

Data Compression

- ▶ fixed-length codes
- ▶ variable-length codes
- ▶ an application
- ▶ adaptive codes

References:

Algorithms 2nd edition, Chapter 22
<http://www.cs.princeton.edu/algs4/65compression>

Algorithms in Java, 4th Edition · Robert Sedgewick and Kevin Wayne · Copyright © 2008 · November 20, 2008 7:52:01 AM

Data compression

Compression reduces the size of a file:

- To save **space** when storing it.
- To save **time** when transmitting it.
- Most files have lots of redundancy.

Who needs compression?

- Moore's law: # transistors on a chip doubles every 18-24 months.
- Parkinson's law: data expands to fill space available.
- Text, images, sound, video, ...

"All of the books in the world contain no more information than is broadcast as video in a single large American city in a single year. Not all bits have equal value." — Carl Sagan

Basic concepts ancient (1950s), best technology recently developed.

2

Applications

Generic file compression.

- Files: GZIP, BZIP, BOA.
- Archivers: PKZIP.
- File systems: NTFS.

Multimedia.

- Images: GIF, JPEG.
- Sound: MP3.
- Video: MPEG, DivX™, HDTV.

Communication.

- ITU-T T4 Group 3 Fax.
- V.42bis modem.

Databases. Google.



3

Encoding and decoding

Message. Binary data M we want to compress.

← uses fewer bits (you hope)

Encoder. Generates a "compressed" representation $C(M)$.

Decoder. Reconstructs original message or some approximation M' .



Compression ratio. Bits in $C(M)$ / bits in M .

Lossless. $M = M'$, 50-75% or better compression ratio. ← this lecture

Ex. Natural language, source code, executables.

Lossy. $M \approx M'$, 10% or better compression ratio.

Ex. Images, sound, video.

4

Food for thought

Data compression has been omnipresent since antiquity:

- Number systems.
- Natural languages.
- Mathematical notation.

has played a central role in communications technology,

- Braille.
- Morse code.
- Telephone system.

and is part of modern life.

- MP3.
- MPEG.

Q. What role will it play in the future?

5

What data can be compressed?

US Patent 5,533,051 on "Methods for Data Compression", which is capable of compression all files.

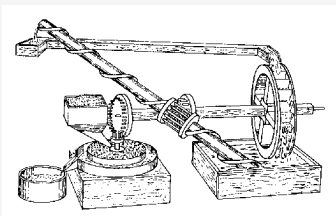
Slashdot reports of the Zero Space Tuner™ and BinaryAccelerator™.

“ZeoSync has announced a breakthrough in data compression that allows for 100:1 lossless compression of random data. If this is true, our bandwidth problems just got a lot smaller....”

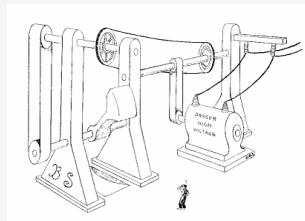
6

Perpetual motion machines

Universal data compression is the analog of perpetual motion.



Closed-cycle mill by Robert Fludd, 1618



Gravity engine by Bob Schadewald

Reference: Museum of Unworkable Devices by Donald E. Simanek
<http://www.lhup.edu/~dsimanek/museum/unwork.htm>

7

What data can be compressed?

Proposition. Impossible to losslessly compress all files.

Pf 1. [by counting]

- Consider all 1,000 bit messages.
- 2^{1000} possible messages.
- Only $2^{999} + 2^{998} + \dots + 1$ can be encoded with ≤ 999 bits.
- Only 1 in 2^{499} can be encoded with ≤ 500 bits!

Pf 2. [by contradiction]

- Given a file M , compress it to get a smaller file M_1 .
- Compress that file to get a still smaller file M_2 .
- Continue until reaching file size 0.
- Implication: all files can be compressed with 0 bits!

8

A difficult file to compress

one million pseudorandom characters (a-p)

```
fclkkaci fobjofmkgdcooicnfmcpojfccabckjamolnihkgobcjbngjiceeelpfgcjjihppene flhglfemdemgahlbpigmlmn  
efnhjelmgjncjcidlhkghoeninidmgnobkleglpnadanfbecoonbiehglmphnkamdfpacjmgogmcaabpcjcecpflbgamlidock  
lhfkmloljdnnoagihaiapaimcnlljngigpeanbmogjkccogpmkmoifioeikefjdbadgdcpehdpfjaeeapdjoefkpddeghidbg  
caiemajllhndigeihbeifemacfdnkhblgincpimindogimeeomgeljfgklkdnghafohongpbmlkapddhmpednckeaiejmekn  
meenjnenmnmnfebdbhpmigbjknjmobimamjjaaffhlhggallbjaijnebidpaieydgogochihodnlhahlhhojdfcanhadhgkfhah  
meaebccacgeogjgkcoapknlomfignamedmajinlompjoaifiaejbcjcdibpkofobmjiobbdhfilfajkhfmppcngdneei nfnfafe  
ladbhhi fechinknplnplamakpkhekogkjpdmjnbngklhibohdfeagmclllmdha fkl dmi ndbp lggbbe jcmh lkjocjjl cngckfp  
fakmpiaanfjdlleinii laenbnikgfnjfcopbhgkhdgmfoehfmbpiaignphogbkelphobonmfgpdmkfedfkchceeldkcoof  
aldinljggafaimaanelmfokokjekefkmegegjjfjcpjpnabldjoaafpbdafifgcoibbcmofbbgqimngfepkmbhbghlbdjngenl  
dhgnfbdcmjdmoflhocgfjoldfjpaokpepnde jmbiealkaofifekdjkgedglgbioacflfjafbaemgpjlagbdgilhcfcdamhmp  
pfoghjphlmheghechgdpkkljpdphfongnanmbmngpphncbieknjhilafkegboila jdpccodpeoddldjfcpi al oal feomj bphkm  
hndmcpkggeaohfmcnecgmibjka jcdcpjcpjgmihnhakihfgiaachfefpfnllcooiciepoapmdjniimfbolchikbkmhbkgconimk  
dchahcnhpfdkiaipikencegcjapkjkljgdlnmncpbakhjidapblodgeekjaoihbnbigmhboengpmedlio fgioof dophelapijce  
gejgldcfodikal ehbccbbcfakkbhmoobdmgdka fbbkjndioikfakjclbchambcpaepfeinmenmpoodadoecbgm fkekeabilaoe  
ogghoekamaibhjibe fmpopbhfbhfapjnodlofeihmjahmeipejlfhloefgmjhjn lomajakhhjpncomippeanbkkhekpcfgbgk  
mklipfbii kdkdcholofhelipbkbjmjfoemppcneae bklimb caddlmdcajpmhhae ddbfpjafcn dianlfcjmm bfnopccccodel dhm  
nbdjmeajmbocikggogjghlohblhgjkhkmlcolhkgjamfmochkchmiadggjhjehflchklfiackbecgjogpbkhlcmfhipflhmmi fppj  
mcoldbeghpcckhgmnahi jpanomnokldjcpbbpcpcjofngmbdcpseeiiclm bmfjkhlanckidhmbeanmlabnncppbhoafajjic  
nfeenppokml dho lndjapbfcabjblboeiaepfmmeoafeflmdcbaodgeahimogpcamnjloebp fmghogfckgmomecdipmodbcemp  
idfnlcggpbf foncajpncomal goiikeolmiglikjkolgol fkdgiijjiooikdihjbbfoioibakad jnedlodeeijkl iicnioima  
blfdpjiafc fineecbafaamhei pegegibioocmlmhjekifkfe fmdmdoaklni fdhckmbonbchfhclecjamjildonj jdpifngboj  
ianpljahpkindkdoanlldcbmlmhj fomifhmcicolj jhebidj dphdpeibfgdonj l jfgifimni ipogockpidamncpi pglafmlmoa  
cjabogbnplejnikdoefccdpfkomimfjgjelocdemnblimfbkfbkkelkpfchoekfofochbmi fleecbglnmfnfnfcjme fnihdco  
eieflemohlfdombdfedmbseebbalggfbajdamp lphdgi imehg lpi kbi pnkceekhilchhhfae fbbfdmjo jfhppongkfdmhpj  
cieofcnjgkpi bchiblfpnjlejkcoppbhophdghljlcookdoahfmlqblkliajbmnkfkookhllelhjhoigainamgabcabefmjndbfn  
ohkjphnkcbhccjgpbada koecbkjcaebanhnfhpnfkbfbfpo hmnkligpgfkjadomdjnhlnfal fpcmnooldjkeohlhkebif fba  
jjpcqlhllmemegncnkmkeoogilijmkmnlbkkabelmodcohppdakbelmljednmbfmcj de befni hne jmnoeeafldabjcgfoa  
ehldcmkbnafpcie fhlopicifadbpmpgngecjhefnkbjml odhelhi cnfoongnemddephokkdjafegnpjedakmbcpmckkbf  
feihpkajginfdolfnlgnade faml foodibhfkiaofeeppcjilndep leihpkkgkphbnkgjjiaolno lnbjpbj dcehgleckbhjila  
fccfipgebpcc...
```

9

A difficult file to compress

```
public class Rand  
{  
    public static void main(String[] args)  
    {  
        int N = Integer.parseInt(args[0]);  
        for (int i = 0; i < N; i++)  
        {  
            char c = 'a';  
            c += (char) (Math.random() * 16);  
            StdOut.print(c);  
        }  
    }  
}
```

231 bytes, but output is hard to compress
(assume random seed is fixed)

```
javac Rand.java  
java Rand 1000000 > temp.txt  
compress -c temp.txt > temp.z  
gzip -c temp.txt > temp.gz  
bzip2 -c temp.txt > temp.bz2  
  
ls -l  
    231 Rand.java  
1000000 temp.txt  
576861 temp.z  
570872 temp.gz  
499329 temp.bz2
```

resulting file sizes (bytes)

10

Rdenudancy in Enlgisih Inagugae

Q. How much redundancy is in the English language?

“... randomising letters in the middle of words [has] little or no effect on the ability of skilled readers to understand the text. This is easy to demntrasote. In a pubiltacion of New Scnieitst you could ramdinose all the letetrs, keipeng the first two and last two the same, and reibadaily would hadrly be aftcfeed. My ansaylis did not come to much beucase the thoery at the time was for shape and senqeuce retigcionon. Saberi's work sugesgts we may have some pofrweul parlrael prsooscers at work. The resaon for this is suerly that idnetiyfing coentnt by paarllel prseocsing speeds up regnicoiton. We only need the first and last two letetrs to spot chganes in menieng.” — Graham Rawlinson

A. Quite a bit.

11

- ▶ fixed-length codes
- ▶ variable length codes
- ▶ an application
- ▶ adaptive codes

12

Fixed-length coding

- Use same number of bits for each symbol.
- k -bit code supports 2^k different symbols.

Ex. 7-bit ASCII code.

char	decimal	codeword
NUL	0	0000000
...	...	
a	97	1100001
b	98	1100010
c	99	1100011
d	100	1100100
...	...	
DEL	127	1111111

a	b	r	a	c	a	d	a	b	r	a	!
1100001	1100010	1110010	1100001	1100011	1100001	1100100	1100001	1100010	1110010	1100001	1111111

12 symbols × 7 bits per symbol = 84 bits in code

13

Fixed-length coding

- Use same number of bits for each symbol.
- k -bit code supports 2^k different symbols.

Ex. 3-bit custom code.

char	codeword
a	000
b	001
c	010
d	011
r	100
!	111

a	b	r	a	c	a	d	a	b	r	a	!
000	001	100	000	010	000	011	000	001	100	000	111

12 symbols × 3 bits per symbol = 36 bits in code

14

Fixed-length coding: general scheme

- Count number of different symbols.
- $\sim \lg M$ bits suffice to support M different symbols.

Ex. Genomic sequences.

- 4 different nucleotides.
- 2 bits suffice.
- Amazing but true: initial databases in 1990s did not use such a code!

$\sim 2N$ bits to encode genome with N nucleotides

a	c	t	a	c	a	g	a	t	g	a	a
00	01	10	00	01	00	11	00	10	11	00	00

2-bit DNA code

char	codeword
a	00
c	01
t	10
g	11

Important detail. Decoder needs to know the code!

15

- › fixed-length codes
- › **variable-length codes**
- › an application
- › adaptive codes

16

Variable-length coding

Use different number of bits to encode different symbols.

Ex. Morse code: `•••---•••`

Issue. Ambiguity.

SOS ?
IAMIE ?
EEWNI ?
V7O ?

In practice. Use a medium gap to separate codewords.

Letters	Numbers
A	1
B	2
C	3
D	4
E	5
F	6
G	7
H	8
I	9
J	0
K	
L	
M	
N	
O	
P	
Q	
R	
S	
T	
U	
V	
W	
X	
Y	
Z	

17

Variable-length coding

Q. How do we avoid ambiguity?

A. Ensure that no codeword is a **prefix** of another.

Ex 1. Fixed-length code.

Ex 2. Append special stop symbol to each codeword.

Ex 3. Custom prefix-free code.

char	codeword
S	•••
E	•
I	••
V	•••-

← prefix of V
← prefix of I, S
← prefix of S

char	codeword
a	0
b	111
c	1011
d	100
r	110
!	1010

a	b	r	a	c	a	d	a	b	r	a	!				
0	1	1	1	1	0	0	1	0	1	1	0	1	0	1	0

28 bits to encode message (vs. 36 bits for fixed-length 3-bit code)

18

Prefix-free code: encoding and decoding

Q. How to represent a prefix-free code?

A. Binary trie.

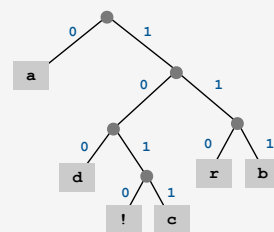
- Symbols are stored in leaves.
- Codeword is path to leaf.

Encoding.

- Method 1: start at leaf; follow path up to the root, and print bits in reverse order.
- Method 2: create ST of symbol-codeword pairs.

Decoding.

- Start at root of tree.
- Go left if bit is 0; go right if 1.
- If leaf node, print symbol and return to root.



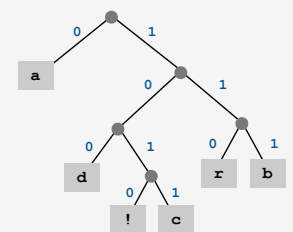
char	codeword
a	0
b	111
c	1011
d	100
r	110
!	1010

19

Representing the code

Q. How to transmit the trie?

A. Send preorder traversal of trie.



* used as sentinel for internal node

```
*a**d**!c*rb
```

← preorder traversal
← # chars to decode
← the message bits (pack 8 to the byte)

char	codeword
a	0
b	111
c	1011
d	100
r	110
!	1010

Note. If message is long, overhead of transmitting trie is small.

20

Prefix-free decoding: Java implementation

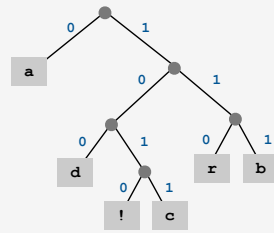
```

public class PrefixFreeDecoder
{
    private Node root = new Node();
    private class Node
    {
        private char ch;
        private Node left, right;
        public Node()
        {
            ch = StdIn.readChar();
            if (ch == '*')
            {
                left = new Node();
                right = new Node();
            }
        }

        public boolean isInternal()
        {
            return left != null && right != null;
        }

        public void decode()
        {
            /* See next slide. */
        }
    }
}

```



build binary trie from preorder traversal
*a**d*c*rb

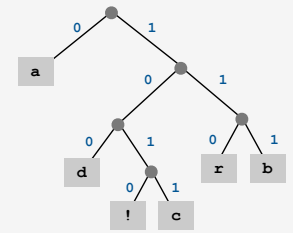
21

Prefix-free decoding: Java implementation (cont)

```

public void decode()
{
    int N = StdIn.readInt();
    for (int i = 0; i < N; i++)
    {
        Node x = root;
        while (x.isInternal())
        {
            char bit = StdIn.readChar();
            if (bit == '0') x = x.left;
            else if (bit == '1') x = x.right;
        }
        StdOut.print(x.ch);
    }
}

```



decode message abracadabra! from bits
12
0111110010110100011111001010

use bits, not chars in actual applications
(see next programming assignment)

22

Huffman coding

Q. What is the best variable-length code for a given message?

A. Huffman code.



David Huffman

To compute Huffman code:

- Count frequency p_s for each symbol s in message.
- Start with one node corresponding to each symbol s (with weight p_s).
- Repeat until single trie formed:
 - select two tries with min weight p_1 and p_2
 - merge into single trie with weight $p_1 + p_2$

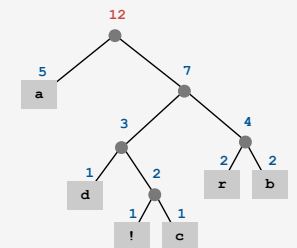
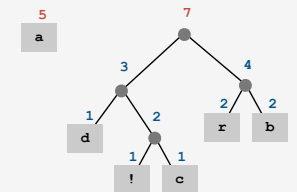
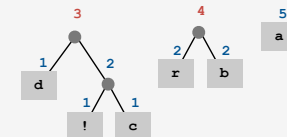
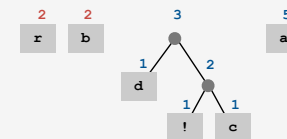
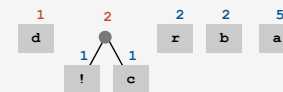
Applications. JPEG, MP3, MPEG, PKZIP, GZIP, ...

23

Huffman coding example

a b r a c a d a b r a !

! c d r b a



24

Huffman trie construction code

```

int[] freq = new int[R];
for (int i = 0; i < input.length(); i++)
    freq[input.charAt(i)]++;

MinPQ<Node> pq = new MinPQ<Node>();
for (int r = 0; r < R; r++)
    if (freq[r] > 0)
        pq.insert(new Node((char) r, freq[r], null, null));

while (pq.size() > 1)
{
    Node x = pq.delMin();
    Node y = pq.delMin();
    Node parent = new Node('*', x.freq + y.freq, x, y);
    pq.insert(parent);
}
root = pq.delMin();
    
```

Annotations:

- tabulate frequencies (points to the first for loop)
- initialize PQ (points to the second for loop)
- merge tries (points to the while loop)
- internal node marker (points to the '*' character in the parent node creation)
- total frequency (points to the `x.freq + y.freq` expression)
- two subtrees (points to the `x, y` arguments in the parent node creation)

25

Huffman encoding summary

Proposition. [Huffman 1950s] Huffman coding is an optimal prefix-free code.
Pf. See textbook.

no prefix-free code uses fewer bits

Implementation.

- Pass 1: tabulate symbol frequencies and build trie.
- Pass 2: encode file by traversing trie or lookup table.

Running time. Using a binary heap $\Rightarrow O(N + R \log R)$.

input symbols (points to N)
 distinct symbols (points to R)

Q. Can we do better? [stay tuned]

26

- › fixed-length codes
- › variable-length codes
- › **an application**
- › adaptive codes

27

An application: compress a bitmap

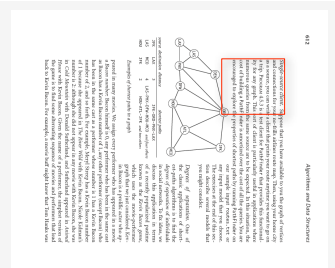
Typical black-and-white-scanned image.

- 300 pixels/inch.
- 8.5-by-11 inches.
- $300 \times 8.5 \times 300 \times 11 = 8.415$ million bits.

Observation. Bits are mostly white.

Typical amount of text on a page.

40 lines * 75 chars per line = 3000 chars.



as a source, you can write a tip. PROGRAM 4.5.3 is a program that computes the Huffman tree for any graph. This program is a good example of how numerous queries for the cost of building a Path are encouraged to explore

b

28

Natural encoding of a bitmap

Natural encoding. $(19 \times 51) + 6 = 975$ bits.

one bit per pixel to encode number of characters per line

```
000000000000000000000000000000000001111111111111111111100000000
000000000000000000000000000000000001111111111111111111110000000
0000000000000000000000000000000000011111111111111111111110000
000000000000000000000000000000000001111111111111111111111110000
000000000000000000000000000000000001111111111111111111111111000
000000000000000000000000000000000001111111111111111111111111100
000000000000000000000000000000000001111111111111111111111111110
000000000000000000000000000000000001111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0110000000000000000000000000000000000000000000000000000000000011
```

19-by-51 raster of letter 'q' lying on its side

Run-length encoding of a bitmap

Run-length encoding. $(63 \times 6) + 6 = 384$ bits.

63 6-bit run lengths

```
0000000000000000000000000000000000011111111111111111111100000000
000000000000000000000000000000000001111111111111111111110000000
0000000000000000000000000000000000011111111111111111111110000
000000000000000000000000000000000001111111111111111111111110000
000000000000000000000000000000000001111111111111111111111111000
000000000000000000000000000000000001111111111111111111111111100
000000000000000000000000000000000001111111111111111111111111110
000000000000000000000000000000000001111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0000000000000000000000000000000000011111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0111111111111111111111111111111111111111111111111111111111111
0110000000000000000000000000000000000000000000000000000000000011
```

51
28 14 9
26 18 7
23 24 4
22 26 3
20 30 1
19 7 18 7
19 5 22 5
19 3 26 3
19 3 26 3
19 3 26 3
19 3 26 3
20 4 23 3 1
22 3 20 3 3
1 50
1 50
1 50
1 50
1 50
1 2 46 2

19-by-51 raster of letter 'q' lying on its side

run-length encoding

Run-length encoding

Goal. Exploit long runs of repeated characters.

Bitmaps Runs alternate between 0 and 1; just output run lengths.

Q. How to encode run lengths? (!)

```
0 0 1 1 1 1 1 1 0 0 1 1 1 1 1 0 0 0 0 0 20-bit message
```

```
2 6 2 5 5 run lengths
```

```
0 1 0 1 1 0 0 1 0 1 0 1 1 0 1 1 run lengths encoded using 3-bit codewords (15 bits)
```

- 001: 1
- 010: 2
- 101: 3
- 100: 4
- 101: 5
- 110: 6
- 111: 7

Note. Runs are long in typical applications (such as black-and-white bitmaps).

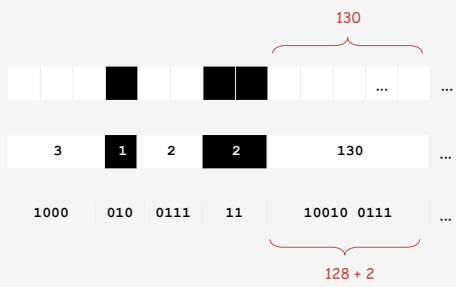
RLE and Huffman codes in the wild

ITU-T T4 Group 3 Fax. [for black-and-white bitmap images]

- Up to 1728 pixels per line.
- Typically mostly white.

Step 1. Use run-length encoding.
Step 2. Encode run lengths using two Huffman codes.

one for white and one for black
based on statistics from huge number of faxes



run	white	black
0	00110101	0000110111
1	000111	010
2	0111	11
3	1000	10
...
63	00110100	000001100111
64+	11011	0000001111
128+	10010	000011001000
...
1728+	010011011	0000001100101

Black and white bitmap compression: another approach

Fax machine (~1980).

- Slow scanner produces lines in sequential order.
- Compress to save time (reduce number of bits to send).

Electronic documents (~2000).

- High-resolution scanners produce huge files.
- Compress to save space (reduce number of bits to save).

Idea.

- use OCR to get back to ASCII (!)
- use Huffman on ASCII string (!)

Bottom line. Any extra information about file can yield dramatic gains.

33

- › fixed-length codes
- › variable-length codes
- › an application
- › **adaptive codes**

34

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.

35

Lempel-Ziv-Welch encoding

LZW encoding.

- Create ST associating a fixed-length codeword with some previous substring.
- When input matches string in ST, output associated codeword.
- Length of strings in ST grows, hence compression.



To send (encode) message M.

- Find longest string s in ST that is a prefix of unsent part of M.
- Send codeword associated with s .
- Add $s \cdot x$ to ST, where x is next char in M.

Ex. ST: a, aa, ab, aba, abb, abaa, **abaab**, abaaa.

- unsent part of M: **abaab**abbb...
- $s = \text{abaab}$, $x = a$.
- Output integer associated with s ; insert **abaaba** into ST.

36

LZW encoding example

input	codeword	add to ST
a	97	ab
b	98	br
r	114	ra
a	97	ac
c	99	ca
a	97	ad
d	100	da
a		
b	128	abr
r		
a	130	rac
c		
a	132	cad
d		
a	134	dab
b		
r	129	bra
a	97	

abracadabracadabra

ASCII

key	value
NUL	0
...	...
a	97
b	98
c	99
d	100
e	101
f	102
g	103
...	...
r	114
...	...
DEL	127

ST

key	value
ab	128
br	129
ra	130
ac	131
ca	132
ad	133
da	134
abr	135
rac	136
cad	137
dab	138
bra	139
...	...

To send (encode) M:

- Find longest string s in ST that is a prefix of unsent part of M .
- Send integer associated with s .
- Add $s \cdot x$ to ST, where x is next char in M .

37

LZW encoding example

input	codeword	add to ST
a	97	ab
b	98	br
r	114	ra
a	97	ac
c	99	ca
a	97	ad
d	100	da
a		
b	128	abr
r		
a	130	rac
c		
a	132	cad
d		
a	134	dab
b		
r	129	bra
a	97	

abracadabracadabra

Message.

- 7-bit ASCII.
- 19 chars.
- 133 bits.

Encoding.

- 8-bit codewords.
- 14 codewords.
- 112 bits.

Key point. Don't need to send ST (!)

Q. Which ST?

38

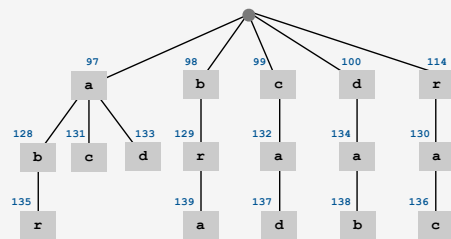
LZW encode ST implementation

Q. How to do longest prefix match?

A. Use a **trie** for the ST.

Encode.

- Lookup string suffix in trie.
- Output ST index at bottom.
- Add new node to bottom of trie.



ASCII

key	value
NUL	0
...	...
a	97
b	98
c	99
d	100
e	101
f	102
g	103
...	...
r	114
...	...
DEL	127

ST

key	value
ab	128
br	129
ra	130
ac	131
ca	132
ad	133
da	134
abr	135
rac	136
cad	137
dab	138
bra	139
...	...

Note. All prefixes of strings are also in ST.

39

Lempel-Ziv-Welch decoding

LZW decoding. Decode by rebuilding ST from code.

To decode received message to M.

- Let s be ST entry associated with received integer.
- Add s to M .
- Add $p \cdot x$ to ST, where x is first char in s , p is previous value of s .

Ex. ST: a, aa, ab, aba, abb, abaa, **abaab**, abaaa.

- undecoded part of M : **abaab**abb...
- $s = \text{abaab}$, $x = a$.
- Output integer associated with s ; insert **abaaba** into ST.

40

LZW decoding example

role of keys and values switched

codeword	output	add to ST
97	a	
98	b	ab
114	r	br
97	a	ra
99	c	ac
97	a	ca
100	d	ad
128	a	da
	b	
130	r	abr
	a	
132	c	rac
	a	
134	d	cad
	a	
129	b	dab
	r	
97	a	bra

key	value
0	NUL
...	...
97	a
98	b
99	c
100	d
101	e
102	f
103	g
...	...
114	r
...	...
127	DEL

value	key
128	ab
129	br
130	ra
131	ac
132	ca
133	ad
134	da
135	abr
136	rac
137	cad
138	dab
139	bra
...	...

use an array to implement ST

To decode received message to M:

- Let s be ST entry associated with received integer.
- Add s to M.
- Add $p \cdot x$ to ST, where x is first char in s and p is previous value of s .

41

LZW decoding example (tricky situation)

input	code	add to ST
a	97	ab
b	98	ba
a		
b	128	aba
a		
b		
a	130	abab
b		
STOP	255	

abababab

key	value
128	ab
129	ba
130	aba
131	abab
...	
255	

codeword	output	add to ST
97	a	
98	b	ab
128	a	ba
	b	
130	a	aba
	b	
	a	
98	b	
255	STOP	

needed before aba added to ST!

To send (encode) M:

- Find longest string s in ST that is a prefix of unsent part of M.
- Send integer associated with s .
- Add $s \cdot x$ to ST, where x is next char in M.

To decode received message to M:

- Let s be ST entry for received integer.
- Add s to M.
- Add $p \cdot x$ to ST where x is first char in s and p is previous value of s .

42

LZW implementation details

How big to make ST?

- How long is message?
- Whole message similar model?
- [many variations have been developed]

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]

43

LZW in the real world

Lempel-Ziv and friends.

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

LZ77 not patented \Rightarrow widely used in open source
LZW patent #4,558,302 expired in US on June 20, 2003
some versions copyrighted

PNG: LZ77.

Winzip, gzip, jar: deflate.

Unix compress: LZW.

Pkzip: LZW + Shannon-Fano.

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

Google: zlib which is based on deflate.

never expands a file

44

Lossless data compression benchmarks

year	scheme	bits / char
1967	ASCII	7.00
1950	Huffman	4.70
1977	LZ77	3.94
1984	LZMW	3.32
1987	LZH	3.30
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

← next programming assignment

data compression using Calgary corpus

45

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.

Practical compression. Use extra knowledge whenever possible.

46