

# Data Compression III

## LZW compression



Mark Matthews PhD

# Statistical methods

---

**Static model.** Same model for all texts.

- Fast.
- Not optimal: different texts have different statistical properties.
- Ex: ASCII, Morse code.

**Dynamic model.** Generate model based on text.

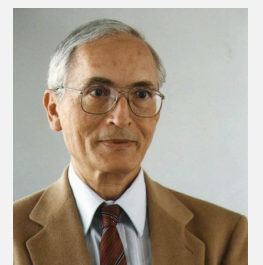
- Preliminary pass needed to generate model.
- Must transmit the model.
- Ex: Huffman code.

**Adaptive model.** Progressively learn and update model as you read text.

- More accurate modeling produces better compression.
- Decoding must start from beginning.
- Ex: LZW.



Abraham Lempel



Jacob Ziv

# LZW compression demo

input	A	B	B	R	C	A	C	A	R	D	A	R	B	E	R	A	B	R	A	B	R	A
matches	A	B		R	A	C		A	D		AB			RA			BR		ABR			A
value	41	42		52	41	43		41	44		81			83			82		88			41 80

LZW compression for A B R A C A D A B R A B R A B R A

key	value	key	value	key	value
:	:	AB	81	DA	87
A	41	BR	82	ABR	88
B	42	RA	83	RAB	89
C	43	AC	84	BRA	8A
D	44	CA	85	ABRA	8B
:	:	AD	86		

codeword table

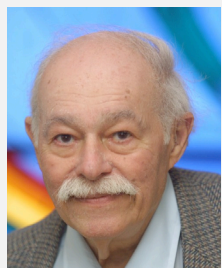
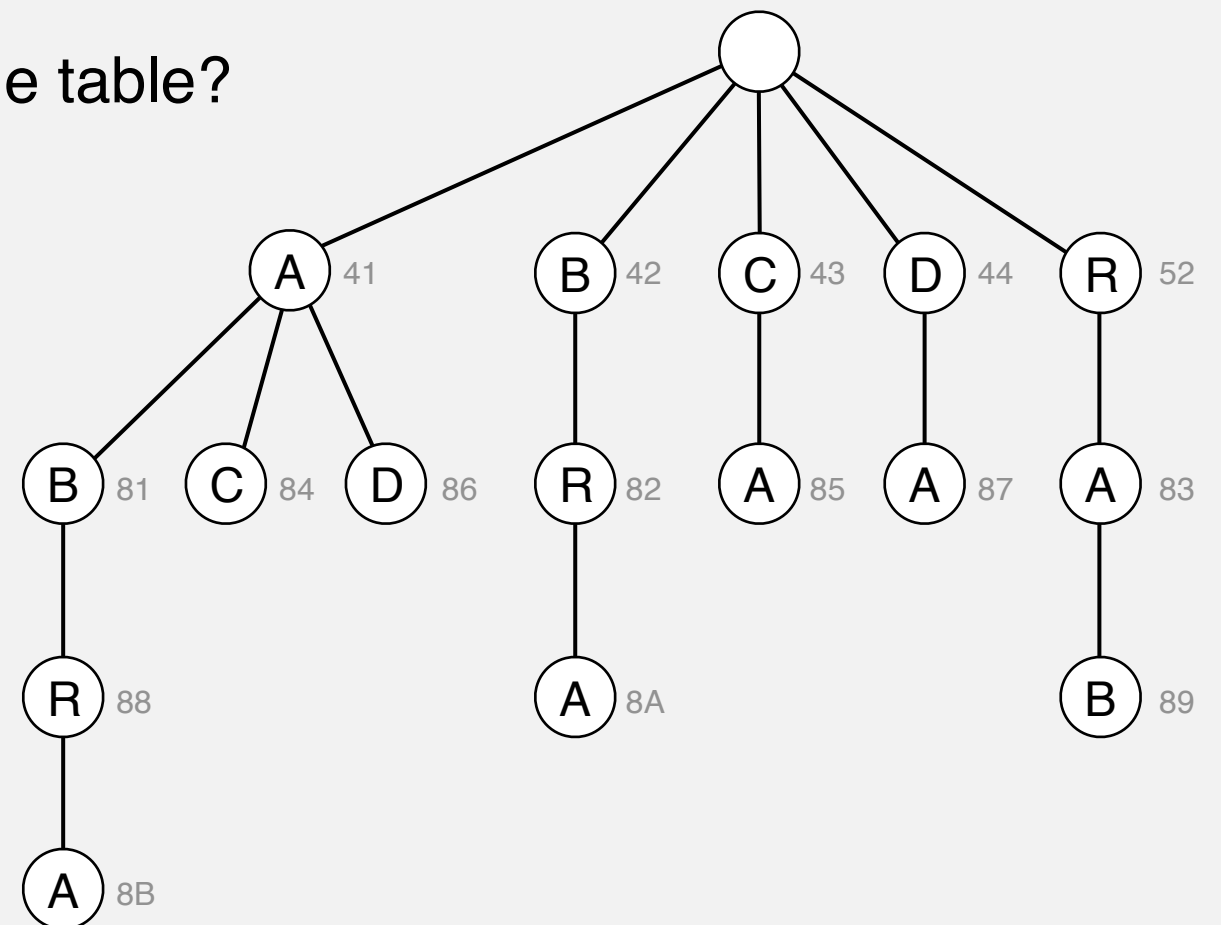
# Lempel-Ziv-Welch compression

## LZW compression.

- Create ST associating  $W$ -bit codewords with string keys.
- Initialize ST with codewords for single-char keys.
- Find longest string  $s$  in ST that is a prefix of unscanned part of input.
- Write the  $W$ -bit codeword associated with  $s$ .
- Add  $s + c$  to ST, where  $c$  is next char in the input.

Q. How to represent LZW compression code table?

A. A trie to support longest prefix match.



Abraham Lempel



Jacob Ziv

# LZW expansion demo

---

value    41   42   52   41   43   41   44   81    83    82    88       41   80

output   A     B     R     A     C     A     D     AB     RA     BR     ABR     A

LZW expansion for 41 42 52 41 43 41 44 81 83 82 88 41 80

key	value	key	value	key	value
⋮	⋮	81	AB	87	DA
41	A	82	BR	88	ABR
42	B	83	RA	89	RAB
43	C	84	AC	8A	BRA
44	D	85	CA	8B	ABRA
⋮	⋮	86	AD		

codeword table

# LZW expansion

---

## LZW expansion.

- Create ST associating string values with  $W$ -bit keys.
- Initialize ST to contain single-char values.
- Read a  $W$ -bit key.
- Find associated string value in ST and write it out.
- Update ST.

**Q.** How to represent LZW expansion code table?

**A.** An array of size  $2^W$ .

key	value
:	:
65	A
66	B
67	C
68	D
:	:
129	AB
130	BR
131	RA
132	AC
133	CA
134	AD
135	DA
136	ABR
137	RAB
138	BRA
139	ABRA
:	:

# LZW tricky case: compression

---

input	A	B	B		A	\	B	A	A	B	A
matches	A	B	AB					ABA			
value	41	42	81					83			80

LZW compression for ABABABA

key	value	key	value
⋮	⋮	AB	81
A	41	BA	82
B	42	ABA	83
C	43		
D	44		
⋮	⋮		

codeword table

# LZW tricky case: expansion

value	41	42	81	83	80
output	A	B	AB	ABA	

need to know which  
key has value 83  
before it is in ST!

LZW expansion for 41 42 81 83 80

key	value	key	value
⋮	⋮	81	AB
41	A	82	BA
42	B	83	ABA
43	C		
44	D		
⋮	⋮		

codeword table



# LZW implementation details

---

## How big to make ST?

- How long is message?
- Whole message similar model?
- [many other variations]

## What to do when ST fills up?

- Throw away and start over. [GIF]
- Throw away when not effective. [Unix compress]
- [many other variations]

## Why not put longer substrings in ST?

- [many variations have been developed]

# LZW in the real world

## Lempel-Ziv and friends.

- LZ77.
- LZ78.
- LZW.
- Deflate / zlib = LZ77 variant + Huffman.

LZ77 not patented  $\Rightarrow$  widely used in open source

LZW patent #4,558,302 expired in U.S. on June 20, 2003

### United States Patent [19] Welch

[11] Patent Number: 4,558,302

[45] Date of Patent: Dec. 10, 1985

[54] HIGH SPEED DATA COMPRESSION AND DECOMPRESSION APPARATUS AND METHOD

[75] Inventor: Terry A. Welch, Concord, Mass.

[73] Assignee: Sperry Corporation, New York, N.Y.

[21] Appl. No.: 505,638

[22] Filed: Jun. 20, 1983

[51] Int. Cl.<sup>4</sup> ..... G06F 5/00

[52] U.S. Cl. .... 340/347 DD; 235/310

[58] Field of Search ..... 340/347 DD; 235/310, 235/311; 364/200, 900

#### [56] References Cited

##### U.S. PATENT DOCUMENTS

4,464,650 8/1984 Eastman ..... 340/347 DD

##### OTHER PUBLICATIONS

Ziv, "IEEE Transactions on Information Theory", IT-24-5, Sep. 1977, pp. 530-537.

Ziv, "IEEE Transactions on Information Theory", IT-23-3, May 1977, pp. 337-343.

Primary Examiner—Charles D. Miller

Attorney, Agent, or Firm—Howard P. Terry; Albert B. Cooper

#### [57] ABSTRACT

A data compressor compresses an input stream of data character signals by storing in a string table strings of data character signals encountered in the input stream. The compressor searches the input stream to determine

the longest match to a stored string. Each stored string comprises a prefix string and an extension character where the extension character is the last character in the string and the prefix string comprises all but the extension character. Each string has a code signal associated therewith and a string is stored in the string table by, at least implicitly, storing the code signal for the string, the code signal for the string prefix and the extension character. When the longest match between the input data character stream and the stored strings is determined, the code signal for the longest match is transmitted as the compressed code signal for the encountered string of characters and an extension string is stored in the string table. The prefix of the extended string is the longest match and the extension character of the extended string is the next input data character signal following the longest match. Searching through the string table and entering extended strings therein is effected by a limited search hashing procedure. Decompression is effected by a decompressor that receives the compressed code signals and generates a string table similar to that constructed by the compressor to effect lookup of received code signals so as to recover the data character signals comprising a stored string. The decompressor string table is updated by storing a string having a prefix in accordance with a prior received code signal and an extension character in accordance with the first character of the currently recovered string.

181 Claims, 9 Drawing Figures



# LZW in the real world

---

## Lempel-Ziv and friends.

- LZ77.
- LZ78.
- LZW.
- Deflate / zlib = LZ77 variant + Huffman.



Unix compress, GIF, TIFF, V.42bis modem: LZW.

zip, 7zip, gzip, jar, png, pdf: deflate / zlib.

iPhone, Sony Playstation 3, Apache HTTP server: deflate / zlib.



# Lossless data compression benchmarks

---

year	scheme	bits / char
1967	ASCII	7
1950	Huffman	4.7
1977	LZ77	3.94
1984	LZMW	3.32
1987	LZH	3.3
1987	move-to-front	3.24
1987	LZB	3.18
1987	gzip	2.71
1988	PPMC	2.48
1994	SAKDC	2.47
1994	PPM	2.34
1995	Burrows-Wheeler	2.29
1997	BOA	1.99
1999	RK	1.89

data compression using Calgary corpus

# Data compression summary

---

## Lossless compression.

- Represent fixed-length symbols with variable-length codes. [Huffman]
- Represent variable-length symbols with fixed-length codes. [LZW]

## Lossy compression. [not covered in this course]

- JPEG, MPEG, MP3, ...
- FFT, wavelets, fractals, ...

## Theoretical limits on compression. Shannon entropy:

$$H(X) = - \sum_i^n p(x_i) \lg p(x_i)$$

## Practical compression. Use extra knowledge whenever possible.