

# Finite Orbits of Language Operations

Émilie Charlier<sup>1</sup> Mike Domaratzki<sup>2</sup> Tero Harju<sup>3</sup> Jeffrey Shallit<sup>1</sup>

<sup>1</sup>University of Waterloo   <sup>2</sup>University of Manitoba   <sup>3</sup>University of Turku

Algorithms and Complexity Seminar  
Waterloo, May 18, 2011

# Closure operations

Let  $x : 2^{\Sigma^*} \rightarrow 2^{\Sigma^*}$  be an operation on languages. Suppose  $x$  satisfies the following three properties:

1.  $L \subseteq x(L)$  (expanding);
2. If  $L \subseteq M$  then  $x(L) \subseteq x(M)$  (inclusion-preserving);
3.  $x(x(L)) = x(L)$  (idempotent).

Then we say that  $x$  is a **closure operation**.

## Example

Kleene closure, positive closure, prefix, suffix, factor, subword.

## Some notation and a first result

If  $x(L) = y(L)$  for all languages  $L$ , then we write  $x \equiv y$ .

We write  $\epsilon(L) = L$  and  $xy = x \circ y$ , that is,  $xy(L) = x(y(L))$ .

Define  $c$  to be the complementation:  $c(L) = \Sigma^* - L$ . In particular, we have  $cc \equiv \epsilon$ .

### Theorem

*Let  $x, y$  be closure operations. Then  $xcycxcy \equiv xcy$ .*

## Corollary (Peleg 1984, Brzozowski-Grant-Shallit 2009)

*Let  $x$  be any closure operation and  $L$  be any language.*

*If  $S = \{x, c\}$ , then the orbit  $\mathcal{O}_S(L) = \{y(L) : y \in S^*\}$  contains at most 14 languages, which are given by the images of  $L$  under the 14 operations*

$\epsilon, x, c, xc, cx, xcxc, cxc, xcxc, cxcx,$   
 $xcxcx, cxcxc, xcxcxc, cxcxcx, cxcxcxc.$

NB: This result is the analogous for languages of Kuratowski-14 sets-theorem for topological spaces.

# Orbits of languages

Given a set  $S$  of operations, we consider the orbit of languages  $\mathcal{O}_S(L) = \{x(L) : x \in S^*\}$  under the monoid generated by  $S$ .

So compositions of operations in  $S$  are considered as “words over the alphabet  $S$ ”.

We are interested in the following questions: When is this monoid finite? Is the cardinality of  $\mathcal{O}_S(L)$  bounded, independently of  $L$ ?

## Operations with infinite orbit

It is possible for the orbit under a single operation to be infinite even if the operation is expanding and inclusion-preserving.

### Example

Consider the operation of fractional exponentiation, defined by

$$n(L) = \{x^\alpha : x \in L \text{ and } \alpha \geq 1 \text{ rational}\} = \bigcup_{x \in L} x^+ p(\{x\}).$$

Let  $M = \{ab\}$ . Then the orbit

$$\mathcal{O}_{\{n\}}(M) = \{M, n(M), n^2(M), n^3(M), \dots\}$$

is infinite, since we have

$$aba^i \in n^i(M) \text{ and } aba^i \notin n^j(M) \text{ for } j < i.$$

## Some notation and definitions

If  $t, x, y, z$  are words with  $t = xyz$ , we say

- ▶  $x$  is a **prefix** of  $t$ ;
- ▶  $z$  is a **suffix** of  $t$ ; and
- ▶  $y$  is a **factor** of  $t$ .

If  $t = x_1y_1x_2y_2 \cdots x_ny_nx_{n+1}$  for some words  $x_i$  and  $y_j$ , we say

- ▶  $y_1 \cdots y_n$  is a **subword** of  $t$ .

Thus a factor is a contiguous block, while a subword can be “scattered”.

Further,  $x^R$  denotes the reverse of the word  $x$ .

## 8 natural operations on languages

$$k: L \mapsto L^*$$

$$s: L \rightarrow \text{suff}(L)$$

$$e: L \mapsto L^+$$

$$f: L \rightarrow \text{fact}(L)$$

$$c: L \mapsto \bar{L} = \Sigma^* - L$$

$$w: L \rightarrow \text{subw}(L)$$

$$p: L \mapsto \text{pref}(L)$$

$$r: L \rightarrow L^R$$

where

$$L^* = \bigcup_{n \geq 0} L^n \quad \text{and} \quad L^+ = \bigcup_{n \geq 1} L^n$$

$$\text{pref}(L) = \{x \in \Sigma^* : x \text{ is a prefix of some } y \in L\}$$

$$\text{suff}(L) = \{x \in \Sigma^* : x \text{ is a suffix of some } y \in L\}$$

$$\text{fact}(L) = \{x \in \Sigma^* : x \text{ is a factor of some } y \in L\}$$

$$\text{subw}(L) = \{x \in \Sigma^* : x \text{ is a subword of some } y \in L\}$$

$$L^R = \{x \in \Sigma^* : x^R \in L\}$$



# Kuratowski identities

We now consider the set  $S = \{k, e, c, p, s, f, w, r\}$ .

## Lemma

*The 14 operations  $k, e, p, s, f, w, kp, ks, kf, kw, ep, es, ef,$  and  $ew$  are closure operations.*

## Theorem (mentioned above)

*Let  $x, y$  be closure operations. Then  $xcycxcy \equiv xcy$ .*

Together, these two results thus give  $196 = 14^2$  separate identities.

## Further identities

### Lemma

Let  $a \in \{k, e\}$  and  $b \in \{p, s, f, w\}$ . Then  $aba \equiv bab \equiv ab$ .

In a similar fashion, we obtain many kinds of Kuratowski-style identities involving the operations  $k$ ,  $e$ ,  $c$ ,  $p$ ,  $s$ ,  $f$ ,  $w$ , and  $r$ .

### Proposition

Let  $a \in \{k, e\}$  and  $b \in \{p, s, f, w\}$ . Then we have the following identities:

- ▶  $abcacaca \equiv abca$
- ▶  $bcbcbcab \equiv bcab$
- ▶  $abcbcabcab \equiv abcab$

## Additional identities (I)

We obtain many additional identities connecting the operations  $k$ ,  $e$ ,  $c$ ,  $p$ ,  $s$ ,  $f$ ,  $w$ , and  $r$ .

### Proposition

*We have the following identities:*

- ▶  $rp \equiv sr$ ;  $rs \equiv pr$
- ▶  $rf \equiv fr$ ;  $rw \equiv wr$ ;  $rc \equiv cr$ ;  $rk \equiv kr$
- ▶  $ps \equiv sp \equiv pf \equiv fp \equiv sf \equiv fs \equiv f$
- ▶  $pw \equiv wp \equiv sw \equiv ws \equiv fw \equiv wf \equiv w$
- ▶  $rkw \equiv kw \equiv wk$
- ▶  $ek \equiv ke \equiv k$
- ▶  $fks \equiv pks$ ;  $fkp \equiv skp$
- ▶  $rkf \equiv skf \equiv pkf \equiv fkf \equiv kf$

## Additional identities (II)

### Proposition

For all languages  $L$ , we have

- ▶  $pcs(L) = \Sigma^*$  or  $\emptyset$ .
- ▶ The same result holds for  $pcf$ ,  $fcs$ ,  $fcf$ ,  $scp$ ,  $scf$ ,  $fcf$ ,  $wcp$ ,  $wcs$ ,  $wcf$ ,  $pcw$ ,  $scw$ ,  $fcw$ , and  $wcw$ .

Let's prove this for  $pcs$ :

If  $s(L) = \Sigma^*$ , then  $cs(L) = \emptyset$  and  $pcs(L) = \emptyset$ .

Otherwise,  $s(L)$  omits some word  $w$ .

Then  $s(L) \cap \Sigma^*w = \emptyset$ .

Then  $\Sigma^*w \subseteq cs(L)$ .

Then  $\Sigma^* = p(\Sigma^*w) \subseteq pcs(L)$ , hence  $pcs(L) = \Sigma^*$ .

# Additional identities (III)

## Proposition

*For all languages  $L$ , we have*

- ▶  $sckp(L) = \Sigma^*$  or  $\emptyset$ .
- ▶ *The same result holds for  $fckp$ ,  $pcks$ ,  $fcks$ ,  $pckf$ ,  $sckf$ ,  $fckf$ ,  $wckp$ ,  $wcks$ ,  $wckf$ ,  $wckw$ ,  $pckw$ ,  $sckw$ ,  $fckw$ .*

## Proposition

*For all languages  $L$ , we have*

- ▶  $scskp(L) = \Sigma^*$  or  $\emptyset$ .
- ▶ *The same result holds for  $pcpks$ .*

## Additional identities (IV)

### Proposition

For all languages  $L$  and for all  $b \in \{p, s, f, w\}$ , we have

- ▶  $kcb(L) = cb(L) \cup \{\epsilon\}$
- ▶  $kckb(L) = ckb(L) \cup \{\epsilon\}$
- ▶  $kckck(L) = ckck(L) \cup \{\epsilon\}$
- ▶  $kbc bckb(L) = bcbckb(L) \cup \{\epsilon\}$ .

Let's prove  $kcp(L) \subseteq cp(L) \cup \{\epsilon\}$ :

Assume  $x \in kcp(L)$  and  $x \neq \epsilon$ .

We have  $x = x_1x_2 \cdots x_n$  for some  $n \geq 1$ , where each  $x_i \in cp(L)$ .

Then  $x_1x_2 \cdots x_n \notin p(L)$ , because if it were, then  $x_1 \in p(L)$ .

Hence  $x \in cp(L)$ .

# Main Result

Theorem (C-Domaratzki-Harju-Shallit 2011)

*Let  $S = \{k, e, c, p, f, s, w, r\}$ . Then for every language  $L$ , the orbit  $\mathcal{O}_S(L)$  contains at most 5676 distinct languages.*

## Sketch of the proof

We used breadth-first search to examine the set

$S^* = \{k, e, c, p, f, s, w, r\}^*$  w.r.t. the radix order with  
 $k < e < c < p < f < s < w < r$ .

As each new word  $x$  is examined, we test it to see if any factor is of the form given by “certain identities”.

If it is, then the corresponding language must be either  $\Sigma^*$ ,  $\emptyset$ ,  $\{\epsilon\}$ , or  $\Sigma^+$ ; furthermore, each descendant language will be of this form. In this case the word  $x$  is discarded.

Otherwise, we use the remaining identities to try to reduce  $x$  to an equivalent word that we have previously encountered. If we succeed, then  $x$  is discarded.

Otherwise we append all the words in  $Sx$  to the end of the queue.



## Sketch of the proof (cont'd)

If the process terminates, then  $\mathcal{O}_S(L)$  is of finite cardinality.

For  $S = \{k, e, c, p, f, s, w, r\}$ , the process terminated with 5672 nodes that could not be simplified using our identities. We did not count  $\emptyset$ ,  $\{\epsilon\}$ ,  $\Sigma^+$ , and  $\Sigma^*$ . The total is thus 5676.

(The longest word examined was *ckcpcpckpckpckpcpcpckckcr*, of length 25, and the same word with *p* replaced by *s*.)

If we use two arbitrary closure operations  $a$  and  $b$  with no relation between them, then the monoid generated by  $\{a, b\}$  is infinite, since any two finite prefixes of  $ababab \cdots$  are distinct.

### Example

Define the exponentiation operation

$$t(L) = \{x^i : x \in L \text{ and } i \text{ is an integer } \geq 1\}.$$

Then  $t$  is a closure operation.

Hence the orbits  $\mathcal{O}_{\{p\}}(L)$  and  $\mathcal{O}_{\{t\}}(L)$  are finite, for all  $L$ .

However, if  $M = \{ab\}$ , then the orbit  $\mathcal{O}_{\{p,t\}}(M)$  is infinite, as

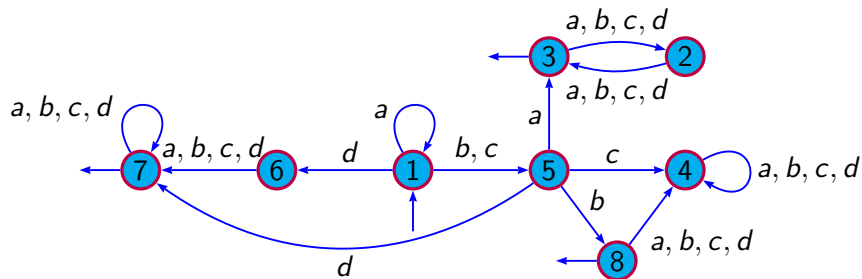
$$aba^i \in (pt)^i(M) \text{ and } aba^i \notin (pt)^j(M) \text{ for } j < i.$$

.

## Prefix and complement

In this case at most 14 distinct languages can be generated.

The bound of 14 can be achieved, e.g., by the regular language over  $\Sigma = \{a, b, c, d\}$  accepted by the following DFA:



The following table gives the appropriate set of final states under the operations.

language	final states	language	final states
$L$	3,7,8	$pcpc(L)$	1,5,6,7
$c(L)$	1,2,4,5,6	$cpcp(L)$	2,3,6,7
$p(L)$	1,2,3,5,6,7,8	$cpcpc(L)$	2,3,4,8
$pc(L)$	1,2,3,4,5,6,8	$pcpcp(L)$	1,2,3,5,6,7
$cp(L)$	4	$pcpcpc(L)$	1,2,3,4,5,8
$cpc(L)$	7	$cpcpcp(L)$	4, 8
$pcp(L)$	1,4,5,8	$cpcpcpc(L)$	6, 7

## Factor, Kleene star, complement

Here breadth-first search gives 78 languages, so our bound is  $78 + 4 = 82$ . We can improve this bound by considering new kinds of arguments.

### Lemma

*There are at most 4 languages distinct from  $\Sigma^*$ ,  $\emptyset$ ,  $\Sigma^+$ , and  $\{\epsilon\}$  in*

$$\mathcal{O}_{\{k,f,kc,fc\}}(f(L)).$$

*These languages are among  $f(L)$ ,  $kf(L)$ ,  $kckf(L)$ , and  $kcf(L)$ .*

### Lemma

*There are at most 2 languages distinct from  $\Sigma^*$ ,  $\emptyset$ ,  $\Sigma^+$ , and  $\{\epsilon\}$  in*

$$\mathcal{O}_{\{k,f,kc,fc\}}(fk(L)) - \mathcal{O}_{\{k,f,kc,fc\}}(f(L)).$$

*These languages are among  $fk(L)$  and  $kcfk(L)$ .*

## Lemma

*For all languages  $L$ , we have either  $f(L) = \Sigma^*$  or  $fc(L) = \Sigma^*$ .*

## Theorem (C-Domaratzki-Harju-Shallit 2011)

*Let  $L$  be an arbitrary language. Then 50 is a tight upper bound for the size of  $\mathcal{O}_{\{k,c,f\}}(L)$ .*

## Sketch of the proof

The languages in  $\mathcal{O}_{\{k,c,f\}}(L)$  that may differ from  $\Sigma^*$ ,  $\emptyset$ ,  $\Sigma^+$ , and  $\{\epsilon\}$  are among the images of  $L$  and  $c(L)$  under the 16 operations

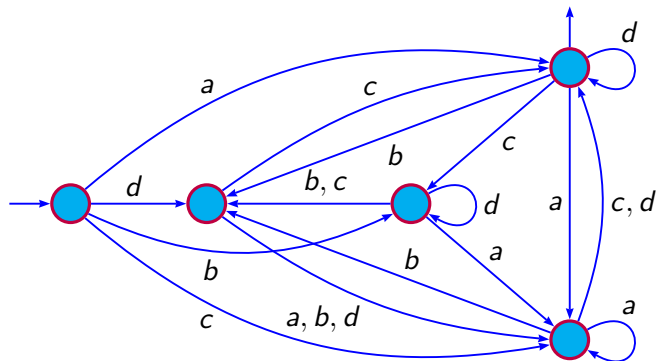
$$\begin{aligned} f, kf, kckf, kcf, fk, kcfk, fck, kfck, kckfck, kcfck, \\ fkck, kcfkck, fckck, kfckck, kckfckck, kcfckck, \end{aligned} \quad (1)$$

the complements of these images, together with the 14 languages in  $\mathcal{O}_{\{k,c\}}(L)$ .

We show that there are at most 32 distinct languages among the  $64 = 16 \cdot 4$  languages given by the images of  $L$  and  $c(L)$  under the 16 operations (1) and their complements.

Adding the 14 languages in  $\mathcal{O}_{\{k,c\}}(L)$ , and  $\Sigma^*$ ,  $\emptyset$ ,  $\Sigma^+$ , and  $\{\epsilon\}$ , we obtain that  $50 = 32 + 14 + 4$  is an upper bound for the size of the orbit of  $\{k, c, f\}$ .

## Sketch of the proof (cont'd)



The DFA made of two copies of this DFA accepts a language  $L$  with orbit size 50 under operations  $k$ ,  $c$ , and  $f$ .



## Kleene star, prefix, suffix, factor

Here there are at most 13 distinct languages, given by the action of

$$\{\epsilon, k, p, s, f, kp, ks, kf, pk, sk, fk, pks, skp\}.$$

The bound of 13 is achieved, for example, by  $L = \{abc\}$ .

## Summary of results

<i>r</i>	<b>2</b>	<i>w</i>	<b>2</b>	<i>f</i>	<b>2</b>
<i>s</i>	<b>2</b>	<i>p</i>	<b>2</b>	<i>c</i>	<b>2</b>
<i>k</i>	<b>2</b>	<i>w, r</i>	<b>4</b>	<i>f, r</i>	<b>4</b>
<i>f, w</i>	<b>3</b>	<i>s, w</i>	<b>3</b>	<i>s, f</i>	<b>3</b>
<i>p, w</i>	<b>3</b>	<i>p, f</i>	<b>3</b>	<i>c, r</i>	<b>4</b>
<i>c, w</i>	<b>6*</b>	<i>c, f</i>	<b>6*</b>	<i>c, s</i>	<b>14</b>
<i>c, p</i>	<b>14</b>	<i>k, r</i>	<b>4</b>	<i>k, w</i>	<b>4</b>
<i>k, f</i>	<b>5</b>	<i>k, s</i>	<b>5</b>	<i>k, p</i>	<b>5</b>
<i>k, c</i>	<b>14</b>	<i>f, w, r</i>	<b>6</b>	<i>s, f, w</i>	<b>4</b>
<i>p, f, w</i>	<b>4</b>	<i>p, s, f</i>	<b>4</b>	<i>c, w, r</i>	<b>10*</b>
<i>c, f, r</i>	<b>10*</b>	<i>c, f, w</i>	<b>8*</b>	<i>c, s, w</i>	<b>16*</b>
<i>c, s, f</i>	<b>16*</b>	<i>c, p, w</i>	<b>16*</b>	<i>c, p, f</i>	<b>16*</b>
<i>k, w, r</i>	<b>7</b>	<i>k, f, r</i>	<b>9</b>	<i>k, f, w</i>	<b>6</b>
<i>k, s, w</i>	<b>7</b>	<i>k, s, f</i>	<b>9</b>	<i>k, p, w</i>	<b>7</b>
<i>k, p, f</i>	<b>9</b>	<i>k, c, r</i>	<b>28</b>	<i>k, c, w</i>	<b>38*</b>
<i>k, c, f</i>	<b>50*</b>	<i>k, c, s</i>	1070	<i>k, c, p</i>	1070

## Summary of results (Cont'd)

$p, s, f, r$	<b>8</b>	$p, s, f, w$	<b>5</b>	$c, f, w, r$	<b>12*</b>
$c, s, f, w$	<b>16*</b>	$c, p, f, w$	<b>16*</b>	$c, p, s, f$	<b>16*</b>
$k, f, w, r$	<b>11</b>	$k, s, f, w$	<b>10</b>	$k, p, f, w$	<b>10</b>
$k, p, s, f$	<b>13</b>	$k, c, w, r$	72*	$k, c, f, r$	<b>84*</b>
$k, c, f, w$	66*	$k, c, s, w$	1114	$k, c, s, f$	1450
$k, c, p, w$	1114	$k, c, p, f$	1450	$p, s, f, w, r$	<b>10</b>
$c, p, s, f, r$	<b>30*</b>	$c, p, s, f, w$	<b>16*</b>	$k, p, s, f, r$	<b>25</b>
$k, p, s, f, w$	<b>14</b>	$k, c, f, w, r$	120*	$k, c, s, f, w$	1474
$k, c, p, f, w$	1474	$k, c, p, s, f$	2818	$c, p, s, f, w, r$	<b>30*</b>
$k, p, s, f, w, r$	<b>27</b>	$k, c, p, s, f, r$	5628	$k, c, p, s, f, w$	2842
$k, c, p, s, f, w, r$	5676				

## Further work

We plan to continue to refine our estimates of the previous tables, and pursue the status of other sets of operations.

For example, if  $t$  is the exponentiation operation, then, using the identities  $kt \equiv tk \equiv k$ , and the inclusion  $t \subseteq k$ , we get the additional Kuratowski-style identities

- ▶  $kctckck \equiv kck$ ,
- ▶  $kckctck \equiv kck$ ,
- ▶  $kctctck \equiv kck$ ,
- ▶  $tctctck \equiv tck$ ,
- ▶  $kctctct \equiv kct$ .

This allows us to prove that  $\mathcal{O}_{\{k,c,t\}}(L)$  is finite and of cardinality at most 126.