

A q -Matrix Encoding Extending the Parikh Matrix Mapping

Ömer Eğecioğlu

Abstract: We introduce a generalization of the Parikh mapping called the *Parikh q -matrix encoding*, which takes its values in matrices with polynomial entries. The encoding represents a word w over a k -letter alphabet as a $(k + 1)$ -dimensional upper-triangular matrix with entries that are nonnegative integral polynomials in variable q . Putting $q = 1$, we obtain the morphism introduced by Mateescu, Salomaa, Salomaa, and Yu [6] which extends the Parikh mapping to $(k + 1)$ -dimensional (numerical) matrices. The Parikh q -matrix encoding however, produces matrices that carry more information about w than the numerical Parikh matrix. In fact it is injective. The entries of the q -matrix image of w under this encoding is constructed by *q -counting* the number of occurrences of certain words as scattered subwords of w . This construction is distinct from the Parikh q -matrix mapping into k -dimensional upper-triangular matrices with integral polynomial entries introduced by Eğecioğlu and Ibarra [2].

Keywords: Parikh mapping, Parikh matrix mapping, scattered subword, q -analogue.

1 Introduction

By Parikh's theorem [7], the commutative image of any context-free language is a semilinear set, and is therefore also the commutative image of some regular set. Consider the (ordered) alphabet $\Sigma_k = \{a_1 < a_2 < \dots < a_k\}$ and for $w \in \Sigma^*$, define by $|w|_{a_i}$ the number of occurrences of a_i in w . The *Parikh mapping* is a morphism

$$\Psi : \Sigma^* \rightarrow \mathbb{N}^k$$

where \mathbb{N} is the set of nonnegative integers and $\Psi(w) = (|w|_{a_1}, |w|_{a_2}, \dots, |w|_{a_k})$ is the Parikh vector.

The Parikh mapping is an important concept in the theory of formal languages. Various languages accepted by automata that are more powerful than pushdown automata have been shown to have effectively computable semilinear sets. For example, it is known that every language accepted by a pushdown automaton augmented with reversal-bounded counters has a semilinear Parikh map [4]. The fact that the emptiness problem for semilinear sets is decidable implies that the emptiness problem for these automata is decidable. This fact has a number of applications in formal languages (e.g., [3]) and formal verification (e.g., [5]).

The *Parikh matrix mapping* introduced in [6] is a morphism $\Psi_{\mathcal{M}_k} : \Sigma^* \rightarrow \mathcal{M}_{k+1}$ where \mathcal{M}_{k+1} is a collection of $(k + 1)$ -dimensional upper-triangular matrices with nonnegative integral entries and unit diagonal. The classical Parikh vector $\Psi(w)$ appears in the image matrix as the second diagonal.

In [2] the *Parikh q -matrix mapping* was introduced. It is a morphism $\Psi_q^k : \Sigma^* \rightarrow \mathcal{M}_k(q)$ where $\mathcal{M}_k(q)$ is a collection of k -dimensional upper-triangular matrices with nonnegative integral polynomials in q as entries. The diagonal entries of $\Psi_q^k(w)$ are

$$(q^{|w|_{a_1}}, q^{|w|_{a_2}}, \dots, q^{|w|_{a_k}})$$

which readily encodes the Parikh vector. In addition, it can be shown [2] that viewing $w \in \Sigma_k$ as a word in Σ_{k+1} with $|w|_{a_{k+1}} = 0$, the Parikh q -matrix $\Psi_q^{k+1}(w)$ evaluated at $q = 1$ is precisely the $(k + 1)$ -dimensional numerical Parikh matrix $\Psi_{\mathcal{M}_k}(w)$.

It is a basic property of the Parikh matrix mapping that two words with the same Parikh matrix have the same Parikh vector, but two words with the same Parikh vector in many cases have different Parikh matrices [1]. Thus, the Parikh matrix gives more information about a word than a Parikh vector. Similarly, two words with the same Parikh q -matrix have the same Parikh matrix (and therefore the same Parikh vector), but there are cases in which two words with the same Parikh matrix have different q -matrices. Thus the Parikh q -matrix gives more information about a word than the Parikh matrix. We show that the Parikh q -matrix encoding introduced here is injective, and for $q = 1$, it also reduces to the $(k + 1)$ -dimensional numerical Parikh matrix mapping.

The basic idea in the construction of the entries of the Parikh q -matrix encoding of w is *q -counting* the number of occurrences of certain words as scattered subwords of w . This is done by keeping track of the indices (positions) of the letters of the scattered subword in w as exponents of monomials in q .

The structure of this paper is as follows. Sections 2 gives some basic notation and definitions. Section 3 recalls the notion of a Parikh matrix mapping introduced in [6] and the fundamental theorem concerning these mappings. Section 4 presents the polynomials which generalize the count of scattered subwords. Section 5 presents our Parikh q -matrix encoding and the main results, including Theorem 2, which gives the main properties of the encoding. Section 6 looks at properties such as injectivity and gives an interpretation of the inverse matrix in terms of mirror images.

2 Definitions

We start with some basic notation and definitions. Most of these are as they appear in references [6, 1, 2]. The set of all nonnegative integers is denoted by \mathbb{N} . We denote by $\mathbb{N}[q]$ the collection of polynomials in the variable q with coefficients from \mathbb{N} . \mathbb{Z} denotes integers, and $\mathbb{Z}[q]$ denotes the ring of polynomials in the variable q with integral coefficients. For an alphabet Σ , we denote the set of all words over Σ by Σ^* and the empty word by λ . We use “ordered” alphabets. An ordered alphabet is an alphabet $\Sigma = \{a_1, a_2, \dots, a_k\}$ with a relation of order (“ $<$ ”) on it. If for instance $a_1 < a_2 < \dots < a_k$, then we use the notation

$$\Sigma = \{a_1 < a_2 < \dots < a_k\}.$$

If $w \in \Sigma^*$ then $|w|$ denotes the length of w . For $a_i \in \Sigma$ and $w \in \Sigma^*$ the number of occurrences of the letter a_i in w is denoted by $|w|_{a_i}$.

Let v, w be words over Σ . As defined in [6], the word v is called a *scattered subword* of w if there exists a word u such that $w \in u \sqcup v$, where \sqcup denotes the shuffle operation. If $v, w \in \Sigma^*$, then the number of occurrences of v in w as a scattered subword is denoted by $|w|_{scatt-v}$. Partially overlapping occurrences of a word as a scattered subword of a word are counted as distinct occurrences. For example, $|acbb|_{scatt-ab} = 2$, $|acba|_{scatt-ab} = 1$.

Notation We also use the symbol $P_{w,v}$ to denote $|w|_{scatt-v}$. Using this notation, $P_{acbb,ab} = 2$, $P_{acba,ab} = 1$, and $P_{w,a_i} = |w|_{a_i}$ for any letter $a_i \in \Sigma$.

Notation Consider the ordered alphabet $\{a_1 < a_2 < \dots < a_k\}$ where $k \geq 1$. As in [6], we denote by the symbol $a_{i,j}$ the word $a_i a_{i+1} \dots a_j$ where $1 \leq i \leq j \leq k$.

For motivation and further issues about the Parikh mapping as well as language-theoretic notions not considered here, we refer the reader to [8].

3 Parikh matrix mapping

We first describe the extension of the Parikh mapping to matrices as originally defined in [6]. The extension involves special types of triangular matrices. These are square matrices $m = (m_{i,j})_{1 \leq i,j \leq k}$ such that $m_{i,j} \in \mathbb{N}$, for all $1 \leq i, j \leq k$, $m_{i,j} = 0$, for all $1 \leq j < i \leq k$, and moreover, $m_{i,i} = 1$, for all $1 \leq i \leq k$. The set of all these matrices of dimension k is denoted by \mathcal{M}_k . Thus \mathcal{M}_k is the collection $k \times k$ upper-triangular matrices with entries from \mathbb{N} and unit diagonal. The set \mathcal{M}_k is a monoid with respect to multiplication of matrices and has a unit which is the identity matrix I_k . The main notion introduced in [6] is as follows:

Definition 3. Let $\Sigma = \{a_1 < a_2 < \dots < a_k\}$ be an ordered alphabet, where $k \geq 1$. The Parikh matrix mapping, denoted by $\Psi_{\mathcal{M}_k}$, is the morphism:

$$\Psi_{\mathcal{M}_k} : \Sigma^* \rightarrow \mathcal{M}_{k+1},$$

defined as follows: If $\Psi_{\mathcal{M}_k}(a_l) = (m_{i,j})_{1 \leq i,j \leq (k+1)}$, then for each $1 \leq i \leq k+1$, $m_{i,i} = 1$, $m_{i,l+1} = 1$ and all other elements of the matrix $\Psi_{\mathcal{M}_k}(a_l)$ are zero.

Example 3.1. Let Σ be the ordered alphabet $\{a < b < c\}$. Then the Parikh matrix mapping $\Psi_{\mathcal{M}_3}$ represents each word over Σ^* as a 4×4 upper triangular matrix with unit diagonal with nonnegative integral entries. As an example $\Psi_{\mathcal{M}_3}(ab^2) = \Psi_{\mathcal{M}_3}(a)\Psi_{\mathcal{M}_3}(b)\Psi_{\mathcal{M}_3}(b)$. Thus

$$\Psi_{\mathcal{M}_3}(ab^2) = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 2 & 0 \\ 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Since the summation in the definition (5) of $P_{w,a_i,j}(q)$ is over all occurrences of a_i,j in w as a scattered subword, the following proposition is immediate:

Proposition 1. Let $\Sigma = \{a_1 < a_2 < \dots < a_k\}$ and $1 \leq i \leq j \leq k$. Then $P_{w,a_i,j}(1) = P_{w,a_i,j} (= |w|_{scatt-a_i,j})$.

Thus the polynomials $P_{w,a_i,j}(q)$ “ q -count” the number of occurrences of $a_i a_{i+1} \dots a_j$ as a scattered subword of w .

5 The Parikh q -matrix encoding

We denote by $\mathcal{M}_k(q)$ the collection of k -dimensional upper-triangular matrices with entries in $\mathbb{N}[q]$. Let I_k denote the identity matrix of dimension k . We define a mapping

$$\bar{\Psi} : \Sigma \times \mathbb{N} \rightarrow \mathcal{M}_{k+1}(q)$$

as follows. The matrix $\bar{\Psi}(a_l, j)$ corresponding to a pair $a_l \in \Sigma = \{a_1 < a_2 < \dots < a_k\}$ and $j \in \mathbb{N}$, is defined as the matrix obtained from I_{k+1} by changing the $(l, l+1)$ -st entry in I_{k+1} to q^j . Thus if $\bar{\Psi}(a_l, j) = (m_{i,j})_{1 \leq i, j \leq k+1}$, then $m_{i,i} = 1$ for $1 \leq i \leq k+1$, $m_{l,l+1} = q^j$, and all other entries of the matrix $\bar{\Psi}(a_l, j)$ are zero.

Example 5.1. When the alphabet is $\Sigma = \{a < b < c\}$,

$$\bar{\Psi}(a, j) = \begin{bmatrix} 1 & q^j & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \bar{\Psi}(b, j) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & q^j & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, \quad \bar{\Psi}(c, j) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & q^j \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Remark 2. The matrix $\bar{\Psi}(a_l, j)$ differs from the matrix $\Psi(a_l)$ introduced in [6] in the definition of the Parikh matrix mapping only in the $(l, l+1)$ -st entry: In $\Psi(a_l)$ this entry is 1, whereas in $\bar{\Psi}(a_l, j)$ it is q^j . We note that the matrices $\bar{\Psi}(a_l, j)$ specialize to $\Psi(a_l)$ for $q = 1$, for any value of j .

We now define a mapping (which we again denote by $\bar{\Psi}$) from Σ^* to $\mathcal{M}_{k+1}(q)$ by setting $\bar{\Psi}(\lambda) = I_{k+1}$ and

$$\bar{\Psi}(w_1 w_2 \dots w_n) = \bar{\Psi}(w_1, 1) \bar{\Psi}(w_2, 2) \dots \bar{\Psi}(w_n, n), \quad w_i \in \Sigma, 1 \leq i \leq n.$$

We refer to $\bar{\Psi}$ as the *Parikh q -matrix encoding*.

Example 5.2. Let Σ be the ordered alphabet $\{a < b < c\}$. Then $\bar{\Psi}(ab^2) = \bar{\Psi}(a, 1) \bar{\Psi}(b, 2) \bar{\Psi}(b, 3)$. Thus

$$\bar{\Psi}(ab^2) = \begin{bmatrix} 1 & q & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & q^2 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & q^3 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & q & q^3 + q^4 & 0 \\ 0 & 1 & q^2 + q^3 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Furthermore

$$\begin{aligned} \bar{\Psi}(ab^2ac^2a) &= \begin{bmatrix} 1 & q & q^3 + q^4 & 0 \\ 0 & 1 & q^2 + q^3 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \bar{\Psi}(a, 4) \bar{\Psi}(c, 5) \bar{\Psi}(c, 6) \bar{\Psi}(a, 7) \\ &= \begin{bmatrix} 1 & q + q^4 + q^7 & q^3 + q^4 & q^8 + 2q^9 + q^{10} \\ 0 & 1 & q^2 + q^3 & q^7 + 2q^8 + q^9 \\ 0 & 0 & 1 & q^5 + q^6 \\ 0 & 0 & 0 & 1 \end{bmatrix} \end{aligned}$$

Note that for $q = 1$, this is the matrix $\Psi_{\mathcal{M}_3}(ab^2ac^2a)$ displayed as (1).

Remark 3. The Parikh q -matrix encoding is not a morphism for arbitrary q . For example over $\Sigma = \{a < b\}$ $\bar{\Psi}(ab) = \bar{\Psi}(a, 1) \bar{\Psi}(b, 2) \neq \bar{\Psi}(a, 1) \bar{\Psi}(b, 1) = \bar{\Psi}(a) \bar{\Psi}(b)$

Remark 4. Neither the Parikh matrix mapping of [6] nor the Parikh q -matrix mapping of [2] is injective. We will later show that the Parikh q -matrix encoding $\overline{\Psi}$ is injective. On the other hand, both Parikh matrix mapping and the Parikh q -matrix mapping are morphisms, whereas Parikh q -matrix encoding is not.

Remark 5. Just as in the case of the Parikh matrix mapping and the Parikh q -matrix mapping, it is not true for the Parikh q -matrix encoding that if \mathcal{L} is a context-free language, then its image is some suitable extension of the notion of semilinearity to matrices over $\mathbb{N}[q]$. This is a direct consequence of Theorem 3 and the similar negative result concerning the Parikh matrix mapping ([6], Remark 3.2).

Now we give a characterization of the entries of the matrix $\overline{\Psi}(w)$ in terms of the polynomials $P_{w,a_i,j}(q)$ introduced in section 4.

Theorem 2. Let $\Sigma = \{a_1 < a_2 < \dots < a_k\}$ be an ordered alphabet, where $k \geq 1$ and assume that $w \in \Sigma^*$. The matrix $\overline{\Psi}(w) = (m_{i,j}(q))_{1 \leq i,j \leq k+1}$, has the following properties

1. $m_{i,j} = 0$, for all $1 \leq j < i \leq k+1$,
2. $m_{i,i} = 1$, for all $1 \leq i \leq k+1$,
3. $m_{i,j+1} = P_{w,a_i,j}(q)$, for all $1 \leq i \leq j < k+1$.

Proof. The proof of the parts 1. and 2. are immediate. We now prove property 3. Assume that $|w| = n$. The proof is by induction on n . If $n \leq 1$, the assertion holds. Assume now that part 3. holds for all words of length n and let w be of length $n+1$. Write $w = w'a_j$ where $|w'| = n$ and $a_j \in \Sigma$. Then

$$\overline{\Psi}(w) = \overline{\Psi}(w')\overline{\Psi}(a_j, n+1)$$

Assume that

$$\overline{\Psi}(w') = \begin{bmatrix} 1 & m'_{1,2} & \dots & \dots & m'_{1,k+1} \\ 0 & 1 & \dots & \dots & m'_{2,k+1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & 1 & m'_{k,k+1} \\ 0 & 0 & \dots & \dots & 1 \end{bmatrix} = M'$$

By the inductive hypothesis the matrix $M' = \overline{\Psi}(w')$ has property 3. We have

$$\overline{\Psi}(a_j, n+1) = \begin{bmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 1 & q^{n+1} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \dots & 1 \end{bmatrix}$$

where the matrix differs from I_{k+1} only in the entry in position $(j, j+1)$ where it is q^{n+1} instead of 0. Let $M = \overline{\Psi}(w)$. Then

$$M = \begin{bmatrix} 1 & m'_{1,2} & \dots & \dots & m'_{1,k+1} \\ 0 & 1 & \dots & \dots & m'_{2,k+1} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & 1 & m'_{k,k+1} \\ 0 & 0 & \dots & \dots & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 1 & q^{n+1} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \dots & 1 \end{bmatrix}$$

If $M = (m_{r,s})_{1 \leq r,s \leq k+1}$, then

$$m_{i,j+1} = m'_{i,j+1} + m'_{i,j} \cdot q^{n+1} \quad \text{for all } 1 \leq i \leq j, \quad (6)$$

and for all other indices, $m_{r,s} = m'_{r,s}$. But these are immediate from the definition of the polynomials $P_{w,a_i,j}(q)$ which satisfy

$$P_{w'a_j,a_i,j}(q) = P_{w',a_i,j}(q) + P_{w',a_i,j-1}(q) \cdot q^{n+1} \quad \text{for } 1 \leq i \leq j \leq k, \quad (7)$$

and the proof follows by induction. \square

Remark 7. The Parikh matrix mapping and the Parikh q -matrix mapping are non-injective morphisms, and the latter is a refinement of the former [2]. The Parikh q -matrix encoding introduced in the present paper is injective, and reduces to the (numerical) Parikh matrix mapping of [6]. It is possible to construct weaker and non-injective versions of the Parikh q -matrix encoding in a number of ways. One is to set $q^{p+1} = q$ for some $p \in \mathbb{N}$ and keep track of the positions of the letters in w modulo p only. The encoding is then with upper-triangular matrices with entries over the quotient ring $\mathbb{Z}[q]/\langle q^{p+1} - q \rangle$.

References

- [1] Atanasiu, A., Martin-Vide, C., Mateescu, A.: On the injectivity of the Parikh matrix mapping. *Fundamenta Informaticae*, 49 (2002) 289-299.
- [2] Egecioglu, Ö., and Ibarra, O.: A Matrix q -Analogue of the Parikh Map. UCSB Technical Report, TR 2004-06, February 2004.
- [3] Harju, T., Ibarra, O., Karhumaki, J., Salomaa, A.: Some decision problems concerning semilinearity and commutation. *J. Computer and System Sciences*, 65 (2002), 278-294.
- [4] Ibarra, O.: Reversal-bounded multicounter machines and their decision problems. *J. Assoc. Comput. Mach.*, 24 (1978) 123-137.
- [5] Ibarra, O., Su, J., Dang, Z., Bultan, T., Kemmerer, R.: Counter machines and verification problems. *Theoretical Computer Science*, 289 (2002) 165-189.
- [6] Mateescu, A., Salomaa, A., Salomaa K., Yu, S.: A sharpening of the Parikh mapping. *Theoretical Informatics and Applications*, 35 (2001) 551-564.
- [7] Parikh, R.J.: On context-free languages. *J. Assoc. Comput. Mach.*, **13** (1966) 570-581.
- [8] Rozenberg, G., Salomaa, A. (eds): *Handbook of Formal Languages*. Springer, Berlin, 1997.

Sabanci University
Faculty of Engineering and Natural Sciences
Tuzla, 34956 Istanbul, Turkey
E-mail: omer@cs.ucsb.edu

On sabbatical leave from
Department of Computer Science, University of California,
Santa Barbara, CA 93106, USA