

CSEP 521
Applied Algorithms
Autumn 2009

Dictionary Coding for Lossless
Compression

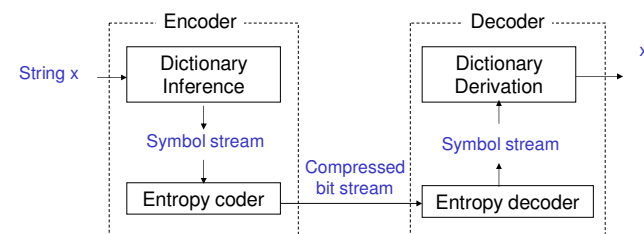
Plan for Tonight

- 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

Dictionary Coding

5

LZW Encoding Example (1)

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

Dictionary Coding

6

LZW Encoding Example (2)

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

Dictionary Coding

7

LZW Encoding Example (3)

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

Dictionary Coding

8

LZW Encoding Example (4)

Dictionary

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

Dictionary Coding

9

LZW Encoding Example (5)

Dictionary

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

Dictionary Coding

10

LZW Encoding Example (6)

Dictionary

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

Dictionary Coding

11

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;
```

Dictionary Coding

12

LZW Decoding Example (1)

Dictionary

0	a	<u>0</u> 1 2 4 3 6
1	b	a
2	a?	

Dictionary Coding

13

LZW Decoding Example (2a)

Dictionary

0	a	<u>0</u> <u>1</u> 2 4 3 6
1	b	a b
2	ab	

Dictionary Coding

14

LZW Decoding Example (2b)

Dictionary

0	a	<u>0</u> <u>1</u> 2 4 3 6
1	b	a b
2	ab	
3	b?	

Dictionary Coding

15

LZW Decoding Example (3a)

Dictionary

0	a	<u>0</u> <u>1</u> <u>2</u> 4 3 6
1	b	a b a
2	ab	
3	ba	

Dictionary Coding

16

LZW Decoding Example (3b)

Dictionary

- 0 a
- 1 b
- 2 ab
- 3 ba
- 4 ab?

0 1 2 4 3 6
a b ab

Dictionary Coding

17

LZW Decoding Example (4a)

Dictionary

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

0 1 2 4 3 6
a b ab a

Dictionary Coding

18

LZW Decoding Example (4b)

Dictionary

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

0 1 2 4 3 6
a b ab aba

Dictionary Coding

19

LZW Decoding Example (5a)

Dictionary

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

0 1 2 4 3 6
a b ab aba b

Dictionary Coding

20

LZW Decoding Example (5b)

Dictionary

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

0 1 2 4 3 6
a b ab aba ba

Dictionary Coding

21

LZW Decoding Example (6a)

Dictionary

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

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

Dictionary Coding

22

LZW Decoding Example (6b)

Dictionary

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

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

Dictionary Coding

23

Decoding Exercise

Base Dictionary

0 a
1 b
2 c
3 d
4 r

0 1 4 0 2 0 3 5 7

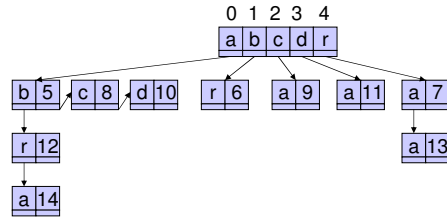
Dictionary Coding

24

Trie Data Structure for Encoder's Dictionary

- Fredkin (1960)

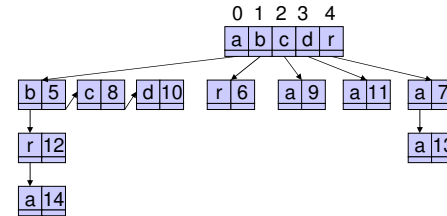
0	a	9	ca
1	b	10	ad
2	c	11	da
3	d	12	abr
4	r	13	raa
5	ab	14	abra
6	br		
7	ra		
8	ac		



Dictionary Coding

25

Encoder Uses a Trie (1)

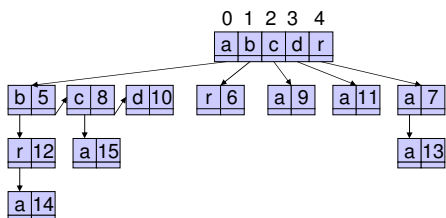


abracadabraabra
0 1 4 0 2 0 3 5 7 12

Dictionary Coding

26

Encoder Uses a Trie (2)



abracadabraabra
0 1 4 0 2 0 3 5 7 12 8

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	abra?
6	br		
7	ra		
8	ac		

0 1 4 0 2 0 3 5 7 12 8 ...
a b r a c a d a b r a a b r

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.

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

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!

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

Dictionary Coding

32

Context-Free Grammar Example

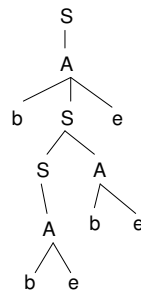
- $S \rightarrow SA$
- $S \rightarrow A$
- $A \rightarrow bSe$
- $A \rightarrow be$

Example: b and e matched as parentheses

derivation of bbebee

S
A
bSe
bSAe
bAAe
bbeAe
bbebee

hierarchical parse tree



Dictionary Coding

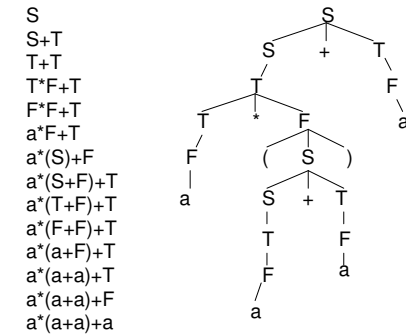
33

Arithmetic Expressions

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

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

parse tree

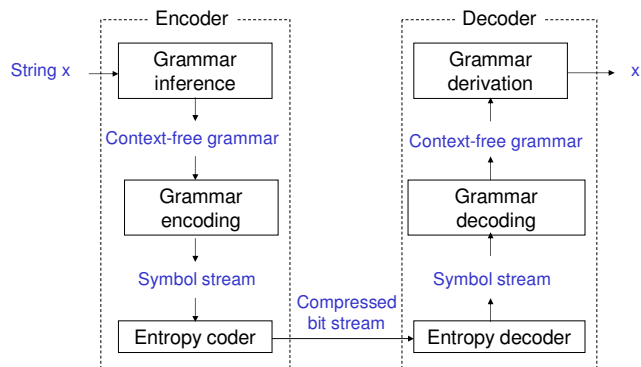


S
S+T
T+T
T*F+T
F*F+T
a*F+T
a*(S)+F
a*(S+F)+T
a*(T+F)+T
a*(F+F)+T
a*(a+F)+T
a*(a+a)+T
a*(a+a)+F
a*(a+a)+a

Dictionary Coding

34

Overview of Grammar Compression



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.

Dictionary Coding

36

Sequitur Example (1)

bbebebebebbebee

S → b

Dictionary Coding

37

Sequitur Example (2)

bbebebebebbebee

S → bb

Dictionary Coding

38

Sequitur Example (3)

bbebebebebbebee

S → bbe

Dictionary Coding

39

Sequitur Example (4)

bbbeebebebbebee

S → bbbe

Dictionary Coding

40

Sequitur Example (5)

bbebeebebebbeee

$S \rightarrow bbebe$

Enforce digram uniqueness.
be occurs twice.
Create new rule $A \rightarrow be$.

Dictionary Coding

41

Sequitur Example (6)

bbebeebebebbeee

$S \rightarrow bAA$

$A \rightarrow be$

Dictionary Coding

42

Sequitur Example (7)

bbebeeebebebbeee

$S \rightarrow bAAe$

$A \rightarrow be$

Dictionary Coding

43

Sequitur Example (8)

bbebeebebebebbeee

$S \rightarrow bAAeb$

$A \rightarrow be$

Dictionary Coding

44

Sequitur Example (9)

bbebebebebbebee

S → bAAe**be**
A → **be**

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

Dictionary Coding

45

Sequitur Example (10)

bbebebebebbebee

S → bAAeA
A → be

Dictionary Coding

46

Sequitur Example (11)

bbebebebbebbebee

S → bAAeAb
A → be

Dictionary Coding

47

Sequitur Example (12)

bbebebebebebbebee

S → bAAeA**be**
A → **be**

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

Dictionary Coding

48

Sequitur Example (13)

bbebebebebbebee

S → bAAeAA
A → be

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

Dictionary Coding

49

Sequitur Example (14)

bbebebebebbebee

S → bBeB
A → be
B → AA

Dictionary Coding

50

Sequitur Example (15)

bbebebebebbbebee

S → bBeBb
A → be
B → AA

Dictionary Coding

51

Sequitur Example (16)

bbebebebebbbbbebee

S → bBeBbb
A → be
B → AA

Dictionary Coding

52

Sequitur Example (17)

bbebeebbebbebee

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

Dictionary Coding

53

Sequitur Example (18)

bbebeebbebbebee

S → bBeBbA
A → be
B → AA

Dictionary Coding

54

Sequitur Example (19)

bbebeebbebbebee

S → bBeBbAb
A → be
B → AA

Dictionary Coding

55

Sequitur Example (20)

bbebeebbebbebee

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

Dictionary Coding

56

Sequitur Example (21)

bbebeebbebbee

$S \rightarrow bBeBbAA$ Enforce digram uniqueness.
 $A \rightarrow be$ AA occurs twice.
 $B \rightarrow AA$ Use existing rule $B \rightarrow AA$.

Dictionary Coding

57

Sequitur Example (22)

bbebeebbebbee

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

Dictionary Coding

58

Sequitur Example (23)

bbebeebbebbee

$S \rightarrow CeBC$
 $A \rightarrow be$
 $B \rightarrow AA$
 $C \rightarrow bB$

Dictionary Coding

59

Sequitur Example (24)

bbebeebbebbee

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

Dictionary Coding

60

Sequitur Example (25)

bbebeebbebbbeee

S → DBD
A → be
B → AA
C → bB
D → Ce

Enforce rule utility.
C occurs only once.
Remove C → bB.

Dictionary Coding

61

Sequitur Example (26)

bbebeebbebbbeee

S → DBD
A → be
B → AA
D → bBe

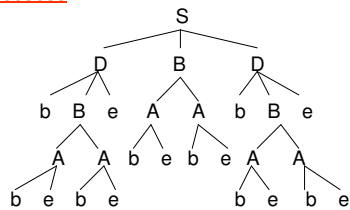
Dictionary Coding

62

The Hierarchy

bbebeebbebbbeee

S → DBD
A → be
B → AA
D → bBe



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

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:

$A \rightarrow \dots B \dots$ \rightarrow $C \rightarrow XY$
 $B \rightarrow X_1 X_2 \dots X_k$ \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

Dictionary Coding

64

Exercise

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

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

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.

Dictionary Coding

67

Sequitur Rule Complexity

- Digram Uniqueness - match an existing rule.

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

- Digram Uniqueness - create a new rule.

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

- Rule Utility - Remove a rule.

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

Dictionary Coding

68

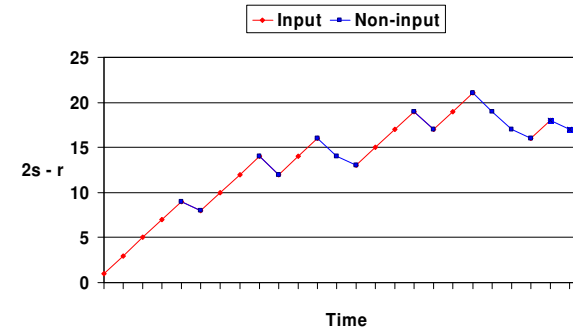
Amortized Complexity Argument

- $2s - r \geq 0$ 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
- $2n - m \geq 0$. $m \leq 2n$.

Dictionary Coding

69

Amortized Complexity Argument



Dictionary Coding

70

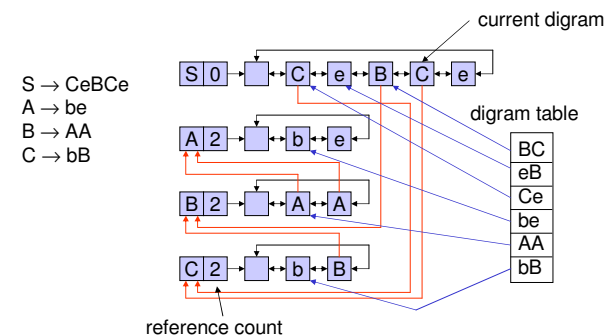
Linear Time Algorithm

- 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.

Dictionary Coding

71

Data Structure Example



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

Dictionary Coding

73

Better Encoding of the Grammar

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

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		B → aE	
C → ab		C → ab	
D → cd		D → cd	
E → bD		E → bD	

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

Dictionary Coding

75

Other Grammar Based Methods

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

Dictionary Coding

76

Notes on Sequitur

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

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 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.

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

Dictionary Coding

79

Example

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

0
0 1 2 3
a b c d

Dictionary Coding

80

Example

- Example: a b a b a a b c c b b c c c b d b c c
0 1

0	1	2	3
a	b	c	d

↓

0	1	2	3
b	a	c	d

Dictionary Coding

81

Example

- Example: a b a b a a b c c b b c c c b d b c c
0 1 1

0	1	2	3
b	a	c	d

↓

0	1	2	3
a	b	c	d

Dictionary Coding

82

Example

- Example: a b a b a a b c c b b c c c b d b c c
0 1 1 1

0	1	2	3
a	b	c	d

↓

0	1	2	3
b	a	c	d

Dictionary Coding

83

Example

- Example: a b a b a a b c c b b c c c b d b c c
0 1 1 1 1

0	1	2	3
b	a	c	d

↓

0	1	2	3
a	b	c	d

Dictionary Coding

84

Example

- Example: a b a b a a b c c b b c c c c b d b c c
0 1 1 1 1 0

0 1 2 3
a b c d

Dictionary Coding

85

Example

- Example: a b a b a a b c c b b c c c c b d b c c
0 1 1 1 1 0 1

0 1 2 3
a b c d
↓
0 1 2 3
b a c d

Dictionary Coding

86

Example

- Example: a b a b a a b c c b b c c c c b d b c c
0 1 1 1 1 0 1 2

0 1 2 3
b a c d
↓
0 1 2 3
c b a d

Dictionary Coding

87

Example

- Example: a b a b a a b c c b b c c c c b d b c c
0 1 1 1 1 0 1 2 0 1 0 1 0 0 1 3 1 2 0

0 1 2 3
c b d a

Dictionary Coding

88

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

Dictionary Coding

93

Encoding Example

3. L = the last column

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

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	L = rdarcaaaaabb
	2	abracadabra	X = 2
	3	acadabraabr	
	4	adabraabrac	
	5	braabracada	
	6	bracadabraa	
	7	cadabraabra	
	8	dabraabraca	
	9	raabracadab	
	10	racadabraab	

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.

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

<table border="0"> <tr><td>A</td><td>0</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>1</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td></tr> <tr><td></td><td>2</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>3</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>4</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td></tr> <tr><td></td><td>5</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td></tr> <tr><td></td><td>6</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td></tr> <tr><td></td><td>7</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>8</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td></tr> <tr><td></td><td>9</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td></tr> <tr><td></td><td>10</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td></tr> </table>	A	0	a	a	b	r	a	c	a	d	a	b	r		1	a	b	r	a	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	a	b	r		4	a	d	a	b	r	a	a	b	r	a	c		5	b	r	a	a	b	r	a	c	a	d	a		6	b	r	a	c	a	d	a	b	r	a	a		7	c	a	d	a	b	r	a	a	b	r	a		8	d	a	b	r	a	a	b	r	a	c	a		9	r	a	a	b	r	a	c	a	d	a	b		10	r	a	c	a	d	a	b	r	a	a	b	<table border="0"> <tr><td>A^s</td><td>0</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td></tr> <tr><td></td><td>1</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td></tr> <tr><td></td><td>2</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>3</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td></tr> <tr><td></td><td>4</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>5</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td></tr> <tr><td></td><td>6</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>7</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>8</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td></tr> <tr><td></td><td>9</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td></tr> <tr><td></td><td>10</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td></tr> </table>	A^s	0	r	a	a	b	r	a	c	a	d	a	b		1	d	a	b	r	a	a	b	r	a	c	a		2	a	a	b	r	a	c	a	d	a	b	r		3	r	a	c	a	d	a	b	r	a	a	b		4	c	a	d	a	b	r	a	a	b	r	a		5	a	b	r	a	a	b	r	a	c	a	d		6	a	b	r	a	c	a	d	a	b	r	a		7	a	c	a	d	a	b	r	a	a	b	r		8	a	d	a	b	r	a	a	b	r	a	c		9	b	r	a	a	b	r	a	c	a	d	a		10	b	r	a	c	a	d	a	b	r	a	a
A	0	a	a	b	r	a	c	a	d	a	b	r																																																																																																																																																																																																																																																																																			
	1	a	b	r	a	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	a	b	r																																																																																																																																																																																																																																																																																			
	4	a	d	a	b	r	a	a	b	r	a	c																																																																																																																																																																																																																																																																																			
	5	b	r	a	a	b	r	a	c	a	d	a																																																																																																																																																																																																																																																																																			
	6	b	r	a	c	a	d	a	b	r	a	a																																																																																																																																																																																																																																																																																			
	7	c	a	d	a	b	r	a	a	b	r	a																																																																																																																																																																																																																																																																																			
	8	d	a	b	r	a	a	b	r	a	c	a																																																																																																																																																																																																																																																																																			
	9	r	a	a	b	r	a	c	a	d	a	b																																																																																																																																																																																																																																																																																			
	10	r	a	c	a	d	a	b	r	a	a	b																																																																																																																																																																																																																																																																																			
A^s	0	r	a	a	b	r	a	c	a	d	a	b																																																																																																																																																																																																																																																																																			
	1	d	a	b	r	a	a	b	r	a	c	a																																																																																																																																																																																																																																																																																			
	2	a	a	b	r	a	c	a	d	a	b	r																																																																																																																																																																																																																																																																																			
	3	r	a	c	a	d	a	b	r	a	a	b																																																																																																																																																																																																																																																																																			
	4	c	a	d	a	b	r	a	a	b	r	a																																																																																																																																																																																																																																																																																			
	5	a	b	r	a	a	b	r	a	c	a	d																																																																																																																																																																																																																																																																																			
	6	a	b	r	a	c	a	d	a	b	r	a																																																																																																																																																																																																																																																																																			
	7	a	c	a	d	a	b	r	a	a	b	r																																																																																																																																																																																																																																																																																			
	8	a	d	a	b	r	a	a	b	r	a	c																																																																																																																																																																																																																																																																																			
	9	b	r	a	a	b	r	a	c	a	d	a																																																																																																																																																																																																																																																																																			
	10	b	r	a	c	a	d	a	b	r	a	a																																																																																																																																																																																																																																																																																			

Dictionary Coding

97

Decoding Example

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

<table border="0"> <tr><td>A</td><td>0</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>1</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td></tr> <tr><td></td><td>2</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>3</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>4</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td></tr> <tr><td></td><td>5</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td></tr> <tr><td></td><td>6</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td></tr> <tr><td></td><td>7</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>8</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td></tr> <tr><td></td><td>9</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td></tr> <tr><td></td><td>10</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td></tr> </table>	A	0	a	a	b	r	a	c	a	d	a	b	r		1	a	b	r	a	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	a	b	r		4	a	d	a	b	r	a	a	b	r	a	c		5	b	r	a	a	b	r	a	c	a	d	a		6	b	r	a	c	a	d	a	b	r	a	a		7	c	a	d	a	b	r	a	a	b	r	a		8	d	a	b	r	a	a	b	r	a	c	a		9	r	a	a	b	r	a	c	a	d	a	b		10	r	a	c	a	d	a	b	r	a	a	b	<table border="0"> <tr><td>A^s</td><td>0</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td></tr> <tr><td></td><td>1</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td></tr> <tr><td></td><td>2</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>3</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td></tr> <tr><td></td><td>4</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>5</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td></tr> <tr><td></td><td>6</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td></tr> <tr><td></td><td>7</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td></tr> <tr><td></td><td>8</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td></tr> <tr><td></td><td>9</td><td>b</td><td>r</td><td>a</td><td>a</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td></tr> <tr><td></td><td>10</td><td>b</td><td>r</td><td>a</td><td>c</td><td>a</td><td>d</td><td>a</td><td>b</td><td>r</td><td>a</td><td>a</td></tr> </table>	A^s	0	r	a	a	b	r	a	c	a	d	a	b		1	d	a	b	r	a	a	b	r	a	c	a		2	a	a	b	r	a	c	a	d	a	b	r		3	r	a	c	a	d	a	b	r	a	a	b		4	c	a	d	a	b	r	a	a	b	r	a		5	a	b	r	a	a	b	r	a	c	a	d		6	a	b	r	a	c	a	d	a	b	r	a		7	a	c	a	d	a	b	r	a	a	b	r		8	a	d	a	b	r	a	a	b	r	a	c		9	b	r	a	a	b	r	a	c	a	d	a		10	b	r	a	c	a	d	a	b	r	a	a
A	0	a	a	b	r	a	c	a	d	a	b	r																																																																																																																																																																																																																																																																																			
	1	a	b	r	a	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	a	b	r																																																																																																																																																																																																																																																																																			
	4	a	d	a	b	r	a	a	b	r	a	c																																																																																																																																																																																																																																																																																			
	5	b	r	a	a	b	r	a	c	a	d	a																																																																																																																																																																																																																																																																																			
	6	b	r	a	c	a	d	a	b	r	a	a																																																																																																																																																																																																																																																																																			
	7	c	a	d	a	b	r	a	a	b	r	a																																																																																																																																																																																																																																																																																			
	8	d	a	b	r	a	a	b	r	a	c	a																																																																																																																																																																																																																																																																																			
	9	r	a	a	b	r	a	c	a	d	a	b																																																																																																																																																																																																																																																																																			
	10	r	a	c	a	d	a	b	r	a	a	b																																																																																																																																																																																																																																																																																			
A^s	0	r	a	a	b	r	a	c	a	d	a	b																																																																																																																																																																																																																																																																																			
	1	d	a	b	r	a	a	b	r	a	c	a																																																																																																																																																																																																																																																																																			
	2	a	a	b	r	a	c	a	d	a	b	r																																																																																																																																																																																																																																																																																			
	3	r	a	c	a	d	a	b	r	a	a	b																																																																																																																																																																																																																																																																																			
	4	c	a	d	a	b	r	a	a	b	r	a																																																																																																																																																																																																																																																																																			
	5	a	b	r	a	a	b	r	a	c	a	d																																																																																																																																																																																																																																																																																			
	6	a	b	r	a	c	a	d	a	b	r	a																																																																																																																																																																																																																																																																																			
	7	a	c	a	d	a	b	r	a	a	b	r																																																																																																																																																																																																																																																																																			
	8	a	d	a	b	r	a	a	b	r	a	c																																																																																																																																																																																																																																																																																			
	9	b	r	a	a	b	r	a	c	a	d	a																																																																																																																																																																																																																																																																																			
	10	b	r	a	c	a	d	a	b	r	a	a																																																																																																																																																																																																																																																																																			

Dictionary Coding

98

Decoding Example

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

$F =$

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

$T =$

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

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

Dictionary Coding

99

Decoding Example

$F =$

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

$T =$

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

a

Dictionary Coding

100

Decoding Example

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

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

ab

Dictionary Coding

101

Decoding Example

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

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

abr

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		T		A^s
0	aabracadabr	2		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

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.

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

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

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

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

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
```

Dictionary Coding

109

BWT Encoding Exercise

Encode the string abababababababab = (ab)⁸

1. Find L and X

Dictionary Coding

110

BWT Decoding Exercise

Decode L = baaaaaba, X = 6

1. First Compute F and T
2. Use those to decode.

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

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

Dictionary Coding

113