

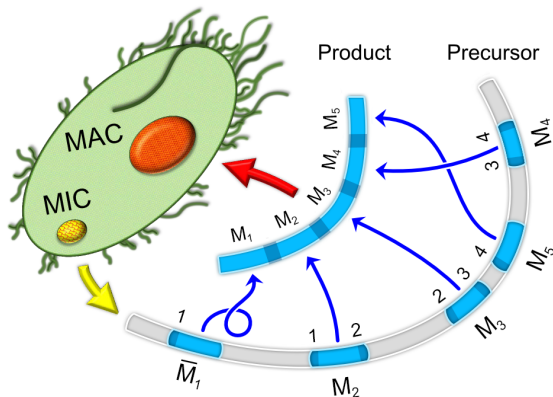
Transformations on Double Occurrence Words Motivated by DNA Rearrangement

Daniel A. Cruz, Margherita Maria Ferrari, Nataša Jonoska, Lukas
Nabergall, and Masahico Saito
University of South Florida

dcruz@mail.usf.edu
December 4, 2018



Motivation: Analysis of DNA Scrambling in Ciliates



$$\begin{array}{cccccc} \bar{M}_1 & M_2 & M_3 & M_5 & M_4 & \\ 1 & 12 & 23 & 4 & 34 &) \quad w = 1123434 \end{array}$$

Jonoska, N. et al. Patterns and Distances in Words Related to DNA Rearrangement. *Fundamenta Informaticae* **154**:1-4 (2017) pp 225-238.

Preliminaries

Given alphabet = $\{0, 1, \dots, 9\}$,

$w = 015164443$ is a word over

The length of w is 9, written $|w| = 9$

$w^R = 344461510$ is the reverse of w

The set of symbols used in w is $\Sigma = \{0, 1, 3, 4, 5, 6\}$

The set of all words over Σ is Σ^* and includes the empty word.

Preliminaries

Given alphabet = $\{0, 1, \dots, 9\}$,

$w = 015164443$ is a word over

The length of w is 9, written $|w| = 9$

$w^R = 344461510$ is the reverse of w

The set of symbols used in w is $[w] = \{0, 1, 3, 4, 5, 6\}$

The set of all words over $[w]$ is $[w]^*$ and includes the empty word.

Definition

The word w is a double occurrence word (DOW) if each symbol in w appears 0 or 2 times in w . The set of all DOWs is DOW .

$11; 1221; 11223434 \in \text{DOW}$

Single occurrence words (SOWs) and SOW are similarly defined.

Definition: Repeat and Return Words

Definition

Given $w \in \Sigma^+$ and $SOW(w) \subseteq \Sigma^+$,

the word uu is a repeat word in w if $w = z_1uz_2uz_3$ for some $z_1; z_2; z_3 \in \Sigma^+$

the word uu^R is a return word in w if $w = z_1uz_2u^Rz_3$ for some $z_1; z_2; z_3 \in \Sigma^+$

w	Repeat words
1123455 <u>236</u> 78876	234234, 2323, 88, etc.

w	Return words
1123455 <u>234678876</u>	678876, 6776, 22, etc.

A repeat word uu or return word uu^R is trivial if $|u| = 1$.

Repeat and Return Words in Ciliate DNA

M_6	M_7	M_8	M_9	M_{11}	M_1	M_3	M_{10}	M_2	M_4	M_5	M_{12}	M_{13}
56	67	78	89	ab	1	23	9a	12	34	45	bc	c

$$w_0 = 56677889ab1239a123445bcc$$

Burns, J. et al. Recurring patterns among scrambled genes in the encrypted genome of the ciliate *Oxytricha trifallax*. *Journal of Theoretical Biology* 410 (2016) pp 171-180.

Repeat and Return Words in Ciliate DNA

M_6	M_7	M_8	M_9	M_{11}	M_1	M_3	M_{10}	M_2	M_4	M_5	M_{12}	M_{13}
56	67	78	89	ab	1	23	9a	12	34	45	bc	c

$$w_0 = 5\overline{66}7\overline{78}89ab1239a123\overline{44}5b\overline{cc}$$

$$w_1 = 59ab1239a1235b$$

Burns, J. et al. Recurring patterns among scrambled genes in the encrypted genome of the ciliate *Oxytricha trifallax*. *Journal of Theoretical Biology* 410 (2016) pp 171-180.

Repeat and Return Words in Ciliate DNA

M ₆	M ₇	M ₈	M ₉	M ₁₁	M ₁	M ₃	M ₁₀	M ₂	M ₄	M ₅	M ₁₂	M ₁₃
56	67	78	89	ab	1	23	9a	12	34	45	bc	c

$$w_0 = 5\overline{66}7\overline{78}89ab1239a123\overline{44}5b\overline{cc}$$

$$w_1 = 59ab1239a1235b$$

$$w_2 = 5b5b$$

Burns, J. et al. Recurring patterns among scrambled genes in the encrypted genome of the ciliate *Oxytricha trifallax*. *Journal of Theoretical Biology* 410 (2016) pp 171-180.

Repeat and Return Words in Ciliate DNA

M_6	M_7	M_8	M_9	M_{11}	M_1	M_3	M_{10}	M_2	M_4	M_5	M_{12}	M_{13}
56	67	78	89	ab	1	23	9a	12	34	45	bc	c

$$w_0 = 5\overline{66}7\overline{78}8\overline{9}ab123\overline{9a}123\overline{44}5\overline{b}cc$$

$$w_1 = 5\overline{9}ab1\overline{23}9\overline{a}1\overline{23}5b$$

$$w_2 = 5\overline{b}5\overline{b}$$

$$w_3 =$$

Burns, J. et al. Recurring patterns among scrambled genes in the encrypted genome of the ciliate *Oxytricha trifallax*. *Journal of Theoretical Biology* 410 (2016) pp 171-180.

Jonoska, N. et al. Patterns and Distances in Words Related to DNA Rearrangement. *Fundamenta Informaticae* 154:1-4 (2017) pp 225-238.

Definition: Repeat and Return Insertions

Definition

Given $w = a_1 a_2 \dots a_n \in \Sigma^*$,

let $1 \leq k \leq n+1$,

let $l \in \Sigma \setminus \Sigma_w$ be a symbol not in w , and

let $u \in \Sigma^*$ be a SOW such that $[u] \setminus [w] = \{l\}$.

Then $I(u; k; l)$ is an insertion on w which acts as follows:

$w \cdot I(u; k; l) = a_1 a_2 \dots a_{k-1} u a_k a_{k+1} \dots a_n$ where

$$u^0 = \begin{cases} u & \text{if } l = \epsilon \\ u^R & \text{if } l = \text{!} \end{cases}$$

1232314554 $\xrightarrow{(abc; 4; 6)}$ 123abc23abc14554

1232314554 $\xrightarrow{(abc; 7; 10)}$ 123231abc4554cba

Definition: Equivalence

Definition

Words $v, w \in \Sigma^*$ are equivalent if there exists a bijection $f: \Sigma \rightarrow \Sigma$ such that $f(v) = w$; in this case, we write $v \sim w$.

Examples of equivalent pairs of words:

123123		1 2 3	1221?	(ab; 3; 3)=	12 abba 21		1 2 a b
		# # #					# # # #
321321		3 2 1	1221?	(ab; 1; 5)=	ab 1221 ba		a b 1 2

The following words are not equivalent:

$$1232314554 \not\sim (abc; 4; 6) = \underline{1}23\underline{abc}23\underline{abc}14554$$

$$1232314554 \not\sim (abc; 7; 10) = \underline{1}2323\underline{1abc}4554\underline{cba}$$

When Do Insertions Yield Equivalent Words?

Let $w = a_1 a_2 \dots a_n$ be given.

- 1 If $w_1 = w ? (u_1; k_1; \text{'}_1)$ and $w_2 = w ? (u_2; k_2; \text{'}_2)$, then is it possible for w_1 and w_2 to be equivalent?

When Do Insertions Yield Equivalent Words?

Let $w = a_1 a_2 \dots a_n$ be given.

- 1 If $w_1 = w ? (u_1; k_1; \cdot_1)$ and $w_2 = w ? (u_2; k_2; \cdot_2)$, then is it possible for w_1 and w_2 to be equivalent? **Yes.**

$$\begin{array}{r} w_1 = 1212 ? (a; 3; 5) = 12a12a \quad | \quad 1 \quad 2 \quad a \\ \quad \quad \quad \quad \quad \quad \quad \quad \quad | \quad \# \quad \# \quad \# \\ w_2 = 1212 ? (a; 1; 3) = a12a12 \quad | \quad a \quad 1 \quad 2 \end{array}$$

But what if $|u_1| = |u_2| \neq 1$?

- 2 In general if $w_1 = w ? |u_1|(u_1; k_1; \cdot_1)$ and $w_2 = w ? |u_2|(u_2; k_2; \cdot_2) = w_2$, what can we say about w if the insertions are "distinct"?

$$\begin{array}{r} 1221 ? (ab; 3; 3) = 12abba21 \quad | \quad 1 \quad 2 \quad a \quad b \\ \quad \quad \quad \quad \quad \quad \quad \quad \quad | \quad \# \quad \# \quad \# \quad \# \\ 1221 ? (ab; 1; 5) = ab1221ba \quad | \quad a \quad b \quad 1 \quad 2 \end{array}$$

Definition: Distinct Insertions

Definition

Two insertions $I_1(u_1; k_1; \ell_1)$ and $I_2(u_2; k_2; \ell_2)$ on w are distinct if at least one of the following holds:

$$(k_1; \ell_1) \neq (k_2; \ell_2); \quad I_1 \neq I_2; \quad \text{or} \quad |u_1| \neq |u_2|$$

If $w_1 = w_2$, then $|u_1| = |u_2|$. What if $w_1 \neq w_2$ but only $I_1 \neq I_2$?

Definition: Distinct Insertions

Definition

Two insertions $I_1(u_1; k_1; \ell_1)$ and $I_2(u_2; k_2; \ell_2)$ on $w \in \text{DOW}$ are distinct if at least one of the following holds:

$$(k_1; \ell_1) \neq (k_2; \ell_2); \quad I_1 \neq I_2; \quad \text{or} \quad |u_1| \neq |u_2|$$

If $w_1 = w_2$, then $|u_1| = |u_2|$. What if $w_1 \neq w_2$ but only $I_1 \neq I_2$?

w_1 u_1 u_1

w_2

u_1 u_2 u_2^R

$|u_1| = |u_2|$ and $|u_1| = |u_2^R| \Rightarrow |u_2| = |u_2^R| = 1$ since $u_2 \in \text{DOW}$ so w

Distinct Insertions and Equivalent DOWs

Without loss of generality, we take $k_1 < k_2$. Suppose that $k_1 = k_2$:



Thus, $k_1 \in k_2$; similarly $\ell_1 \in \ell_2$. We have three cases:



Interleaving

$(k_1 < k_2 \quad \ell_1 < \ell_2)$



Nested

$(k_1 < k_2 \quad \ell_2 < \ell_1)$



Sequential

$(k_1 \quad \ell_1 < k_2 \quad \ell_2)$

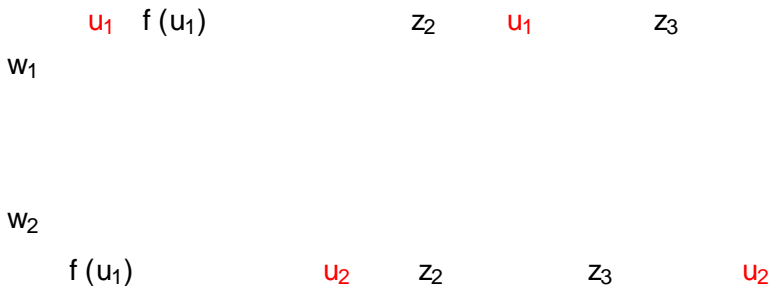
Interleaving Insertions $k_1 < k_2$ ($\backslash_1 < \backslash_2$)

w_1 u_1 z_1 z_2 u_1 z_3

w_2 z_1 u_2 z_2 z_3 u_2

Note that $u_1 z_1 z_1 u_2$. We consider $l_1 = l_2 =$ to start.

Interleaving Insertions $k_1 < k_2$ $\backslash_1 < \backslash_2$



Interleaving Insertions $k_1 < k_2$ $\backslash_1 < \backslash_2$

$$w_1 \quad u_1 \quad f(u_1) \quad f^h(u_1) \quad z_2 \quad u_1 \quad f(u_1) \quad f^h(u_1)$$

$$w_2 \quad f(u_1) f^2(u_1) \quad u_2 \quad z_2 \quad f(u_1) \quad u_2$$

We adapt a result by Lyndon and Schützenberger:

Lemma

If $xz = zy$ and $x \notin \langle z \rangle$, then $x = st$, $z = (st)^h s$, and $y = ts$ for some $s, t \neq \epsilon$ and $h \geq 0$.

Lyndon, R.C., and Schützenberger, M.-P. "The equation $a^m = b^n c^p$ in a free group." The Michigan Mathematical Journal 9:4 (1962) pp 289-298.

Interleaving Insertions $(k_1 < k_2, \ell_1 < \ell_2)$

w_1 u_1 $f(u_1)$ $f^h(u_1)$ z_2 u_1^0 $f(u_1^0)$ $f^h(u_1^0)$

w_2 $f(u_1) f^2(u_1)$ u_2 z_2 $f(u_1^0)$ u_2^0

Proposition (Interleaving)

If $l_1 = l_2 = \dots$, then $z_1 z_3$ is a repeat word.

If $l_1 = l_2 = \dots$, then $z_1 z_3 \in R(k_2, k_1; j u_{1j})$.

$R(h; q) = x_1 x_2 \dots x_h x_1^R x_2^R \dots x_h^R$ where each $x_i x_i^R$ is a return word and $|x_i| = q$ for $1 \leq i \leq h$.

Nested Insertions $k_1 < k_2$ $\backslash_2 < \backslash_1$

w_1 u_1 $f(u_1)$ $f^h(u_1)$ z_2 u_1^0

w_2
 $f(u_1) f^2(u_1)$ u_2 z_2 u_2^0

Nested Insertions $(k_1 < k_2 \quad \ell_2 < \ell_1)$

w_1 $u_1 \quad f(u_1) \quad f^h(u_1) \quad z_2 \quad f^h(u_1^0) \quad f(u_1^0) \quad u_1^0$

w_2 $f(u_1) f^2(u_1) \quad u_2 \quad z_2 \quad u_2^0 \quad f^2(u_1^0) f(u_1^0)$

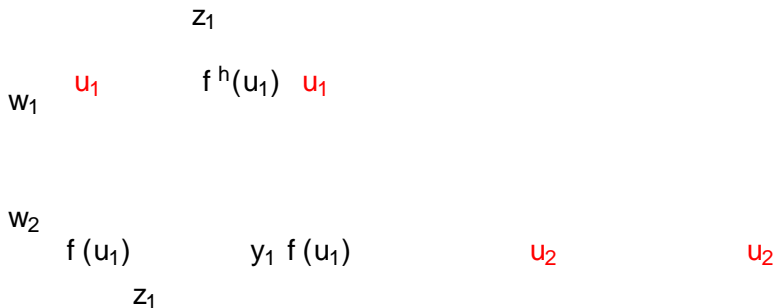
Proposition (Nested)

If $l_1 = l_2 = \dots$, then $z_1 z_3 \dots T(k_2 \quad k_1; j u_{1j})$.

If $l_1 = l_2 = \dots$, then $z_1 z_3$ is a return word.

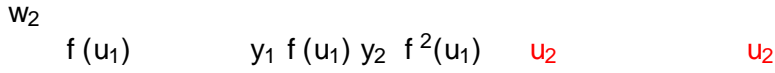
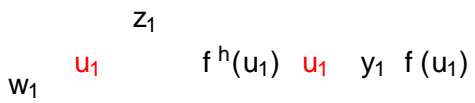
$T(h; q) = x_1 x_2 \dots x_{h-1} x_h x_{h-1} x_h x_{h-2} x_h x_{h-1} \dots x_2 x_1$ where each $x_i x_i$ is a repeat word and $|x_i| = q$ for $1 \leq i \leq h$.

Sequential Insertions $k_1 < k_2$



Note that $|y_{1j}| = |u_{1j}|$. We consider $l_1 = l_2 =$ to start.

Sequential Insertions $k_1 < k_2$



Sequential Insertions $k_1 < k_2$

$$u_1 \quad f^h(u_1) \quad u_1 \quad y_1 f(u_1) \quad y_p f^p(u_1) \quad y_p$$

 w_2

$$f(u_1) \quad y_1 f(u_1) \quad y_2 \quad u_2 \quad u_2$$

$$w = f(u_1) \quad f^h(u_1) y_1 f(u_1) y_2 f(u_2) \quad y_p f^p(u_1) \quad y_p$$

For example:

$$123456127812345678 \text{ (ab; 1; 5)} \quad 123456127812345678 \text{ (ab; 13; 17)}$$

$$w = \begin{array}{cccccccc} 12 & 34 & 56 & 12 & 78 & 12 & 34 & 56 & 78 \\ x_1 & x_2 & y_1 & x_1 & y_2 & x_2 & x_2 & y_1 & y_2 \end{array}$$

Sequential Insertions $k_1 \leq k_2 \leq \dots \leq k_n$

Consider the following words:

$$\begin{array}{ll}
 v_0 = 123123 & |v_0| = 6 \\
 v_1 = 123a123a & = v_0 ? \text{ (a; } |v_0| + 1) \\
 v_2 = 123a1b23ab & = v_1 ? \text{ (b; } |v_1| + 1) \\
 v_3 = 123a1b2c3abc & = v_2 ? \text{ (c; } |v_2| + 1) \\
 v_4 = 123a1b2c3dabcd & = v_3 ? \text{ (d; } |v_3| + 1)
 \end{array}$$

Words v_i are generalized i -tangled cords denoted $C(m; q; i)$ with $m = 3$ and $q = 1$. Tangled cords $C(1; 1; i)$, were introduced in:

Burns, J. et al. Four-regular graphs with rigid vertices associated to DNA recombination. Discrete Applied Mathematics, 161:10-11 (2013) pp 1378-1394.

Sequential Insertions $k_1 \leq k_2$

$$u_1 \quad f^h(u_1) \quad u_1^0 \quad y_1 f(u_1^0) \quad y_p f^p(u_1^0) \quad y_p^0$$

$$w_2 \quad f(u_1) \quad y_1 f(u_1^0) \quad y_2 \quad u_2 \quad u_2^0$$

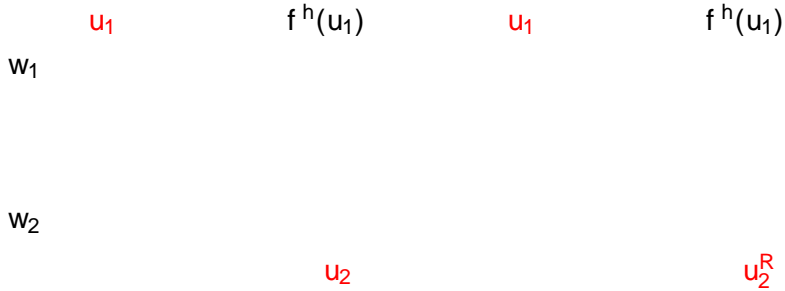
Proposition (Sequential)

If $l_1 = l_2 = \dots$, then $z_1 z_2 z_3 \dots \in C_{k_1; j}^{u_1; \frac{k_2 - l_1}{2j}}$.

If $l_1 = l_2 = \dots$, then $z_1 z_2 z_3 \dots \in C_{k_1; j}^{u_1; \frac{k_2 - l_1}{2j}}$.

Generalized l -tangled cord $C(m; q; j)$ is defined similarly.

Repeat and Return Insertions



Lemma

If uu and vv^R are repeat and return words in Σ^2 such that $[u] \setminus [v] \neq \emptyset$, then $|juj| = 1$ or $|jvj| = 1$.

Proposition (Repeat and Return)

Suppose that $u_1 \in I_2$. If $w_1 = w_2$, then $|ju_1j| = |ju_2j| = 1$.

Future Work: Graph of Words

[1a1a]

[1aa1]

(A)

(B)

[x]

[11]

[w]

(C)

[y]

[11aa]

