

CSEP 590  
Data Compression  
Autumn 2007

Dictionary Coding  
LZW, LZ77

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,
- Applications: Unix Compress, gzip, GIF

CSEP 590 - Lecture 4 - Autumn 2007

2

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

CSEP 590 - Lecture 4 - Autumn 2007

3

LZW Encoding Example (1)

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

CSEP 590 - Lecture 4 - Autumn 2007

4

LZW Encoding Example (2)

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

CSEP 590 - Lecture 4 - Autumn 2007

5

LZW Encoding Example (3)

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

CSEP 590 - Lecture 4 - Autumn 2007

6

### LZW Encoding Example (4)

Dictionary

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

### LZW Encoding Example (5)

Dictionary

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

### LZW Encoding Example (6)

Dictionary

0	a	a	b	a	b	a	b	a	b	a
1	b	0	1	2	4	3				
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

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

### LZW Decoding Example (2a)

Dictionary

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

### LZW Decoding Example (2b)

Dictionary

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

CSEP 590 - Lecture 4 - Autumn 2007

13

### LZW Decoding Example (3a)

Dictionary

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

CSEP 590 - Lecture 4 - Autumn 2007

14

### LZW Decoding Example (3b)

Dictionary

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

CSEP 590 - Lecture 4 - Autumn 2007

15

### LZW Decoding Example (4a)

Dictionary

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

CSEP 590 - Lecture 4 - Autumn 2007

16

### 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?						

CSEP 590 - Lecture 4 - Autumn 2007

17

### LZW Decoding Example (5a)

Dictionary

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

CSEP 590 - Lecture 4 - Autumn 2007

18

### LZW Decoding Example (5b)

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

CSEP 590 - Lecture 4 - Autumn 2007

19

### LZW Decoding Example (6a)

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

CSEP 590 - Lecture 4 - Autumn 2007

20

### LZW Decoding Example (6b)

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

CSEP 590 - Lecture 4 - Autumn 2007

21

### Decoding Exercise

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

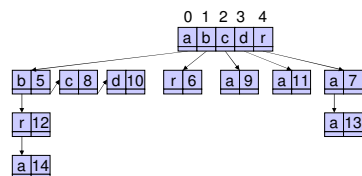
CSEP 590 - Lecture 4 - Autumn 2007

22

### Trie Data Structure for Encoder's Dictionary

- Fredkin (1960)

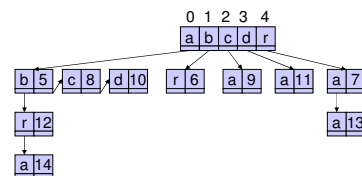
0	a	9	ca
1	b	10	ad
2	c	11	da
3	d	12	abr
4	r	13	raa
5	ab	14	abra
6	br		
7	ra		
8	ac		



CSEP 590 - Lecture 4 - Autumn 2007

23

### Encoder Uses a Trie (1)

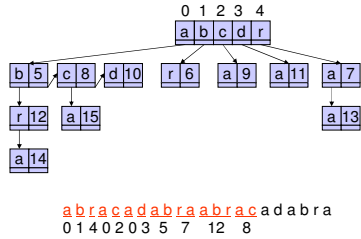


abracadabraabra  
0 1 4 0 2 0 3 5 7 12

CSEP 590 - Lecture 4 - Autumn 2007

24

## Encoder Uses a Trie (2)



CSEP 590 - Lecture 4 - Autumn 2007

25

## 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

CSEP 590 - Lecture 4 - Autumn 2007

26

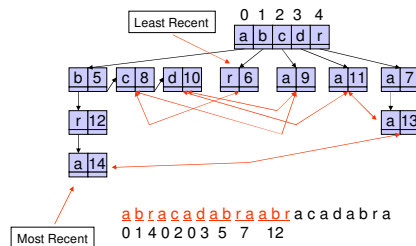
## Bounded Size Dictionary

- 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.
  - Don't add more, just use what is there.
  - Throw it away and start a new dictionary.
  - Double the dictionary, adding one more bit to indices.
  - Throw out the least recently visited entry to make room for the new entry.

CSEP 590 - Lecture 4 - Autumn 2007

27

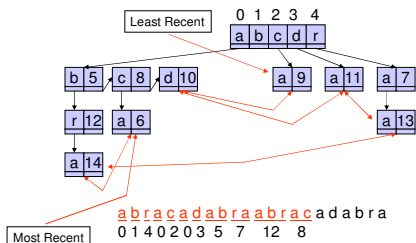
## Implementing the LRV Strategy



CSEP 590 - Lecture 4 - Autumn 2007

28

## Implementing the LRV Strategy



CSEP 590 - Lecture 4 - Autumn 2007

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

CSEP 590 - Lecture 4 - Autumn 2007

30

### The Dictionary is Implicit

- Ziv and Lempel, 1977
- Use the string coded so far as a dictionary.
- Given that  $x_1x_2\dots x_n$  has been coded we want to code  $x_{n+1}x_{n+2}\dots x_{n+k}$  for the largest  $k$  possible.

CSEP 590 - Lecture 4 - Autumn 2007 31

### Solution A

- If  $x_{n+1}x_{n+2}\dots x_{n+k}$  is a substring of  $x_1x_2\dots x_n$  then  $x_{n+1}x_{n+2}\dots x_{n+k}$  can be coded by  $\langle j,k \rangle$  where  $j$  is the beginning of the match.
- Example

```

ababababa bababababababab....
  coded
ababababa babababa babababab....
                <2,8>
  
```

CSEP 590 - Lecture 4 - Autumn 2007 32

### Solution A Problem

- What if there is no match at all in the dictionary?

```

ababababa cababababababab....
  coded
  
```

- Solution B. Send tuples  $\langle j,k,x \rangle$  where
  - If  $k = 0$  then  $x$  is the unmatched symbol
  - If  $k > 0$  then the match starts at  $j$  and is  $k$  long and the unmatched symbol is  $x$ .

CSEP 590 - Lecture 4 - Autumn 2007 33

### Solution B

- If  $x_{n+1}x_{n+2}\dots x_{n+k}$  is a substring of  $x_1x_2\dots x_n$  and  $x_{n+1}x_{n+2}\dots x_{n+k}x_{n+k+1}$  is not then  $x_{n+1}x_{n+2}\dots x_{n+k}x_{n+k+1}$  can be coded by  $\langle j,k, x_{n+k+1} \rangle$  where  $j$  is the beginning of the match.
- Examples

```

ababababa cababababababab....
ababababa c abababab ababab....
<0,0,c> <1,9,b>
  
```

CSEP 590 - Lecture 4 - Autumn 2007 34

### Solution B Example

```

a bababababababababab....
<0,0,a>
a b ababababababababab....
<0,0,b>
a b aba bababababababab....
<1,2,a>
a b aba babab ababababab....
<2,4,b>
a b aba babab abababababab....
<1,10,a>
  
```

CSEP 590 - Lecture 4 - Autumn 2007 35

### Surprise Code!

```

a bababababababababab$
<0,0,a>
a b ababababababababab$
<0,0,b>
a b ababababababababab$
<1,22,$>
  
```

CSEP 590 - Lecture 4 - Autumn 2007 36

## Surprise Decoding

$\langle 0,0,a \rangle \langle 0,0,b \rangle \langle 1,22,\$ \rangle$

$\langle 0,0,a \rangle$     a  
 $\langle 0,0,b \rangle$     b  
 $\langle 1,22,\$ \rangle$     a  
 $\langle 2,21,\$ \rangle$     b  
 $\langle 3,20,\$ \rangle$     a  
 $\langle 4,19,\$ \rangle$     b  
 ...  
 $\langle 22,1,\$ \rangle$     b  
 $\langle 23,0,\$ \rangle$     \$

CSEP 590 - Lecture 4 - Autumn 2007

37

## Surprise Decoding

$\langle 0,0,a \rangle \langle 0,0,b \rangle \langle 1,22,\$ \rangle$

$\langle 0,0,a \rangle$     a  
 $\langle 0,0,b \rangle$     b  
 $\langle 1,22,\$ \rangle$     a  
 $\langle 2,21,\$ \rangle$     b  
 $\langle 3,20,\$ \rangle$     a  
 $\langle 4,19,\$ \rangle$     b  
 ...  
 $\langle 22,1,\$ \rangle$     b  
 $\langle 23,0,\$ \rangle$     \$

CSEP 590 - Lecture 4 - Autumn 2007

38

## Solution C

- The matching string can include part of itself!
- If  $x_{n+1}x_{n+2}\dots x_{n+k}$  is a substring of  $x_1x_2\dots x_n x_{n+1}x_{n+2}\dots x_{n+k}$  that begins at  $j \leq n$  and  $x_{n+1}x_{n+2}\dots x_{n+k}x_{n+k+1}$  is not then  $x_{n+1}x_{n+2}\dots x_{n+k}x_{n+k+1}$  can be coded by  $\langle j,k, x_{n+k+1} \rangle$

CSEP 590 - Lecture 4 - Autumn 2007

39

## In Class Exercise

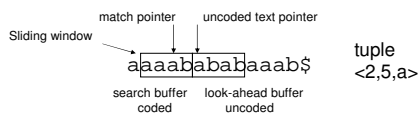
- Use Solution C to code the string
  - abaabaaabaaab\$
  - aaaabaaababab\$

CSEP 590 - Lecture 4 - Autumn 2007

40

## Bounded Buffer – Sliding Window

- We want the triples  $\langle j,k,x \rangle$  to be of bounded size. To achieve this we use bounded buffers.
  - Search buffer of size  $s$  is the symbols  $x_{n-s+1}\dots x_n$   $j$  is then the offset into the buffer.
  - Look-ahead buffer of size  $t$  is the symbols  $x_{n+1}\dots x_{n+t}$
- Match pointer can start in search buffer and go into the look-ahead buffer but no farther.



CSEP 590 - Lecture 4 - Autumn 2007

41

## Search in the Sliding Window

	offset	length	
aaaabababaaab\$	1	0	
aaaabababaaab\$	2	1	
aaaabababaaab\$	2	2	
aaaabababaaab\$	2	3	
aaaabababaaab\$	2	4	
aaaabababaaab\$	2	5	tuple <2,5,a>

CSEP 590 - Lecture 4 - Autumn 2007

42

## Coding Example

$s = 4, t = 4, a = 3$

	tuple
aaaabababaaab\$	<0, 0, a>
aaaabababaaab\$	<1, 3, b>
aaaabababaaab\$	<2, 5, a>
aaaabababaaab\$	<4, 2, \$>

CSEP 590 - Lecture 4 - Autumn 2007

43

## Coding the Tuples

- Simple fixed length code

$$\lceil \log_2(s+1) \rceil + \lceil \log_2(s+t+1) \rceil + \lceil \log_2 a \rceil$$

$s = 4, t = 4, a = 3$       tuple      fixed code  
 <2,5,a>      010 0101 00

- Variable length code using adaptive Huffman or arithmetic code on Tuples
  - Two passes, first to create the tuples, second to code the tuples
  - One pass, by pipelining tuples into a variable length coder

CSEP 590 - Lecture 4 - Autumn 2007

44

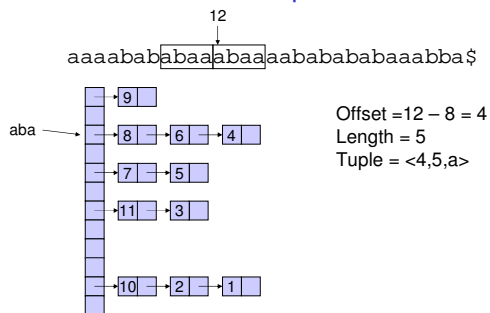
## Zip and Gzip

- Search Window
  - Search buffer 32KB
  - Look-ahead buffer 258 Bytes
- How to store such a large dictionary
  - Hash table that stores the starting positions for all three byte sequences.
  - Hash table uses chaining with newest entries at the beginning of the chain. Stale entries can be ignored.
- Second pass for Huffman coding of tuples.
- Coding done in blocks to avoid disk accesses.

CSEP 590 - Lecture 4 - Autumn 2007

45

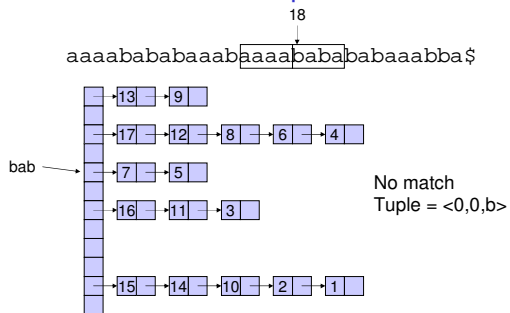
## Example



CSEP 590 - Lecture 4 - Autumn 2007

46

## Example



CSEP 590 - Lecture 4 - Autumn 2007

47

## Notes on LZ77

- Very popular especially in unix world
- Many variants and implementations
  - Zip, Gzip, PNG, PKZip, Lharc, ARJ
- Tends to work better than LZW
  - LZW has dictionary entries that are never used
  - LZW has past strings that are not in the dictionary
  - LZ77 has an implicit dictionary. Common tuples are coded with few bits.

CSEP 590 - Lecture 4 - Autumn 2007

48