Space Efficient Linear Time Lempel-Ziv Factorization for Small Alphabet

Kyushu University Keisuke Goto, Hideo Bannai

- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies the following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before position $i, f_j = T[i]$
 - 2) otherwise, f_j is the longest substring that occurs at *i* and somewhere before *i*



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies the following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before position $i, f_j = T[i]$
 - 2) otherwise, f_j is the longest substring that occurs at *i* and somewhere before *i*



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$



- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$

2) otherwise, f_j is the longest substring that occurs at *i* and somewhere before *i*

Example $f_1 f_2$ f_3 f_4 f_5 f_6 f_7 f_8 T =a b a b a b a b a b a b a a a a a b b a b a b

- For a string *T* of length *N*, the LZ77 factorization of *T* is a factorization $T = f_1 f_2 \dots f_z$ such that each factor satisfies following conditions
- + For a factor f_j , let $i = |f_1 f_2 ... f_{j-1}| + 1$,
 - 1) if T[i] does not occur before, $f_j = T[i]$

2) otherwise, f_j is the longest substring that occurs at *i* and somewhere before *i*

Algorithm by using PrevOcc, LPF Arrays

• Consider the arrays that store PrevOcc(i) and LPF(i) for all postions i(1) For a factor f_j starting at i, output PrevOcc(i) and LPF(i)

(2) Update i = i + LPF(i), and repeat (1) when i < N



Previous Works Before 2013

* Most previous works focus on how to efficiently compute the *LPF* and *PrevOcc* arrays

	Sp	bace	required Integer Arrays							
Algorithm	Stack	# of Integer Arrays of	LCP	LPF	Prev Occ	SA	PSV	NSV	SA-1	
Crachemora & Ilia 2008		5	~		~	V	V	V		
Crochemore & me, 2008	~	4	V	V	V	V				
	V	4	V	V	V	V				
Chen+, 2008	~	3	V	V	~					
	~	2	V		V					
Crochemore+, 2008	V	4	V	V	V	V				
Crochemore+, 2009	V	4	V	V	V	V				
Ohlenbusch & Gog, 2011		3		V	V	V				

Previous Works in 2013

+ In 2013, algorithms not based on *LPF* and *PrevOcc* were proposed.

• Karkkainen+ reduced integer arrays of length N from three to two.

	5	required Integer Arrays								
Algorithm	Algorithm Name		# of Integer Arrays of length N	LPF	Prev Occ	SA	PSV	NSV	Φ	SA-1
Ohlenbusch & Gog,2011	OG		3	V	V	V				
Goto & Ronnoi 2012	BG4	V	4			~	V	V		V
Goto & Dannai, 2015	BG3		3			V		V	V	
V = 1 + 1 + 2 + 2 + 2 + 2 + 2 + 2 + 2 + 2 +	KKP3		3			V	V	~		
	KKP2		2			V		V	/	

Our Contribution

We propose a linear time algorithm that uses a single integer array of length N in addition to $O(\sigma \log N)$ bits, where σ is alphabet size

	(Space	required Integer Arrays									
Algorithm	Name	Stack	# of Integer Arrays of length N	LPF	Prev Occ	SA	PSV	NSV	Φ	SA-1		
Ohlenbusch & Gog,2011	OG		3	く	V	く						
Cata & Dannai 2012	BG4	V	4			V	~	~		V		
Ooto & Dannai, 2015	BG3		3			V		V	V			
Körkköinon⊥ 2012	KKP3		3			V	V	~				
Karkkannen ⁺ , 2013	KKP2		2			く		V	/			
Cata & Damai 2014	BG2		2		V		/	~				
Outo & Dannai, 2014	BG1		1*					V	/			

* require additional space of $O(\sigma \log N)$ bits $\frac{14}{46}$

KKP3 / BG3

compute LPF(i), PrevOcc(i) by naive character comparison between T[i..N] and T[k..N] for all k < i



compute LPF(i), PrevOcc(i) by naive character comparison between T[i..N] and T[k..N] for all k < i



compute LPF(i), PrevOcc(i) by naive character comparison between T[i..N] and T[k..N] for all k < i



compute LPF(i), PrevOcc(i) by naive character comparison between T[i..N] and T[k..N] for all k < i



Idea of KKP3 / BG3

Lemma [Crochemore and Ilie, 2008]

The candidates of *PrevOcc*(*i*) can be reduced to 2 positions; the lexicographic predecessor and successor of suffix $i_{\underline{sep}}$ in suffixes that start at previous positions



Idea of KKP3 / BG3

Lemma [Crochemore and Ilie, 2008]

PrevOcc(SA[i]) = PSV(SA[i]) OR NSV(SA[i]) $PSV(SA[i]) = \max\{SA[j] \mid j < i, SA[j] < SA[i]\}$ $NSV(SA[i]) = \min\{SA[j] \mid j > i, SA[j] < SA[i]\}$



Idea of KKP3 / BG3

Lemma [Crochemore and Ilie, 2008]

PrevOcc(SA[i]) = PSV(SA[i]) OR NSV(SA[i]) $PSV(SA[i]) = \max\{SA[j] \mid j < i, SA[j] < SA[i]\}$ $NSV(SA[i]) = \min\{SA[j] \mid j > i, SA[j] < SA[i]\}$



For a factor f_j , the number of comparison is $2 |f_j|$ Total # of comparison is $\Sigma 2 |f_j| = O(N)$

Linear Time Computation of PSV and NSV Arrays

Lemma [Crochemore and Ilie, 2008]

All *PSV* and *NSV* values can be obtained in linear time by sequentially scanning *SA* from left to right OR right to left



Overview of KKP3

KKP3 run in linear time and 3N log N bits space



KKP2

 Φ Array

[Kärkkäinen+, 2009]

 Φ array is an array that stores each lexicographic predecessor of each suffix in text order.

 $\Phi[SA[i]] = SA[i-1] (i > 1), \ \Phi[SA[1]] = N$

NOTE: all SA[i] can be obtained from right to left by Φ

i	$\Phi[i]$	SA[i]	T[SA[i]]
1	4	-3	aaacatat
2	7	- 4	aacatat
3	0	1	acaaacatat
4	3	5	acatat
5	1	9	at
6	2	7	atat
7	9	2	caaacatat
8	10	6	catat
9	5	10	t
10	6	€8♥	tat

T = acaaacatat

Lemma [Kärkkäinen+, 2013]

 Φ array can be constructed in linear time and in-place from *NSV(PSV)* array by sequentially scanning from left to right

NOTE: for each step *i*, <u>*PSV(i)*</u>, <u>*NSV(i)*</u> can be obtained



Overview of KKP2

KKP2 runs in linear time and 2N log N bits space



31/46

BG2

Idea of BG2

If we can rewrite *SA* to *PSV* array, LZ77 factorization can be computed by two integer arrays



Idea of BG2

It can be accomplished through Φ array



34/46

Observation

- By scanning *SA*[*i*] from left to right, *PSV*(*SA*[*i*]), *NSV*(*SA*[*i*]) can be obtained
- * *SA*[*i*] never be read after read once

Is it possible to store *NSV*[*SA*[*i*]] to *SA*[*i*]?

It seems IMPOSSIBLE because they have a different ordering

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
SA	9	11	2	8	13	1	14	6	3	10	16	5	12	7	15	4



Observation

+ All SA[i] can be obtained from right to left by Φ

+ $\Phi[SA[i]]$ never be used after read once

PSV and *NSV* arrays can be computed by Φ







Overview of BG2

BG2 runs in linear time and 2N log N bits space



BG1

Idea of BG1

• If given Φ not SA, LZ77 factorization can be computed in $N \log N$ bits space by using ideas of BG2 and KKP2



In-place computation of Φ array

Lemma

For a string *T* of length *N*, the Φ array of *T* can be computed in linear time and $O(\sigma \log N)$ bits additional space

* Idea: simulating Nong's suffix array construction algorithm on Φ array

Theorem [Nong, 2013]

For a string *T* of length *N*, the *SA* of *T* can be computed in linear time and $O(\sigma \log N)$ bits additional space

Overview of BG1

BG1 runs in linear time and $N \log N + O(\sigma \log N)$ bits space



Computational Experiments

- Nevertheless BG1 uses a third space than KKP3,
 the runtime of BG1 does not over 2.5 times than that of KKP3
- Nevertheless BG1 uses a half space than KKP2,
 the runtime of BG1 does not over 2 times than that of KKP2



corpus: http://www.cas.mcmaster.ca/~bill/strings 43/46

DNA 100MB $\sigma = 16$



Summary

We proposed a space efficient linear time LZ77 factorization algorithm for small alphabets, which uses $N \log N$ bits + $O(\sigma \log N)$ bits of space

Future Perspective

•Can we reduce the space $O(\sigma \log N)$ bits of BG1?

It seems difficult since we have to compute Φ array truly in-place.
We should perhaps shift our focus to practically fast non-linear time algorithms which use less space than BG1.