we use a reduction to a lemma on first-order queries which is quite similar to the
 90 one outlined in [5, §6.1]; but then, instead of using the heavy machinery of factorization
 91 forests⁴ as in [4] to prove our lemma, we explain how to quickly deduce it from a simple
 92 known property of \mathbb{N} -weighted automata. Let us remark an immediate corollary of our
 93 result: if a first-order polyregular function admits an MSO interpretation of dimension k ,
 94 it also admits an FO interpretation of dimension at most k .
- 95 ■ In Section 4, we construct a sequence $(f_k)_{k \geq 1}$ of quadratic polyregular functions such
 96 that f_k requires $k + 1$ pebbles. More than that: the output language of f_k differs from
 97 that of any k -fold composition of macro tree transducers. (This composition hierarchy is
 98 quite canonical and well-studied, see §4.1.) Again, our proof is quite short,⁵ and even
 99 arguably easier to check than our ad-hoc argument for `innsq`, though it is less elementary
 100 since its “trusted base” is larger: we use a powerful “bridge theorem” on MTTs from [23].

101 **Subclasses of polyregular functions** Some restrictions on pebble transducers ensure that:

- 102 ■ the computed function sits at a low level of the aforementioned composition hierarchy;
- 103 ■ a pebble minimization property holds.

104 This is for instance the case for *k-marble transducers* – they can compute precisely the
 105 same string-to-string functions as MTTs⁶ with growth⁷ $O(n^k)$ [15, §5] – or for *blind pebble*
 106 *transducers*, which define *polyblind*⁸ *functions* (cf. Theorem 2.8, taken from [40]). The present
 107 paper makes a small contribution to the study of polyblind functions by refuting their
 108 conjectured logical characterization (Corollary 3.6; the counterexample is the aforementioned
 109 inner squaring). Furthermore, Douéneau-Tabot has recently proved pebble minimization
 110 for “last pebble” transducers [14], a model that subsumes both marble and blind pebble
 111 transducers; and a function computed by any such device can be obtained by a composition
 112 of *two* MTTs (a consequence of [21, Theorem 53 (in §15)]).

⁴ The factorization forest theorem admits a slight weakening that applies to the aperiodic case. This should enable in principle an adaptation of the proof of [4] to FO interpretations, so we do not claim that our ability to handle this case is a technical advantage of our approach. However, by avoiding factorization forests, our proof strategy should extend to MSO interpretations on trees (see Remark 3.14).

⁵ And mostly unoriginal: it consists of little adjustments to an argument by Engelfriet & Maneth [24, §4]. In a followup paper [18], Engelfriet mentions and corrects a mistake in [24, §3], but it does not affect the section which is relevant for our purposes.

⁶ Using the fact that string-to-string MTTs are syntactically isomorphic, up to insignificant details, to copyful streaming string transducers.

⁷ Interestingly, the proof of marble minimization uses the same property of \mathbb{N} -automata that we apply to dimension minimization of FO interpretations.

⁸ A name given by Douéneau-Tabot [12, 13, 14] to what Nguyễn, Noûs and Pradic [40] originally called the “comparison-free” subclass of polyregular functions.

Tito — see also [49]

113

114 **Contributions (2): unary inputs** For general polyregular functions with a *unary output*
 115 *alphabet*, pebble minimization also works, because every such function can be computed
 116 by a marble transducer [12, §4] (see also Theorem 5.20). What about unary *inputs* (with
 117 arbitrary outputs)? Section 5 focuses on this case, that is, on polyregular *word sequences*
 118 $\mathbb{N} \rightarrow \Sigma^*$ (interpreting input words as unary numerals). We establish pebble minimization
 119 in this special case by showing that polyregular sequences⁹ can be expressed as “ultimately
 120 periodic combinations” of word sequences generated by some simple operations. We believe
 121 this latter characterization might be of independent interest: it is inspired by an analogous
 122 property of polyblind word sequences [40, §9] that has been applied to prove a separation
 123 result [40, Theorem 8.3(i)].

124 Finally, still in Section 5, we give simple characterizations for polyregular and polyblind
 125 *integer sequences* (unary input and output), as well as their first-order (i.e. aperiodic [48])
 126 subclasses. This leads us to a somewhat surprising observation: $n \mapsto n(n-1)/2$ is both
 127 polyblind and first-order polyregular, but not first-order polyblind.

128 **Organization of the paper** A plan of the sections has already been outlined above; each of
 129 them comes with its own preliminaries subsections. Sections 2 and 3 are independent from
 130 each other, since they use different characterizations of polyregular functions, while Section 4
 131 reuses some of the preliminaries (but not the new results) of Section 3. Section 5 is mostly
 132 independent from the rest of the paper.

133 **Notations.** The set of natural numbers is $\mathbb{N} = \{0, 1, \dots\}$. Alphabets are always finite sets.

134 Let Σ be an alphabet. We write ε for the empty string and Σ^* for the set of strings (or
 135 words) with letters in Σ , i.e. the free monoid over Σ ; the Kleene star $(-)^*$ will also be applied
 136 to languages $L \subseteq \Sigma^*$ as part of usual regular expression syntax. We denote by $\underline{\Sigma}$ a disjoint
 137 copy of Σ , made of “underlined letters”, so that $a \in \Sigma \mapsto \underline{a} \in \underline{\Sigma}$ is a bijection.

138 Let $w = w_1 \dots w_n \in \Sigma^*$ ($w_i \in \Sigma$ for $i \in \{1, \dots, n\}$). We use the notations $w[i] = w_i$ and
 139 $w[i \dots j] = w_i \dots w_j$ for $1 \leq i \leq j \leq n$. The length $|w|$ of w is n , and $|w|_c$ refers to the number
 140 of occurrences of $c \in \Sigma$ in w . We also write $w \dot{\downarrow} i$ for $w_1 \dots w_{i-1} \underline{w_i} w_{i+1} \dots w_n \in (\Sigma \cup \underline{\Sigma})^*$.

141 The *output language* of a function $f: X \rightarrow \Sigma^*$ is $f(X) \subseteq \Sigma^*$, also denoted by $\text{Im}(f)$.

142 **In sections 2 and 3, all our (poly)regular functions map ε to ε .** This convention goes
 143 against the usual practice of allowing the image of the empty string to be any string over the
 144 output alphabet. However, we make this choice to avoid inessential inconveniences that arise
 145 when dealing with ε in origin semantics and FO interpretations. This makes no difference
 146 concerning the strength of our results, since ε could always be dealt with as a special case.

147 **2 A simple proof that inner squaring requires 3 pebbles**

148 While the inner squaring function (Example 1.1) has quadratic growth, we reprove that:

149 ► **Theorem 2.1** ([5, Theorem 6.3]). $\text{innsq} \in \text{Pebble}_3 \setminus \text{Pebble}_2$.

150 Note however that the restriction of innsq to inputs in $\{a, \#\}^*$ is in Pebble_2 .

⁹ In fact, our results apply to a notion of polyregular function $\mathbb{N}^d \rightarrow \Sigma^*$.

151 ▶ Remark 2.2. Bojańczyk’s proof sketch for [5, Theorem 6.3] actually shows that the function
 152 $a_1 \dots a_n \in \mathbb{A}^* \mapsto (a_1)^n \dots (a_n)^n$ cannot be computed by a 2-pebble *atom-oblivious* transducer
 153 (here the alphabet \mathbb{A} is an *infinite* set of *atoms*). Combining this with the Deatomization
 154 Theorem from [4] – whose proof is rather complicated – shows that no $f: \{\langle, \rangle, \bullet\}^* \rightarrow \{\langle, \rangle, \bullet\}^*$
 155 in Pebble_2 can satisfy $f(\langle \bullet^{p_1} \rangle \dots \langle \bullet^{p_n} \rangle) = (\langle \bullet^{p_1} \rangle)^n \dots (\langle \bullet^{p_n} \rangle)^n$, from which Theorem 2.1 can
 156 easily be deduced using the composition properties of pebble transducers [18].

157 Instead of introducing the hierarchy $(\text{Pebble}_k)_{k \in \mathbb{N}}$ by a concrete machine model (k -pebble
 158 transducers), Section 2.2 defines it using combinators (operators on functions). This abstract
 159 presentation depends on the *origin semantics* [2] (see also [39, §5]) of regular functions,
 160 about which we say a few words in Section 2.1. We also recall in §2.2 a similar definition of
 161 polyblind functions; their only use in this section will be to state and prove Corollary 2.9,
 162 but they will appear again in later sections. After all this, Section 2.3 proves Theorem 2.1.

163 2.1 Regular functions with origin information

164 We will avoid explicit manipulations of automata computing regular functions. For our
 165 purposes here, we just need to know that they are the functions computed by *two-way*
 166 *transducers* (2DFTs), since we use an old result (Lemma 2.13) from the literature on 2DFTs
 167 to prove Theorem 2.1. For more on regular functions and their origin semantics, we refer to
 168 the survey [39]. A typical example of a regular string-to-string function is:

$$169 \quad a^{m_0} \# \dots \# a^{m_n} \in \{a, \#\}^* \mapsto a^{m_0} b^{m_0} \# \dots \# a^{m_n} b^{m_n}$$

170 There are several natural ways to lift it to a regular function with origins, which reflect
 171 different ways to compute it with a 2DFT: for example $aaa\#aa$ could be mapped to

$$172 \quad \begin{array}{ccc} a a a b b b \# a a b b & \text{or} & a a a b b b \# a a b b \\ 173 \quad 1 2 3 1 2 3 4 5 6 5 6 & & 1 2 3 3 2 1 4 5 6 6 5 \end{array}$$

174 The second component indicates, for each output letter, which input position it “comes from”.

175 We shall write $f^\circ, g^\circ, \dots: \Gamma^* \rightarrow (\Sigma \times \mathbb{N})^*$ for regular functions with origin information
 176 and $f, g, \dots: \Gamma^* \rightarrow \Sigma^*$ for the corresponding regular string-to-string functions. Thus, we have
 177 $f = (\pi_1)^* \circ f^\circ$ and $(\pi_2)^* \circ f(w) \in \{1, \dots, |w|\}^*$ for any $w \in \Gamma^*$, where π_i is the projection on
 178 the i -th component of a pair ($i \in \{1, 2\}$).

179 2.2 Pebble transducers in an abstract style and polyblind functions

180 Recall the notations $\underline{\Gamma}$ and $w \dot{\downarrow} i$ from the end of the introduction.

181 ▶ **Definition 2.3.** Let $f^\circ: \Gamma^* \rightarrow (I \times \mathbb{N})^*$ be a regular function with origins (so $f: \Gamma^* \rightarrow I^*$
 182 is regular) and, for $i \in I$, let $g_i: (\Gamma \cup \underline{\Gamma})^* \rightarrow \Sigma^*$ and $h_i: \Gamma^* \rightarrow \Sigma^*$. For $w \in \Gamma^*$, define

$$183 \quad \text{pebble}(f^\circ, (g_i)_{i \in I})(w) = g_{i_1}(w \dot{\downarrow} j_1) \cdot \dots \cdot g_{i_n}(w \dot{\downarrow} j_n) \quad \text{where} \quad f^\circ(w) = (i_1, j_1) \dots (i_n, j_n)$$

$$184 \quad \text{blind}(f, (h_i)_{i \in I})(w) = h_{i_1}(w) \cdot \dots \cdot h_{i_n}(w)$$

186 Using these combinators, we define the hierarchies of function classes Pebble_n and Blind_n
 187 inductively. $\text{Pebble}_0 = \text{Blind}_0$ is the class of string-to-string functions with finite range (or
 188 equivalently, whose output has bounded length), and

$$189 \quad \forall n \in \mathbb{N}, \quad \text{Pebble}_{k+1} = \{\text{pebble}(f^\circ, (g_i)) \mid f^\circ \text{ regular}, g_i \in \text{Pebble}_k\}$$

$$190 \quad \text{Blind}_{k+1} = \{\text{blind}(f, (g_i)) \mid f \text{ regular}, g_i \in \text{Blind}_k\}$$

6 Revisiting the growth of polyregular functions

192 The following properties can be proved routinely.

193 ► **Proposition 2.4.** Pebble_1 and Blind_1 are both equal to the class of regular functions.
 194 Furthermore, $\text{Blind}_n \subseteq \text{Pebble}_n$ for all $n \in \mathbb{N}$. Finally, the hierarchies are monotone:
 195 $\text{Pebble}_n \subset \text{Pebble}_{n+1}$ and $\text{Blind}_n \subset \text{Blind}_{n+1}$ for all $n \in \mathbb{N}$.

196 ► **Example 2.5.** Let $f^\circ: \{a, b, \#\}^* \rightarrow (\{\bullet, \#\} \times \mathbb{N})^*$ be defined by

$$197 \quad f^\circ(\underbrace{w[1] \dots w[i_1 - 1]}_{\text{each block is in } \{a,b\}^*} \# \dots \# w[i_m + 1] \dots w[n]) = \begin{matrix} \bullet \# & \bullet & \# & \bullet \\ 1 & i_1 & i_1 + 1 & \dots & i_m & i_m + 1 \end{matrix}$$

each \bullet is omitted if the corresponding input block is empty

198 and $h(\dots \# c w \# \dots) = c w$ for $c \in \{a, b\}$ and $w \in \{a, b\}^*$. Then

$$199 \quad \underbrace{\text{innsq}}_{\text{Example 1.1}} = \underbrace{\text{pebble}(f^\circ, (g_i)_{i \in \{\bullet, \#\}})}_{\in \text{Pebble}_3} \quad g_\bullet = \underbrace{\text{blind}(w \mapsto \bullet^{|w|}, (h)_{j \in \{\bullet\}})}_{\in \text{Blind}_2 \subset \text{Pebble}_2} \quad g_\# : \underbrace{w \mapsto \#}_{\in \text{Pebble}_0}$$

regular with origin information both regular i.e. in Pebble_1

200 For $k \geq 1$, the class Blind_k corresponds *by definition* to the polyblind functions of rank
 201 at most¹⁰ k . Indeed, the combinator blind is precisely the “composition by substitution”
 202 operation used to define polyblind functions in [40, §4]. It is then proved in [40, §5]
 203 that this coincides with the expressive power of a machine model, the comparison-free (or
 204 blind [12, 13, 14]) k -pebble transducers. Similarly, we have:

205 ► **Proposition 2.6.** Pebble_k is exactly the class of string-to-string functions computed by
 206 k -pebble transducers [5, §2] for any $k \geq 1$.

207 **Proof idea.** By a straightforward adaptation of the proof of [40, Corollary 5.7]; therefore, we
 208 refer the reader to the long version of [40] for details. To translate Pebble_k into transducers,
 209 observe that $\text{pebble}(f^\circ, (g_i)_{i \in I})$ is definable by a k -pebble transducer whenever some two-way
 210 transducer (2DFT) computes f° and each g_i is computed by some $(k - 1)$ -pebble transducer.
 211 The converse relies on the fact that the functions computed by 2DFTs *with regular lookahead*
 212 are regular [20, Lemma 4]. Note that a similar “origin-based” definition of “last pebble”
 213 transducers is stated in [14, Definition 3.3], without explicit proof that they indeed coincide
 214 with the expected machine model. ◀

215 Tito — Actually we should probably give the proof in an appendix or something, and also prove the thing said to be tedious in the following remark

216 ► **Remark 2.7.** By replacing “regular” by “first-order regular” in the definition of Pebble_k , we
 217 would get a natural definition of the FO k -pebble hierarchy. However, it is less immediate
 218 than for the above proposition to relate it to the first-order k -pebble transducers as defined
 219 in [5, Theorem 7.1(2)]. While we expect that the two should indeed coincide, and that this
 220 should not involve any conceptual difficulty, we did not attempt a rigorous proof since it
 221 would probably be tedious. That said, we shall characterize in Section 5.1 the FO polyregular
 222 functions using a similarly abstracted version of *for-transducers* [5, §1].

223 Finally, let us recall that unlike general pebble transducers (as we shall see soon), blind
 224 pebble transducers enjoy the pebble minimization property [40, Theorem 7.1]. For an *effective*
 225 minimization proof and a generalization to “last pebble” transducers, see [14].

¹⁰ Actually, in [40, Definition 4.6] they are called the comparison-free polyregular functions of rank at most $k - 1$, so there is an offset of 1 with respect to our numbering of the hierarchy.

226 ► **Theorem 2.8** ([40, 14]). For all $k \in \mathbb{N}$, $\text{Blind}_k = \left\{ f \in \bigcup_{\ell \in \mathbb{N}} \text{Blind}_\ell \mid |f(w)| = O(|w|^k) \right\}$.

227 ► **Corollary 2.9.** The inner squaring function (Example 1.1) is not polyblind.

228 **Proof.** Since it has quadratic growth, if it were polyblind, it would be in $\text{Blind}_2 \subset \text{Pebble}_2$.
 229 This would contradict the main result of this section (Theorem 2.1). ◀

230 2.3 Proof of Theorem 2.1

231 We have seen in Example 2.5 that $\text{innsq} \in \text{Pebble}_3$. Our approach to prove $\text{innsq} \notin \text{Pebble}_2$
 232 goes through the output languages of regular functions.

233 ► **Definition 2.10.** Call a language $L \subseteq \Sigma^*$ a regular image if there exists a regular function
 234 f with codomain Σ^* such that $L = \text{Im}(f)$.

235 Not all regular images are regular languages, or even context-free languages: consider for
 236 instance $\{a^n b^n c^n d^n \mid n \in \mathbb{N}\}$. The usual closure properties of regular functions entail the
 237 following for their images:

238 ► **Proposition 2.11.** The class of regular images contains all regular languages L with $\varepsilon \in L$.
 239 If f is a regular function and L, L' are regular images, so are $f(L)$ and $L \cup L'$.

240 We will show that no function in Pebble_2 can coincide with innsq on the subset of inputs

$$241 (a^*b\#)^*\#\# = \{\text{strings with the shape } a \dots ab\# \dots \#a \dots ab\#\# \dots \#\}$$

242 For the sake of contradiction, assume the opposite.

243 ▷ **Claim 2.12.** Under this assumption, there exists a language $L \subseteq b\{a, b\}^*b$

- 244 ■ which is a regular image
- 245 ■ that contains only factors of words from $\text{innsq}((a^*b\#)^*\#\#)$
- 246 ■ and such that, for any $N \in \mathbb{N}$, some word in L has the shape

$$247 \overbrace{ba \dots ab \dots b a \dots a b}^{\text{there are at least } N \text{ bs}}$$

every block of a s has the same length $n \geq N$

248 **Proof.** We start by an analysis of some output factors arising in the computation of innsq
 249 by a hypothetical 2-pebble transducer, which will later inform the construction of a suitable
 250 language L . Unfolding our formal definition of Pebble_2 , our assumption means that innsq
 251 coincides on $(a^*b\#)^*\#\#$ with some function $\text{pebble}(g^\circ, (h_i)_{i \in I})$ where

- 252 ■ $g^\circ: \{a, b, \#\}^* \rightarrow (\{a, b, \#\} \times \mathbb{N})^*$ is a regular function with origins;
- 253 ■ for $i \in I$, the function $h_i: \{a, b, \#, \underline{a}, \underline{b}, \underline{\#}\}^* \rightarrow \{a, b, \#\}^*$ is in Pebble_1 , so it is regular.

254 The functions g° and h_i , being regular, have linear growth: for some $C \in \mathbb{N}$, we have
 255 $|g^\circ(u)| \leq C|u|$ and $|h_i(v)| \leq C|v|$ for all nonempty input words u, v and $i \in I$.

256 Consider the input word $u = (a^n b\#)^n \#^{nm}$ for some $n, m \in \mathbb{N}$ (which shall be taken large
 257 enough to satisfy constraints that will arise during the proof). By definition,

$$258 \text{innsq}(u) = \underbrace{h_{i_1}(u \downarrow j_1)}_{:=v_1} \cdot \dots \cdot \underbrace{h_{i_k}(u \downarrow j_k)}_{:=v_k} \quad \text{where} \quad g^\circ(w) = (i_1, j_1) \dots (i_k, j_k)$$

8 Revisiting the growth of polyregular functions

259 For large enough n and m , we have that for all $\ell \in \{1, \dots, k\}$,

$$260 \quad |v_\ell| \leq C|u \dot{\downarrow} j_\ell| = Cn(n+2+m) < (n+1)n(1+m) = \underbrace{|a^n b| \times |u|_\#}_{\text{distance between two consecutive non-adjacent occurrences of } \# \text{ in } \mathbf{innsq}(u)}$$

261 Thanks to the above, if we call p the largest integer such that $|v_1 \dots v_p|_\# < n$, then:

262 ■ $|v_\ell|_\# \leq 1$ for every $\ell \in \{1, \dots, p\}$;

263 ■ $v_1 \dots v_p$ is a prefix of $\mathbf{innsq}(u)$ that contains all the $n-1$ first occurrences of $\#$.

264 As a consequence of the latter item, $|v_1 \dots v_p|_b \geq |((a^n b)^{n(1+m)})\#|^{n-1}$. Meanwhile, we also
265 have $p < k = |g^\circ(u)| \leq C|u|$. Hence

$$266 \quad \frac{|v_1 \dots v_p|_b}{p} \geq \frac{n(1+m)(n-1)}{Cn(n+2+m)} \xrightarrow{n,m \rightarrow +\infty} +\infty$$

267 Thus, by the pigeonhole principle, for an input u parameterized by large enough values of n
268 and m , there is at least one output factor v_ℓ containing $2N$ or more occurrences of b . Since
269 $|v_\ell|_\# \leq 1$, either its largest $\#$ -free prefix or its largest $\#$ -free suffix (or both) contains N
270 occurrences of b . After applying the regular function $a \dots abwba \dots a \mapsto bw b$ (for $w \in \{a, b\}^*$)
271 to trim this prefix or suffix, we get a factor of v_ℓ – and therefore of $\mathbf{innsq}(u)$ – in $b(a^n b)^*$,
272 and we may take $n \geq N$ to ensure that the blocks of as have length at least N .

273 These observations justify the following definition: let L consist of all words that can
274 be obtained by trimming (in the above sense) the largest $\#$ -free prefix or suffix of some
275 $h_i(u \dot{\downarrow} j)$ where $u \in (a^* b \#)^* \#^*$ and (i, j) appears in $g^\circ(u)$. We just argued above that for
276 every N , there are some $n, r \geq N$ such that $b(a^n b)^r \in L$. To conclude, the closure properties
277 of Proposition 2.11 can be used to show that L is a regular image. \triangleleft

278 Tito — Perhaps expand that last sentence to give more details (also in Proposition 2.11), because it does not have a proof

279 We shall now use the following old pumping lemma (see also Remark 2.14):

280 ► **Lemma 2.13** (Rozoy [44, §4.1]). *If L is a regular image, then for some $k, K \in \mathbb{N}$, every*
281 *$w \in L$ with $|w| \geq K$ has a decomposition $w = u_0 v_1 \dots u_{k-1} v_k u_k$ with*

282 ■ $\exists i \in \{1, \dots, k\}: v_i \neq \varepsilon$

283 ■ $\forall i \in \{1, \dots, k\}, |v_i| \leq K$

284 ■ $\{u_0(v_1)^n \dots u_{k-1}(v_k)^n u_k \mid n \in \mathbb{N}\} \subseteq L$

285 Let us apply this lemma to the language L given by Claim 2.12, yielding two constants
286 $k, K \in \mathbb{N}$. One of the properties stated in Claim 2.12 is that there exist $n \geq K$ and $r \geq 2k+1$
287 such that $b(a^n b)^r \in L$. This string has length greater than K , so it is pumpable:

$$288 \quad b(a^n b)^r = u_0 v_1 \dots u_{k-1} v_k u_k \quad \text{and} \quad w = u_0 (v_1)^2 \dots u_{k-1} (v_k)^2 u_k \in L$$

289 Let us now show that $ba^n b$ is a factor of w . Since $|v_i| \leq K \leq n$ for every $i \in \{1, \dots, k\}$, and
290 any two occurrences of b in $b(a^n b)^r$ are separated by a^n , each factor v_i can contain at most
291 one b . So $|v_1|_b + \dots + |v_k|_b \leq k$ which leads to $|u_0|_b + \dots + |u_k|_b \geq r+1-k \geq k+2$. One of
292 the u_j must then contain two occurrences of b , and this u_j is also a factor of w .

293 Recall that L is comprised exclusively of $\#$ -free factors of words in $\mathbf{innsq}((a^* b \#)^* \#^*)$,
294 so *all the “blocks of as ” in w must have the same size*. Having seen above that this size must
295 be n , we can now conclude by a case analysis, reaching a contradiction in both cases:

- 296 ■ First assume that $|v_i|_b \geq 1$ for some $i \in \{1, \dots, k\}$. We may write $v_i = a^\ell b \dots b a^m$ (where
 297 the first and last b can coincide) in this case. Then w contains $(v_i)^2$ which in turn contains
 298 a factor $ba^{\ell+m}b$. So $\ell + m = n$ since it is the size of a block of a s in w ; but at the same
 299 time, using the second item of Lemma 2.13, $\ell + m < |v_i| \leq K \leq n$.
- 300 ■ Otherwise, $v_i \in \{a\}^*$ for all i , so pumping does not increase the number of b s. Therefore
 301 w has as many blocks of a s as $b(a^n b)^r$, and we have also seen that its blocks have size n ,
 302 so $w = b(a^n b)^r$. But since at least one v_i is nonempty, $|w| > |b(a^n b)^r|$.
- 303 ► Remark 2.14. Rozoy proves Lemma 2.13 for two-way transducers (cf. §2.1). She also
 304 shows [44, §4.2] that the output languages of *nondeterministic* two-way transducers enjoy a
 305 weaker version of the lemma without the bound on the length of the v_i (this bound is refuted
 306 by the example $\{(w\#)^n \mid w \in \Sigma^*, n \in \mathbb{N}\}$); in general, the languages that satisfy this weaker
 307 pumping lemma are called *k-iterative* in the literature (see e.g. [47, 34]). For more on regular
 308 images, see [16, 19]; several references are also given in [32, p. 18].¹¹

309 3 Dimension minimization for first-order interpretations

310 We first recall the definition of FO interpretations in Section 3.1 and give several examples.
 311 Then Section 3.2 states the dimension minimization theorem, and explains how to derive it
 312 as a consequence of a lemma on queries defined by first-order formulas. To establish this
 313 lemma, we first recall some background on \mathbb{N} -weighted automata in §3.3 and then use it to
 314 give a proof in §3.4.

315 3.1 First-order interpretations

316 We assume basic familiarity with first-order logic (FO).

317 Let Σ be an alphabet. A word $w \in \Sigma^*$ can be seen as a structure whose domain is the
 318 set of *positions* $\{1, \dots, |w|\}$, over the relational vocabulary consisting of:

- 319 ■ for each $c \in \Sigma$, a unary symbol c where $c(i)$ is interpreted as true whenever $w[i] = c$,
 320 ■ and a binary relation symbol \leq , interpreted as the total order on positions.

321 As an example, for $\Sigma = \{a, b\}$, let $F(x) = b(x) \vee \forall y. (a(y) \vee x \leq y)$. A string $w \in \Sigma^*$
 322 satisfies $F(i)$ for a given $i \in \{1, \dots, w\}$ – notation: $w \models F(i)$ – when either $w[i] = b$, or the
 323 position i contains an a which occurs before (at the left of) all the b s in the word w . Thus,
 324 the formula $\forall x. F(x)$ evaluates to true exactly over the words in a^*b^* . (It is precisely the
 325 star-free languages that can be defined this way, see [48] for a history of this seminal result.)

326 This model-theoretic perspective also leads to a way to specify string-to-string functions.
 327 (Note that FO interpretations make sense for – and have been intensively studied in the
 328 context of – more general relational structures such as graphs, see e.g. [10, 11].)

329 ► **Definition 3.1** (variant of [5, §5.1]¹²). *Let $k \geq 1$ and Γ, Σ be alphabets.*

330 *A k -dimensional first-order interpretation \mathcal{I} from Γ^* to Σ^* consists of*

- 331 ■ *a finite set of components $\mathcal{I}_{\text{comp}}$;*
 332 ■ *an FO formula $\mathcal{I}_c^\alpha(x_1, \dots, x_k)$ over Γ for each $c \in \Sigma$ and $\alpha \in \mathcal{I}_{\text{comp}}$;*

¹¹ Instead of Latteux’s unobtainable technical report cited by [32], see [36, Prop. I.2] & compare with [42].

¹² Bojańczyk’s definition in [5, §5.1] uses components with varying arities $\leq k$. Those are essentially equivalent to our definition except that the arity 0 components provide a way to encode non-empty outputs for the empty input string – a case that we have excluded by convention (as mentioned at the end of the introduction). The definition used by Bojańczyk, Kiefer & Lhote in [6] is equivalent to our single-component interpretations; but see Remark 3.8.

10 Revisiting the growth of polyregular functions

333 ■ an FO formula $\mathcal{I}_{\leq}^{\alpha, \beta}(x_1, \dots, x_k, y_1, \dots, y_k)$ over Γ for each $(\alpha, \beta) \in (\mathcal{I}_{\text{comp}})^2$.

334 It defines the function $u \in \Gamma^* \mapsto \begin{cases} v & \text{if } O_u \cong \text{the structure corresponding to } v \\ \varepsilon & \text{when}^{13} \text{ there is no such } v \in \Sigma^* \end{cases}$

335 where O_u is the structure

336 ■ with domain $\{ \langle \alpha, i_1, \dots, i_k \rangle \in \times \{1, \dots, |u|\}^k \mid \exists c \in \Sigma: u \models \mathcal{I}_c^\alpha(i_1, \dots, i_k) \}$;

337 ■ where $c(\langle \alpha, i_1, \dots, i_k \rangle)$ iff $u \models \mathcal{I}_c(i_1, \dots, i_k)$ for $c \in \Sigma$ and $\alpha \in \mathcal{I}_{\text{comp}}$;

338 ■ where $\langle \alpha, i_1, \dots, i_k \rangle \leq \langle \beta, j_1, \dots, j_k \rangle$ iff $u \models \mathcal{I}_{\leq}^{\alpha, \beta}(i_1, \dots, i_k, j_1, \dots, j_k)$.

339 \mathcal{I} is said to be a single-component FO interpretation when $|\mathcal{I}_{\text{comp}}| = 1$. In that case we
340 usually omit the components α, β, \dots .

341 ► Remark 3.2. It follows directly from the definitions that if a function f is defined by a
342 first-order interpretation with N components, then $|f(u)| \leq N|u|^k$.

343 ► Example 3.3. Consider the following 1-dimensional FO interpretation: $\mathcal{I}_{\text{comp}} = \{ \lambda, \rho \}$ and

344 ■ $\mathcal{I}_a^\lambda(x) = \mathcal{I}_b^\rho(x) = \text{true}$ whereas $\mathcal{I}_b^\lambda(x) = \mathcal{I}_a^\rho(x) = \text{false}$

345 ■ $\mathcal{I}_{\leq}^\lambda(x, y) = \mathcal{I}_{\leq}^\rho(x, y) = (x \leq y)$ and $\mathcal{I}^{\lambda, \rho}(x, y) = \text{true}$ and $\mathcal{I}^{\rho, \lambda}(x, y) = \text{false}$

346 The function that it defines is $w \in \Gamma^* \mapsto a^{|w|}b^{|w|}$. Note that we could get the same function
347 by taking instead $\mathcal{I}_{\leq}^\rho(x, y) = (y \leq x)$ for example.

348 ► Example 3.4. $a^n \in \{a\}^* \mapsto (a^{n-1}b)^{n-1}$ is defined by a single-component 2-dimensional
349 FO interpretation \mathcal{I} such that \mathcal{I}_{\leq} is the lexicographic order over pairs and

$$350 \quad \mathcal{I}_a(x_1, x_2) = \neg \underbrace{\max(x_1)}_{\text{i.e. } \forall y. y \leq x_1} \wedge \neg \max(x_2) \quad \mathcal{I}_b(x_1, x_2) = \neg \max(x_1) \wedge \max(x_2)$$

351 ► Example 3.5. As a more subtle example, the inner squaring function (Example 1.1) admits
352 a single-component 2-dimensional FO interpretation that we intuitively illustrate as follows:

$$353 \quad \begin{array}{ccc} & \text{baba\#abba\#b} & \\ & \vdots & \\ & \dots & \\ \# & \text{baba abba b} & \rightsquigarrow \text{babababa\#abbaabba\#bb} \\ & \vdots & \\ & \dots & \\ \# & \text{baba abba b} & \\ b & \# \# & \end{array}$$

354 ■ $\mathcal{I}_a(x_1, x_2) = a(x_1) \wedge \#(x_2)$ and $\mathcal{I}_b(x_1, x_2) = b(x_1) \wedge \#(x_2)$

355 ■ $\mathcal{I}_{\#}(x_1, x_2) = \#(x_1) \wedge \forall y_2. y_2 \leq x_2$

356 ■ $\mathcal{I}_{\leq}(x_1, x_2, y_1, y_2) = \text{either } \#(y_1) \wedge x_1 \leq y_1, \text{ or there exist } x_1, y_1 \text{ such that}$

357 ■ $x_3 \leq x_1$, and there are no $\#$ s strictly in-between x_3 and x_1 ;

358 ■ $y_3 \leq y_1$, and there are no $\#$ s strictly in-between y_3 and y_1 ;

359 ■ neither x_3 nor y_3 has an immediate predecessor which is an a or a b ;

360 ■ $(x_3, x_2, x_1) \leq (y_3, y_2, y_1)$ for the lexicographic order on 3-tuples.

361 ► Corollary 3.6. The characterization of polyblind functions conjectured in [40, §10] is wrong.

¹³This is purely for convenience, to avoid having to distinguish between partial and total functions. The language of input words u for which such a v exists is star-free/FO-definable, so a partial FO interpretation can always be completed to a total one.

362 **Proof.** For purely syntactic reasons, the 2-dimensional FO interpretations such that:
 363 ■ in every formula defining the interpretation, each variable can be given a sort in $\{1, 2\}$ in
 364 such a way that variables of different sorts are never compared (for equality or ordering)
 365 ■ and in $\mathcal{I}(x_1, x_2, y_1, y_2)$ and $\mathcal{I}_c(x_1, x_2)$ ($c \in \Sigma$), x_i and y_i have sort $i \in \{1, 2\}$
 366 can be seen as a special case of the logical interpretations considered in [40, Conjecture 10.1].
 367 Example 3.5 fits these criteria (with x_3 and y_3 having sort 1). If the conjecture were true,
 368 **innsq** would therefore be polyblind, contradicting Corollary 2.9. ◀

369 ▶ **Theorem 3.7** (Bojańczyk, Kiefer & Lhote [6]). *The functions specified by FO interpretations*
 370 *are exactly the first-order polyregular functions.*

371 ▶ **Remark 3.8.** In [6] it is explained that for any first-order polyregular function, there is a
 372 single-component FO interpretation which coincides with it on every input of length ≥ 2 .
 373 But dimension minimization requires multiple components.

374 3.2 Reducing dimension minimization to a lemma on FO queries

375 We wish to prove a sort of converse to Remark 3.2:

376 ▶ **Theorem 3.9.** *Let f be a first-order polyregular function, and let $k \in \mathbb{N}$ be the minimal*
 377 *exponent such that $|f(w)| = O(|w|^k)$. Then f is specified by some FO interpretation \mathcal{J} of*
 378 *dimension at most k . Moreover, both k and \mathcal{J} can be effectively computed from a given*
 379 *interpretation of possibly non-optimal dimension defining f .*

380 As mentioned in the introduction, the counterpart to this theorem for MSO interpretations
 381 and general polyregular functions has been proved by Bojańczyk [4, §2]. The survey [5]
 382 gives a reduction [5, proof of Theorem 6.1] to a “Seed Lemma” [5, Lemma 6.2] concerning
 383 arbitrary MSO formulas. We analogously reduce Theorem 3.9 to a similar (but weaker, to
 384 make our job easier) lemma on FO formulas.

385 ▶ **Lemma 3.10.** *Let $F(x_1, \dots, x_\ell)$ be a first-order formula over Γ^* . One can compute:*
 386 ■ *the least $k \in \mathbb{N}$ such that $|\{(i_1, \dots, i_\ell) \mid w \models F(i_1, \dots, i_\ell)\}| = O(|w|^k)$ (so $k \leq \ell$);*
 387 ■ *a formula $G(x_1, \dots, x_\ell, z_1, \dots, z_k)$ and a bound $B \in \mathbb{N}$ such that for every $w \in \Gamma^*$,*
 388 ■ $\forall j_1, \dots, j_k, |\{(i_1, \dots, i_\ell) \mid w \models G(i_1, \dots, i_\ell, j_1, \dots, j_k)\}| \leq B$;
 389 ■ $\forall i_1, \dots, i_\ell, w \models F(i_1, \dots, i_\ell) \implies |\{(j_1, \dots, j_k) \mid w \models G(i_1, \dots, i_\ell, j_1, \dots, j_k)\}| = 1$.

390 **Proof of Theorem 3.9.** Let \mathcal{I} be an ℓ -dimensional FO interpretation from Γ^* to Σ^* . Let

$$391 F(x_1, \dots, x_\ell) = \bigvee_{\alpha \in \mathcal{I}_{\text{comp}}, c \in \Sigma} \mathcal{I}_c^\alpha(x_1, \dots, x_\ell)$$

392 By Lemma 3.10, one can compute k , G and B verifying several properties. Note that we
 393 have $|\mathcal{I}(w)| = \Theta(|\{(i_1, \dots, i_\ell) \mid w \models F(i_1, \dots, i_\ell)\}|)$ (their ratio is between 1 and $|\mathcal{I}_{\text{comp}}|$);
 394 therefore, k is the smallest integer such that $|\mathcal{I}(w)| = O(|w|^k)$.

395 For $m \in \{1, \dots, B\}$, there is an FO formula $H_m(x_1, \dots, x_\ell, z_1, \dots, z_k)$ whose meaning is:
 396 among the tuples (y_1, \dots, y_ℓ) such that $G(y_1, \dots, y_\ell, z_1, \dots, z_k)$, the tuple (x_1, \dots, x_ℓ) is the
 397 m -th one (for some fixed first-order definable total order, e.g. lexicographic). (H_m is false
 398 whenever there are strictly less than m such tuples, or (x_1, \dots, x_ℓ) is not one of those tuples.)
 399 The following k -dimensional interpretation \mathcal{J} then computes the same function as \mathcal{I} :

- 400 ■ $\mathcal{J}_{\text{comp}} = \mathcal{I}_{\text{comp}} \times \{1, \dots, B\}$
 401 ■ $\mathcal{J}_c^{(\alpha, m)}(z_1, \dots, z_k) = \exists x_1. \dots \exists x_\ell. H_m(x_1, \dots, x_\ell, z_1, \dots, z_k) \wedge \mathcal{I}_c^\alpha(x_1, \dots, x_\ell)$
 402 ■ $\mathcal{J}_{\leq}^{(\alpha, m), (\beta, p)}(\vec{z}, \vec{z}') = \exists \vec{x}. \exists \vec{x}'. H_m(\vec{x}, \vec{z}) \wedge H_p(\vec{x}', \vec{z}') \wedge \mathcal{I}_{\leq}^{\alpha, \beta}(\vec{x}, \vec{x}')$. ◀

403 To conclude our proof, we must show Lemma 3.10. Towards this aim, we now recall some
 404 tools used in our approach.

405 **3.3 \mathbb{N} -weighted automata**

406 The classical notion of weighted automaton over a semiring is presented for instance in the
 407 reference book [45, Chapter III]. We consider only the case where the weights are in \mathbb{N} , and
 408 follow the notations of [15, §5.2].

409 ► **Definition 3.11.** *An \mathbb{N} -automaton \mathcal{A} over the alphabet Γ consists of:*

- 410 ■ *a finite set Q of states;*
- 411 ■ *an initial row vector $\mathbf{i} \in \mathbb{N}^Q$ and final column vector $\mathbf{f} \in \mathbb{N}^Q$;*
- 412 ■ *a monoid morphism $\mu: \Gamma^* \rightarrow \mathbb{N}^{Q \times Q}$ (with matrix multiplication).*
- 413 *It computes the function $w \in \Gamma^* \mapsto \mathcal{A}(w) = \mathbf{i} \cdot \mu(w) \cdot \mathbf{f} \in \mathbb{N}$.*

414 ▷ **Claim 3.12.** *The usual notion of deterministic finite automaton (DFA) is equivalent to an*
 415 *\mathbb{N} -automaton such that:*

- 416 ■ *\mathbf{i} has a single coordinate equal to 1, and the others are equal to 0;*
- 417 ■ *the same property applies to every row of every matrix $\mu(c)$ for $c \in \Gamma$;*
- 418 ■ *$\mathbf{f} \in \{0, 1\}^{\mathbb{N}}$.*

419 This implies $\mathcal{A}(w) \in \{0, 1\}$ for every $w \in \Gamma^*$, so \mathcal{A} seen as an \mathbb{N} -automaton computes the
 420 characteristic function of the regular language recognized by \mathcal{A} seen as an usual DFA.

421 Throughout the rest of this section we shall therefore identify DFA with \mathbb{N} -automata
 422 of this shape. We will also use \mathbb{N} -automata that take unbounded – indeed, polynomially
 423 growing – values. Their main property of interest for us is:

424 ► **Lemma 3.13** ([15, §5.2]). *Let $\mathcal{A} = (Q, \mathbf{i}, \mathbf{f}, \mu)$ be an \mathbb{N} -automaton over Γ . Suppose that:*

- 425 ■ *\mathcal{A} is trim, i.e. $\forall q \in Q, \exists u, v \in \Gamma^*: (\mathbf{i} \cdot \mu(u))(q) \geq 1 \wedge (\mu(v) \cdot \mathbf{f})(q) \geq 1$;*
- 426 ■ *$\mathcal{A}(w)$ grows at most polynomially¹⁴ in $|w|$.*

427 *Then one can compute from \mathcal{A} the smallest $k \in \mathbb{N}$ such that $\mathcal{A}(w) = O(|w|^k)$. Furthermore,*
 428 *there is a computable partition $Q = Q_0 \sqcup \dots \sqcup Q_k$ of the states such that:*

- 429 ■ *$\mu(w)(q, q') = 0$ for any $i > j, q \in Q_i, q' \in Q_j$ and $w \in \Gamma^*$;*
- 430 ■ *$\mu(w)(q, q') = O(1)$ for any fixed $i \in \{0, \dots, k\}$ and $q, q' \in Q_i$ – the bound depends on \mathcal{A} ,*
 431 *and is computable.*

432 ► **Remark 3.14.** *This is very similar to some results in the literature on polynomially ambiguous*
 433 *weighted automata (over arbitrary semirings) such as [35, Theorem 6.2]. But the extensional*
 434 *notion of growth of the function computed by an \mathbb{N} -automaton fits better our intended*
 435 *application than the intensional degree of ambiguity of automata.*

436 We also note that polynomially ambiguous *tree* automata admit a decomposition into
 437 simpler analyzable parts [41, Theorems 9 and 12] which “morally” seems to entails an
 438 extension of Bojańczyk’s “Seed Lemma” [5, Lemma 6.2] to trees (though formally we could
 439 not deduce the latter from the precise statement given for the former). Thus, we expect that
 440 MSO tree-to-tree interpretations should also admit a dimension minimization theorem. It
 441 is unclear whether dimension minimization for trees could be proved using the techniques
 442 of [4], since factorization forests do not work on trees.

¹⁴ Actually, subexponential growth implies polynomial growth for \mathbb{N} -automata, but we will not need this.

3.4 Proof of Lemma 3.10

First recall that the celebrated correspondence between first-order logic and star-free languages (cf. [48]) also extends to formulas with free variables in the following way, that provides a useful “logic-free” perspective on queries defined by first-order formulas:

- **Theorem 3.15.** *Let $\ell \in \mathbb{N}$ and $\Xi \subseteq \bigcup_{w \in \Gamma^*} \{w\} \times \{1, \dots, |w|\}^\ell$. The following are equivalent:*
- *there is an FO formula $F(x_1, \dots, x_\ell)$ such that $(w, i_1, \dots, i_\ell) \in \Xi$ iff $w \models F(i_1, \dots, i_\ell)$;*
 - *there is a star-free language $L \subseteq (\Gamma \times \{0, 1\}^\ell)^*$ such that $(w, i_1, \dots, i_\ell) \in \Xi$ iff L contains the word $\text{marked}(w, i_1, \dots, i_\ell)$ of length $|w|$ whose j -th letter*
 - *has $w[j]$ as its first coordinate;*
 - *has 1 as its $(m+1)$ -th coordinate if $i_m = j$, and 0 otherwise, for $m \in \{1, \dots, \ell\}$;*
 - *the “language of marked words” $\{\text{marked}(\xi) \mid \xi \in \Xi\}$ is star-free.*

Let $F(x_1, \dots, x_\ell)$ be an FO formula over Γ^* , and let $\mathcal{A} = (Q, \mathbf{i}, \mathbf{f}, \mu)$ be the *minimal* DFA over $(\Gamma \times \{0, 1\}^\ell)^*$ (seen as an \mathbb{N} -automaton) computing its language of marked words. We define another \mathbb{N} -automaton $\widehat{\mathcal{A}} = (Q, \mathbf{i}, \mathbf{f}, \widehat{\mu})$ over Γ^* by taking for each $c \in \Gamma$:

$$\widehat{\mu}(c) = \sum_{\substack{\vec{b} \in \{0, 1\}^\ell \\ \text{a 1-letter word}}} \mu\left(\underbrace{(c, \vec{b})}_{\text{a 1-letter word}}\right) \quad \text{so that} \quad \forall w \in \Gamma^*, \widehat{\mathcal{A}}(w) = |\{(i_1, \dots, i_\ell) \mid w \models F(i_1, \dots, i_\ell)\}|$$

Since \mathcal{A} is minimal, it is trim, so $\widehat{\mathcal{A}}$ is also trim. Apply Lemma 3.13 to $\widehat{\mathcal{A}}$ to get $k \in \mathbb{N}$ – the degree of growth of $\widehat{\mathcal{A}}$ and therefore of F – and a partition $Q = Q_0 \sqcup \dots \sqcup Q_k$. The fact that $\widehat{\mathcal{A}}$ is trim and the conditions on the partition imply that the initial state of \mathcal{A} is in Q_0 .

We use this partition to define a new DFA $\widetilde{\mathcal{A}} = (Q, \mathbf{i}, \mathbf{f}, \widetilde{\mu})$ over $(\Gamma \times \{0, 1\}^{\ell+k})^*$: for $c \in \Gamma$, $\vec{b} \in \{0, 1\}^\ell$, $\vec{b}' \in \{0, 1\}^k$, $i, j \in \{0, \dots, k\}$, $q \in Q_i$ and $q' \in Q_j$,

$$\widetilde{\mu}\left(\underbrace{(c, \vec{b}, \vec{b}')}_{\text{a 1-letter word}}\right)(q, q') = \begin{cases} \mu\left(\underbrace{(c, \vec{b})}_{\text{a 1-letter word}}\right)(q, q') & \text{if } \vec{b}' = (0, \dots, 0, \underbrace{1, \dots, 1}_{\text{from the } (i+1)\text{-th to the } j\text{-th coordinate}}, 0, \dots, 0) \\ & \text{so } \vec{b}' = \vec{0} \text{ if } i=j \\ 0 & \text{otherwise} \end{cases}$$

$\widetilde{\mathcal{A}}$ is a counter-free DFA. Indeed, for any $w \notin (\Gamma \times \{0, 1\}^\ell \times \{0\}^k)^*$, the matrix $\widetilde{\mu}(ww)$ is zero (due to the first item in the conclusions of Lemma 3.13), so checking that $\widetilde{\mathcal{A}}$ is counter-free reduces to the fact that \mathcal{A} is counter-free – which is the case since \mathcal{A} is the minimal DFA of a star-free language. So $\widetilde{\mathcal{A}}$ recognizes a star-free language. Because (again) of the first item of Lemma 3.13, this language is of the form $\{\text{marked}(\xi) \mid \xi \in \widetilde{\Xi}\}$ where $\widetilde{\Xi}$ contains tuples of the form $(w, i_1, \dots, i_\ell, j_1, \dots, j_m)$ for $m \in \{0, \dots, k\}$ (the corresponding marked word reaches a final state in Q_m when read by $\widetilde{\mathcal{A}}$). Let

$$\widetilde{\Xi} = \left\{ (w, i_1, \dots, i_\ell, \underbrace{j_1, \dots, j_m}_{k-m \text{ times}}, \underbrace{1, \dots, 1}_{k-m \text{ times}}) \mid (w, i_1, \dots, j_m) \in \Xi \right\} \subseteq \bigcup_{w \in \Gamma^*} \{w\} \times \{1, \dots, |w|\}^{\ell+k}$$

Then the language of marked words for $\widetilde{\Xi}$ is still star-free, so $\widetilde{\Xi}$ corresponds to some first-order formula $G(x_1, \dots, x_\ell, z_1, \dots, z_k)$. Using the second item of Lemma 3.13, one can see that G verifies the conclusions of the Lemma 3.10 that we wanted to prove.

Tito — more details?

476 **4 Quadratic polyregular functions vs macro tree transducers**

477 Our goal now is to show the following:

- 478 ► **Theorem 4.1.** *For any $k \geq 1$, there exists a string-to-string function f_k such that:*
- 479 ■ *f_k is computed by a single-component 2-dimensional first-order interpretation (cf. §3.1);*
 - 480 ■ *for any k -tuple of functions (g_1, \dots, g_k) such that each g_i is computed by some macro*
 - 481 *tree transducer, $\text{Im}(g_1 \circ \dots \circ g_k) \neq \text{Im}(f_k)$.*

482 In the second item, the domain of g_i must be equal to the codomain of g_{i+1} for the
 483 composition to be well-defined; and to make the comparison of output languages meaningful,
 484 the codomains of g_1 and f_k should be equal.

485 But g_1 outputs ranked trees, whereas f_1 outputs strings. To make sense of this, as usual,
 486 we identify Σ^* with the set of trees over the ranked alphabet that consists of a letter \hat{c} of
 487 rank 1 for each $c \in \Sigma$, plus a single letter $\hat{\epsilon}$ of rank 0. Through this identification, the domain
 488 of g_k and f_k may also be equal; in that case we may conclude that $g_1 \circ \dots \circ g_k \neq f_k$.

489 The relevance of Theorem 4.1 to the question of pebble minimization comes from:

- 490 ► **Theorem 4.2** (Engelfriet & Maneth [25, item (2) of the abstract]). *Any tree-to-tree function*
 491 *computed by some k -pebble¹⁵ tree transducer can also be expressed as a k -fold composition of*
 492 *macro tree transducers. (But the converse is false.)*

493 A k -pebble string transducer is none other than a k -pebble tree transducer working
 494 on the encodings of strings as unary trees described above. Thus it follows directly that
 495 $f_k \notin \text{Pebble}_k$ (without having to recall the definition of Pebble_k from §2.2), and even that
 496 $\text{Im}(f_k) \neq \text{Im}(h)$ for any $h \in \text{Pebble}_k$.

497 On the other hand, any 2-dimensional FO interpretation is polyregular (Theorem 3.7)
 498 with quadratic growth (Remark 3.2). So the f_k are indeed counterexamples to pebble
 499 minimization.

500 ► **Remark 4.3.** Our main reason for using first-order interpretations is that they provide a nice
 501 notion of multidimensional origin semantics, which we can use in our inductive construction
 502 of the f_k . But we do not truly depend on the difficult translation of FO interpretations
 503 to pebble transducers [6] to refute pebble minimization: it would be possible to show by a
 504 simpler ad-hoc argument that each f_k is computable by some ℓ -pebble transducer.

505 We recall the required properties of macro tree transducers in Section 4.1, then we prove
 506 Theorem 4.1 in Section 4.2.

507 **4.1 Compositions of macro tree transducers (MTTs)**

508 We do not formally define MTTs here, but only recall useful facts for our purposes. In this
 509 paper, we only consider *total deterministic* MTTs. Let MTT^k be the class of tree-to-tree
 510 functions computed by some composition of $k \geq 1$ macro tree transducers.

511 ► **Remark 4.4.** MTT^k can be characterized equivalently as the class of functions computed
 512 by k -iterated pushdown transducers [27], or by “level- k ” tree transducers [28]. For any k ,
 513 the functions in MTT^k with linear growth are precisely the regular tree functions [22]¹⁶ –

¹⁵As explained in Footnote 1 of the introduction, the indexing convention for the pebble transducer hierarchy in [25] is off by one compared to ours.

¹⁶The paper [22] talks about compositions of tree-walking transducers (TWT), but MTT^1 is included in the class of functions obtained by composing 3 TWTs [25, Lemma 37].

514 combined with Theorem 4.2, this proves pebble minimization for polyregular functions of
515 linear growth,¹⁷ so our counterexamples with quadratic growth are in some sense minimal.

516 For a class of tree-to-tree functions \mathcal{C} , we also write $\text{SO}(\mathcal{C})$ for the subclass of functions
517 that output strings (via the identification described at the beginning of §4). The literature
518 on tree transducers often refers to the classes $\text{SO}(\text{MTT}^k)$ by equivalent descriptions involving
519 a *yield* operation¹⁸ that maps trees to strings:

520 \triangleright **Claim 4.5.** $\text{SO}(\text{MTT}^{k+1}) = \{\text{yield} \circ f \mid f \in \text{MTT}^k\}$ for any $k \geq 1$, whereas $\text{SO}(\text{MTT}^1)$
521 consists of the functions $\text{yield} \circ (\text{some top-down tree transducer})$.

522 *Explanation.* The case $k = 1$ of the first claim – which directly entails the general case – is a
523 consequence of the work of Engelfriet and Vogler [27, items (a) + (d) of the introduction].
524 To understand why requires a terminological clarification: while, according to old definitions,

525 X tree-to-string transducer = $\text{yield} \circ (X$ tree transducer) for $X \in \{\text{macro, top-down}\}$

526 (for the top-down case see e.g. [37, 46]), this does not apply to $X = \text{pushdown}^n$. Actually,
527 pushdown^n tree-to-string transducers = $\text{SO}(\text{pushdown}^n$ tree transducers) in [27].

528 The claim about $\text{SO}(\text{MTT}^1)$ is well-known. For instance it is stated by Maneth as follows:
529 “Note that macro tree transducers with monadic output alphabet are essentially the same as
530 top-down tree-to-string transducers” [37, end of §5]. \triangleleft

531 This allows us to rephrase a key “bridge theorem” by Engelfriet and Maneth [23] in a form
532 that suits us better. For a class of string-valued functions \mathcal{C}' , let $\text{Im}(\mathcal{C}') = \{\text{Im}(f) \mid f \in \mathcal{C}'\}$.

533 \blacktriangleright **Theorem 4.6** ([23, Theorem 18] + Claim 4.5). *Let $k \geq 1$ and L, L' be string languages. If*
534 *L' is d -complete (cf. below) for L and $L' \in \text{Im}(\text{SO}(\text{MTT}^{k+1}))$, then $L \in \text{Im}(\text{SO}(\text{MTT}^k))$.*

535 \blacktriangleright **Definition 4.7.** *Let $L \subseteq \Sigma^*$ and $L' \subseteq (\Sigma \cup \Delta)^*$ with $\Sigma \cap \Delta = \emptyset$. We say that L' is:*

- 536 \blacksquare δ -complete for L [23, §5] *when for every $u \in L$ there exist $w_0, \dots, w_n \in \Delta^*$ such that*
 - 537 \blacksquare *all the w_i for $i \in \{1, \dots, n-1\}$ are pairwise distinct words;*
 - 538 \blacksquare *$n = |u|$ and $w_0 u[1] w_1 \dots u[n] w_n \in L'$;*
- 539 \blacksquare d -complete for L [24, §4] *when it is δ -complete for L and “conversely”, by erasing the*
540 *letters from Δ in the words in L' one gets exactly L , i.e. $\varphi(L') = L$ where φ is the monoid*
541 *morphism mapping each letter in Σ to itself and Δ to ε .*

542 4.2 Proof of Theorem 4.1

543 Our strategy is a minor adaptation of Engelfriet and Maneth’s proof [24, §4] that the hierarchy
544 of output languages of k -pebble string transducers is strict. (See also [25, Theorem 41] for a
545 similar result on tree languages.) Writing $\text{Im}(\mathcal{I}) = \text{Im}(f)$ when the first-order interpretation
546 \mathcal{I} defines the string function f , we show that:

547 \blacktriangleright **Lemma 4.8.** *From any single-component 2D FO interpretation \mathcal{I} , one can build another*
548 *single-component 2D FO interpretation $\Psi(\mathcal{I})$ such that $\text{Im}(\Psi(\mathcal{I}))$ is d -complete for $\text{Im}(\mathcal{I})$.*

549 \blacktriangleright **Remark 4.9.** Our construction preserves the class of interpretations that verify the property
550 described in the proof of Corollary 3.6.

¹⁷This special case is also a consequence of dimension minimization for MSO interpretations.

¹⁸ $\text{yield}(t)$ is the word formed by listing the labels of the leaves of t in infix order.

16 Revisiting the growth of polyregular functions

551 Before proving the lemma, we explain how it leads to Theorem 4.1. For $k \geq 1$, let

$$552 \quad \mathcal{I}_k = \Psi^{k-1}(\text{the FO interpretation of Example 3.4})$$

553 \mathcal{I}_1 maps a^n to $(a^{n-1}b)^{n-1}$ for $n \geq 1$, so we have¹⁹ $\text{Im}(\mathcal{I}_1) \notin \text{Im}(\text{SO}(\text{MTT}^1))$ by immediate
554 application of an old result of Engelfriet [17, Theorem 3.16].²⁰ Together, Lemma 4.8 and
555 the contrapositive of Theorem 4.6 entail that $\text{Im}(\mathcal{I}_k) \notin \text{Im}(\text{SO}(\text{MTT}^k))$ for all k – which is
556 precisely what we want.

557 **Proof of Lemma 4.8.** The idea is inspired by the coding of atom-oblivious functions used
558 in the Deatomization Theorem of [4] (cf. Remark 2.2). Let us fix our set of *atoms* to be
559 $\mathbb{A} = \mathbb{N} \setminus \{0\}$. Let \mathcal{I} be a single-component 2D FO interpretation specifying a function
560 $f: \Gamma^* \rightarrow \Sigma^*$. The interpretation canonically defines a function $f^{(\mathbb{A})}: (\Gamma \times \mathbb{A})^* \rightarrow (\Sigma \times \mathbb{A}^2)^*$
561 in the following way: reusing the notation O_u associated to \mathcal{I} in Definition 3.1,

$$562 \quad f^{(\mathbb{A})} \left(\begin{array}{ccc} u_1 & \dots & u_n \\ a_1 & \dots & a_n \end{array} \right) = \begin{array}{ccc} v_1 & \dots & v_m \\ a_{i_1, j_1} & \dots & a_{i_m, j_m} \end{array} \quad \text{where} \quad \begin{array}{l} f(u) = v_1 \dots v_m \\ O_u = \{\langle i_1, j_1 \rangle \leq \dots \leq \langle i_m, j_m \rangle\} \end{array}$$

563 (this is morally a form of “origin semantics” as in §2.1). For instance, for the interpretation
564 of Example 3.4 defining $f_1: a^n \mapsto (a^{n-1}b)^{n-1}$, we would have

$$565 \quad f_1(aaa) = aabaab \quad f_1^{(\mathbb{A})} \left(\begin{array}{ccc} a & a & a \\ 3 & 1 & 2 \end{array} \right) = \begin{array}{ccc} a & a & b \\ 3, 3 & 3, 1 & 3, 2 \end{array} \quad \begin{array}{ccc} a & a & b \\ 1, 3 & 1, 1 & 1, 2 \end{array}$$

566 Next, choose $\bullet \notin \Gamma$ and $\square, \diamond \notin \Sigma$ with $\square \neq \diamond$. We define the codings

$$567 \quad \text{enc}_1: \begin{array}{c} c \\ n \end{array} \in \Gamma \times \mathbb{A} \mapsto c \bullet^n \in (\Gamma \cup \{\bullet\})^* \quad \text{enc}_2: \begin{array}{c} c' \\ p, q \end{array} \in \Sigma \times \mathbb{A}^2 \mapsto c' \square^p \diamond^q \in (\Sigma \cup \{\square, \diamond\})^*$$

568 which we extend to input strings as morphisms of free monoids.

569 \triangleright **Claim 4.10.** There is a function f' specified by some 2D FO interpretation $\Psi(\mathcal{I})$ such
570 that $f' \circ \text{enc}_1 = \text{enc}_2 \circ f^{(\mathbb{A})}$ and $f'(\bullet w) = f'(w)$ for every input word w .

571 For example, $\mathcal{I}_2 = \Psi(\text{Example 3.4})$ can be illustrated (similarly to Example 3.5) by

$$572 \quad \begin{array}{ccc} a & \bullet & \bullet & \bullet & a & \bullet & a & \bullet & \bullet \\ a & a & \diamond & \diamond & \diamond & a & \diamond & b & \diamond & \diamond \\ \bullet & \square & & & \square & & \square & & & \\ \bullet & \square & & & \square & & \square & & & \\ \bullet & \square & & & \square & & \square & & & \\ a & a & \diamond & \diamond & \diamond & a & \diamond & b & \diamond & \diamond \\ \bullet & \square & & & \square & & \square & & & \\ a & & & & & & & & & \\ \bullet & & & & & & & & & \\ \bullet & & & & & & & & & \end{array} \rightsquigarrow a \square \square \square \diamond \diamond \diamond a \square \square \square \diamond \diamond b \square \square \square \diamond \diamond a \square \diamond \diamond \diamond a \square \diamond b \square \diamond \diamond \\ = \text{enc}_2 \circ f_1^{(\mathbb{A})} \left(\begin{array}{ccc} a & a & a \\ 3 & 1 & 2 \end{array} \right)$$

573 **Proof.** Given a formula F over the relational signature for Γ^* , let F^{R} be the relativized formula
574 over $(\Gamma \cup \{\bullet\})$ where all quantifiers $\forall z.(\dots)$ and $\exists z.(\dots)$ are replaced by $\forall z. \neg \bullet(z) \Rightarrow (\dots)$
575 and $\exists z. \neg \bullet(z) \wedge (\dots)$ respectively. Let

$$576 \quad P(x, y) = y \leq x \wedge \neg \bullet(y) \wedge \forall z. (\neg(z \leq y) \wedge z \leq x) \Rightarrow \bullet(z)$$

577 We take the following definition for the interpretation $\Psi(\mathcal{I})$:

¹⁹This implies $(a^n \mapsto (a^{n-1}b)^{n-1}) \notin \text{MTT}^1$, a fact that was later reproved in [40, Theorem 8.1(a)].

²⁰The theorem says that some superclass of $\text{Im}(\text{SO}(\text{MTT}^1))$ does not contain any language of the form $\{(a^n b)^{f(n)} \mid n \in X\}$ where $X \subseteq \mathbb{N} \setminus \{0\}$ is infinite and $f: X \rightarrow \mathbb{N} \setminus \{0\}$ is injective.

- 578 ■ $\Psi(\mathcal{I})_c(x_1, x_2) = \neg \bullet(x_1) \wedge \neg \bullet(x_2) \wedge \mathcal{I}_c^R(x_1, x_2)$ for $c \in \Sigma$
- 579 ■ $\Psi(\mathcal{I})_{\square}(x_1, x_2) = \bullet(x_1) \wedge \neg \bullet(x_2) \wedge \exists y. P(y, x_2) \wedge \bigvee_{c \in \Sigma} \mathcal{I}_c^R(y, x_2)$
- 580 ■ $\Psi(\mathcal{I})_{\diamond}(x_1, x_2) = \neg \bullet(x_1) \wedge \bullet(x_2) \wedge \exists y. P(x_1, y) \wedge \bigvee_{c \in \Sigma} \mathcal{I}_c^R(x_1, y)$
- 581 ■ $\Psi(\mathcal{I})_{\leq}(x_1, x_2, y_1, y_2) = \exists \hat{x}_1, \hat{x}_2, \hat{y}_1, \hat{y}_2. \bigwedge_{i=1,2} P(x_i, \hat{x}_i) \wedge P(y_i, \hat{y}_i)$ and
- 582 ■ either $\neg \mathcal{I}_{\leq}^R(\hat{y}_1, \hat{y}_2, \hat{x}_1, \hat{x}_2)$
- 583 ■ or $(\hat{x}_1, \hat{x}_2) = (\hat{x}_2, \hat{y}_2)$ and $(x_2, x_1) \leq (y_2, y_1)$ lexicographically. ◁

584 Let f' be given by Claim 4.10. We must now check that $\text{Im}(f')$ is d-complete for $\text{Im}(f)$.

- 585 ■ The part concerning the erasing morphism is immediate.
- 586 ■ So what remains is to verify δ -completeness (Definition 4.7). Let $u \in \text{Im}(f)$, i.e. $u = f(v)$
- 587 for some $v \in \Gamma^*$, and $m = |v|$. Then $f'(v[1] \bullet v[2] \bullet \dots \bullet v[m] \bullet^m) \in \text{Im}(f')$ is of the required
- 588 form. (It is here that we use the fact that our FO interpretation is single-component to
- 589 ensure that the maximal factors of the output in $\square^* \diamond^*$ are pairwise distinct.) ◀

5 Polyregular word sequences

5.1 For-transducers in an abstract style

Tito — $\prod_{i=a}^b$ in non-commutative monoid

592 Recall the notations $\underline{\Sigma}$ and $w \not\downarrow i$ from the end of the introduction.

594 ▶ **Definition 5.1.** Let M be any monoid.

595 Define the higher-order functions $\text{for}^\uparrow, \text{for}^\downarrow: [(\Sigma \cup \underline{\Sigma})^* \rightarrow M] \rightarrow [\Sigma^* \rightarrow M]$ by

$$\begin{array}{ccc} \text{for}^\uparrow(f): \Sigma^* & \rightarrow & M \\ w & \mapsto & \prod_{i=1}^{|w|} f(w \not\downarrow i) \end{array} \quad \begin{array}{ccc} \text{for}^\downarrow(f): \Sigma^* & \rightarrow & M \\ w & \mapsto & \prod_{j=0}^{|w|-1} f(w \not\downarrow |w| - j) \end{array}$$

597 ▶ **Proposition 5.2.** The class of (string-to-string) polyregular functions can be characterized

598 as the smallest class containing constant functions, closed under concatenation, regular

599 conditionals and the for^\uparrow and for^\downarrow operators. For first-order polyregular functions, an

600 analogous statement holds by replacing “regular conditional” by “star-free conditional”.

601 **Proof idea.** It is clear that all such functions are polyregular. The converse corresponds to

602 converting the mutable finite state of a for-transducer [5, §1] into regular conditionals. ◀

603 Call the “product rank” of a polyregular function the minimum number of nestings of

604 for^\uparrow and for^\downarrow operators necessary to define it.

605 ▷ Claim 5.3 (looks obvious, but is it?). If f is polyregular and g is a concatenation of rational

606 functions, then the product rank of $f \circ g$ is at most that of f .

Tito — prove this

607 An example of concatenation of rational functions is given, for any fixed $n \in \mathbb{N}$ and

608 $\sigma \in \mathfrak{S}_n$, by $\tilde{\sigma}: a^{x_1} \# \dots \# a^{x_n} \mapsto a^{x_{\sigma(1)}} \# \dots \# a^{x_{\sigma(n)}}$.

18 Revisiting the growth of polyregular functions

610 ► **Definition 5.4.** A function $f: \mathbb{N}^n \rightarrow \Sigma^*$ is polyregular if there exists a polyregular function
 611 $\hat{f}: \{a, \#\}^* \rightarrow \Sigma^*$ such that $f(x_1, \dots, x_n) = \hat{f}(a^{x_1} \# \dots \# a^{x_n})$ for all $x_1, \dots, x_n \in \mathbb{N}$. We
 612 then say that \hat{h} represents f (or is a representative of f).

613 Our goal is to prove the following, which for $n = 1$ characterizes the polyregular functions
 614 with unary input.

615 ► **Theorem 5.5.** The polyregular functions $\mathbb{N}^n \rightarrow \Sigma^*$ for $n \geq 1$ can be characterized
 616 inductively as follows:

- 617 ■ Constant functions are polyregular.
- 618 ■ The concatenation of two polyregular functions (with the same arity) is polyregular.
- 619 ■ If $f: \mathbb{N}^n \rightarrow \Sigma^*$ is polyregular and $\sigma \in \mathfrak{S}_n$ then $f \circ \tilde{\sigma}$ is polyregular.
- 620 ■ Indexed product: if $f: \mathbb{N}^{n+1} \rightarrow \Sigma^*$ is polyregular, then so is

$$621 \quad (x_1, \dots, x_n) \mapsto \prod_{y=0}^{x_1-1} f(y, x_1 - 1 - y, x_2, \dots, x_n)$$

- 622 ■ Ultimately periodic combination: if g_i and h_j are polyregular, $k \geq 0$ and $p \geq 1$ then so is

$$623 \quad (x_1, \dots, x_n) \mapsto \begin{cases} g_{x_1}(x_2, \dots, x_n) & \text{when } x_1 < k \\ h_{(x_1-k) \bmod p}(\lfloor (x_1-k)/p \rfloor, x_2, \dots, x_n) & \text{when } x_1 \geq k \end{cases}$$

624 By restricting to $p = 1$ in the ultimately periodic combinations, we get a characterization of
 625 the FO-polyregular functions $\mathbb{N}^n \rightarrow \Sigma^*$.

626 **Proof.** The fact that all functions inductively generated this way are polyregular is routine.
 627 We therefore focus on the converse, by induction first on $r \in \mathbb{N}$ and then on Proposition 5.2
 628 restricted to the functions with representatives of product rank $\leq r$. One key case is to
 629 handle representatives defined by regular conditionals, which is done by nesting n ultimately
 630 periodic combinations (one for each variable) thanks to an idempotent power argument.
 631 Another key case is as follows.

632 Let $f: \{a, \#, \underline{a}, \#\} \rightarrow \Sigma^*$ be of product rank $\leq r$, and let \hat{g} be the function defined by
 633 the indexed product in Proposition 5.2. Let $g: \mathbb{N}^n \rightarrow \Sigma^*$ be represented by \hat{g} . We want to
 634 show that g fits our inductive definition.

635 Let $\chi_i(x_1, \dots, x_{i-1}, y, z, x_{i+1}, \dots, x_n) = a^{x_1} \# \dots \# a^{x_{i-1}} \# a^y \underline{a} a^z \# a^{x_{i+1}} \# \dots \# a^{x_n}$ and
 636 $\chi'_i(x_1, \dots, x_n) = a^{x_1} \# \dots \# a^{x_i} \underline{\#} a^{x_{i+1}} \# \dots \# a^{x_n}$. We then have, for $x_1, \dots, x_n \in \mathbb{N}$,

$$637 \quad g(x_1, \dots, x_n) = \left(\prod_{i=1}^{n-1} u_i \cdot v_i \right) \cdot u_n$$

638 where for each i , we have $v_i = f \circ \chi'_i(x_1, \dots, x_n)$ and, assuming $x_i \geq 1$,

$$639 \quad u_i = \prod_{y=0}^{x_i-1} f \circ \chi_i(\dots, x_{i-1}, y, x_i - 1 - y, x_{i+1}, \dots)$$

640 Since χ_i and χ'_i have rational representatives, according to Claim 5.3, $f \circ \chi_i$ and $f \circ \chi'_i$
 641 are represented by polyregular functions with product rank $\leq r$, allowing us to invoke the
 642 induction hypothesis, and we are then done. (Well, almost, because we need to handle the
 643 edge case $x_i = 0$ and get rid of the -1 to fit with our actual inductive definition – but this
 644 can be done using an ultimately periodic combination.) ◀

5.2 Minimizing the number of loops for polyregular sequences

Warning: Very rough and sketchy draft; we ignore the permutations of variables throughout for convenience

Cécilia — ↑

Now we will turn to showing a pebble minimization theorem for polyregular word sequences. Theorem 5.5 gives us an inductive characterization of such sequence as being closed under a number of operators that include ultimately periodic combination. Now we shall show that those operators can always be brought up at the toplevel.

► **Definition 5.6.** Call simple polyregular word matrices $\mathbb{N}^n \rightarrow \Sigma^*$ for $n \geq 1$ the smallest class of functions containing constant functions and closed under concatenation, permutation of input variables and indexed products.

► **Theorem 5.7.** Polyregular functions $\mathbb{N}^n \rightarrow \Sigma^*$ correspond exactly to ultimately periodic combinations of simple polyregular word matrices.

Proof. The proof goes by an induction over the syntax corresponding to the inductive description given in Theorem 5.5, observing beforehand a couple of facts

Now to set up the induction, call a *polyregular AST* (abstract syntax tree) a syntactic representation of a polyregular function as per the description given in Theorem 5.5. Given an AST T , assign a multiset of natural numbers as follows:

- the AST corresponding to constant functions get assigned the empty multiset.
- given two ASTs T_1 and T_2 , the multiset corresponding to the concatenation AST $T_1 \cdot T_2$ will be $\mu(T_1 \cdot T_2) = \mu(T_1)^\uparrow + \mu(T_2)^\uparrow$, where $X^\uparrow = \{i + 1 \mid i \in X\}$.
- given an AST T of an indexed product, if T' is the AST corresponding to the inner component, the resulting AST will be $\mu(T) = \mu(T')^\uparrow$
- the AST corresponding to a ultimately periodic combination $\text{upc}(x_1; T'_1, \dots, T'_{k-1}; T''_0, \dots, T''_{p-1})$ will have measure $\mu(T) = \max\{T'_i, T''_j \mid 0 < i < k, j \leq p\}$

We show that any AST associated to a multi-set containing a non-zero element can be rewritten to an AST associated to a strictly small multiset in the canonical well-founded order. Such rewrites are terminating and will give us the desired result. To do so, we need a few auxiliary lemmas.

► **Lemma 5.8.** If $f: \mathbb{N}^{n+1} \rightarrow \Sigma^*$ is a simple numerical polyregular function, for any $k \in \mathbb{N}$, the function $g: (x_2, \dots, x_n) \mapsto f(k, x_2, \dots, x_n)$ is also a simple polyregular word sequence. Further, if T_f is an AST representing f , we can build an AST T_g representing g such that $\mu(T_g) \leq \mu(T_f)$.

Proof. Turn indexed product involving the x_1 component into a finitary concatenation. ◀

► **Lemma 5.9.** If $f: \mathbb{N}^{n+1} \rightarrow \Sigma^*$ is a simple numerical polyregular function, for any $a, b \in \mathbb{N}$, the function $g: (x_1, x_2, \dots, x_n) \mapsto f(ax_1 + b, x_2, \dots, x_n)$ is also a simple polyregular word sequence. Further, if T_f is an AST representing f with $\mu(T_f)$ containing only 0s, we can build an AST T_g representing g such that $\mu(T_g) = \mu(T_f)$.

Proof. By induction over the input AST, unrolling a times each loop containing the variable x_1 as a bound and adding b initial executions. ◀

Now, let's carry on with the induction. Fix an AST and consider a critical maximally deep ultimately periodic combination therein; it is necessarily preceded by either a concatenation or an indexed product node.

687 ■ Suppose that we have a concatenation involving an ultimately periodic combination

$$688 \quad \text{upc}(x_1; f_0, f_1, \dots, f_{k-1}; g_0, \dots, g_{p-1}) \cdot h$$

689 where x_1 is the variable we are making the combination over, the f_i representing the
690 functions corresponding to the offset and the g_j the functions corresponding to the period
691 of the first component. For $i < k$, define h'_i so that the corresponding function is

$$692 \quad (x_2, \dots, x_n) \mapsto h(i, x_2, \dots, x_n)$$

693 and h''_j for $i < p$ by

$$694 \quad (x_1, x_2, \dots, x_n) \mapsto h(k + px_1, x_2, \dots, x_n)$$

695 By Lemma 5.8 and Lemma 5.9, those two functions can be represented by ASTs with
696 smaller measure than h . We can then represent the original function with the AST

$$697 \quad \text{upc}(x_1; f_0 h'_0, \dots, f_{k-1} h'_{k-1}; g_0 \cdot h''_0, \dots, g_{p-1} \cdot h''_{p-1})$$

698 which has strictly smaller measure.

699 ■ The case where the ultimately periodic combination as the second component of a
700 concatenation is handled in the same way.

701 ■ The key case is the one of the indexed product. There are two subcases

- 702 ■ Either the combination operates according to a variable that is not iterated over in the
703 indexed product; in such a case, the combination trivially commutes with the indexed
704 product.
- 705 ■ Otherwise, assume that the combination is according to the first variable created by
706 the indexed product, i.e., that the function F under consideration is

$$707 \quad (x_1, \dots, x_n) \mapsto \prod_{y=0}^{x_1-1} f(y, x_1 - 1 - y, \dots)$$

708 with $k \geq 0$, $p \geq 1$ such that $f(y, z, \dots) = g_y(z, \dots)$ for $y < k$ and $f(y, z, \dots) =$
709 $h_{(y-k) \bmod p}(\lfloor (y-k)/p \rfloor, z, \dots)$ with g_y and the h_i simple. We express it as a new
710 combination using offset $2k$ but keeping the period p using a sort of loop unrolling.
711 The idea For $x_1 < k$, we can consider the functions l_{x_1}

$$712 \quad l_{x_1} : (x_2, \dots, x_n) \mapsto \prod_{z=0}^{x_1-1} g_z(x_1 - 1 - z, \dots)$$

713 and for $k \geq x_1 < 2k$

$$714 \quad l_{x_1} : (x_2, \dots, x_n) \mapsto \left(\prod_{z=0}^{k-1} g_z(x_1 - 1 - z, \dots) \right) \cdot \left(\prod_{z=k}^{x_1} h_{(z-k) \bmod p}(\lfloor (z-k)/p \rfloor, \dots, x_n) \right)$$

715 which can all be represented by ASTs of smaller measure thanks to our lemmas. For
716 the period, we consider the functions m_j for $j < p$ defined as

$$717 \quad m_j : (x_1, \dots, x_n) \mapsto l_{k+j}(x_2, \dots, x_n) \cdot \prod_{y=0}^{x_1-1} \theta_j(y, x_1 - 1 - y, \dots)$$

718 where

$$719 \quad \theta_j : (y, z, x_2, \dots, x_n) \mapsto \prod_{k=j}^{p+j-1} h_{k \bmod p}(y, \dots, x_n)$$

720 All of those also correspond to AST of smaller measure, so overall the ultimately
721 periodic combination has strictly smaller measure than before.

722 ■ The case where the combination is according to the second variable created by the
723 indexed product is handled in a similar way. ◀

724 Once we have converted a polyregular word sequence into a combination of simple
725 polyregular matrices, the asymptotic growth is necessarily bounded by the growth of a
726 dominant simple polyregular matrix in the decomposition. Further, it is easy to normalize
727 an expression corresponding to a simple polyregular matrix by removing subexpressions
728 containing only empty outputs so that the asymptotic growth of the associated function is
729 $O(n^d)$ where d is the number of imbricated loops. To obtain the minimization theorem, it
730 suffices to show such combinations of simple polyregular matrices can be turned into d -pebble
731 transducers.

732 Cécilia — the paragraph above should be expanded somewhat.

733 Tito — mention that it entails pebble minimization in the usual sense

734 5.3 Polyregular integer sequences (unary output too)

735 We may identify polynomials with the functions they denote in the sequel: $A[X] \subset A^A$.
736 (Recall that when A is an infinite integral domain, this identification is injective.)

737 ► **Remark 5.10 (Related work).** The specialization of [40, Theorem 9.2] to unary outputs says
738 that a sequence $\mathbb{N} \rightarrow \mathbb{N}$ is polyblind of rank at most $k \in \mathbb{N}$ if and only if it is an ultimately
739 periodic combination of *polynomials* in $\mathbb{N}[X]$ of degree at most $k+1$. The regular case ($k=0$)
740 was first shown in [8].

741 This should be included in the *poly-rational sequences* $\mathbb{N} \rightarrow \mathbb{Q}$ introduced recently in [1]
742 (counter-intuitively those don't have polynomial growth). Poly-rational sequences form a
743 subclass of linear recursive sequences (LRS) over \mathbb{Q} (see [1, Theorem 1] for several equivalent
744 definitions), while the notion of “catenative recurrence” which generalizes LRS over \mathbb{N} to
745 non-unary inputs characterizes \mathbb{N} -rational series (i.e. HDTOL with unary output) [29].

746 For bigger classes of automata-theoretic sequences $\mathbb{N} \rightarrow \mathbb{N}$, see [31, 7]; and for polyregular
747 sequences with \mathbb{Z} -valued output see [9].

748 **Non-a-periodic case** We refer to Section 5 for the notion of ultimately periodic combination.

749 ► **Theorem 5.11.** *Polyregular, polyblind and marble (Definition 5.13) integer sequences all*
750 *correspond to ultimately periodic combination of integer polynomials.*

751 (This collapse of the three classes in the case $\mathbb{N} \rightarrow \mathbb{N}$ has been independently proved in [13,
752 §4] as a corollary of a much stronger result.)

753 By [40, Theorem 9.2], we know that every sequence of $\text{UPC}(\mathbb{N}[X])$ is polyblind, and a
754 fortiori polyregular. By [12, Corollary 17], we know that polyregular and marble sequences
755 coincide. So it suffices to show the following.

756 ► **Lemma 5.12.** *Every marble sequence is in $\text{UPC}(\mathbb{N}[X])$.*

757 To do so, let us employ the definition of marble sequence in terms of bimachines of [12].
 758 Restricted to the setting of input over unary alphabets, the definition is equivalent to the
 759 following.

760 ► **Definition 5.13.** Let $S \subseteq \mathbb{N}^{\mathbb{N}}$. A marble S -bimachine (also called marble bimachine with
 761 external functions in S) is a tuple $\mathcal{M} = (M, x, \mathfrak{F}, \lambda)$ with

- 762 ■ M is a finite monoid and $x \in M$
- 763 ■ \mathfrak{F} is a finite subset of S
- 764 ■ λ is a function $M \times M \rightarrow \mathfrak{F}$

765 The output function associated with \mathcal{M} (written abusively $n \mapsto \mathcal{M}(n)$) is defined as

$$766 \quad \mathcal{M}(n) = \sum_{i=0}^n \lambda(x^i, x^{n-i})(i)$$

767 Define inductively the notion of n -marble bimachine and the associated set of integer sequences
 768 Marble_n :

- 769 ■ A 0-marble bimachine is a bimachine with external functions which are constant.
- 770 ■ A $n + 1$ -marble bimachine is a marble Marble_n -bimachine.

771 Lemma 5.12 thus states that $\text{Marble}_n \subseteq \text{UPC}(\mathbb{N}[X])$. This is proven by induction over n .
 772 The base case is trivial, and the inductive case can be reduced to the following lemma.

773 ► **Lemma 5.14.** Every $\text{UPC}(\mathbb{N}[X])$ -bimachine computes a function in $\text{UPC}(\mathbb{N}[X])$.

774 **Proof.** Let $\mathcal{M} = (M, x, \mathfrak{F}, \lambda)$ be the bimachine under consideration. Let $m \in \mathbb{N} \setminus \{0\}$ be
 775 such that $x^m = x^{2m}$ and, for every $f \in \mathfrak{F}$, there exist polynomials $(P_i)_{i=0}^{m-1}$ such that
 776 $f((n+1)m+k) = P_k(n)$ for every $k < m$.

777 To conclude, it suffices to show that the functions $g_k: n \mapsto \mathcal{M}((n+2)m+k)$ are in
 778 $\text{UPC}(\mathbb{N}[X])$ for every $k < m$. We have

$$779 \quad g_k(n) = \sum_{i=0}^{(n+2)m+k} \lambda(x^i, x^{(n+2)m+k-i})(i)$$

$$780 \quad = \sum_{i=0}^{n+2} \sum_{j=0}^{m-1} \lambda(x^{im+j}, x^{(n+2)m+k-(im+j)})(im+j) + \sum_{i=(n+2)m}^{(n+2)m+k} \lambda(x^i, x^{(n+2)m+k-i})(i)$$

$$781$$

782 Exploiting idempotency of x^m , we thus have

$$783 \quad g_k(n) = \sum_{i=1}^n \sum_{j=0}^{m-1} \lambda(x^{m+j}, x^{2m+k-j})(im+j) + g_k(m-1)$$

$$784 \quad + \sum_{j=0}^{m+k} \lambda(x^{m+j}, x^{m+k-j})((n+1)m+j)$$

$$785$$

786 Define $Q_k(n)$ as $g_k(m-1) + \sum_{j=0}^{m+k} \lambda(x^{m+j}, x^{m+k-j})((n+1)m+j)$, so that we have

$$787 \quad g_k(n) = Q_k(n) + \sum_{i=1}^n \sum_{j=0}^{m-1} \lambda(x^{m+j}, x^{2m+k-j})(im+j)$$

$$788$$

789 Recall that we chose m so that $\lambda(x, y)((n + 1)m + k) = P_{x, y, k}(n)$ for some polynomial
 790 $P_{x, y, k} \in \mathbb{N}[X]$. Hence, we have that Q_k is an integer polynomial and that there are
 791 polynomials $P_{k, j} \in \mathbb{N}[X]$ for $j < m$ such that

$$792 \quad g_k(n) = Q_k(n) + \sum_{i=1}^n \sum_{j=0}^{m-1} P_{k, j}(i)$$

794 Taking P_k to be the sum $\sum_{j=0}^{m-1} P_{k, j}$, we thus have $g_k(n) = Q_k(n) + \sum_{i=1}^n P_k(i)$.

795 Now we can conclude that $g_k \in \text{UPC}(\mathbb{N}[X])$ using the two following facts. The first follows
 796 from Faulhaber's formula for sums of powers.

797 \triangleright **Claim 5.15.** For every polynomial $P \in \mathbb{N}[X]$, there is $Q \in \mathbb{Q}[X]$ such that $\sum_{i=0}^n P(i) =$
 798 $Q(n)$.

799 The second is easy to check.

800 \triangleright **Claim 5.16.** $\{f: \mathbb{N} \rightarrow \mathbb{N} \mid \exists Q \in \mathbb{Q}[X]. \forall n \in \mathbb{N}. f(n) = Q(n)\} \subseteq \text{UPC}(\mathbb{N}[X])$

801 **Proof.** We can decompose the proof in two steps. First we show that polynomials $Q \in \mathbb{Q}[X]$
 802 that preserve positive integers are periodic combination of polynomials $P \in \mathbb{Z}[X]$ that
 803 preserve integers, and then than any such integer polynomial composed with $X + N$ for a
 804 sufficiently large N yields a polynomial of $\mathbb{N}[X]$. We can then conclude using the fact that
 805 UPC is a closure operator.

806 First, assume that $Q \in \mathbb{Q}[X]$ with $Q(n) \in \mathbb{N}$ for every $n \in \mathbb{N}$. Then, take the least
 807 common multiple of the denominators of the coefficients of Q and call it m . For $k < m$,
 808 consider $Q_k(X) = Q(mX + k)$. The non-constant coefficients can be seen to be integers by
 809 expanding the expressions. The constant coefficient of Q_k is equal to $Q(k) \in \mathbb{N}$. Clearly, all
 810 the Q_k preserve positive integers, so we may conclude that $Q \in \text{UPC}(\{f \mid \exists P \in \mathbb{Z}[X]. \forall n \in$
 811 $\mathbb{N}. f(n) = P(n)\})$.

812 Now suppose that $P \in \mathbb{Z}[X] \setminus \{0\}$ and that it preserves positive integers (the case
 813 $P = 0$ is trivial as $0 \in \mathbb{N}[X]$). Let d be the degree of P and $P(X) = \sum_{i=0}^d p_i X^i$. Take
 814 $N = d! \max_{i \leq d} |p_i|$ and consider the polynomial $R(X) = P(X + N)$. R and P have the same
 815 leading coefficient p_d , which is necessarily positive because $\lim_{n \rightarrow +\infty} P(n) > 0$. The k th
 816 coefficient of R is $\sum_{i=k}^d \binom{d}{i} N^{i-k} p_i$, which is positive because

$$817 \quad \binom{d}{k} N^{d-k} p_d \geq N^{d-k} \geq N^{d-k} (d-1-k) \binom{d-1}{k} \max_{i \leq d} |p_i| \geq N^{d-k-1} \max_{i \leq d} |p_i| \sum_{i=k}^{d-1} \binom{i}{k}$$

$$818 \quad \geq \sum_{i=k}^{d-1} \binom{i}{k} N^{i-k} p_i$$

820 Whence $P \in \text{UPC}(\mathbb{N}[X])$. \triangleleft

821 By combining the g_k for $k < m$, we have $\mathcal{M} \in \text{UPC}(\text{UPC}(\mathbb{N}[X])) = \text{UPC}(\mathbb{N}[X])$. \blacktriangleleft

822 First-order case

823 \blacktriangleright **Theorem 5.17.** A function $f: \mathbb{N} \rightarrow \mathbb{N}$ is FO polyblind of rank $\leq k$ if and only if there
 824 exists $P \in \mathbb{Z}[X]$ of degree $\leq k + 1$ such that for all $n \in \mathbb{N}$ except finitely many, we have
 825 $f(n) = P(n)$.

826 **Proof.** By induction on the rank. ◀

827 ▶ **Proposition 5.18** (useful for the “if” direction above). *Let $P \in \mathbb{Z}[X]$ and assume that its*
828 *leading coefficient is positive. Then $P(X + n) \in \mathbb{N}[X]$ for all $n \in \mathbb{N}$ except finitely many.*

829 **Proof.** Write the Taylor formula for $P(X + n)$; the k -th derivatives of P for $k \leq \deg(P)$
830 have a positive leading coefficient and therefore take positive values for large enough n . ◀

831 ▶ **Corollary 5.19** (A surprising phenomenon). *There exists a function $\mathbb{N} \rightarrow \mathbb{N}$ which is both*
832 *polyblind and first-order polyregular, but not FO polyblind.*

833 **Proof.** Take $n \mapsto n(n + 1)/2$. ◀

834 ▶ **Theorem 5.20.** *For unary outputs, FO-marble = FO-polyregular.*

835 **Proof.** A sum of n FO queries can be translated into the sum of $n \times k!$ FO k -marble
836 bimachines, each computing the part of the output produced when the variables used in the
837 query are in a certain ordering relative to one another. (This is a different construction from
838 the one used in [12] for the non-FO case which uses streaming string transducers. Note also
839 that it means converting an logical interpretation to a pebble transducer is trivial in the
840 unary output case.) ◀

841 ▶ **Theorem 5.21.** *A function $f: \mathbb{N} \rightarrow \mathbb{N}$ is FO-polyregular if and only if there exists $P \in \mathbb{Q}[X]$*
842 *such that for almost all $n \in \mathbb{N}$, we have $f(n) = P(n)$.*

843 **Proof.** Note that this forces P to be integer-valued.

844 **“if”** Use Newton’s interpolation formula

$$845 \quad P(X + a) = \sum_{i=0}^{\deg P} \Delta^i(P)(a) \binom{X}{i} \quad \text{where} \quad \Delta(P) = P(X + 1) - P(X)$$

846 and note that just like the k -th derivative, $\Delta^k(P)$ also takes positive values for large enough
847 arguments. Computing binomial coefficients with FO-marble transducers is trivial.

848 **“only if”** Similar argument to Lemma 5.12, using the previous theorem. ◀

849 — References —

- 850 1 Corentin Barloy, Nathanaël Fijalkow, Nathan Lhote, and Filip Mazowiecki. A robust class of
851 linear recurrence sequences. *Information and Computation*, 2022. doi:10.1016/j.ic.2022.
852 104964.
- 853 2 Mikołaj Bojańczyk. Transducers with origin information. In Javier Esparza, Pierre Fraigniaud,
854 Thore Husfeldt, and Elias Koutsoupias, editors, *Automata, Languages, and Programming - 41st*
855 *International Colloquium, ICALP 2014, Copenhagen, Denmark, July 8-11, 2014, Proceedings,*
856 *Part II*, volume 8573 of *Lecture Notes in Computer Science*, pages 26–37. Springer, 2014.
857 doi:10.1007/978-3-662-43951-7_3.
- 858 3 Mikołaj Bojańczyk. Polyregular functions, 2018. arXiv:1810.08760.
- 859 4 Mikołaj Bojańczyk. On the growth rate of polyregular functions, 2022. arXiv:2212.11631.
- 860 5 Mikołaj Bojańczyk. Transducers of polynomial growth. In Christel Baier and Dana Fisman,
861 editors, *LICS '22: 37th Annual ACM/IEEE Symposium on Logic in Computer Science, Haifa,*
862 *Israel, August 2 - 5, 2022*, pages 1:1–1:27. ACM, 2022. doi:10.1145/3531130.3533326.

- 863 6 Mikołaj Bojańczyk, Sandra Kiefer, and Nathan Lhote. String-to-string interpretations with
864 polynomial-size output. In *46th International Colloquium on Automata, Languages, and*
865 *Programming, ICALP 2019, July 9-12, 2019, Patras, Greece*, pages 106:1–106:14, 2019.
866 doi:10.4230/LIPIcs.ICALP.2019.106.
- 867 7 Michaël Cadilhac, Filip Mazowiecki, Charles Paperman, Michał Pilipczuk, and Géraud Séniz-
868 ergues. On polynomial recursive sequences. *Theory of Computing Systems*, June 2021.
869 doi:10.1007/s00224-021-10046-9.
- 870 8 Christian Hoffrut and Bruno Guillon. An algebraic characterization of unary two-way
871 transducers. In Erzsébet Csuhaj-Varjú, Martin Dietzfelbinger, and Zoltán Ésik, editors,
872 *Mathematical Foundations of Computer Science 2014 - 39th International Symposium, MFCS*
873 *2014, Budapest, Hungary, August 25-29, 2014. Proceedings, Part I*, volume 8634 of *Lecture*
874 *Notes in Computer Science*, pages 196–207. Springer, 2014. doi:10.1007/978-3-662-44522-8_
875 17.
- 876 9 Thomas Colcombet, Gaëtan Douéneau-Tabot, and Aliaume Lopez. Z-polyregular functions,
877 2022. arXiv:2207.07450.
- 878 10 Bruno Courcelle and Joost Engelfriet. *Graph structure and monadic second-order logic. A*
879 *language-theoretic approach*. Encyclopedia of Mathematics and its applications, Vol. 138. Cam-
880 bridge University Press, June 2012. Collection Encyclopedia of Mathematics and Applications,
881 Vol. 138. URL: <https://hal.archives-ouvertes.fr/hal-00646514>.
- 882 11 Patrice Ossona de Mendez. First-order transductions of graphs (invited talk). In Markus
883 Bläser and Benjamin Monmege, editors, *38th International Symposium on Theoretical Aspects*
884 *of Computer Science, STACS 2021, March 16-19, 2021, Saarbrücken, Germany (Virtual*
885 *Conference)*, volume 187 of *LIPIcs*, pages 2:1–2:7. Schloss Dagstuhl - Leibniz-Zentrum für
886 Informatik, 2021. doi:10.4230/LIPIcs.STACS.2021.2.
- 887 12 Gaëtan Douéneau-Tabot. Pebble Transducers with Unary Output. In Filippo Bonchi and
888 Simon J. Puglisi, editors, *46th International Symposium on Mathematical Foundations of*
889 *Computer Science (MFCS 2021)*, volume 202 of *Leibniz International Proceedings in Informatics*
890 *(LIPIcs)*, pages 40:1–40:17, Dagstuhl, Germany, 2021. Schloss Dagstuhl – Leibniz-Zentrum für
891 Informatik. doi:10.4230/LIPIcs.MFCS.2021.40.
- 892 13 Gaëtan Douéneau-Tabot. Hiding Pebbles When the Output Alphabet Is Unary. In Mikołaj
893 Bojańczyk, Emanuela Merelli, and David P. Woodruff, editors, *49th International Colloquium*
894 *on Automata, Languages, and Programming (ICALP 2022)*, volume 229 of *Leibniz International*
895 *Proceedings in Informatics (LIPIcs)*, pages 120:1–120:17, Dagstuhl, Germany, 2022. Schloss
896 Dagstuhl – Leibniz-Zentrum für Informatik. doi:10.4230/LIPIcs.ICALP.2022.120.
- 897 14 Gaëtan Douéneau-Tabot. Pebble minimization: the last theorems, 2022. arXiv:2210.02426.
- 898 15 Gaëtan Douéneau-Tabot, Emmanuel Filiot, and Paul Gastin. Register Transducers Are Marble
899 Transducers. In Javier Esparza and Daniel Král, editors, *45th International Symposium*
900 *on Mathematical Foundations of Computer Science (MFCS 2020)*, volume 170 of *Leibniz*
901 *International Proceedings in Informatics (LIPIcs)*, pages 29:1–29:14, Dagstuhl, Germany, 2020.
902 Schloss Dagstuhl–Leibniz-Zentrum für Informatik. doi:10.4230/LIPIcs.MFCS.2020.29.
- 903 16 Roger W. Ehrich and Stephen S. Yau. Two-way sequential transductions and stack automata.
904 *Information and Control*, 18(5):404–446, 1971. doi:10.1016/S0019-9958(71)90483-9.
- 905 17 Joost Engelfriet. Three hierarchies of transducers. *Mathematical Systems Theory*, 15(2):95–125,
906 1982. doi:10.1007/BF01786975.
- 907 18 Joost Engelfriet. Two-way pebble transducers for partial functions and their composition.
908 *Acta Informatica*, 52(7-8):559–571, 2015. doi:10.1007/s00236-015-0224-3.
- 909 19 Joost Engelfriet and Linda Heyker. The string generating power of context-free hypergraph
910 grammars. *Journal of Computer and System Sciences*, 43(2):328–360, 1991. doi:10.1016/
911 0022-0000(91)90018-Z.
- 912 20 Joost Engelfriet and Hendrik Jan Hoogeboom. MSO definable string transductions and
913 two-way finite-state transducers. *ACM Transactions on Computational Logic*, 2(2):216–254,
914 April 2001. doi:10.1145/371316.371512.

- 915 21 Joost Engelfriet, Hendrik Jan Hoogeboom, and Bart Samwel. XML navigation and transforma-
 916 tion by tree-walking automata and transducers with visible and invisible pebbles. *Theoretical*
 917 *Computer Science*, 850:40–97, January 2021. doi:10.1016/j.tcs.2020.10.030.
- 918 22 Joost Engelfriet, Kazuhiro Inaba, and Sebastian Maneth. Linear-bounded composition of tree-
 919 walking tree transducers: linear size increase and complexity. *Acta Informatica*, 58(1-2):95–152,
 920 2021. doi:10.1007/s00236-019-00360-8.
- 921 23 Joost Engelfriet and Sebastian Maneth. Output string languages of compositions of determin-
 922 istic macro tree transducers. *Journal of Computer and System Sciences*, 64(2):350–395, 2002.
 923 doi:10.1006/jcss.2001.1816.
- 924 24 Joost Engelfriet and Sebastian Maneth. Two-way finite state transducers with nested pebbles.
 925 In Krzysztof Diks and Wojciech Rytter, editors, *Mathematical Foundations of Computer*
 926 *Science 2002, 27th International Symposium, MFCS 2002, Warsaw, Poland, August 26-30,*
 927 *2002, Proceedings*, volume 2420 of *Lecture Notes in Computer Science*, pages 234–244. Springer,
 928 2002. doi:10.1007/3-540-45687-2_19.
- 929 25 Joost Engelfriet and Sebastian Maneth. A comparison of pebble tree transducers with macro
 930 tree transducers. *Acta Informatica*, 39(9):613–698, 2003. doi:10.1007/s00236-003-0120-0.
- 931 26 Joost Engelfriet and Heiko Vogler. Macro tree transducers. *Journal of Computer and System*
 932 *Sciences*, 31(1):71–146, 1985. doi:10.1016/0022-0000(85)90066-2.
- 933 27 Joost Engelfriet and Heiko Vogler. Pushdown machines for the macro tree transducer.
 934 *Theoretical Computer Science*, 42:251–368, 1986. doi:10.1016/0304-3975(86)90052-6.
- 935 28 Joost Engelfriet and Heiko Vogler. High level tree transducers and iterated pushdown tree
 936 transducers. *Acta Informatica*, 26(1/2):131–192, 1988. doi:10.1007/BF02915449.
- 937 29 Julien Ferté, Nathalie Marin, and Géraud Sénizergues. Word-Mappings of Level 2. *Theory of*
 938 *Computing Systems*, 54(1):111–148, January 2014. doi:10.1007/s00224-013-9489-5.
- 939 30 Emmanuel Filiot and Pierre-Alain Reynier. Copyful streaming string transducers. *Fundamenta*
 940 *Informaticae*, 178(1-2):59–76, January 2021. doi:10.3233/FI-2021-1998.
- 941 31 Séverine Fratani and Géraud Sénizergues. Iterated pushdown automata and sequences of
 942 rational numbers. *Annals of Pure and Applied Logic*, 141(3):363–411, September 2006. doi:
 943 10.1016/j.apal.2005.12.004.
- 944 32 Olivier Gauwin. *Transductions: resources and characterizations*. Habilitation à diriger des
 945 recherches, Université de Bordeaux, October 2020. URL: [https://tel.archives-ouvertes.](https://tel.archives-ouvertes.fr/tel-03118919)
 946 [fr/tel-03118919](https://tel.archives-ouvertes.fr/tel-03118919).
- 947 33 Liam Jordon. *An Investigation of Feasible Logical Depth and Complexity Measures via Automata*
 948 *and Compression Algorithms*. PhD thesis, National University of Ireland Maynooth, 2022.
 949 URL: <https://mural.maynoothuniversity.ie/16566/>.
- 950 34 Makoto Kanazawa, Gregory M. Kobele, Jens Michaelis, Sylvain Salvati, and Ryo Yoshinaka.
 951 The failure of the strong pumping lemma for multiple context-free languages. *Theory of*
 952 *Computing Systems*, 55(1):250–278, 2014. doi:10.1007/s00224-014-9534-z.
- 953 35 Stephan Kreutzer and Cristian Riveros. Quantitative monadic second-order logic. In *28th*
 954 *Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2013, New Orleans, LA,*
 955 *USA, June 25-28, 2013*, pages 113–122. IEEE Computer Society, 2013. doi:10.1109/LICS.
 956 2013.16.
- 957 36 Michel Latteux. Substitutions dans les EDT0L systèmes ultralinéaires. *Information and*
 958 *Control*, 42(2):194–260, 1979. doi:10.1016/S0019-9958(79)90641-7.
- 959 37 Sebastian Maneth. A survey on decidable equivalence problems for tree transducers. *In-*
 960 *ternational Journal of Foundations of Computer Science*, 26(8):1069–1100, 2015. doi:
 961 10.1142/S0129054115400134.
- 962 38 Tova Milo, Dan Suciu, and Victor Vianu. Typechecking for XML transformers. *Journal of*
 963 *Computer and System Sciences*, 66(1):66–97, 2003. Journal version of a PODS 2000 paper.
 964 doi:10.1016/S0022-0000(02)00030-2.
- 965 39 Anca Muscholl and Gabriele Puppis. The Many Facets of String Transducers. In Rolf
 966 Niedermeier and Christophe Paul, editors, *36th International Symposium on Theoretical*

- 967 *Aspects of Computer Science (STACS 2019)*, volume 126 of *Leibniz International Proceedings*
968 *in Informatics (LIPIcs)*, pages 2:1–2:21. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik,
969 2019. doi:10.4230/LIPIcs.STACS.2019.2.
- 970 40 Lê Thành Dũng Nguyễn, Camille Noûs, and Cécilia Pradic. Comparison-Free Polyregular
971 Functions. In Nikhil Bansal, Emanuela Merelli, and James Worrell, editors, *48th International*
972 *Colloquium on Automata, Languages, and Programming (ICALP 2021)*, volume 198 of *Leibniz*
973 *International Proceedings in Informatics (LIPIcs)*, pages 139:1–139:20. Schloss Dagstuhl –
974 Leibniz-Zentrum für Informatik, 2021. doi:10.4230/LIPIcs.ICALP.2021.139.
- 975 41 Erik Paul. On finite and polynomial ambiguity of weighted tree automata. In Srečko Brlek
976 and Christophe Reutenauer, editors, *Developments in Language Theory - 20th International*
977 *Conference, DLT 2016, Montréal, Canada, July 25-28, 2016, Proceedings*, volume 9840 of
978 *Lecture Notes in Computer Science*, pages 368–379. Springer, 2016. Full proofs available in
979 the diploma thesis https://www.informatik.uni-leipzig.de/~epaul/Diplom_Thesis.pdf
980 titled *Weighted Tree Automata and Quantitative Logics with a Focus on Ambiguity*. doi:
981 10.1007/978-3-662-53132-7_30.
- 982 42 Václav Rajlich. Absolutely parallel grammars and two-way finite state transducers. *Journal of*
983 *Computer and System Sciences*, 6(4):324–342, 1972. doi:10.1016/S0022-0000(72)80025-4.
- 984 43 Jonathan Rawski, Hossep Dolatian, Jeffrey Heinz, and Eric Raimy. Regular and polyregular
985 theories of reduplication. *Glossa: a journal of general linguistics*, 8(1), 2023. doi:10.16995/
986 glossa.8885.
- 987 44 Brigitte Rozoy. Outils et résultats pour les transducteurs boustrophédons. *RAIRO – Theoretical*
988 *Informatics and Applications*, 20(3):221–249, 1986. doi:10.1051/ita/1986200302211.
- 989 45 Jacques Sakarovitch. *Elements of Automata Theory*. Cambridge University Press, 2009.
990 Translated by Reuben Thomas. doi:10.1017/CB09781139195218.
- 991 46 Helmut Seidl, Sebastian Maneth, and Gregor Kemper. Equivalence of deterministic top-
992 down tree-to-string transducers is decidable. *Journal of the ACM*, 65(4):21:1–21:30, 2018.
993 doi:10.1145/3182653.
- 994 47 Tim Smith. A pumping lemma for two-way finite transducers. In Erzsébet Csuhaaj-Varjú,
995 Martin Dietzfelbinger, and Zoltán Ésik, editors, *Mathematical Foundations of Computer*
996 *Science 2014 - 39th International Symposium, MFCS 2014, Budapest, Hungary, August 25-29,*
997 *2014. Proceedings, Part I*, volume 8634 of *Lecture Notes in Computer Science*, pages 523–534.
998 Springer, 2014. doi:10.1007/978-3-662-44522-8_44.
- 999 48 Howard Straubing. First-order logic and aperiodic languages: a revisionist history. *ACM*
1000 *SIGLOG News*, 5(3):4–20, 2018. doi:10.1145/3242953.3242956.
- 1001 49 Géraud Sénizergues. Word-mappings of level 3, 2023. arXiv:2301.09966.