

# 15-583: Algorithms in the Real World

Data Compression III  
 Lempel-Ziv algorithms  
 Burrows-Wheeler  
 Introduction to Lossy Compression

## Lempel-Ziv Algorithms

### LZ77 (Sliding Window)

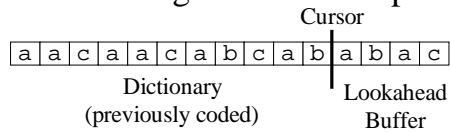
- Variants: LZSS (Lempel-Ziv-Storer-Szymanski)
- Applications: gzip, Squeeze, LHA, PKZIP, ZOO

### LZ78 (Dictionary Based)

- Variants: LZW (Lempel-Ziv-Welch), LZC (Lempel-Ziv-Compress)
- Applications: compress, GIF, CCITT (modems), ARC, PAK

Traditionally LZ77 was better but slower, but the gzip version is almost as fast as any LZ78.

## LZ77: Sliding Window Lempel-Ziv

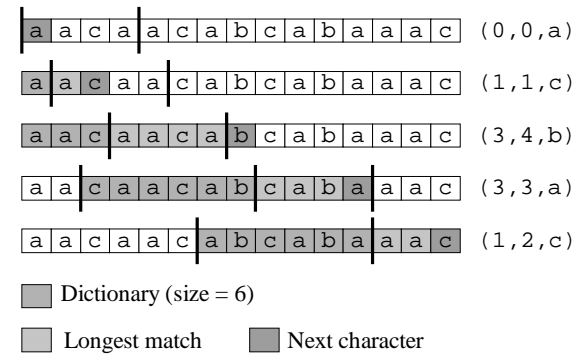


Dictionary and buffer “windows” are fixed length and slide with the cursor

On each step:

- Output (p,l,c)  
 p = relative position of the longest match in the dictionary  
 l = length of longest match  
 c = next char in buffer beyond longest match
- Advance window by l + 1

## LZ77: Example



## LZ77 Decoding

Decoder keeps same dictionary window as encoder.

- For each message it looks it up in the dictionary and inserts a copy

What if  $l > p$ ? (only part of the message is in the dictionary.)

- E.g. dict = abcd, codeword = ( 2 , 9 , e )
- Simply copy starting at the cursor  

```
for (i = 0; i < length; i++)
    out[cursor+i] = out[cursor-offset+i]
```
- Out = abcdcdcdcdcd

15-853

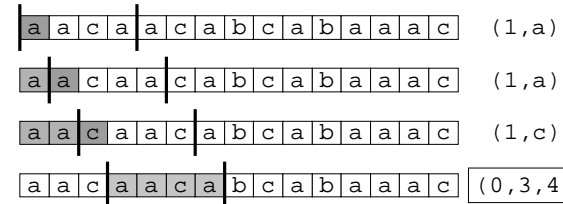
Page5

## LZ77 Optimizations used by gzip

LZSS: Output one of the following formats

(0, position, length) or (1, char)

Typically use the second format if length < 3.



15-853

Page6

## Optimizations used by gzip (cont.)

- Huffman code the positions, lengths and chars
- Non greedy: possibly use shorter match so that next match is better
- Use hash table to store dictionary.
  - Hash is based on strings of length 3.
  - Find the longest match within the correct hash bucket.
  - Limit on length of search.
  - Store within bucket in order of position

15-853

Page7

## Theory behind LZ77

Sliding Window LZ is Asymptotically Optimal  
 [Wyner-Ziv,94]

Will compress long enough strings to the source entropy as the window size goes to infinity.

$$H_n = \sum_{X \in A^n} p(X) \log \frac{1}{p(X)}$$

$$H = \lim_{n \rightarrow \infty} H_n$$

Uses logarithmic code for position.

Problem: “long enough” is really **really** long.

15-853

Page8

## LZ78: Dictionary Lempel-Ziv

Basic algorithm:

- Keep dictionary of words with integer *id* for each entry (e.g. keep it as a trie).
- Coding loop
  - find the longest match *S* in the dictionary
  - Output the entry *id* of the match and the next character past the match from the input (*id,c*)
  - Add the string *Sc* to the dictionary
- Decoding keeps same dictionary and looks up *ids*

15-853

Page9

## LZ78: Coding Example

	Output	Dict.
a a b a a c a b c a b c b	(0, a)	1 = a
a a b a a c a b c a b c b	(1, b)	2 = ab
a a b a a c a b c a b c b	(1, a)	3 = aa
a a b a a c a b c a b c b	(0, c)	4 = c
a a b a a c a b c a b c b	(2, c)	5 = abc
a a b a a c a b c a b c b	(5, b)	6 = abcb

15-853

Page10

## LZ78: Decoding Example

Input	Dict.
(0, a) a	1 = a
(1, b) a a b	2 = ab
(1, a) a a b a a	3 = aa
(0, c) a a b a a c	4 = c
(2, c) a a b a a c a b c	5 = abc
(5, b) a a b a a c a b c a b c b	6 = abcb

15-853

Page11

## LZW (Lempel-Ziv-Welch) [Welch84]

Don't send extra character *c*, but still add *Sc* to the dictionary.

The dictionary is initialized with byte values being the first 256 entries (e.g. a = 112, ascii), otherwise there is no way to start it up.

The decoder is one step behind the coder since it does not know *c*

- Only an issue for strings of the form *SSc* where *S[0] = c*, and these are handled specially

15-853

Page12

## LZW: Encoding Example

	Output	Dict.
a a b a a c a b c a b c b	112	256=aa
a a b a a c a b c a b c b	112	257=ab
a a b a a c a b c a b c b	113	258=ba
a a b a a c a b c a b c b	256	259=aac
a a b a a c a b c a b c b	114	260=ca
a a b a a c a b c a b c b	257	261=abc
a a b a a c a b c a b c b	260	262=cab

15-853

Page13

## LZ78 and LZW issues

What happens when the dictionary gets too large?

- Throw the dictionary away when it reaches a certain size (used in GIF)
- Throw the dictionary away when it is no longer effective at compressing (used in unix compress)
- Throw the least-recently-used (LRU) entry away when it reaches a certain size (used in BTLZ, the British Telecom standard)

15-853

Page14

## LZW: Decoding Example

Input		Dict
112	a a b a a c a b c a b c b	
112	a a b a a c a b c a b c b	256=aa
113	a a b a a c a b c a b c b	257=ab
256	a a b a a c a b c a b c b	258=ba
114	a a b a a c a b c a b c b	259=aac
257	a a b a a c a b c a b c b	260=ca
260	a a b a a c a b c a b c b	261=abc

15-853

Page15

## Lempel-Ziv Algorithms Summary

Both LZ77 and LZ78 and their variants keep a “dictionary” of recent strings that have been seen.

The differences are:

- How the dictionary is stored
- How it is extended
- How it is indexed
- How elements are removed

15-853

Page16

## Lempel-Ziv Algorithms Summary (II)

Adapt well to changes in the file (e.g. a Tar file with many file types within it).

Initial algorithms did not use probability coding and perform very poorly in terms of compression (e.g. 4.5 bits/char for English)

More modern versions (e.g. gzip) do use probability coding as “second pass” and compress much better.

Are becoming outdated, but ideas are used in many of the newer algorithms.

15-853

Page17

## Burrows -Wheeler

Currently best “balanced” algorithm for text

### Basic Idea:

- Sort the characters by their full context (typically done in blocks). This is called the **block sorting transform**.
- Use move-to-front encoding to encode the sorted characters.

The ingenious observation is that the decoder only needs the sorted characters and a pointer to the first character of the original sequence.

15-853


Page18

## Burrows Wheeler: Example

Let's encode:  $d_1e_2c_3o_4d_5e_6$

We've numbered the characters to distinguish them.

Context “wraps” around.

<u>Context</u>	<u>Char</u>		<u>Context</u>	<u>Output</u>
ecode <sub>6</sub>	d <sub>1</sub>	<b>Sort Context</b> 	dedec <sub>3</sub>	o <sub>4</sub>
coded <sub>1</sub>	e <sub>2</sub>		coded <sub>1</sub>	e <sub>2</sub>
odede <sub>2</sub>	c <sub>3</sub>		<u>decod<sub>5</sub></u>	<u>e<sub>6</sub></u>
dedec <sub>3</sub>	o <sub>4</sub>		odede <sub>2</sub>	c <sub>3</sub>
edeco <sub>4</sub>	d <sub>5</sub>		ecode <sub>6</sub>	d <sub>1</sub> ←
decod <sub>5</sub>	e <sub>6</sub>		edeco <sub>4</sub>	d <sub>5</sub>

15-853

Page19

## Burrows-Wheeler (Continued)

**Theorem:** After sorting, equal valued characters appear in the same order in the output as in the most significant position of the context.

**Proof:** Since the chars have equal value in the most-sig-position of the context, they will be ordered by the rest of the context, *i.e.* the previous chars. This is also the order of the output since it is sorted by the previous characters.

<u>Context</u>	<u>Output</u>
dedec <sub>3</sub>	o <sub>4</sub>
<b>coded<sub>1</sub></b>	e <sub>2</sub>
<b>decod<sub>5</sub></b>	e <sub>6</sub>
odede <sub>2</sub>	c <sub>3</sub>
<b>ecode<sub>6</sub></b>	d <sub>1</sub>
<b>edeco<sub>4</sub></b>	d <sub>5</sub>

15-853

Page20

## Burrows-Wheeler Decode

```

Function BW_Decode(In, Start, n)
  S = MoveToFrontDecode(In,n)
  R = Rank(S)
  j = Start
  for i=1 to n do
    Out[i] = S[j]
    j = R[i]

```

Rank gives position of each char in sorted order.

15-853

Page21

## Decode Example

<u>S</u>	<u>Sort(S)</u>	<u>Rank(S)</u>	<u>Out</u>
o <sub>4</sub>	c <sub>3</sub>	6	e <sub>6</sub> d <sub>1</sub> ←
e <sub>2</sub>	d <sub>1</sub>	4	d <sub>1</sub> e <sub>2</sub>
e <sub>6</sub>	d <sub>5</sub>	5	e <sub>2</sub> c <sub>3</sub>
c <sub>3</sub>	e <sub>2</sub>	1	c <sub>3</sub> o <sub>4</sub>
d <sub>1</sub>	e <sub>6</sub>	2 ←	o <sub>4</sub> d <sub>5</sub>
d <sub>5</sub>	o <sub>4</sub>	3	d <sub>5</sub> e <sub>6</sub>

15-853

Page22

## ACB (Associate Coder of Buyanovsky)

Currently the second best compression for text  
(BOA, based on PPM, is marginally better)

Keep dictionary sorted by context

- Find best match for context     Context Contents
- Find best match for contents     dec ode
- Code     d ecode
- Distance between matches     decod e
- Length of contents match     de code

Has aspects of Burrows-Wheeler,  
residual coding, and LZ77

deco de

15-853

Page23

## Other Ideas?

15-853

Page24

More Ideas?

15-853

Page25