On the Number of Non-equivalent Parameterized Squares in a String

presented at SPIRE 2024

Rikuya Hamai Kazushi Taketsugu <u>Yuto Nakashima</u> Shunsuke Inenaga Hideo Bannai



Square = string represented by *xx*

Q. How many distinct squares can a string of length *n* contain?

caabaabcaabaabcc

of distinct squares = 5
of occ. of squares = 9

Considering the distinctness is more interesting than the occurrences because the maximum number of occ. of squares is $\Theta(n^2)$ for a string a^n .

Distinct square conjecture [Fraenkel and Simpson, 1998] The maximum number of distinct squares is less than *n*.

< 2 <i>n</i>	[Fraenkel and Simpson, 1998]	< 11 <i>n</i> /6	[Deza et al., 2015]
$< 2n - \Theta(\log n)$	[llie, 2007]	< 1.5 <i>n</i>	[Thierry, 2020]
< 95 <i>n</i> /48	[Lam, 2013]	$< n - \sigma + 1$	[Brlek & Li, 2022]

Squares in other matching models

Square = concatenation of two equal strings

Q. What happens if "equal" is changed from exact matching to another matching model (equivalence)?



Parameterized equivalence [Baker, 1996]

Definition. Two strings *x* and *y* of length *k* are said to be **parameterized** equivalent if there is a bijection *f* on Σ such that f(x[i]) = y[i] for all $1 \le i \le k$. We write $x \equiv y$ if two strings *x* and *y* are parameterized equivalent.



Parameterized square (P-square)

Definition. A string w is called a parameterized square when w = xy for strings x, y and x and y are parameterized equivalent.



We sometimes use "P" to denote the term "parameterized" (e.g., P-square).

There are two types of the distinctness.

A) # of p-squares that are distinct as strings.

B) # of p-squares that are non-equivalent in the equivalence model.



There are two types of the distinctness.

A) # of p-squares that are distinct as strings.

B) # of p-squares that are non-equivalent in the equivalence model.

	distinct as strings		non-equivalent	
	upper bound	lower bound	upper bound	lower bound
Parameterized	$\leq 2(\sigma!)^2n$ [Kociumaka et al., 2016]	_	$\leq 2\sigma!n$ [Kociumaka et al., 2016]	_
Order-preserving	$\mathrm{O}(\sigma n)$ [Gawrychowski et al., 2023]	$\Omega(\sigma n)$ [Gawrychowski et al., 2023]	$O(\sigma n)$ [Gawrychowski et al., 2023]	
Abelian	$\mathrm{O}(\sigma n)$ [Kociumaka et al., 2016]	$\Omega(\sigma n)$ [Kociumaka et al., 2016]	\leq $(n\!\!-\!\!1)^{11/6}$ [Kociumaka et al., 2016]	$\Omega(n^{1.5})$ [Kociumaka et al., 2016]

n : string length, σ : alphabet size

There are two types of the distinctness.

- A) # of p-squares that are distinct as strings.
- B) # of p-squares that are non-equivalent in the equivalence model.

	distinct as strings		non-equivalent	
	upper bound	lower bound	upper bound	lower bound
Parameterized	$\leq 2(\sigma!)^2n$ [Kociumaka et al., 2016]	_	$\leq 2\sigma! n$ [Kociumaka et al., 2016]	—
Order-preserving	$\mathrm{O}(\sigma n)$ [Gawrychowski et al., 2023]	$\Omega(\sigma n)$ [Gawrychowski et al., 2023]	$O(\sigma n)$ [Gawrychowski et al., 2023]	_
Abelian	$\mathrm{O}(\sigma n)$ [Kociumaka et al., 2016]	$\Omega(\sigma n)$ [Kociumaka et al., 2016]	\leq $(n-1)^{11/6}$ [Kociumaka et al., 2016]	$\Omega(n^{1.5})$ [Kociumaka et al., 2016]

Theorem 1. Any string of length *n* that contains σ distinct characters can contain at most σn non-equivalent parameterized squares.

Upper bound (main theorem)

Theorem 2. For any string *s* that contains σ distinct characters, there can be at most σ prefixes of *s* that are P-squares and have no other parameterized occurrences in *s*.

at most σ parameterized squares for each of *n* positions

Theorem 1. Any string of length *n* that contains σ distinct characters can contain at most σn non-equivalent parameterized squares.

We consider the periodicity in the parameterized equivalent model.

Standard period

Definition. An integer *p* is said to be a period of string *w* 1. if w[i] = w[i+p] for all *i* satisfying $1 \le i \le |w| - p$, or equivalently, 2. if w[1..|w|-p] = w[p+1..|w|].

E.g.



Standard period

Definition. An integer *p* is said to be a period of string *w* 1. if w[i] = w[i+p] for all *i* satisfying $1 \le i \le |w| - p$, or equivalently, 2. if w[1 + |w| + p] = w[p+1 + |w|]

2. if
$$w[1..|w|-p] = w[p+1..|w|]$$
.



Parameterized period (P-period)

Definition. An integer *p* is said to be a P-period of string *w* 1. if f(w[i]) = w[i+p] for all *i* satisfying $1 \le i \le |w| - p$ with a bijection *f*, or equivalently,

2. if $w[1..|w|-p] \equiv w[p+1..|w|]$.



 We sometimes use p ||_f w to denote the fact "p is a P-period of w with a bijection f" (we just write p || w when no confusions occur). To prove our main result, we use the following lemmas about P-periods. <u>Periodicity lemma</u>

Lemma 1. Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 2) \le |s|$, then gcd(p, q) is a P-period of *s*.

Period connection lemma

Lemma 2. Let *x*, *y*, *z* be strings s.t. *p* is a P-period of *xy* with bijection *f* and a P-period of *yz* with bijection *g*. If $|y| \ge p \cdot (|\Sigma_{xyz}| - 1) + 1$, then *p* is a P-period of *xyz* with bijection f(=g).

Period extension lemma

Lemma 3. Let *s* be a string that has a P-period *q*, and *s* ' be a substring of *s* that has a P-period *p*. If $p \cdot (|\Sigma_s| - 2) + q + 1 \le |s'|$ and q = kp for some integer *k*, then *p* is a P-period of *s*.

Properties on P-period

To prove our main result, we use the follo

Periodicity lemma

Lemma 1. Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 2) \le |s|$, then gcd(p, q) is a P-period of *s*.

Period connection lemma

Lemma 2. Let *x*, *y*, *z* be strings s.t. *p* is a P-period of *xy* with bijection *f* and a P-period of *yz* with bijection *g*. If $|y| \ge p \cdot (|\Sigma_{xyz}| - 1) + 1$, then *p* is a P-period of *xyz* with bijection f(=g).

Period extension lemma

Lemma 3. Let *s* be a string that has a P-period *q*, and *s* ' be a substring of *s* that has a P-period *p*. If $p \cdot (|\Sigma_s| - 2) + q + 1 \le |s'|$ and q = kp for some integer *k*, then *p* is a P-period of *s*.

Periodicity lemma for standard period

Lemma. [Fine & Wilf, 1965] Let *s* be a string that has periods *p* and *q*. If $p + q - \text{gcd}(p, q) \le |s|$, then gcd(p, q) is a period of *s*.

(Previous) Periodicity lemmas for P-period

Lemma. [Apostolico & Giancarlo, 2007] Let *s* be a string that has a P-period of *p* with bijection *f* and a P-period *q* with bijection *g*. If $p + q \le |s|$ and $f \circ g = g \circ f$, then gcd(p, q) is a P-period of *s*.

Lemma. [Ideguchi et al., 2023] Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 1) \le |s|$, then gcd(p, q) is a P-period of *s*. **Lemma.** [Ideguchi et al., 2023] Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 1) \le |s|$, then gcd(p, q) is a P-period of *s*.

We give an improved version of the above lemma.

Lemma 1. Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 2) \le |s|$, then gcd(p, q) is a P-period of *s*.

■ Our result gives a tight bound when $\sigma = 2$ [Apostolico & Giancarlo, 2007]. ▶ When $\sigma = 2$, *f* and *g* always commute.

Tighter versions

There exist tighter versions for each alphabet size.

Alphabet size	Lower bound of s		
$\sigma = 2$	$\geq p + q$ tight!	[Apostolico & Giancarlo, 2007]	
$\sigma = 3$	$\geq p + q + 1$ tight!	[Our work (not published)]	
$\sigma = 4$	$\geq p + q + [\min(p, q) / 2]$ tight!	[Our work (not published)]	
$\sigma = 5$	$\geq p + q + [\min(p, q) / 2] + 1$ tight!	[Our work (not published)]	
$\sigma \ge 6$	$\geq p + q + \min(p, q) \cdot (\sigma - 3)$	[Our work (not published)]	

Lemma 1. Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 2) \le |s|$, then gcd(p, q) is a P-period of *s*.

In this work, we only use this general version.

Parameterized period connection lemma

Lemma 2. Let *x*, *y*, *z* be strings s.t. *p* is a P-period of *xy* with bijection *f* and a P-period of *yz* with bijection *g*. If $|y| \ge p \cdot (|\Sigma_{xyz}| - 1) + 1$, then *p* is a P-period of *xyz* with bijection f(=g).

Almost all the characters occur in y.



Intuitively, if the overlapping part y is sufficiently long, the whole string has the same period p and they use the same bijection.

Parameterized period extension lemma

Lemma 3. Let *s* be a string that has a P-period *q*, and *s* ' be a substring of *s* that has a P-period *p*. If $p \cdot (|\Sigma_s| - 2) + q + 1 \le |s'|$ and q = kp for some integer *k*, then *p* is a P-period of *s*.



Intuitively, a P-period of a substring which is sufficiently long can extend to the whole string, if the whole string has a P-period that is a multiple of the substring's P-period.

Theorem 2. For any string *s* that contains σ distinct characters, there can be at most σ prefixes of *s* that are P-squares and have no other parameterized occurrences in *s*.



- Assume on the contrary that there exists a string *s* that contains σ distinct characters s.t. the number of parameterized squares of the prefixes of *s* that have no parameterized occurrences other than the prefixes is $\sigma + 1$.
- $\sigma+1 \text{ prefix squares}$



Since x_1x_1 ' has no parameterized equivalent occurrences other than the prefix, then $|x_1x_1'| > |x_{\sigma+1}|$. (otherwise, $x_{\sigma+1}$ ' has a prefix that is parameterized equivalent to x_1x_1' .)



Since x_1x_1 ' has no parameterized equivalent occurrences other than the prefix, then $|x_1x_1'| > |x_{\sigma+1}|$. (otherwise, $x_{\sigma+1}$ ' has a prefix that is parameterized equivalent to x_1x_1' .)



Let $r_i = |x_{i+1}| - |x_i|$.

■ $r_i \parallel x_i$, $(r_i \text{ is a period of } x_i)$ because $x_i \equiv x_i \equiv x_{i+1} [1../x_i] \equiv x_{i+1} [1../x_i]$.



Hence, $r_i \parallel x_i$ also holds for any *i*.

This implies that $r_i \parallel x_1$ holds for any *i* since $r_i < |x_1|$ and x_1 is a prefix of x_i .

Hence, $r_i \parallel x_i$ also holds for any *i*.

This implies that $r_i \parallel x_1$ holds for any *i* since $r_i < |x_1|$ and x_1 is a prefix of x_i .



Claim 1. The smallest P-period $p(x_1)$ is a period of x_{σ} .

- Since the condition of the periodicity, we have
 - $p(x_1) \cdot (\sigma 1) + r_{\sigma} \le r_1 + \dots + r_{\sigma} \le |x_1'| = |x_1| \quad \Leftrightarrow \quad p(x_1) \cdot (\sigma 2) + p(x_1) + r_{\sigma} \le |x_1|.$



• $p(x_1) \leq r_i$

26

Claim 1. The smallest P-period $p(x_1)$ is a period of x_{σ} .

Since the condition of the periodicity, we have

$$p(x_1) \cdot (\sigma - 1) + r_{\sigma} \le r_1 + \dots + r_{\sigma} \le |x_1'| = |x_1| \quad \Leftrightarrow \left[p(x_1) \cdot (\sigma - 2) + p(x_1) + r_{\sigma} \le |x_1| \right]$$

Lemma 1. Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 2) \le |s|$, then gcd(p, q) is a P-period of *s*.

- By combining with Lemma 1, we have $gcd(p(x_1), r_{\sigma}) \parallel x_1$.
- Since $p(x_1)$ is the smallest P-period of x_1 , $r_{\sigma} = c \cdot p(x_1)$ for some integer $c \ge 1$.

• $r_i \parallel x_1$

• $p(x_1) \leq r_i$

Claim 1. The smallest P-period $p(x_1)$ is a period of x_{σ} .

Since the condition of the periodicity, we have

 $p(x_1) \cdot (\sigma - 1) + r_{\sigma} \le r_1 + \dots + r_{\sigma} \le |x_1'| = |x_1| \quad \Leftrightarrow \quad p(x_1) \cdot (\sigma - 2) + p(x_1) + r_{\sigma} \le |x_1|.$

This inequation also implies that $p(x_1) \cdot (\sigma - 2) + r_{\sigma} + 1 \le |x_1|$.

Lemma 3. Let *s* be a string that has a P-period *q*, and *s*' be a substring of *s* that has a P-period *p*. If $p \cdot (|\Sigma_s| - 2) + q + 1 \le |s'|$ and q = kp for some integer *k*, then *p* is a P-period of *s*.

By combining with Lemma 3 and $r_{\sigma} = c \cdot p(x_1)$, we have $p(x_1) \parallel x_{\sigma}$ (Claim 1 holds).

• $r_i \parallel x_1$

• $p(x_1) \leq r_i$

Claim 1. The smallest P-period $p(x_1)$ is a period of x_{σ} .

- Since the condition of the periodicity, we have
 - $p(x_1) \cdot (\sigma 1) + r_{\sigma} \le r_1 + \dots + r_{\sigma} \le |x_1'| = |x_1| \quad \Leftrightarrow \quad p(x_1) \cdot (\sigma 2) + p(x_1) + r_{\sigma} \le |x_1|.$



• $p(x_1) \leq r_i$

29

Claim 2. The smallest P-period $p(x_1)$ is a period of $x_{\sigma+1}$.

- $r_i \parallel x_1$ • $p(x_1) \le r_i$
- Since the condition of the periodicity and $|x_1| \parallel x_{\sigma}$, we have
 - $p(x_1) \cdot (\sigma 1) + |x_1| \le r_1 + \dots + r_{\sigma 1} + |x_1| = |x_{\sigma}| \iff p(x_1) \cdot (\sigma 2) + p(x_1) + |x_1| \le |x_{\sigma}|.$



Claim 2. The smallest P-period $p(x_1)$ is a period of $x_{\sigma+1}$.

•
$$r_i \parallel x_1$$

 $p(x_1) \leq r_i$

Since the condition of the periodicity and $|x_1| \parallel x_{\sigma}$, we have

 $p(x_1) \cdot (\sigma - 1) + |x_1| \le r_1 + \dots + r_{\sigma - 1} + |x_1| = |x_{\sigma}| \iff p(x_1) \cdot (\sigma - 2) + p(x_1) + |x_1| \le |x_{\sigma}|.$

Lemma 1. Let *s* be a string that has P-periods *p* and *q* and satisfies $|\Sigma_s| \ge 2$. If $p + q + \min(p, q) \cdot (|\Sigma_s| - 2) \le |s|$, then gcd(p, q) is a P-period of *s*.

- By combining with Lemma 1, we have $gcd(p(x_1), |x_1|) \parallel x_{\sigma}$.
- Since $p(x_1)$ is the smallest P-period of x_1 , $|x_1| = d \cdot p(x_1)$ for some integer $d \ge 1$.

Claim 2. The smallest P-period $p(x_1)$ is a period of $x_{\sigma+1}$.

- $r_i \parallel x_1$ • $p(x_1) \le r_i$
- Since the condition of the periodicity and $|x_1| \parallel x_{\sigma}$, we have

 $p(x_1) \cdot (\sigma - 1) + |x_1| \le r_1 + \dots + r_{\sigma - 1} + |x_1| = |x_{\sigma}| \iff p(x_1) \cdot (\sigma - 2) + p(x_1) + |x_1| \le |x_{\sigma}|.$

This inequation also implies that $p(x_1) \cdot (\sigma - 2) + |x_1| + 1 \le |x_{\sigma}|$.

Lemma 3. Let *s* be a string that has a P-period *q*, and *s* ' be a substring of *s* that has a P-period *p*. If $p \cdot (|\Sigma_s| - 2) + q + 1 \le |s'|$ and q = kp for some integer *k*, then *p* is a P-period of *s*.

By combining with Lemma 3 and $|x_1| = d \cdot p(x_1)$, we have $p(x_1) \parallel x_{\sigma+1}$ (Claim 2 holds).

Claim 2. The smallest P-period $p(x_1)$ is a period of $x_{\sigma+1}$.

Since the condition of the periodicity and $|x_1| \parallel x_{\sigma}$, we have

$$p(x_1) \cdot (\sigma - 1) + |x_1| \le r_1 + \dots + r_{\sigma - 1} + |x_1| = |x_\sigma| \iff p(x_1) \cdot (\sigma - 2) + p(x_1) + |x_1| \le |x_\sigma|.$$



• $r_i \parallel x_1$

• $p(x_1) \leq r_i$

We can use Lemma 2 for $x_{\sigma+1}$ and x_1 ' because

• $x_{\sigma+1}$ and x_1 ' have a period $p(x_1)$ (by Claim 2 & definition).

 \blacksquare $x_{\sigma+1}$ and x_1 ' overlap each other, and the length is $r_1 + \cdots + r_{\sigma}$.



We can use Lemma 2 for $x_{\sigma+1}$ and x_1 ' because

• $x_{\sigma+1}$ and x_1 ' have a period $p(x_1)$ (by Claim 2 & definition),

• $x_{\sigma+1}$ and x_1 ' overlap each other, and

• the overlap length is $r_1 + \cdots + r_{\sigma} \ge p(x_1) \cdot (\sigma - 1) + 1$.

Lemma 2. Let *x*, *y*, *z* be strings s.t. *p* is a P-period of *xy* with bijection *f* and a P-period of *yz* with bijection *g*. If $|y| \ge p \cdot (|\Sigma_{xyz}| - 1) + 1$, then *p* is a P-period of *xyz* with bijection *f* (= *g*).

By applying Lemma 2 for $x_{\sigma+1}$ and x_1 ', we have $p(x_1) \parallel x_1 x_1$ '.



By a similar argument for x_1x_1 ' and x_2 ', we have $p(x_1) \parallel x_2x_2$ '.

Then, x_1x_1 ' has another parameterized equivalent occurrence in x_2x_2 '.



Theorem 2. For any string *s* that contains σ distinct characters, there can be at most σ prefixes of *s* that are P-squares and have no other parameterized occurrences in *s*.

Open questions

Theorem 1. Any string of length *n* that contains σ distinct characters can contain at most σn non-equivalent parameterized squares.

- Is this upper bound is tight?
 - How can we improve this bound?
 (There is a lower bound for the σ prefixes property.)
 - How can we construct a family of strings that gives a good lower bound?
- Can we improve bounds in another distinctness (distinct as words)?

Example of σ **prefixes**

Theorem. For any σ , there exists a string with alphabet size σ such that there are σ prefixes that are parameterized squares and have no other parameterized occurrence.

 $(c_1 \cdots c_{\sigma-2} c_{\sigma-1}) (c_1 \cdots c_{\sigma-2} c_{\sigma-1}) c_{\sigma} (c_2 \cdots c_{\sigma-2} c_{\sigma-1} c_{\sigma}) (c_2 \cdots c_{\sigma-2} c_{\sigma-1} c_{\sigma}) c_1$ 123412345234523451 $\sigma = 5$