

CSEP 521  
Applied Algorithms  
Spring 2005

Dictionary Coding

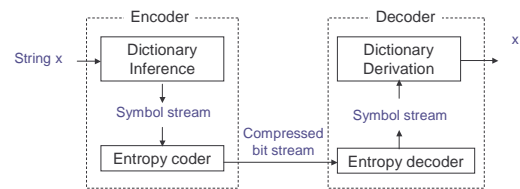
Plan for Tonight

- Overview
- LZW
- Sequitur
- Move-to-front coding
- Burrows-Wheeler Transform

Dictionary Coding

- Does not use statistical knowledge of data.
- Encoder: As the input is processed develop a dictionary and transmit the index of strings found in the dictionary.
- Decoder: As the code is processed reconstruct the dictionary to invert the process of encoding.
- Examples: LZW, LZ77, Sequitur, Burrows-Wheeler
- Applications: Unix Compress, gzip, bzip, GIF

Overview of Dictionary Compression



LZW Dictionary Inference Algorithm

Repeat  
 find the longest match w in the dictionary  
 output the index of w  
 put wa in the dictionary where a was the  
 unmatched symbol

LZW Encoding Example (1)

Dictionary                      a b a b a b a b a  
 0 a  
 1 b

### LZW Encoding Example (2)

Dictionary                      a b a b a b a b a  
0 a                                      0  
1 b  
2 ab

### LZW Encoding Example (3)

Dictionary                      a b a b a b a b a  
0 a                                      0 1  
1 b  
2 ab  
3 ba

### LZW Encoding Example (4)

Dictionary                      a b a b a b a b a  
0 a                                      0 1 2  
1 b  
2 ab  
3 ba  
4 aba

### LZW Encoding Example (5)

Dictionary                      a b a b a b a b a  
0 a                                      0 1 2 4  
1 b  
2 ab  
3 ba  
4 aba  
5 abab

### LZW Encoding Example (6)

Dictionary                      a b a b a b a b a  
0 a                                      0 1 2 4 3  
1 b  
2 ab  
3 ba  
4 aba  
5 abab

### LZW Dictionary Derivation Algorithm

- Emulate the encoder in building the dictionary.  
Decoder is slightly behind the encoder.

```
initialize dictionary;  
decode first index to w;  
put w? in dictionary;  
repeat  
  decode the first symbol s of the index;  
  complete the previous dictionary entry with s;  
  finish decoding the remainder of the index;  
  put w? in the dictionary where w was just decoded;
```

### LZW Decoding Example (1)

Dictionary                    0 1 2 4 3 6  
0 a                            a  
1 b  
2 a?

Lecture 6 - Dictionary Coding

13

### LZW Decoding Example (2a)

Dictionary                    0 1 2 4 3 6  
0 a                            a b  
1 b  
2 ab

Lecture 6 - Dictionary Coding

14

### LZW Decoding Example (2b)

Dictionary                    0 1 2 4 3 6  
0 a                            a b  
1 b  
2 ab  
3 b?

Lecture 6 - Dictionary Coding

15

### LZW Decoding Example (3a)

Dictionary                    0 1 2 4 3 6  
0 a                            a b a  
1 b  
2 ab  
3 ba

Lecture 6 - Dictionary Coding

16

### LZW Decoding Example (3b)

Dictionary                    0 1 2 4 3 6  
0 a                            a b ab  
1 b  
2 ab  
3 ba  
4 ab?

Lecture 6 - Dictionary Coding

17

### LZW Decoding Example (4a)

Dictionary                    0 1 2 4 3 6  
0 a                            a b ab a  
1 b  
2 ab  
3 ba  
4 aba

Lecture 6 - Dictionary Coding

18

### LZW Decoding Example (4b)

Dictionary

0	a	0 1 2 4 3 6
1	b	a b ab aba
2	ab	
3	ba	
4	aba	
5	aba?	

Lecture 6 - Dictionary Coding

19

### LZW Decoding Example (5a)

Dictionary

0	a	0 1 2 4 3 6
1	b	a b ab aba b
2	ab	
3	ba	
4	aba	
5	abab	

Lecture 6 - Dictionary Coding

20

### LZW Decoding Example (5b)

Dictionary

0	a	0 1 2 4 3 6
1	b	a b ab aba ba
2	ab	
3	ba	
4	aba	
5	abab	
6	ba?	

Lecture 6 - Dictionary Coding

21

### LZW Decoding Example (6a)

Dictionary

0	a	0 1 2 4 3 6
1	b	a b ab aba ba b
2	ab	
3	ba	
4	aba	
5	abab	
6	bab	

Lecture 6 - Dictionary Coding

22

### LZW Decoding Example (6b)

Dictionary

0	a	0 1 2 4 3 6
1	b	a b ab aba ba bab
2	ab	
3	ba	
4	aba	
5	abab	
6	bab	
7	bab?	

Lecture 6 - Dictionary Coding

23

### Decoding Exercise

Base Dictionary

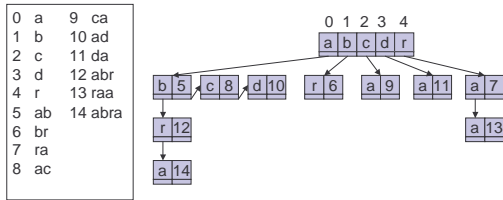
	0 1 4 0 2 0 3 5 7
0	a
1	b
2	c
3	d
4	r

Lecture 6 - Dictionary Coding

24

## Trie Data Structure for Encoder's Dictionary

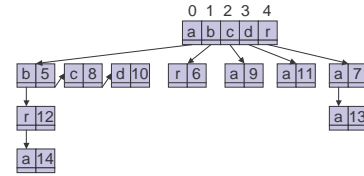
- Fredkin (1960)



Lecture 6 - Dictionary Coding

25

## Encoder Uses a Trie (1)

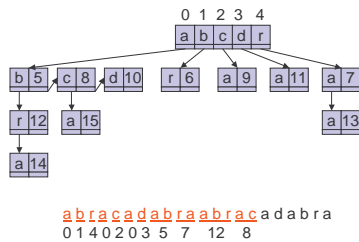


abra cadabra abra cadabra  
0 1 4 0 2 0 3 5 7 12

Lecture 6 - Dictionary Coding

26

## Encoder Uses a Trie (2)



abra cadabra abra cadabra  
0 1 4 0 2 0 3 5 7 12 8

Lecture 6 - Dictionary Coding

27

## Decoder's Data Structure

- Simply an array of strings

0	a	9	ca
1	b	10	ad
2	c	11	da
3	d	12	abr
4	r	13	raa
5	ab	14	abr?
6	br		
7	ra		
8	ac		

0 1 4 0 2 0 3 5 7 12 8 ...  
abra cadabra abra

Lecture 6 - Dictionary Coding

28

## Bounded Size Dictionary

- Bounded Size Dictionary
  - n bits of index allows a dictionary of size  $2^n$
  - Doubtful that long entries in the dictionary will be useful.
- Strategies when the dictionary reaches its limit.
  1. Don't add more, just use what is there.
  2. Throw it away and start a new dictionary.
  3. Double the dictionary, adding one more bit to indices.
  4. Throw out the least recently visited entry to make room for the new entry.

Lecture 6 - Dictionary Coding

29

## Notes on LZW

- Extremely effective when there are repeated patterns in the data that are widely spread.
- Negative: Creates entries in the dictionary that may never be used.
- Applications:
  - Unix compress, GIF, V.42 bis modem standard

Lecture 6 - Dictionary Coding

30

## Sequitur

- Nevill-Manning and Witten, 1996.
- Uses a context-free grammar (without recursion) to represent a string.
- The grammar is inferred from the string.
- If there is structure and repetition in the string then the grammar may be very small compared to the original string.
- Clever encoding of the grammar yields impressive compression ratios.
- Compression plus structure!

Lecture 6 - Dictionary Coding

31

## Context-Free Grammars

- Invented by Chomsky in 1959 to explain the grammar of natural languages.
- Also invented by Backus in 1959 to generate and parse Fortran.
- Example:
  - terminals: b, e
  - non-terminals: S, A
  - Production Rules:  
 $S \rightarrow SA, S \rightarrow A, A \rightarrow bSe, A \rightarrow be$
  - S is the start symbol

Lecture 6 - Dictionary Coding

32

## Context-Free Grammar Example

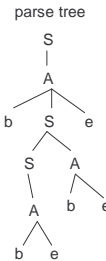
- $S \rightarrow SA$
- $S \rightarrow A$
- $A \rightarrow bSe$
- $A \rightarrow be$

Example: b and e matched as parentheses

derivation of bbebee

S  
A  
bSe  
bSAe  
bAAe  
bbeAe  
bbebee

hierarchical parse tree



Lecture 6 - Dictionary Coding

33

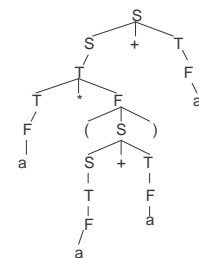
## Arithmetic Expressions

- $S \rightarrow S + T$
- $S \rightarrow T$
- $T \rightarrow T * F$
- $T \rightarrow F$
- $F \rightarrow a$
- $F \rightarrow (S)$

derivation of  $a * (a + a) + a$

S  
S+T  
T+T  
T\*F+T  
F\*F+T  
a\*F+T  
a\*(S)+F  
a\*(S+F)+T  
a\*(T+F)+T  
a\*(F+F)+T  
a\*(a+F)+T  
a\*(a+a)+F  
a\*(a+a)+F  
a\*(a+a)+a

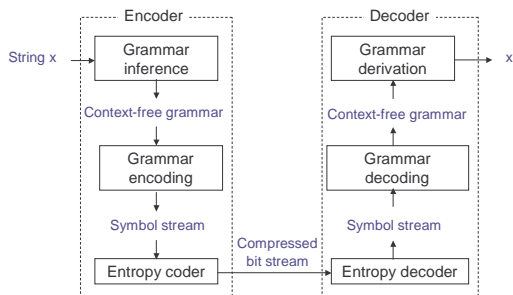
parse tree



Lecture 6 - Dictionary Coding

34

## Overview of Grammar Compression



Lecture 6 - Dictionary Coding

35

## Sequitur Principles

- Digram Uniqueness:
  - no pair of adjacent symbols (digram) appears more than once in the grammar.
- Rule Utility:
  - Every production rule is used more than once.
- These two principles are maintained as an invariant while inferring a grammar for the input string.

Lecture 6 - Dictionary Coding

36

### Sequitur Example (1)

bbeeebebebebeee

S → b

### Sequitur Example (2)

bbebeeebebebeee

S → bb

### Sequitur Example (3)

bbeeebebebebeee

S → bbe

### Sequitur Example (4)

bbbeeebebebebeee

S → bbbe

### Sequitur Example (5)

bbebeebebebebeee

S → bbebe

Enforce digram uniqueness.  
be occurs twice.  
Create new rule A → be.

### Sequitur Example (6)

bbebeebebebebeee

S → bAA  
A → be

### Sequitur Example (7)

bbebeebebbebee

S → bAAe  
A → be

### Sequitur Example (8)

bbeeebbebbebee

S → bAAeb  
A → be

### Sequitur Example (9)

bbeeebebebbebee

S → bAAebe      Enforce digram uniqueness.  
A → be          be occurs twice.  
Use existing rule A → be.

### Sequitur Example (10)

bbeeebebebbebee

S → bAAeA  
A → be

### Sequitur Example (11)

bbeeebebebbbebee

S → bAAeAb  
A → be

### Sequitur Example (12)

bbeeebebeebbbebee

S → bAAeAbe      Enforce digram uniqueness.  
A → be          be occurs twice.  
Use existing rule A → be.



### Sequitur Example (13)

bbeeebebebbebee

S → bAAeAA  
A → be

Enforce digram uniqueness  
AA occurs twice.  
Create new rule B → AA.

### Sequitur Example (14)

bbeeebebebbebee

S → bBeB  
A → be  
B → AA

### Sequitur Example (15)

bbeeebebebbbebee

S → bBeBb  
A → be  
B → AA

### Sequitur Example (16)

bbeeebebebbbebee

S → bBeBbb  
A → be  
B → AA

### Sequitur Example (17)

bbeeebebebbbee

S → bBeBbbe  
A → be  
B → AA

Enforce digram uniqueness.  
be occurs twice.  
Use existing rule A → be.

### Sequitur Example (18)

bbeeebebebbbee

S → bBeBbA  
A → be  
B → AA

### Sequitur Example (19)

bbeeebebebebeee

S → bBeBbAb  
A → be  
B → AA

Lecture 6 - Dictionary Coding

55

### Sequitur Example (20)

bbeeebebebebeee

S → bBeBbAbe      Enforce digram uniqueness.  
A → be              be occurs twice.  
B → AA              Use existing rule A → be.

Lecture 6 - Dictionary Coding

56

### Sequitur Example (21)

bbeeebebebebeee

S → bBeBbAA      Enforce digram uniqueness.  
A → be              AA occurs twice.  
B → AA              Use existing rule B → AA.

Lecture 6 - Dictionary Coding

57

### Sequitur Example (22)

bbeeebebebebeee

S → bBeBbB      Enforce digram uniqueness.  
A → be              bB occurs twice.  
B → AA              Create new rule C → bB.

Lecture 6 - Dictionary Coding

58

### Sequitur Example (23)

bbeeebebebebeee

S → CeBC  
A → be  
B → AA  
C → bB

Lecture 6 - Dictionary Coding

59

### Sequitur Example (24)

bbeeebebebebeee

S → CeBCe      Enforce digram uniqueness.  
A → be              Ce occurs twice.  
B → AA              Create new rule D → Ce.  
C → bB

Lecture 6 - Dictionary Coding

60

## Sequitur Example (25)

bbebeebbebebe

$S \rightarrow DBD$   
 $A \rightarrow be$   
 $B \rightarrow AA$   
 $C \rightarrow bB$   
 $D \rightarrow Ce$

Enforce rule utility.  
 C occurs only once.  
 Remove  $C \rightarrow bB$ .

Lecture 6 - Dictionary Coding

61

## Sequitur Example (26)

bbebeebbebebe

$S \rightarrow DBD$   
 $A \rightarrow be$   
 $B \rightarrow AA$   
 $D \rightarrow bBe$

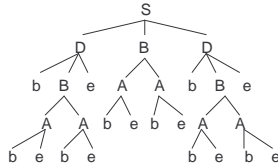
Lecture 6 - Dictionary Coding

62

## The Hierarchy

bbebeebbebebe

$S \rightarrow DBD$   
 $A \rightarrow be$   
 $B \rightarrow AA$   
 $D \rightarrow bBe$



Is there compression? In this small example, probably not.

Lecture 6 - Dictionary Coding

63

## Sequitur Algorithm

Input the first symbol  $s$  to create the production  $S \rightarrow s$ ;  
 repeat  
   match an existing rule:  
      $A \rightarrow \dots XY \dots$      $\rightarrow$      $A \rightarrow \dots B \dots$   
      $B \rightarrow XY$                      $\rightarrow$      $B \rightarrow XY$   
   create a new rule:  
      $A \rightarrow \dots XY \dots$              $\rightarrow$      $A \rightarrow \dots C \dots$   
      $B \rightarrow \dots XY \dots$              $\rightarrow$      $B \rightarrow \dots C \dots$   
   remove a rule:  
      $A \rightarrow \dots B \dots$   
      $B \rightarrow X_1 X_2 \dots X_k$          $\rightarrow$      $A \rightarrow \dots X_1 X_2 \dots X_k \dots$   
   input a new symbol:  
      $S \rightarrow X_1 X_2 \dots X_k$          $\rightarrow$      $S \rightarrow X_1 X_2 \dots X_k s$   
 until no symbols left

Lecture 6 - Dictionary Coding

64

## Exercise

Use Sequitur to construct a grammar for  $aaaaaaaaa = a^{10}$

Lecture 6 - Dictionary Coding

65

## Complexity

- The number of non-input sequitur operations applied  $< 2n$  where  $n$  is the input length.
- Since each operation takes constant time, sequitur is a linear time algorithm

Lecture 6 - Dictionary Coding

66

## Amortized Complexity Argument

- Let  $m$  = # of non-input sequitur operations.  
Let  $n$  = input length. Show  $m \leq 2n$ .
- Let  $s$  = the sum of the right hand sides of all the production rules. Let  $r$  = the number of rules.
- We evaluate  $2s - r$ .
- Initially  $2s - r = 1$  because  $s = 1$  and  $r = 1$ .
- $2s - r > 0$  at all times because each rule has at least 1 symbol on the right hand side.

Lecture 6 - Dictionary Coding

67

## Sequitur Rule Complexity

- Digram Uniqueness - match an existing rule.

$$\begin{array}{l} A \rightarrow \dots XY \dots \\ B \rightarrow XY \end{array} \rightarrow \begin{array}{l} A \rightarrow \dots B \dots \\ B \rightarrow XY \end{array} \quad \begin{array}{ll} s & r \\ -1 & 0 \end{array} \quad \begin{array}{l} 2s - r \\ -2 \end{array}$$

- Digram Uniqueness - create a new rule.

$$\begin{array}{l} A \rightarrow \dots XY \dots \\ B \rightarrow \dots XY \dots \end{array} \rightarrow \begin{array}{l} A \rightarrow \dots C \dots \\ B \rightarrow \dots C \dots \\ C \rightarrow XY \end{array} \quad \begin{array}{ll} s & r \\ 0 & 1 \end{array} \quad \begin{array}{l} 2s - r \\ -1 \end{array}$$

- Rule Utility - Remove a rule.

$$\begin{array}{l} A \rightarrow \dots B \dots \\ B \rightarrow X_1 X_2 \dots X_k \end{array} \rightarrow \begin{array}{l} A \rightarrow \dots X_1 X_2 \dots X_k \dots \end{array} \quad \begin{array}{ll} s & r \\ -1 & -1 \end{array} \quad \begin{array}{l} 2s - r \\ -1 \end{array}$$

Lecture 6 - Dictionary Coding

68

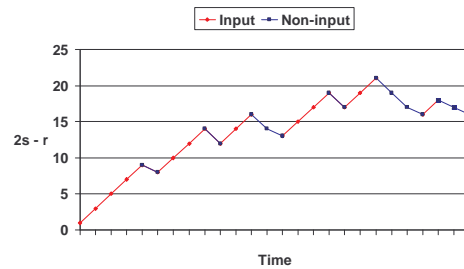
## Amortized Complexity Argument

- $2s - r \geq 0$  at all times because each rule has at least 1 symbol on the right hand side.
- $2s - r$  increases by 2 for every input operation.
- $2s - r$  decreases by at least 1 for each non-input sequitur rule applied.
- $n$  = number of input symbols  
 $m$  = number of non-input operations
- $2n - m \geq 0$ .  $m \leq 2n$ .

Lecture 6 - Dictionary Coding

69

## Amortized Complexity Argument



Lecture 6 - Dictionary Coding

70

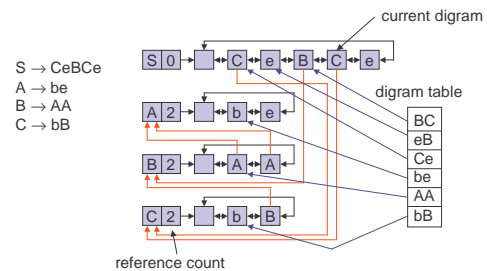
## Linear Time Algorithm

- There is a data structure to implement all the sequitur operations in constant time.
  - Production rules in an array of doubly linked lists.
  - Each production rule has reference count of the number of times used.
  - Each nonterminal points to its production rule.
  - Digrams stored in a hash table for quick lookup.

Lecture 6 - Dictionary Coding

71

## Data Structure Example



Lecture 6 - Dictionary Coding

72

## Basic Encoding a Grammar

Grammar	S → DBD	Symbol Code	b 000	No code for S needed
	A → be		e 001	
	B → AA		A 010	
	D → bBe		B 011	
			D 100	
			# 101	

### Grammar Code

D B D # b e # A A # b B e  
 100 011 100 101 000 001 101 010 010 101 000 011 001 39 bits

$$|\text{Grammar Code}| = (s + r - 1) \lceil \log_2(r + a) \rceil$$

$r$  = number of rules

$s$  = sum of right hand sides

$a$  = number in original symbol alphabet

Lecture 6 - Dictionary Coding

73

## Better Encoding of the Grammar

- Nevill-Manning and Witten suggest a more efficient encoding of the grammar that uses LZ77 ideas.

Lecture 6 - Dictionary Coding

74

## Kieffer-Yang Improvement

- Kieffer and Yang
  - Eliminate rules that are redundant
  - KY is universal; it achieves entropy in the limit
- Add to sequitur Reduction Rule 5:

S → AB	⇒	S → AA	Adding this constraint makes sequitur universal.
A → CD		A → CD	
B → aE		<del>B → aE</del>	
C → ab		C → ab	
D → cd		D → cd	
E → bD		<del>E → bD</del>	

$\langle A \rangle = \langle B \rangle = abcd$

Lecture 6 - Dictionary Coding

75

## Other Grammar Based Methods

- Longest Match
- Most frequent digram
- Match producing the best compression

Lecture 6 - Dictionary Coding

76

## Notes on Sequitur

- Yields compression and hierarchical structure simultaneously.
- With clever encoding is competitive with the best of the standards.

Lecture 6 - Dictionary Coding

77

## Move-to-Front Coding

- Non-numerical data
- The data have a relatively small working set that changes over the sequence.
- Example: a b a b a b c c b b c c c c b d b c c
- Move-to-front coding allows data with a small working set to be transformed to data with with better statistics for entropy coding.

Lecture 6 - Dictionary Coding

78

## Move-to-Front Algorithm

- Move-to-Front
  - Symbols are kept in a list indexed 0 to m-1
  - To code a symbol output its index and move the symbol to the front of the list
  - The index stream is entropy coded using arithmetic coding or some other statistical technique

Lecture 6 - Dictionary Coding

79

## Example

- Example: a b a b a b c c b b c c c c b d b c c

0  
a b c d

Lecture 6 - Dictionary Coding

80

## Example

- Example: a b a b a b a b c c b b c c c c b d b c c

0 1  
a b c d  
↓  
0 1 2 3  
b a c d

Lecture 6 - Dictionary Coding

81

## Example

- Example: a b a b a b a b c c b b c c c c b d b c c

0 1 1  
b a c d  
↓  
0 1 2 3  
a b c d

Lecture 6 - Dictionary Coding

82

## Example

- Example: a b a b a b a b a b c c b b c c c c b d b c c

0 1 1 1  
a b c d  
↓  
0 1 2 3  
b a c d

Lecture 6 - Dictionary Coding

83

## Example

- Example: a b a b a a b a b a b c c b b c c c c b d b c c

0 1 1 1 1  
b a c d  
↓  
0 1 2 3  
a b c d

Lecture 6 - Dictionary Coding

84



## Burrows-Wheeler Transform

- Burrows-Wheeler, 1994
- BW Transform creates a representation of the data which has a small working set.
- The transformed data is compressed with move to front compression.
- The decoder is quite different from the encoder.
- The algorithm requires processing the entire string at once (it is not on-line).
- It is a remarkably good compression method.

Lecture 6 - Dictionary Coding

91

## Encoding Example

- abracadabra
- 1. Create all cyclic shifts of the string.

```

0 abracadabra
1 bracadabraa
2 racadabraab
3 acadabraabr
4 cadabraabra
5 adabraabrac
6 dabraabraca
7 abraabracad
8 braabracada
9 raabracadab
10 aabracadabr
    
```

Lecture 6 - Dictionary Coding

92

## Encoding Example

2. Sort the strings alphabetically in to array A

0	abracadabra	A	0	aabracadabr
1	bracadabraa		1	abraabracad
2	racadabraab		2	abracadabra
3	acadabraabr		3	acadabraabr
4	cadabraabra	→	4	adabraabrac
5	adabraabrac		5	braabracada
6	dabraabraca		6	bracadabraa
7	abraabracad		7	cadabraabra
8	braabracada		8	dabraabraca
9	raabracadab		9	raabracadab
10	aabracadabr		10	racadabraab

Lecture 6 - Dictionary Coding

93

## Encoding Example

3. L = the last column

A	0	aabracadabr		
	1	abraabracad		
	2	abracadabra	L =	rdarcaaabb
	3	acadabraabr		
	4	adabraabrac		
	5	braabracada		
	6	bracadabraa		
	7	cadabraabra		
	8	dabraabraca		
	9	raabracadab		
	10	racadabraab		

Lecture 6 - Dictionary Coding

94

## Encoding Example

4. Transmit X the index of the input in A and L (using a predictive coding scheme).

A	0	aabracadabr		
	1	abraabracad		
	2	abracadabra	L =	rdarcaaabb
	3	acadabraabr	X =	2
	4	adabraabrac		
	5	braabracada		
	6	bracadabraa		
	7	cadabraabra		
	8	dabraabraca		
	9	raabracadab		
	10	racadabraab		

Lecture 6 - Dictionary Coding

95

## Why BW Works

- Ignore decoding for the moment.
- The prefix of each shifted string is a context for the last symbol.
  - The last symbol appears just before the prefix in the original.
- By sorting similar contexts are adjacent.
  - This means that the predicted last symbols are similar.

Lecture 6 - Dictionary Coding

96



### Decoding Example

- We first decode assuming some information. We then show how compute the information.
- Let  $A^s$  be  $A$  shifted by 1

A	0	a	a	b	r	a	c	a	d	a	b	r
	1	a	b	r	a	b	r	a	c	a	d	
	2	a	b	r	a	c	a	d	a	b	r	a
	3	a	c	a	d	a	b	r	a	b	r	
	4	a	d	a	b	r	a	b	r	a	c	a
	5	b	r	a	b	r	a	c	a	d	a	
	6	b	r	a	c	a	d	a	b	r	a	
	7	c	a	d	a	b	r	a	b	r	a	
	8	d	a	b	r	a	c	a	b	r	a	
	9	r	a	b	r	a	c	a	d	a	b	
	10	r	a	c	a	d	a	b	r	a	b	

Lecture 6 - Dictionary Coding

97

### Decoding Example

- Assume we know the mapping  $T[i]$  is the index in  $A^s$  of the string  $i$  in  $A$ .
- $T = [2\ 5\ 6\ 7\ 8\ 9\ 10\ 4\ 1\ 0\ 3]$

A	0	a	a	b	r	a	c	a	d	a	b	r
	1	a	b	r	a	b	r	a	c	a	d	
	2	a	b	r	a	c	a	d	a	b	r	a
	3	a	c	a	d	a	b	r	a	b	r	
	4	a	d	a	b	r	a	b	r	a	c	a
	5	b	r	a	b	r	a	c	a	d	a	
	6	b	r	a	c	a	d	a	b	r	a	
	7	c	a	d	a	b	r	a	b	r	a	
	8	d	a	b	r	a	c	a	b	r	a	
	9	r	a	b	r	a	c	a	d	a	b	
	10	r	a	c	a	d	a	b	r	a	b	

Lecture 6 - Dictionary Coding

98

### Decoding Example

- Let  $F$  be the first column of  $A$ , it is just  $L$ , sorted.

$$F = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ a & a & a & a & a & b & b & c & d & r & r \end{matrix}$$

$$T = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ 2 & 5 & 6 & 7 & 8 & 9 & 10 & 4 & 1 & 0 & 3 \end{matrix}$$

- Follow the pointers in  $T$  in  $F$  to recover the input starting with  $X$ .

Lecture 6 - Dictionary Coding

99

### Decoding Example

$$F = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ a & a & a & a & a & b & b & c & d & r & r \end{matrix}$$

$$T = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ 2 & 5 & 6 & 7 & 8 & 9 & 10 & 4 & 1 & 0 & 3 \end{matrix}$$

a

Lecture 6 - Dictionary Coding

100

### Decoding Example

$$F = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ a & a & a & a & a & b & b & c & d & r & r \end{matrix}$$

$$T = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ 2 & 5 & 6 & 7 & 8 & 9 & 10 & 4 & 1 & 0 & 3 \end{matrix}$$

ab

Lecture 6 - Dictionary Coding

101

### Decoding Example

$$F = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ a & a & a & a & a & b & b & c & d & r & r \end{matrix}$$

$$T = \begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ 2 & 5 & 6 & 7 & 8 & 9 & 10 & 4 & 1 & 0 & 3 \end{matrix}$$

abr

Lecture 6 - Dictionary Coding

102

## Decoding Example

- Why does this work?
- The first symbol of  $A[T[i]]$  is the second symbol of  $A[i]$  because  $A^s[T[i]] = A[i]$ .

A	0	aabracadabr	T	2	A <sup>s</sup>	0	raabracadab
	1	abraabracad		5		1	dabraabraca
	2	abracadabra		6		2	aabracadabr
	3	acadabraabr		7		3	racadabraab
	4	adabraabrac		8		4	cadabraabra
	5	braabracada		9		5	abraabracad
	6	bracadabraa		10		6	abracadabra
	7	cadabraabra		4		7	acadabraabr
	8	dabraabraca		1		8	adabraabrac
	9	raabracadab		0		9	braabracada
	10	racadabraab		3		10	bracadabraa

Lecture 6 - Dictionary Coding

103

## Decoding Example

- How do we compute F and T from L and X?  
F is just L sorted

	0	1	2	3	4	5	6	7	8	9	10
F =	a	a	a	a	a	b	b	c	d	r	r
L =	r	d	a	r	c	a	a	a	a	b	b

Note that L is the first column of  $A^s$  and  $A^s$  is in the same order as A.

If  $i$  is the  $k$ -th  $x$  in F then  $T[i]$  is the  $k$ -th  $x$  in L.

Lecture 6 - Dictionary Coding

104

## Decoding Example

	0	1	2	3	4	5	6	7	8	9	10
F =	a	a	a	a	a	b	b	c	d	r	r
L =	r	d	a	r	c	a	a	a	a	b	b

T =	0	1	2	3	4	5	6	7	8	9	10
	2	5	6	7	8						

Lecture 6 - Dictionary Coding

105

## Decoding Example

	0	1	2	3	4	5	6	7	8	9	10
F =	a	a	a	a	a	b	b	c	d	r	r
L =	r	d	a	r	c	a	a	a	a	b	b

T =	0	1	2	3	4	5	6	7	8	9	10
	2	5	6	7	8	9	10				

Lecture 6 - Dictionary Coding

106

## Decoding Example

	0	1	2	3	4	5	6	7	8	9	10
F =	a	a	a	a	a	b	b	c	d	r	r
L =	r	d	a	r	c	a	a	a	a	b	b

T =	0	1	2	3	4	5	6	7	8	9	10
	2	5	6	7	8	9	10	4			

Lecture 6 - Dictionary Coding

107

## Decoding Example

	0	1	2	3	4	5	6	7	8	9	10
F =	a	a	a	a	a	b	b	c	d	r	r
L =	r	d	a	r	c	a	a	a	a	b	b

T =	0	1	2	3	4	5	6	7	8	9	10
	2	5	6	7	8	9	10	4	1		

Lecture 6 - Dictionary Coding

108

## Decoding Example

```

    0 1 2 3 4 5 6 7 8 9 10
F = a a a a a b b c d r r
    / \
L = r d a r c a a a a b b

T = 0 1 2 3 4 5 6 7 8 9 10
    2 5 6 7 8 9 10 4 1 0 3

```

Lecture 6 - Dictionary Coding

109

## BWT Encoding Exercise

Encode the string abababababababab = (ab)<sup>8</sup>  
 1. Find L and X

Lecture 6 - Dictionary Coding

110

## BWT Decoding Exercise

Decode L = baaaaaba, X = 6  
 1. First Compute F and T  
 2. Use those to decode.

Lecture 6 - Dictionary Coding

111

## Notes on BW

- Alphabetic sorting does not need the entire cyclic shifted inputs.
  - Sort the indices of the string
  - Most significant symbols first radix sort works
- There are high quality practical implementations
  - Bzip
  - Bzip2

Lecture 6 - Dictionary Coding

112

## Compression Quality

	size	comp	gzip	sequitur	PPMC	bzip2
bib	111261	3.35	2.51	2.48	2.12	1.98
book	768771	3.46	3.35	2.82	2.52	2.42
geo	102400	6.08	5.34	4.74	5.01	4.45
obj2	246814	4.17	2.63	2.68	2.77	2.48
pic	513216	0.97	0.82	0.90	0.98	0.78
progc	38611	3.87	2.68	2.83	2.49	2.53

■ = First; ■ = Second; ■ = Third.

Files from the Calgary Corpus  
 Units in bits per character (8 bits)  
 compress - based on LZW  
 gzip - based on LZ77  
 PPMC - adaptive arithmetic coding with context  
 bzip2 - Burrows-Wheeler block sorting

Lecture 6 - Dictionary Coding

113