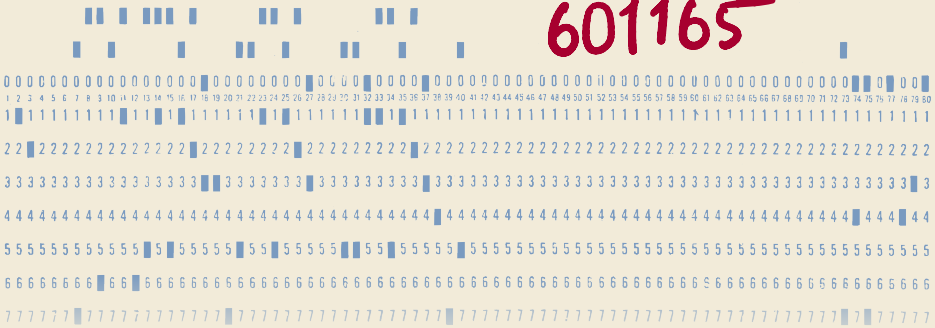


601165



Lecture 15: Data Compression III

COMP526: Efficient Algorithms

Updated: November 21, 2024

Will Rosenbaum
University of Liverpool

Announcements

1. Programming Assignment 2 posted
 - Due 29 November
2. Quiz 6 due Friday
 - Covers Lecture 13 material
 - 1 Question, Short Answer ←
 - Usual rules apply
3. Attendance Code:

601165

Meeting Goals

1. Give a recap of Lempel-Ziv-Welch (LZW) encoding.
2. Describe LZW decoding procedure
3. Introduce text transformations (non-compressing):
 - Move-to-front (MTF) transform
 - Burrows-Wheeler transform

Lempel-Ziv- Welch Encoding

From Last Time



Lempel-Ziv-Welch encoding idea:

- Fixed length codewords, size k
 - 2^k possible codewords
- Each codeword represents a *string* of text
 - initial codewords correspond to source text alphabet Σ_S (strings of length 1)
- Store a dictionary D that maps strings to codewords
- Encoding scheme:
 - scan source text S sequentially
 - if we see a string xc where string x is in dictionary but xc is not
 - append $D[x]$ to encoded string C
 - add xc to D

already has associated codeword

Example

Consider $S = \text{N A N A S B A N A N A S}$

x

$xc = \boxed{\text{NA}}$

$C = 0010$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	
0110	
0111	
1000	
1001	
1010	
1011	

Example

Consider $S = \underline{N} \boxed{AN} A S B A N A N A S$


$C = 0010 \quad 0000$

→

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	
0111	
1000	
1001	
1010	
1011	


Example

Consider $S = \text{N A N A S B A N A N A S}$



$C = 0010 \ 0000$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	
0111	
1000	
1001	
1010	
1011	



Example

Consider $S = N A N A S B A N A N A S$



$C = 0010 \ 0000 \ 0100$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	NAS
0111	
1000	
1001	
1010	
1011	

Example

Consider $S = \text{N A N A S B A N A N A S}$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	NAS
0111	SB
1000	
1001	
1010	
1011	

$C = 0010\ 0000\ 0100\ 0011$

Example

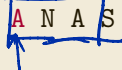
Consider $S = N A N A S B \bar{A} N A N A S$
 $\uparrow \uparrow \uparrow$

$C = 0010 \ 0000 \ 0100 \ 0011 \ 0001 \ 0101$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
→ 0101	AN
0110	NAS
0111	SB
1000	BA
1001	ANA
1010	
1011	

Example

Consider $S = N A N A S B A N A N A S$



$C = 0010 \ 0000 \ 0100 \ 0011 \ 0001 \ 0101 \ 1001$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	NAS
0111	SB
1000	BA
1001	ANA
1010	ANAS
1011	

Example

Consider $S = N A N A S B A N A N A S$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	NAS
0111	SB
1000	BA
1001	ANA
1010	ANAS
1011	

$C = 0010\ 0000\ 0100\ 0011\ 0001\ 0101\ 1001$

Example

Consider $S = \text{N A N A S B A N A N A S}$

code	string
0000	A
0001	B
0010	N
0011	S
0100	NA
0101	AN
0110	NAS
0111	SB
1000	BA
1001	ANA
1010	ANAS
1011	

$C = 0010 \mid 0000 \mid 0100 \mid 0011 \mid 0001 \mid 0101 \mid 1001 \mid 0011$

LZW in Pseudocode

```
1: procedure LZWENCODE( $S[0..n]$ )
2:    $x \leftarrow \varepsilon$ 
3:    $C \leftarrow \varepsilon$ 
4:    $D \leftarrow$  all  $c \in \Sigma_S$ 
5:    $k \leftarrow |\Sigma_S|$ 
6:   for  $i = 0, 1, \dots, n-1$  do
7:      $\rightarrow c \leftarrow S[i]$  cur char
8:     if  $D.CONTAINSKEY(xc)$  then
9:        $x \leftarrow xc$ 
10:    else
11:       $C \leftarrow CD.GET(x)$ 
12:       $D.PUT(xc, k)$ 
13:       $k \leftarrow k + 1$ ,  $x \leftarrow c$ 
14:    end if
15:  end for
16:   $\rightarrow C \leftarrow CD.GET(x)$ 
17: end procedure
```

- ▷ previous word, initially empty
- ▷ output, initially empty
- ▷ dictionary of codewords
- ▷ next free codeword

*iterate over
chars of S*

- ▷ append codeword for x

Decoding LZW

Encoding S according to LZW is reasonably simple...

...but how do we **decode** C ?

Decoding LZW

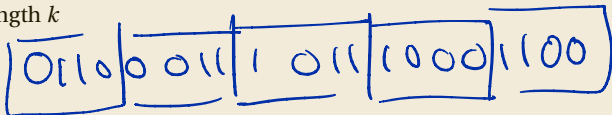
Encoding S according to LZW is reasonably simple...

...but how do we **decode** C ?

Observations.

- All codewords have length k

$$k=4$$



Decoding LZW

Encoding S according to LZW is reasonably simple...

...but how do we **decode** C ?

Observations.

- All codewords have length k
- *Given* the dictionary D , we can compute the inverse map
 - technically, this depends on the representation of D
 - representing D as a **trie** data structure makes this efficient

Decoding LZW

Encoding S according to LZW is reasonably simple...

...but how do we **decode** C ?

Observations.

- All codewords have length k
- *Given* the dictionary D , we can compute the inverse map
 - technically, this depends on the representation of D
 - representing D as a **trie** data structure makes this efficient

Question. How efficient is it to store D ?

Decoding LZW

Encoding S according to LZW is reasonably simple...

...but how do we **decode** C ?

Observations.

- All codewords have length k
- *Given* the dictionary D , we can compute the inverse map
 - technically, this depends on the representation of D
 - representing D as a **trie** data structure makes this efficient

Question. How efficient is it to store D ?

PollEverywhere Question

How many phrases (words) will LZW create on input $S = a^n$ (a run of n as)?

1. $\sim n$
2. $\sim n/2$
3. $\Theta(n/\log n)$
4. $\Theta(\sqrt{n})$
5. $\Theta(\log n)$
6. $\Theta(1)$



polllev.com/comp526

Decoding LZW

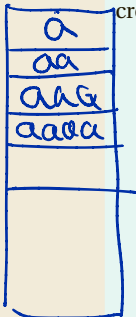
Question. How efficient is it to store D ?

~~a~~ a a a a a a a a a a a a

1
2
3
4
5
⋮

add these together to get n

D :



$k = \text{dict size}$ smallest k ?

$$1 + 2 + 3 + \dots + k \geq n$$

PollEverywhere Question

How many phrases (words) will LZW create on input $S = a^n$ (a run of n as)?

1. $\sim n$
2. $\sim n/2$
3. $\Theta(n/\log n)$
4. $\Theta(\sqrt{n})$
5. $\Theta(\log n)$
6. $\Theta(1)$



pollev.com/comp526

$k^2 \approx \sqrt{2n}$

$$\Theta(\sqrt{n})$$

$$\frac{k(k+1)}{2} \approx \frac{k^2}{2} \approx n$$

Decoding LZW

Question. How efficient is it to store D ?

Important question. Given D could be larger than S for very predictable strings, can we avoid storing D altogether?

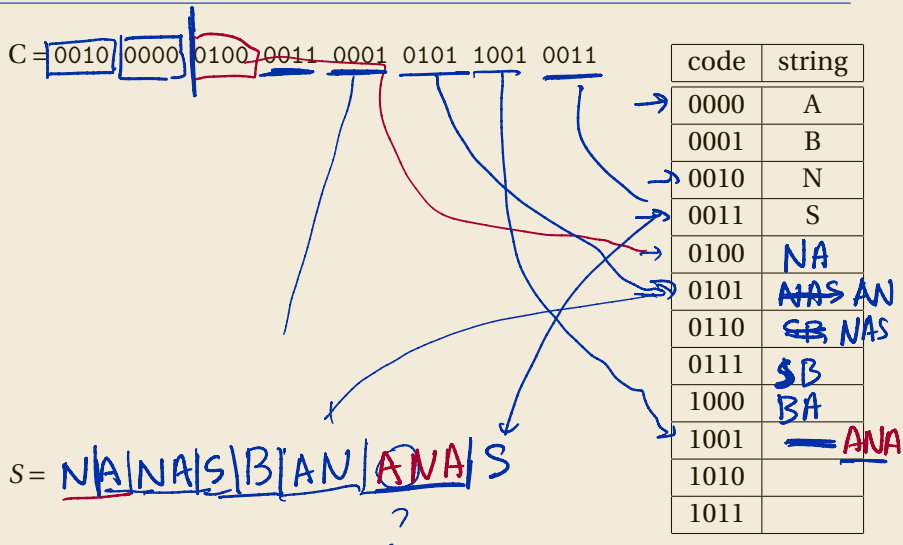
Decoding LZW

Question. How efficient is it to store D ?

Important question. Given D could be larger than S for very predictable strings, can we avoid storing D altogether?

- **Try:** given C and the start of D (containing only Σ), reconstruct D as we decode

Decoding LZW Example



Decoding LZW Pseudocode

This always works!

```
1: procedure LZWDECOMPRESS( $C, \Sigma$ )
2:    $\rightarrow D \leftarrow$  dictionary, initialized with codes for  $c \in \Sigma_S$ 
3:    $\rightarrow k \leftarrow |\Sigma_S|, q \leftarrow C[0]$ 
4:    $y \leftarrow D[q]$ 
5:    $S \leftarrow y$ 
6:   for  $j = 1, 2, \dots, |C| - 1$  do
7:      $x \leftarrow y$ 
8:      $q \leftarrow C[j]$ 
9:     if  $q = k$  then
10:       $y \leftarrow xx[0]$ 
11:     else
12:       $y \leftarrow D[q]$ 
13:     end if
14:      $S \leftarrow Sy$ 
15:      $D[k] \leftarrow xy[0]$ 
16:      $k \leftarrow k + 1$ 
17:   end for
18:   return  $S$ 
19: end procedure
```

Decoded string

tricky case

consider text starts aaa (get stuck after decoding first char.)

\triangleright stored as array

\triangleright first "new" codeword; first codeword

\triangleright the first phrase (single character)

\triangleright output the first phrase

\triangleright remember decoded phrase

\triangleright read next codeword

\triangleright bootstrap case

\triangleright append decoded phrase

\triangleright store new phrase

LZW Discussion

Some Details.

- Implementing the dictionary
 - use a trie data structure (more later...)
- Coded alphabet $\Sigma_C = \{0..2^k\}$, others are possible (e.g., Huffman)
- What happens when dictionary is full? ←
 - start using longer codewords?
 - flush dictionary and start from scratch? F
- Encoding and decoding both run in linear time! (assuming $|\Sigma_S| = O(1)$) ←

Appraisal.

- Fast encoding and decoding
- Works in **streaming model**
- Significant compression for many types of data (45% for English text)
- [Only captures local repetitions
 - e.g., not maximally helpful for $S = \underline{TT}$

Compression Summary

Huffman codes	Run-length encoding	Lempel-Ziv-Welch
fixed-to-variable	variable-to-variable	variable-to-fixed →
2-pass	1-pass	1-pass
must send dictionary	can be worse than ASCII	can be worse than ASCII ←
60% compression on English text	bad on text	45% compression on English text
optimal binary character encoding	good on long runs (e.g., pictures)	good on English text
rarely used directly	rarely used directly	frequently used
part of pkzip, JPEG, MP3	fax machines, old picture- formats	GIF, part of PDF, Unix compress

Text Transforms

Why Settle for What You're Given?

So far we've used the following pipeline:

source text → encoded text

a single algorithm (Huffman, RLE, LZW) to encode the text.

Why Settle for What You're Given?

So far we've used the following pipeline:

source text → encoded text

a single algorithm (Huffman, RLE, LZW) to encode the text.

Inefficiencies:

- Huffman, RLE, and LZW are **limited** in patterns they exploit
 - Huffman: character encoding, changing probabilities
 - RLE: only compresses long runs
 - LZW: only small local repetitions identified

Why Settle for What You're Given?

So far we've used the following pipeline:

source text → encoded text

a single algorithm (Huffman, RLE, LZW) to encode the text.

Inefficiencies:

- Huffman, RLE, and LZW are **limited** in patterns they exploit
 - Huffman: character encoding, changing probabilities
 - RLE: only compresses long runs
 - LZW: only small local repetitions identified

Idea: Text Transformations

- *Reversible* function of text
- does not by itself compress text
- makes text more compressible to algorithms above

Pipeline: source → transformed text → compressed transformed text

reversible

Move to Front Transform

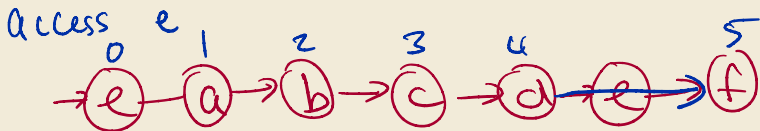
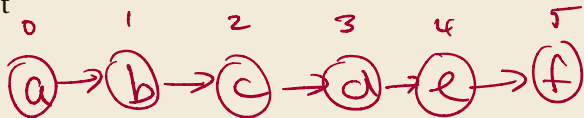
Move to Front Heuristic

Simple Observation. In many texts, a recently used character is likely to be reused

- not just globally frequent characters (i.e., Huffman compressible)
- have local runs with repetition

Idea. Start with indexed alphabet Σ (e.g., linked list)

- Whenever a character c is read, move that character to the front of the list



Self Adjusting Lists

Lists

- elements have sequential indexes (like arrays)
- elements can be removed/inserted
 - other elements' indices are shifted

Self-adjustment

- when element c is accessed, move c to the front of list

Self Adjusting Lists

Lists

- elements have sequential indexes (like arrays)
- elements can be removed/inserted
 - other elements' indices are shifted

Self-adjustment

- when element c is accessed, move c to the front of list

XYZABC
A XYZBC ✓

PollEverywhere Question

Suppose we start with a list containing XYZABC and the next access is to the character A. What is the state of the list after the access?



pollev.com/comp526

Self Adjusting Lists

Lists

- elements have sequential indexes (like arrays)
- elements can be removed/inserted
 - other elements' indices are shifted

Self-adjustment

- when element c is accessed, move c to the front of list

PollEverywhere Question

Suppose we start with a list containing XYZABC and the next access is to the character A. What is the state of the list after the access?



pollev.com/comp526

Puzzle. How to implement lists *efficiently*?

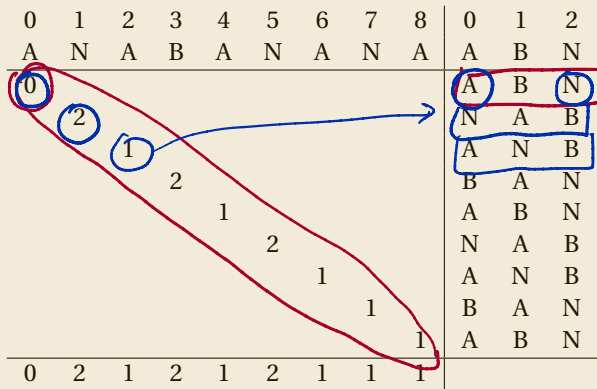
Move to Front Simple Example

0	1	2	3	4	5	6	7	8	0	1	2
A	N	A	B	A	N	A	N	A	A	B	N
0									A	B	N
2									N	A	B
1									A	N	B
2									B	A	N
1									A	B	N
2									N	A	B
1									A	N	B
1									N	A	B
1									A	N	B

Move to Front Simple Example

0	1	2	3	4	5	6	7	8	0	1	2
A	N	A	B	A	N	A	N	A	A	B	N
<hr/>									A	B	N
0									N	A	B
	2								A	N	B
		1							B	A	N
			2						A	B	N
				1					N	A	B
					2				A	N	B
						1			B	A	N
							1		A	B	N
								1	<hr/>		
0	2	1	2	1	2	1	1	1			

Move to Front Simple Example



A
N
A

Observations.

- This process is reasonably efficient
 - store alphabet in a linked list
- This process is **reversible**
 - How?

MTF Heuristic

S	A	N	A	B	A	N	A	N	A
C	0	2	1	2	1	2	1	1	1

Questions.

- What does a *run* in *S* correspond to in *C*?

$a a a a a \dots a \mapsto ? 0 0 0 \dots 0$

MTF Heuristic

S	A	N	A	B	A	N	A	N	A
C	0	2	1	2	1	2	1	1	1

Questions.

- What does a *run* in S correspond to in C ?
 - run in S of length k gives a run of 0s in C of length $k-1$

- What does a run in C correspond to in S ?

C : 11111111

C : 222222



lablablablab
abcabcabc

MTF Code

```
1: procedure MTFENCODE( $S[0..n]$ )
2:  $L \leftarrow$  list containing  $\Sigma_S$  (sorted)
3:  $C \leftarrow \varepsilon$ 
4: for  $i = 0, 1, \dots, n-1$  do
5:    $c \leftarrow S[i]$ 
6:    $p \leftarrow$  position of  $c$  in  $L$ 
7:    $C \leftarrow Cp$ 
8:   Move  $c$  to front of  $L$ 
9: end for
10: return  $C$ 
11: end procedure
```

Handwritten annotations:
- Red arrows point from the word "index" to the variable i in line 4 and the variable i in line 5.
- A red box highlights the word "position" in line 6.
- Red arrows point from the word "position" to the variable c in line 6 and the variable p in line 6.

MTF Code

```
1: procedure MTFENCODE( $S[0..n]$ )
2:    $L \leftarrow$  list containing  $\Sigma_S$  (sorted)
3:    $C \leftarrow \varepsilon$ 
4:   for  $i = 0, 1, \dots, n-1$  do
5:      $c \leftarrow S[i]$ 
6:      $p \leftarrow$  position of  $c$  in  $L$ 
7:      $C \leftarrow Cp$ 
8:     Move  $c$  to front of  $L$ 
9:   end for
10:  return  $C$ 
11: end procedure
```

```
1: procedure MTFDECODE( $C[0..m]$ )
2:    $L \leftarrow$  list containing  $\Sigma_S$  (sorted)
3:    $S \leftarrow \varepsilon$ 
4:   for  $j = 0, 1, \dots, m-1$  do
5:      $p \leftarrow S[j]$ 
6:      $c \leftarrow$  char at position  $p$  in  $L$ 
7:      $S \leftarrow Sc$ 
8:     Move  $c$  to front of  $L$ 
9:   end for
10:  return  $S$ 
11: end procedure
```

MTF Code

```
1: procedure MTFENCODE( $S[0..n]$ )
2:    $L \leftarrow$  list containing  $\Sigma_S$  (sorted)
3:    $C \leftarrow \varepsilon$ 
4:   for  $i = 0, 1, \dots, n-1$  do
5:      $c \leftarrow S[i]$ 
6:      $p \leftarrow$  position of  $c$  in  $L$ 
7:      $C \leftarrow Cp$ 
8:     Move  $c$  to front of  $L$ 
9:   end for
10:  return  $C$ 
11: end procedure
```

```
1: procedure MTFDECODE( $C[0..m]$ )
2:    $L \leftarrow$  list containing  $\Sigma_S$  (sorted)
3:    $S \leftarrow \varepsilon$ 
4:   for  $j = 0, 1, \dots, m-1$  do
5:      $p \leftarrow S[j]$ 
6:      $c \leftarrow$  char at position  $p$  in  $L$ 
7:      $S \leftarrow Sc$ 
8:     Move  $c$  to front of  $L$ 
9:   end for
10:  return  $S$ 
11: end procedure
```

Key observation. Both encode and decode procedure perform same sequence of accesses to L

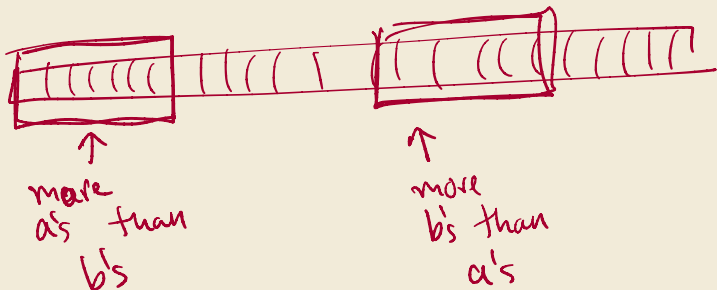
\Rightarrow the decoded letter is always the same as encoded (induction)

MTF Discussion

- MTF does not compress texts (assuming we use fixed-length codewords)
 - MTF is used as part of a longer pipeline
 - For many texts smaller codeword values are more likely after MTF

MTF Discussion

- MTF does not compress texts (assuming we use fixed-length codewords)
 - MTF is used as part of a longer pipeline
 - For many texts smaller codeword values are more likely after MTF
- Intuitive effect: MTF converts patterns with low local entropy to texts with small *global entropy*
 - Huffman is more effective after MTF transform applied



MTF Discussion

- MTF does not compress texts (assuming we use fixed-length codewords)
 - MTF is used as part of a longer pipeline
 - For many texts smaller codeword values are more likely after MTF
- Intuitive effect: MTF converts patterns with low *local entropy* to texts with small *global entropy*
 - Huffman is more effective after MTF transform applied
- Still, many natural texts do not have low local entropy
...but we can try to transform texts so that they *do* have low local entropy!

Burrows- Wheeler Transform

Burrows-Wheeler Transform

A Sophisticated Transform.

- Coded text has same characters as source, but in a different order
- Coded text is typically more compressible (local character frequencies)

Burrows-Wheeler Transform

A Sophisticated Transform.

- Coded text has same characters as source, but in a different order
- Coded text is typically more compressible (local character frequencies)

A Concession

- Encoding requires access to *all* of S
 - not streaming like LZW, or “two pass” like Huffman

Burrows-Wheeler Transform

A Sophisticated Transform.

- Coded text has same characters as source, but in a different order
- Coded text is typically more compressible (local character frequencies)

A Concession

- Encoding requires access to *all* of S
 - not streaming like LZW, or “two pass” like Huffman

An Effective Pipeline: BWT → MTF → RLE → Huffman

- Used by bzip2 compression program:

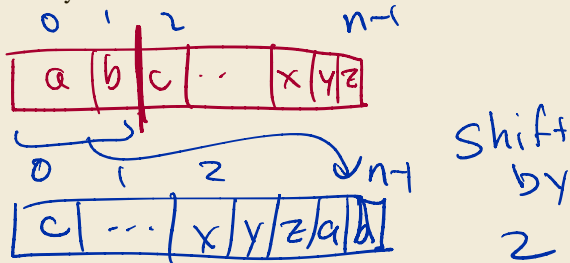
→	5458199	21 Nov 12:06	shakespeare-original.txt	
5.5 MB	1479261	21 Nov 12:04	shakespeare.txt.bz2	- 1.5 MB
	2024091	21 Nov 12:08	shakespeare.txt.gz	
2 MB	2022556	21 Nov 12:06	shakespeare.txt.zip	

↑
bzip2
utility

Cyclic Shifts

Definition. A **cyclic shift** of a string S is a re-indexing of S by some fixed offset with wraparound

- add terminating character \$ to show original end of S
- \$ is alphabetically before all other characters



Cyclic Shifts

Definition. A **cyclic shift** of a string S is a re-indexing of S by some fixed offset with wraparound

- add terminating character $\$$ to show original end of S
- $\$$ is alphabetically before all other characters

Example. All cyclic shifts of $S = \text{alf_eats_alfalfa}\$$:

```
alf_eats_alfalfa$
lf_eats_alfalfa$a
f_eats_alfalfa$a
_eats_alfalfa$a
eats_alfalfa$a
ats_alfalfa$a
ts_alfalfa$a
s_alfalfa$a
_alfalfa$a
alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
alf_eats_alfalfa$a
```

Burrows Wheeler Transform

A Simple Idea.

1. Form all cyclic shifts of S

```
alf_eats_alfalfa$  
lf_eats_alfalfa$a  
f_eats_alfalfa$al  
_eats_alfalfa$alf  
eats_alfalfa$alf_  
ats_alfalfa$alf_e  
ts_alfalfa$alf_ea  
s_alfalfa$alf_eat  
_alfalfa$alf_eats  
alfalfa$alf_eats_  
lfalfa$alf_eats_a  
falfa$alf_eats_alf  
alfa$alf_eats_alf  
lfa$alf_eats_alf  
fa$alf_eats_alfal  
a$alf_eats_alfalf  
$alf_eats_alfalfa
```

Burrows Wheeler Transform

A Simple Idea.

1. Form all cyclic shifts of S
2. Sort the shifts alphabetically

```
alf_eats_alfalfa$  
lf_eats_alfalfa$a  
f_eats_alfalfa$a  
_eats_alfalfa$a  
eats_alfalfa$a  
ats_alfalfa$a  
ts_alfalfa$a  
s_alfalfa$a  
_alfalfa$a  
alfalfa$a  
lfa$a  
falfa$a  
alfalfa$a  
lfa$a  
falfa$a  
alfalfa$a  
lfa$a  
falfa$a
```

sort
rows

alphabetical order

```
$alf_eats_alfalfa  
_alfalfa$a  
eats_alfalfa$a  
a$alf_eats_alfalf  
alf_eats_alfalfa$  
alf$a  
alfalfa$a  
eats_alfalfa$a  
f_eats_alfalfa$a  
fa$a  
falfa$a  
alfalfa$a  
lfa$a  
falfa$a  
alfalfa$a  
lfa$a  
falfa$a
```


Features of BWT

BWT:

1. Form all cyclic shifts of S
2. Sort the shifts alphabetically
3. Return the last **column** of the table

$S = \text{alf_eats_alfalfa\$}$

$B = \text{asff\$f_e_lllaata}$

Remarkably:

- This procedure can be computed in $O(n)$ time
 - “naive” algorithm/analysis is $O(n^2 \log n)$
- This procedure is reverseable!
 - we will see this soon

```
$alf_eats_alfalfa$
_alfalfa$alf_eat
_eats_alfalfa$al
a$alf_eats_alfal
alf_eats_alfalfa$
alfalfa$alf_eats_
alfalfa$alf_eats_
ats_alfalfa$alf_
eats_alfalfa$alf_
f_eats_alfalfa$al
fa$alf_eats_alfal
falfa$alf_eats_alf
lf_eats_alfalfa$a
lfa$alf_eats_alf_
lfalfa$alf_eats_
s_alfalfa$alf_eat
ts_alfalfa$alf_ea
```

What does BWT do?

	r		$\downarrow L[r]$	
alf_eats_alfalfa\$	0	\$alf_eats_alfalfa	a	16
lf_eats_alfalfa\$a	1	_alfalfa\$alf_eats	s	8
f_eats_alfalfa\$a_	2	_eats_alfalfa\$a_	f	3
_eats_alfalfa\$a_	3	a\$alf_eats_alfalf	f	15
eats_alfalfa\$a_	4	alf_eats_alfalfa\$		0
ats_alfalfa\$a_	5	alfa\$a_alf_eats_alf		12
ts_alfalfa\$a_	6	alfalfa\$a_alf_eats_		9
s_alfalfa\$a_	7	ats_alfalfa\$a_	e	5
alfalfa\$a	8	eats_alfalfa\$a_		4
alfalfa\$a_	9	f_eats_alfalfa\$a_		2
lfalfa\$a_	10	fa\$a_alf_eats_alfal		14
falfa\$a_	11	falfa\$a_alf_eats_alf		11
alfa\$a_	12	lf_eats_alfalfa\$a_		1
lfa\$a_	13	lfa\$a_alf_eats_alfal		13
fa\$a_	14	lfalfa\$a_alf_eats_alf		10
a\$a_	15	s_alfalfa\$a_alf_eat		7
\$alf_eats_alfalfa	16	ts_alfalfa\$a_alf_ea		6

orig text, written vertically

Observation 1. BWT contains the same characters as S

- characters are reordered

What does BWT do?

Shift
a to
end



	r		$\downarrow L[r]$
alf_eats_alfalfa\$	0	\$alf_eats_alfalfa	a 16
lf_eats_alfalfa\$a	1	_alfalfa\$alf_eats	s 8
f_eats_alfalfa\$a_	2	_eats_alfalfa\$a_	f 3
_eats_alfalfa\$a_	3	a\$alf_eats_alfalf	f 15
eats_alfalfa\$a_	4	alf_eats_alfalfa\$	0
ats_alfalfa\$a_	5	alfa\$alf_eats_alf	f 12
ts_alfalfa\$a_	6	alfalfa\$alf_eats_	9
s_alfalfa\$a_	7	ats_alfalfa\$a_	e 5
alfalfa\$a	8	eats_alfalfa\$a_	_ 4
alfalfa\$a_	9	f_eats_alfalfa\$a_	2
lfalfa\$a_	10	fa\$alf_eats_alfal	14
falfa\$a_	11	falfa\$alf_eats_alf	11
alfa\$a_	12	lf_eats_alfalfa\$a	a 1
lfa\$a_	13	lfa\$alf_eats_alfi	a 13
fa\$a_	14	lfalfa\$alf_eats_	a 10
a\$a_	15	s_alfalfa\$a_	e 7
\$alf_eats_alfalfa	16	ts_alfalfa\$a_	e 6

Observation 2. BWT groups characters by what follows

- repeated substrings give rise to runs in B
- a always followed by lf $\implies B$ contains a run of as

What does BWT do?

	r		$\downarrow L[r]$
alf_eats_alfalfa\$	0	\$alf_eats_alfalfa	a 16
lf_eats_alfalfa\$a	1	_alfalfa\$alf_eats	s 8
f_eats_alfalfa\$al	2	_eats_alfalfa\$alf	f 3
_eats_alfalfa\$alf	3	a\$alf_eats_alfalf	f 15
eats_alfalfa\$alf_	4	alf_eats_alfalfa\$	0
ats_alfalfa\$alf_e	5	alfa\$alf_eats_alf	f 12
ts_alfalfa\$alf_ea	6	alfalfa\$alf_eats_	9
s_alfalfa\$alf_eat	7	ats_alfalfa\$alf_e	5
_alfalfa\$alf_eats	8	eats_alfalfa\$alf_	4
alfalfa\$alf_eats_	9	f_eats_alfalfa\$al	2
lfalfa\$alf_eats_a	10	fa\$alf_eats_alfal	14
falfa\$alf_eats_al	11	falfa\$alf_eats_al	11
alfa\$alf_eats_alf	12	lf_eats_alfalfa\$a	1
lfa\$alf_eats_alfa	13	lfa\$alf_eats_alfa	13
fa\$alf_eats_alfal	14	lfalfa\$alf_eats_a	10
a\$alf_eats_alfalf	15	s_alfalfa\$alf_eat	7
\$alf_eats_alfalfa	16	ts_alfalfa\$alf_ea	6

Observation 3. BWT outputs are typically amenable to compression after applying MTF

Inverting BWT

Observation. Compression means nothing unless we can invert the procedure!

- What if we just sorted the characters of S directly and compressed the sorted string?

Inverting BWT

Observation. Compression means nothing unless we can invert the procedure!

A Magic Trick.

transformed
text

1. Create an array D of pairs $(B[r], r)$
2. Sort D stably with respect to first entry
3. Interpret D as linked list (follow links!)


Inverting BWT


Observation. Compression means nothing unless we can invert the procedure!

A Magic Trick.

1. Create an array D of pairs $(B[r], r)$
2. Sort D *stably* with respect to first entry
3. Interpret D as linked list (follow links!)

Example.

$B = \text{ard}\$r\text{caaaabb}$ 

- 
- 1: $(a, 1)$
 - 2: $(r, 2)$
 - 3: $(d, 3)$
 - 4: $(\$, 4)$
 - 5: $(r, 5)$
 - 6: $(c, 6)$
 - 7: $(a, 7)$
 - 8: $(a, 8)$
 - 9: $(a, 9)$
 - 10: $(a, 10)$
 - 11: $(b, 11)$
 - 12: $(b, 12)$

Inverting BWT

Observation. Compression means nothing unless we can invert the procedure!

A Magic Trick.

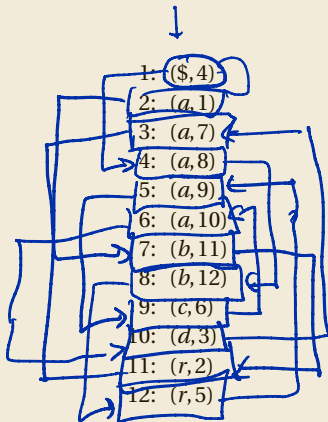
1. Create an array D of pairs $(B[r], r)$
2. Sort D *stably* with respect to first entry
3. Interpret D as linked list (follow links!)

Example.

$B = \text{ard\$rcaaaabb}$

$S = \text{abracadabra}$

- | | |
|-----|-------------|
| 1: | ($a, 1$) |
| 2: | ($r, 2$) |
| 3: | ($d, 3$) |
| 4: | ($\$, 4$) |
| 5: | ($r, 5$) |
| 6: | ($c, 6$) |
| 7: | ($a, 7$) |
| 8: | ($a, 8$) |
| 9: | ($a, 9$) |
| 10: | ($a, 10$) |
| 11: | ($b, 11$) |
| 12: | ($b, 12$) |



Inverting BWT

Observation. Compression means nothing unless we can invert the procedure!

A Magic Trick.

1. Create an array D of pairs $(B[r], r)$
2. Sort D *stably* with respect to first entry
3. Interpret D as linked list (follow links!)

1: $(a, 1)$	1: $(\$, 4)$
2: $(r, 2)$	2: $(a, 1)$
3: $(d, 3)$	3: $(a, 7)$
4: $(\$, 4)$	4: $(a, 8)$
5: $(r, 5)$	5: $(a, 9)$
6: $(c, 6)$	6: $(a, 10)$
7: $(a, 7)$	7: $(b, 11)$
8: $(a, 8)$	8: $(b, 12)$
9: $(a, 9)$	9: $(c, 6)$
10: $(a, 10)$	10: $(d, 3)$
11: $(b, 11)$	11: $(r, 2)$
12: $(b, 12)$	12: $(r, 5)$

Example.

$B = \text{ard}\$r\text{caaaabb}$

For Next Time. Convince yourself this inverts the BWT!

BWT Discussion

What do we know about the Burrows-Wheeler Transform?

- Running time $\Theta(n)$
 - encoding uses **suffix sorting** (future reference)
 - decoding can be done in $\Theta(n)$ time with counting sort
 - decoding is simpler/faster!
- Typically slower than other methods
- Needs access to entire text (or apply to smaller blocks)
- WBT \rightarrow MTF \rightarrow RLE \rightarrow Huffman has great compression!

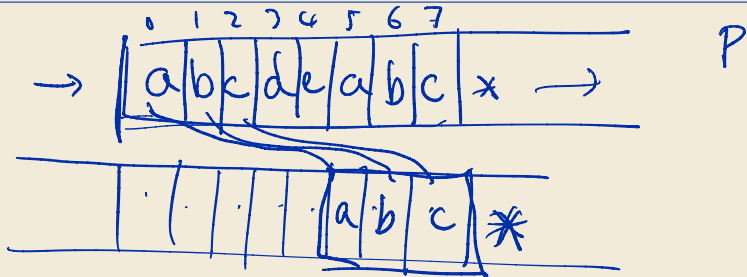
Summary of Compression

- Huffman** Variable-width, single-character (optimal in this case)
- RLE** Variable-width, multiple-character encoding
- LZW** Adaptive, fixed-width, multiple-character encoding
Augments dictionary with repeated substrings
- MTF** Adaptive, transforms to smaller integers
should be followed by variable-width integer encoding
- BWT** Block compression method, should be followed by MTF

Next Time

Error Correcting Codes

Scratch Notes



fail[8]