

CSE 490 GZ
Introduction to Data Compression
Winter 2006

Dictionary Coding
LZ77

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.

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 babababababababababab....
<0,0,a>

a b abababababababababab....
<0,0,b>

a b aba babababababababab....
<1,2,a>

a b aba babab abababababab....
<2,4,b>

a b aba babab abababababa bab....
<1,10,a>

Surprise Code!

```

a bababababababababababab$
<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,$>   b
<23,0,$>   $
    
```

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,$>   b
<23,0,$>   $
    
```

Solution C

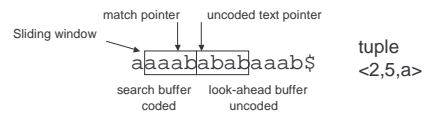
- The matching string can include part of itself!
- If $x_{n+1}x_{n+2}\dots x_{n+k}$ is a substring of $x_1x_2\dots x_n x_{n+1}x_{n+2}\dots x_{n+k}$ that begins at $j \leq n$ and $x_{n+1}x_{n+2}\dots x_{n+k}x_{n+k+1}$ is not then $x_{n+1}x_{n+2}\dots x_{n+k} x_{n+k+1}$ can be coded by $\langle j,k, x_{n+k+1}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

In Class Exercise

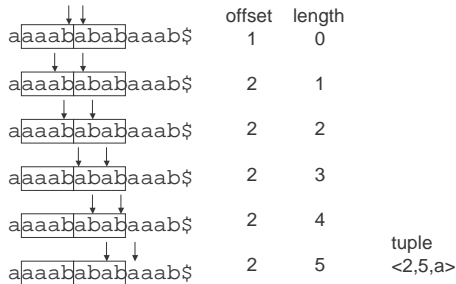
- Use Solution C to code the string
 - abaabaaabaab\$
 - aaaabaaabaab\$

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.



Search in the Sliding Window

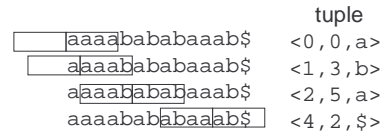


CSE 490g - Lecture 8 - Winter 2006

13

Coding Example

s = 4, t = 4, a = 3



CSE 490g - Lecture 8 - Winter 2006

14

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

CSE 490g - Lecture 8 - Winter 2006

15

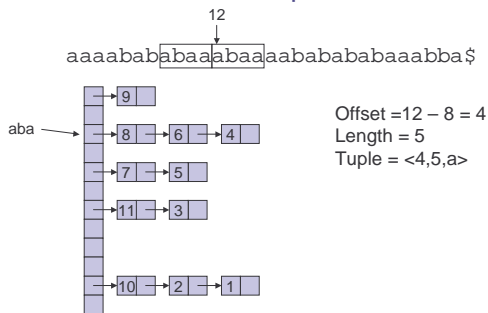
Zip and Gzip

- 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.

CSE 490g - Lecture 8 - Winter 2006

16

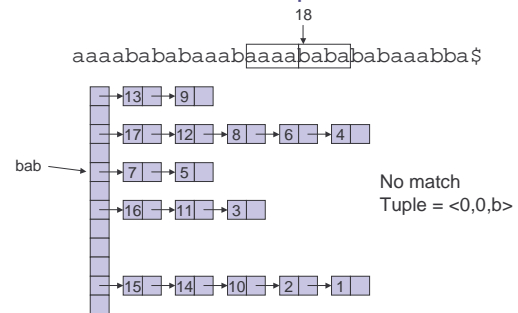
Example



CSE 490g - Lecture 8 - Winter 2006

17

Example



CSE 490g - Lecture 8 - Winter 2006

18

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.