## Computing *q*-gram Non-overlapping Frequencies on SLP Compressed Texts

Kyushu University OKeisuke Goto, Hideo Bannai, Shunsuke Inenaga, Masayuki Takeda

\* Large scale string data is usually stored in compressed form.

# Large Scale String

\* Large scale string data is usually stored in compressed form.



Large scale string data is usually stored in compressed form.
In order to process such data, we usually need to decompress, which requires a lot of space and time.



Large scale string data is usually stored in compressed form.
In order to process such data, we usually need to decompress, which requires a lot of space and time.

# Large Scale String

string processing, e.g. Pattern Matching, Edit Distance

compressed

string

decompress

Large scale string data is usually stored in compressed form.
In order to process such data, we usually need to decompress, which requires a lot of space and time.



- \* Large scale string data is usually stored in compressed form.
- In order to process such data, we usually need to decompress, which requires a lot of space and time.
- One solution is to process compressed strings without explicit decompression.



[2]

[1]

4

[1] [Lifshits, Y. 2007][2] [Hermelin, D. 2009]

- \* Large scale string data is usually stored in compressed form.
- In order to process such data, we usually need to decompress, which requires a lot of space and time.
- One solution is to process compressed strings without explicit decompression.



- \* Large scale string data is usually stored in compressed form.
- In order to process such data, we usually need to decompress, which requires a lot of space and time.
- One solution is to process compressed strings without explicit decompression.



#### Occurrence and Frequency

• Occurrence of pattern *P* in text *T*:  $Occ(T, P) = \{ i | T[i..(i+|P|-1)] = P, 0 \le i < |T| - |P| \}$ 

+ Frequency of P in T: |Occ(T, P)|

Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.

Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



 $Occ(T, aaa) = \{3, 4, 5\}$ 

Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



 $Occ(T, aba) = \{1, 10, 12, 14\}$ 

Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



non-overlapping frequency nOcc(T, P): the size of largest subset of Occ(T, P) where any two occrences in the set are non-overlapping

Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



non-overlapping frequency nOcc(T, P): the size of largest subset of Occ(T, P) where any two occrences in the set are non-overlapping

Any two occurrences  $k_1$ ,  $k_2 \in Occ(T, P)$  with  $k_1 < k_2$ are said to be overlapping if  $k_2 < k_1 + |P|$ , and non-overlapping otherwise.



non-overlapping frequency nOcc(T, P): the size of largest subset of Occ(T, P) where any two occrences in the set are non-overlapping

nOcc(T, aba) = 3

*q*-gram Non-overlapping Frequencies Problem

Input : string T, positive integer q

Output : nOcc(T, x) for  $\forall x \in \Sigma^q$  s.t. nOcc(T, x) > 0

q-gram Non-overlapping Frequencies Problem

Input : string T, positive integer q

Output : nOcc(T, x) for  $\forall x \in \Sigma^q$  s.t. nOcc(T, x) > 0

Example q = 3T = abaababaab

3-gram	nOcc
aab	
aba	
baa	
bab	

q-gram Non-overlapping Frequencies Problem

Input : string T, positive integer q

Output : nOcc(T, x) for  $\forall x \in \Sigma^q$  s.t. nOcc(T, x) > 0

Example	q = 3
	T = abaababaab

3-gram	nOcc
aab	2
aba	
baa	
bab	

q-gram Non-overlapping Frequencies Problem

Input : string T, positive integer q

Output : nOcc(T, x) for  $\forall x \in \Sigma^q$  s.t. nOcc(T, x) > 0

Example	q = 3
	T = abaababaab

3-gram	nOcc
aab	2
aba	2
baa	
bab	

q-gram Non-overlapping Frequencies Problem

Input : string T, positive integer q

Output : nOcc(T, x) for  $\forall x \in \Sigma^q$  s.t. nOcc(T, x) > 0

Example	q = 3
	T = abaababaab

3-gram	nOcc
aab	2
aba	2
baa	2
bab	1

We want to solve this problem when string T is given in compressed form.

#### Straight Line Program (SLP)



- + SLP can represent the output of well-known compression schemes.
  - + e.g. RE-PAIR SEQUITUR, LZ78, LZW
- + LZ77, LZSS of size m can be quickly transformed to SLP of size  $O(m \log |T|)$ . [Rytter, 2003]

#### Example SLP



Length of decompressed string can be  $\Theta(2^n)$ 

#### Example SLP



Length of decompressed string can be  $\Theta(2^n)$ 

#### Previous Work

	Uncompressed String	SLP (Previous Work)
<i>q</i> -gram Freq	O( T )time and space	O(qn) <sup>[2]</sup> time and space
2-gram non- overlapping Freq	O( T )time and space	$O(n^4 \log n) \text{ time} \\ O(n^3) \text{ space}$
<i>q</i> -gram non- overlapping Freq	O( T )time and space	_

[1] [Inenaga and Bannai, 2009][2] [Goto et al, 2011]

#### Previous Work

	Uncompressed	SLP	SLP
	String	(Previous Work)	(This Work)
<i>q</i> -gram Freq	O( T )time and space	O(qn) <sup>[2]</sup> time and space	_
2-gram non- overlapping Freq	O( T ) time and space	$O(n^4 \log n) \text{ time} \\ O(n^3) \text{ space}$	<i>O(n)</i> time and space
<i>q</i> -gram non- overlapping Freq	O( T )time and space	_	$O(q^2n)$ time O(qn) space

[1] [Inenaga and Bannai, 2009][2] [Goto et al, 2011]

## O(qn) algorithm for q-gram Frequencies Problem on SLP

[Goto et al., 2011]

#### Occ\*: Crossing Occurrences

- + For  $X_i = X_l X_r$ , an occurrence of *P* crosses  $X_i$ , if *P* starts in  $X_l$  and ends in  $X_r$ .
- \* Denote the set of such occurrences by  $Occ^*(X_i, P)$ .



vOcc



 $vOcc(X_i)$ : the number of occurrence  $X_i$  in the derivation tree of  $X_n$ 

Example  $vOcc(X_4) = 3$ 

Lemma 1  
$$|Occ(T,P)| = \sum_{i=1}^{n} |Occ^*(X_i,P)| \cdot vOcc(X_i)$$

Lemma 1  
$$|Occ(T,P)| = \sum_{i=1}^{n} |Occ^*(X_i,P)| \cdot vOcc(X_i)$$

For any <u>occurrence</u> of substring P in T, there uniquely exists  $X_i$  such that P crosses  $X_i$ .



Lemma 1  
$$|Occ(T,P)| = \sum_{i=1}^{n} |Occ^*(X_i,P)| \cdot vOcc(X_i)$$

T

For any <u>occurrence</u> of substring P in T, there uniquely exists  $X_i$  such that P crosses  $X_i$ .

For each variable  $X_i$ ,

- the number of times *P* crosses  $X_i$  is  $|Occ^*(X_i, P)|$ .
- $X_i$  appears in the derivation tree  $vOcc(X_i)$  times.



Lemma 1  
$$|Occ(T,P)| = \sum_{i=1}^{n} |Occ^*(X_i,P)| \cdot vOcc(X_i)$$

 $X_i$ 

For any <u>occurrence</u> of substring P in T, there uniquely exists  $X_i$  such that P crosses  $X_i$ .

For each variable  $X_i$ ,

- the number of times *P* crosses  $X_i$  is  $|Occ^*(X_i, P)|$ .
- $X_i$  appears in the derivation tree  $vOcc(X_i)$  times.

 $vOcc(X_{i_1}) = 2, vOcc(X_{i_2}) = 1$  $|Occ^*(X_{i_1}, P)| = 2, |Occ^*(X_{i_2}, P)| = 1$ Frequency of  $P = 2 \cdot 2 + 1 = 5$   $X_n$ 

 $X_{i_1}$ 

 $X_{i\gamma}$ 

#### Computing the Number of Crossing Occurrences

For all  $X_i = X_l X_r$ , let  $t_i = s_l p_r$  where  $s_l$  is the length-(q-1) suffix of  $X_l$  $p_r$  is the length-(q-1) prefix of  $X_r$ 

the number of occurrences of *q*-gram in *t<sub>i</sub>* equals to that of crossing occurrences in *X<sub>i</sub>*.  $|Occ^*(X_i, t_i[j ...j+q-1])| = |Occ(t_i, t_i[j ...j+q-1])|$ where  $1 \le j \le |t_i|-q+1$ 

for 
$$P \in \Sigma^q$$
,  $|Occ(T,P)| = \sum_{i=1}^n |Occ(t_i,P)| \cdot vOcc(X_i)$ 

#### Algorithm for SLP q-gram Frequencies Problem

Theorem 1

SLP q-gram Frequencies Problem can be solved in O(qn) time and space.

#### Idea:

Using the suffix array of the concatenation of all  $t_i$ 's, plus a linear-time construction algorithm of suffix arrays, we can solve the SLP *q*-gram frequency problem in O(qn) time.

## Algorithm for *q*-gram Non-overlapping Frequencies Problem

#### How to deal with crossing occurrences

+ A crossing occurrence of P in  $X_i$  may overlap a crossing occurrence of P in another variable  $X_l$ .



#### How to deal with crossing occurrences

- + A crossing occurrence of P in  $X_i$  may overlap a crossing occurrence of P in another variable  $X_l$ .
- our strategy is that for string P which crosses  $X_i$ ,
  - we compute the interval of maximum chain of overlapping occurrences of *P*.
  - we count *nOcc* in its interval only if the interval cannot be extended either to the left or right.



![](_page_40_Figure_1.jpeg)

27

• Notice that  $loc_q$  in the variable  $X_i$  may or may not extend in  $X_i$ 's ancestor.

![](_page_41_Figure_2.jpeg)

![](_page_42_Figure_1.jpeg)

29

\* Let (b, e) is a pair of beginning and ending point of  $loc_q$  in the variable  $X_i$ , we obtain this lemma

Lemma 4  

$$nOcc(T,P) = \sum_{i=1}^{n} nOcc(X_i[b:e], P) \cdot vOcc(X_i)$$
where  $(b, e)$  is  $loc_q$  of  $P$  in  $T$ , and  $[b-q+1:e+q-1]$  crosses  $X_i$   
 $(0 < b-q+1 \le |X_i| \le e+q-1 \le |X_i|)$ 

#### Lemma 5

For all variables  $X_i = X_l X_r$ , and integer  $|X_l| - (2q-1) \le j \le |X_l| + q-1$ ,  $loc_q(X_i, j) = (b, e)$  and  $nOcc(X_i[b:e], X_i[j:j+q-1]]$  can be computed in  $O(q^2n)$  time and O(qn) space.

#### Theorem 2

Non-overlapping frequencies problem on SLP can be solved in  $O(q^2n)$  time and O(qn) space.

#### Summary

	Uncompressed	SLP	SLP
	String	(Previous Work)	(This Work)
<i>q</i> -gram Freq	O( T )time and space	O(qn) [2] time and space	_
2-gram non-	O( T )	$O(n^4 \log n) \text{ time}$ $O(n^3) \text{ space}$	O(n)
	time and space		time and space
<i>q</i> -gram non- overlapping Freq	$\begin{array}{c c} O( T ) \\ \text{time and space} \end{array}$	_	$O(q^2n)$ time $O(qn)$ space

[1] [Inenaga and Bannai, 2009][2] [Goto et al, 2011]

Thank you for your altention

![](_page_47_Figure_0.jpeg)

$$nOcc \text{ within } loc = nOcc \text{ within } (2) + nOcc \text{ within } (3) + nOcc \text{ within } (1)$$

### no q-gram頻度の重み付け

 すべての変数で対応する *loc<sub>q</sub>*の範囲が確定した *q*-gram の 重み付けを行う

![](_page_48_Figure_2.jpeg)

![](_page_48_Figure_3.jpeg)

*locq*の範囲が確定

なにもしない

重み0

*locq*の範囲が確定 対応するq-gramがまたがる

重み  $P: vOcc(X_i) *$ *locq*のno頻度

重み  $P: vOcc(X_i) * loc_qのno頻度$ 

対応するq-gramがまたがらない

### no q-gram頻度の重み付け

各変数に対しまたがるq-gramは  $O(q^2)$ で $loc_q$ のno頻度を求めることが可能

![](_page_49_Figure_2.jpeg)

*loc*のno頻度 = *lsc*のno頻度 + *lpc*のno頻度 + ①の範囲のno頻度

全ての変数について $loc_q$ のno頻度を求めることで SLP no頻度問題を $O(q^2n)$ で解くことが出来る

![](_page_50_Figure_0.jpeg)

*nOcc* within *loc* = *nOcc* within *loc* + *nOcc* within *loc* + *nOcc* within range of ①

#### Crossing Occurrence

• For all variables  $X_i$ , we count non-overlapping frequencies of loc which crosses  $X_i$  rather than q-gram which crosses  $X_i$ 

![](_page_51_Figure_2.jpeg)

#### Crossing Occurrence

• For all variables  $X_i$ , we count non-overlapping frequencies of loc which crosses  $X_i$  rather than q-gram which crosses  $X_i$ 

![](_page_52_Figure_2.jpeg)

If b < q or  $|X_i|-q < e$ , *loc* can be overlaped the *q*-gram crossing in  $X_i$ 's ancestor

We count the crossing interval [b-q+1:e+q-1]

#### Difference of Algorithm for Overlapping Frequencies Problem

• overlapping algorithm doesn't care whether or not a q-gram which crosses in  $X_i$  is overlapping the same q-gram in  $X_l$  or  $X_r$ 

![](_page_53_Figure_2.jpeg)

non-overlapping frequency of *P* is 2, but algorithm count 3 time  $X_i, X_l, X_r$ 

#### Difference of Algorithm for Overlapping Frequencies Problem

 Important Key is that non-overlapping frequency depend on the way of overlapping in the chain of overlapping

![](_page_54_Figure_2.jpeg)

We want to know

• the interval of the maximum chain of

overlapping occurrences of *q*-gramnon-overlapping frequency in the interval

longest overlapping cover

![](_page_55_Picture_3.jpeg)

longest overlapping cover

![](_page_56_Figure_3.jpeg)

longest overlapping cover

![](_page_57_Figure_3.jpeg)

longest overlapping cover

![](_page_58_Figure_3.jpeg)

#### Difference of Algorithm for Overlapping Frequencies Problem

• overlapping algorithm doesn't care whether or not a q-gram which crosses in  $X_i$  is overlapping the same q-gram in  $X_l$  or  $X_r$ 

![](_page_59_Figure_2.jpeg)

non-overlapping frequency of *P* is 2, but algorithm count 3 time  $X_i, X_l, X_r$ 

#### Difference of Algorithm for Overlapping Frequencies Problem

 Important Key is that non-overlapping frequency depend on the way of overlapping in the chain of overlapping

![](_page_60_Figure_2.jpeg)

We want to know

• the interval of the maximum chain of

overlapping occurrences of *q*-gramnon-overlapping frequency in the interval