

CSE 589  
Applied Algorithms  
Spring 1999

Dictionary Coding  
Sequitur

### LZW Encoding 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

CSE 589 - Lecture 13 - Spring 1999

2

### LZW Encoding Example (1)

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

CSE 589 - Lecture 13 - Spring 1999

3

### LZW Encoding Example (2)

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

CSE 589 - Lecture 13 - Spring 1999

4

### LZW Encoding Example (3)

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

CSE 589 - Lecture 13 - Spring 1999

5

### LZW Encoding Example (4)

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

CSE 589 - Lecture 13 - Spring 1999

6

### LZW Encoding Example (5)

Dictionary	<u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u>
0 a	0 1 2 4
1 b	
2 ab	
3 ba	
4 aba	
5 abab	

### LZW Encoding Example (6)

Dictionary	<u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u> <u>b</u> <u>a</u>
0 a	0 1 2 4 3
1 b	
2 ab	
3 ba	
4 aba	
5 abab	

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

### LZW Decoding Example (2a)

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

### LZW Decoding Example (2b)

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

### LZW Decoding Example (3a)

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

### LZW Decoding Example (3b)

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

### 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

### LZW Decoding Example (4b)

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

### LZW Decoding Example (5a)

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

### LZW Decoding Example (5b)

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

### LZW Decoding Example (6a)

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

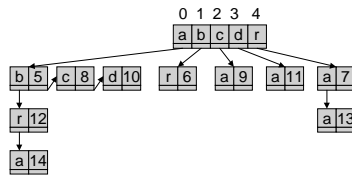
### LZW Decoding Example (6b)

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

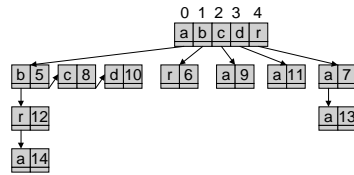
### Trie Data Structure for Dictionary

- 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		

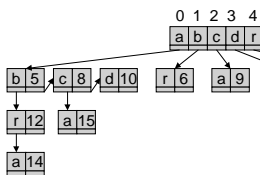


### Encoder Uses a Trie (1)



abracadabraabracadabra  
0 1 4 0 2 0 3 5 7 12

### Encoder Uses a Trie (2)



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

### 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 ...  
a b r a c a d a b r a a b r

## Notes on Dictionary Coding

- Extremely effective when there are repeated patterns in the data that are widely spread. Where local context is not as significant.
  - text
  - some graphics
  - program sources or binaries
- Variants of LZW are pervasive.
  - Unix compress
  - GIF

CSE 589 - Lecture 13 - Spring 1999

25

## 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!

CSE 589 - Lecture 13 - Spring 1999

26

## 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
  - nonterminals: S, A
  - Production Rules:  $S \rightarrow SA$ ,  $S \rightarrow A$ ,  $A \rightarrow bSe$ ,  $A \rightarrow be$
  - S is the start symbol

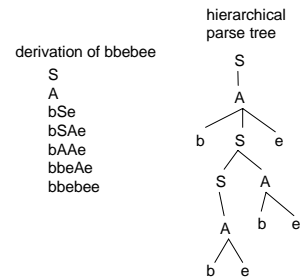
CSE 589 - Lecture 13 - Spring 1999

27

## Context-Free Grammar Example

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

Example: b and e matched as parentheses



CSE 589 - Lecture 13 - Spring 1999

28

## Arithmetic Expressions

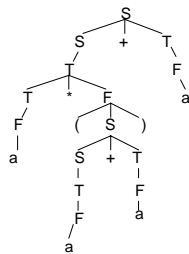
- $S \rightarrow S + T$
- $S \rightarrow T$
- $T \rightarrow T * F$
- $T \rightarrow F$
- $F \rightarrow a$
- $F \rightarrow (S)$

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
    
```



CSE 589 - Lecture 13 - Spring 1999

29

## 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.

CSE 589 - Lecture 13 - Spring 1999

30

### Sequitur Example (1)

bbebebebebebe

S -> b

### Sequitur Example (2)

bbebebebebebe

S -> bb

### Sequitur Example (3)

bbeebebebebebe

S -> bbe

### Sequitur Example (4)

bbebeebebebebe

S -> bbeb

### Sequitur Example (5)

bbebeebebebebe

S -> bbebe

Enforce digram uniqueness.  
be occurs twice.  
Create new rule A -> be.

### Sequitur Example (6)

bbebeebebebebe

S -> bAA

A -> be

### Sequitur Example (7)

bb**e**e**e**e**e**b**e**b**e**e**e**

S -> bAAe  
A -> be

### Sequitur Example (8)

bb**e**e**e**e**e**b**e**b**e**e**e**

S -> bAAeb  
A -> be

### Sequitur Example (9)

bb**e**e**e**e**e**b**e**b**e**e**e**

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

### Sequitur Example (10)

bb**e**e**e**e**e**b**e**b**e**e**e**

S -> bAAeA  
A -> be

### Sequitur Example (11)

bb**e**e**e**e**e**b**e**b**e**e**e**

S -> bAAeAb  
A -> be

### Sequitur Example (12)

bb**e**e**e**e**e**b**e**b**e**e**e**

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

### Sequitur Example (13)

bbebebebebebe

S -> bAAeAA      Enforce digram uniqueness  
A -> be            AA occurs twice.  
                      Create new rule B -> AA.

### Sequitur Example (14)

bbebebebebebe

S -> bBeB  
A -> be  
B -> AA

### Sequitur Example (15)

bbebebebebebe

S -> bBeBb  
A -> be  
B -> AA

### Sequitur Example (16)

bbebebebebebe

S -> bBeBbb  
A -> be  
B -> AA

### Sequitur Example (17)

bbebebebebebe

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

### Sequitur Example (18)

bbebebebebebe

S -> bBeBbA  
A -> be  
B -> AA



### Sequitur Example (19)

bbebeebbebebe

S -> bBeBbAb  
A -> be  
B -> AA

### Sequitur Example (20)

bbebeebbebebe

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

### Sequitur Example (21)

bbebeebbebebe

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

### Sequitur Example (22)

bbebeebbebebe

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

### Sequitur Example (23)

bbebeebbebebe

S -> CeBC  
A -> be  
B -> AA  
C -> bB

### Sequitur Example (24)

bbebeebbebebe

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

## Sequitur Example (25)

bbebeebbebebe

S -> DBD                      Enforce rule utility.  
 A -> be                        C occurs only once.  
 B -> AA                        Remove C -> bB.  
 C -> bB  
 D -> Ce

## Sequitur Example (26)

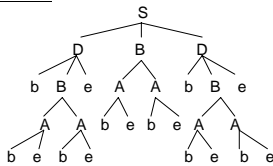
bbebeebbebebe

S -> DBD  
 A -> be  
 B -> AA  
 D -> bBe

## The Hierarchy

bbebeebbebebe

S -> DBD  
 A -> be  
 B -> AA  
 D -> bBe



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

## Sequitur Algorithm

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

## Complexity

- The number of non-input sequitur operations applied  $\leq 2n$  where  $n$  is the input length.
- Amortized Complexity Argument
  - Let  $s$  = the sum of the right hand sides of all the production rules. Let  $r$  = the number of rules.
  - We evaluate  $s - r/2$ .
  - Initially  $s - r/2 = 1/2$  because  $s = 1$  and  $r = 1$ .
  - $s - r/2 \geq 0$  at all times because each rule has at least 1 symbol on the right hand side.
  - $s - r/2$  increases by 1 for every input operation.
  - $s - r/2$  decreases by at least  $1/2$  for each non-input sequitur rule applied.

## Sequitur Rule Complexity

- digram Uniqueness - match an existing rule.

$A \rightarrow \dots XY \dots$	$\rightarrow$	$A \rightarrow \dots B \dots$	$s$	$r$	$s - r/2$
$B \rightarrow XY$		$B \rightarrow XY$	-1	0	-1

- digram Uniqueness - create a new rule.

$A \rightarrow \dots XY \dots$	$\rightarrow$	$A \rightarrow \dots C \dots$	$s$	$r$	$s - r/2$
$B \rightarrow \dots XY \dots$		$B \rightarrow \dots C \dots$	0	1	-1/2
		$C \rightarrow XY$			

- Rule Utility - Remove a rule.

$A \rightarrow \dots B \dots$	$\rightarrow$	$A \rightarrow \dots X_1 X_2 \dots X_k \dots$	$s$	$r$	$s - r/2$
$B \rightarrow X_1 X_2 \dots X_k$		$B \rightarrow X_1 X_2 \dots X_k$	-1	-1	-1/2

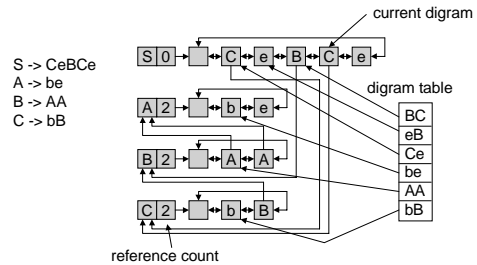
## Linear Time Algorithm

- 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.

CSE 589 - Lecture 13 - Spring 1999

61

## Data Structure Example



CSE 589 - Lecture 13 - Spring 1999

62

## Basic Encoding a Grammar

Grammar	$S \rightarrow DBD$	Symbol Code	S 010
	$A \rightarrow be$		A 011
	$B \rightarrow AA$		B 100
	$D \rightarrow bBe$		D 101
			# 110

Grammar Code

D B D # b e # A A # b B e 39 bits

101 100 101 110 000 001 110 011 011 110 000 100 001

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

$r$  = number of rules  
 $s$  = sum of right hand sides  
 $a$  = number in original symbol alphabet

CSE 589 - Lecture 13 - Spring 1999

63

## Better Encoding of the Grammar

- Nevill-Manning and Witten suggest a more efficient encoding of the grammar that resembles LZ77.
  - The first time a nonterminal is sent, its right hand side is transmitted instead.
  - The second time a nonterminal is sent the new production rule is established with a pointer to the previous occurrence sent along with the length of the rule.
  - Subsequently, the nonterminal is represented by the index of the production rule.

CSE 589 - Lecture 13 - Spring 1999

64

## Compression Quality

- Neville-Manning and Witten 1997

	size	compress	gzip	sequitur	PPMC
bitb	111261	3.35	2.51	2.48	2.12
book	768771	3.46	3.35	2.82	2.52
geo	102400	6.08	5.34	4.74	5.01
obj2	246814	4.17	2.63	2.68	2.77
pic	513216	0.97	0.82	0.90	0.98
progc	38611	3.87	2.68	2.83	2.49

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

CSE 589 - Lecture 13 - Spring 1999

65

## Notes on Sequitur

- Very new and different from the standards.
- Yields compression and structure simultaneously.
- With clever encoding is competitive with the best of the standards.
- Practical linear time encoding and decoding.

CSE 589 - Lecture 13 - Spring 1999

66