

## Dictionaries for Data Compression

CSE 326  
Autumn 2005  
Lecture 19

## 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

Dictionary Data Compression - Lecture 19

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

Dictionary Data Compression - Lecture 19

3

## LZW Encoding Example (1)

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

Dictionary Data Compression - Lecture 19

4

## LZW Encoding Example (2)

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

Dictionary Data Compression - Lecture 19

5

## LZW Encoding Example (3)

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

Dictionary Data Compression - Lecture 19

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?						

### LZW Decoding Example (3a)

Dictionary

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

### LZW Decoding Example (3b)

Dictionary

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

### 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						

### 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?						

### 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						

### LZW Decoding Example (5b)

Dictionary

0	a	012436
1	b	a b ab aba ba
2	ab	
3	ba	
4	aba	
5	abab	
6	ba?	

### LZW Decoding Example (6a)

Dictionary

0	a	012436
1	b	a b ab aba ba b
2	ab	
3	ba	
4	aba	
5	abab	
6	bab	

### LZW Decoding Example (6b)

Dictionary

0	a	012436
1	b	a b ab aba ba bab
2	ab	
3	ba	
4	aba	
5	abab	
6	bab	
7	bab?	

### Decoding Exercise

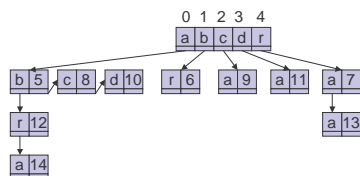
Base Dictionary

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

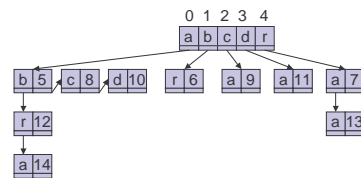
### Trie Data Structure for Encoder's Dictionary

• 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		

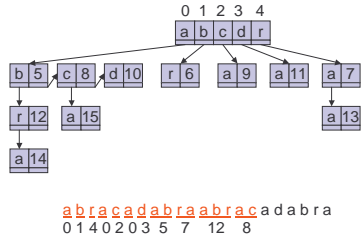


### Encoder Uses a Trie (1)



abra cadabra abra ca da bra  
0 1 4 0 2 0 3 5 7 12

## Encoder Uses a Trie (2)



Dictionary Data Compression - Lecture 19

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 ...  
abraca  
abraca  
abraca

Dictionary Data Compression - Lecture 19

26

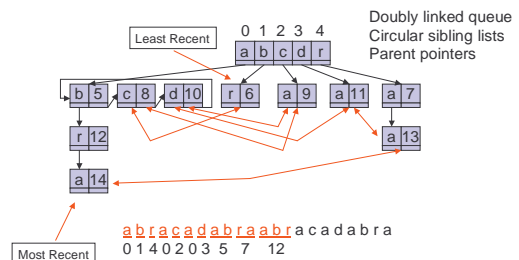
## Bounded Size Dictionary

- 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.

Dictionary Data Compression - Lecture 19

27

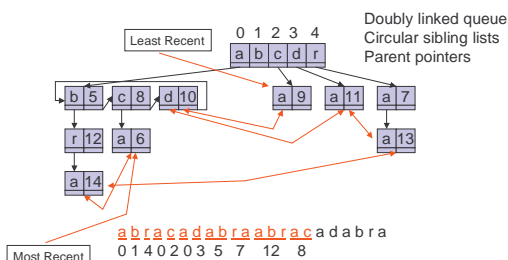
## Implementing the LRV Strategy



Dictionary Data Compression - Lecture 19

28

## Implementing the LRV Strategy



Dictionary Data Compression - Lecture 19

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 Data Compression - Lecture 19

30

## LZ77

- Ziv and Lempel, 1977
- Dictionary is implicit
- 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.

## 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>
```

## 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  $x$ .

## 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>
```

## 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>
```

## Surprise Code!

```
a babababababababababab$
<0,0,a>
a b abababababababababab$
<0,0,b>
a b abababababababababab$
<1,22,$>
```

## Surprise Decoding

<0,0,a><0,0,b><1,22,\$>

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

Dictionary Data Compression - Lecture 19

37

## Surprise Decoding

<0,0,a><0,0,b><1,22,\$>

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

Dictionary Data Compression - Lecture 19

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  $\langle j, k, x_{n+k+1} \rangle$

Dictionary Data Compression - Lecture 19

39

## In Class Exercise

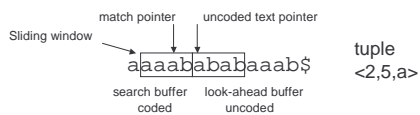
- Use Solution C to code the string  
– abaabaaabaaab\$

Dictionary Data Compression - Lecture 19

40

## Bounded Buffer – Sliding Window

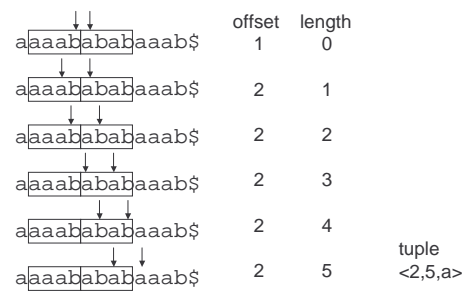
- We want the triples  $\langle j, k, x \rangle$  to be of bounded size. To achieve this we use bounded buffers.
  - Search buffer of size  $s$  is the symbols  $x_{n-s+1}\dots x_n$   
 $j$  is then the offset into the buffer.
  - Look-ahead buffer of size  $t$  is the symbols  $x_{n+1}\dots x_{n+t}$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



Dictionary Data Compression - Lecture 19

41

## Search in the Sliding Window



Dictionary Data Compression - Lecture 19

42

## Coding Example

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

	tuple
aaaa <b>bab</b> abaaab\$	<0, 0, a>
aaaab <b>bab</b> abaaab\$	<1, 3, b>
aaaa <b>bab</b> abaaab\$	<2, 5, a>
aaaabab <b>baa</b> ab\$	<4, 2, \$>

Dictionary Data Compression - Lecture 19

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

Dictionary Data Compression - Lecture 19

44

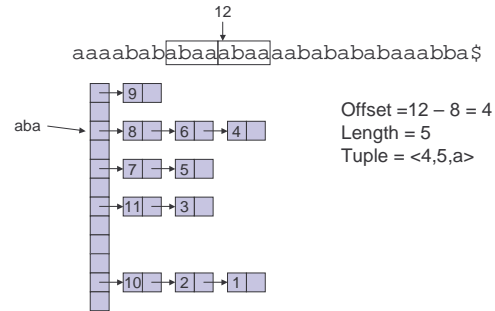
## Zip and Gzip

- 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.

Dictionary Data Compression - Lecture 19

45

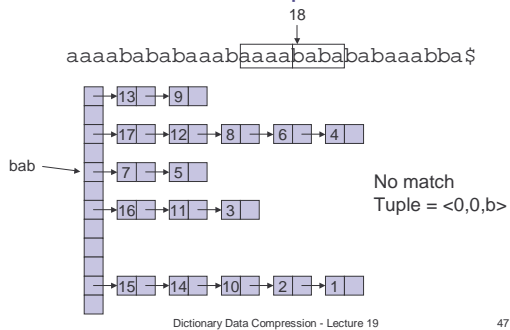
## Example



Dictionary Data Compression - Lecture 19

46

## Example



Dictionary Data Compression - Lecture 19

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.

Dictionary Data Compression - Lecture 19

48