

Data Compression

Reference: Chapter 22, *Algorithms in C, 2nd Edition*, Robert Sedgewick.
 Reference: *Introduction to Data Compression*, Guy Blelloch.

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.

Applications of Data Compression

Generic file compression.

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



Multimedia.

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



Communication.

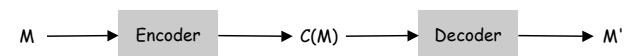
- ITU-T T4 Group 3 Fax.
- V.42bis modem.

Databases. Google.



Encoding and Decoding

Message. Binary data M we want to compress. ↙ hopefully uses fewer bits
Encode. Generate a "compressed" representation $C(M)$.
Decode. Reconstruct original message or some approximation M' .



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

Lossless. $M = M'$, 50-75% or lower.

Ex. Natural language, source code, executables.

Lossy. $M \approx M'$, 10% or lower.

Ex. Images, sound, video.

Fixed Length Coding

Fixed length encoding.

- Use same number of bits for each symbol.
- N symbols \Rightarrow $\lceil \log N \rceil$ bits per symbol.

7-bit ASCII encoding

char	dec	encoding
NUL	0	0000000
...
a	97	1100001
b	98	1100010
c	99	1100011
d	100	1100100
...
~	126	1111110
DEL	127	1111111

a	b	r	a	c	a	d	a	b	r	a
1100001	1100010	1110010	1100001	1100011	1100001	1110100	1100001	1100010	1110010	1100001

$7 \times 11 = 77$ bits

Fixed Length Coding

Fixed length encoding.

- Use same number of bits for each symbol.
- N symbols \Rightarrow $\lceil \log N \rceil$ bits per symbol.

3-bit abracadabra encoding

char	encoding
a	000
b	001
c	010
d	011
r	100

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

$3 \times 11 = 33$ bits

9

10

Variable Length Encoding

Variable-length encoding. Use different number of bits to encode different characters.

Ex. Morse code.

Ambiguity. $\dots - - - \dots$

- SOS
- IAMIE
- EEWNI
- T7O

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

Variable Length Encoding

Variable-length encoding. Use different number of bits to encode different characters.

Q. How do we avoid ambiguity?

- Append special stop symbol to each codeword.
- Ensure no encoding is a **prefix** of another.

101 is a prefix of 1011

char	encoding
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	0	1	1	1	1	0	0	1	0	1	0

variable length coding: 28 bits

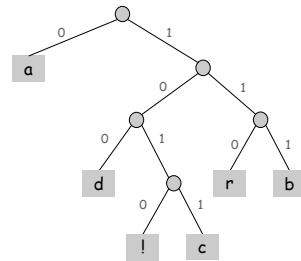
11

12

Prefix-Free Code: Encoding and Decoding

How to represent? Use a binary trie.

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



Encoding.

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

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

Decoding.

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

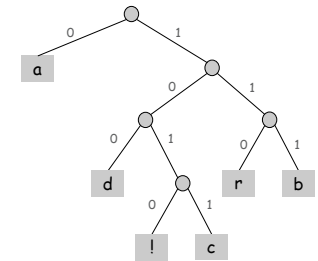
13

How to Transmit the Trie

How to transmit the trie?

- Send preorder traversal of trie.
 - we use * as sentinel for internal nodes
 - what if there is no sentinel?
- Send number of characters to decode.
- Send bits (packed 8 to the byte).

```
*a**d*!c*rb
12
0111110010110100011111001010
```



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

- If message is long, overhead of sending trie is small.

14

Prefix-Free Decoding Implementation

```
public class HuffmanDecoder {
    private Node root = new Node();

    private class Node {
        char ch;
        Node left, right;

        Node() {
            ch = StdIn.readChar();
            if (ch == '*') {
                left = new Node();
                right = new Node();
            }
        }

        boolean isInternal() { }
    }

    // build tree from preorder traversal
    // *a**d*!c*rb
}
```

15

Prefix-Free Decoding Implementation

```
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;
        }
        System.out.print(x.ch);
    }
}

// use bits in real applications instead of chars
12
0111110010110100011111001010
```

16

Huffman Coding

- Q. How to construct a good prefix-free code?
 A. Huffman code. [David Huffman 1950]



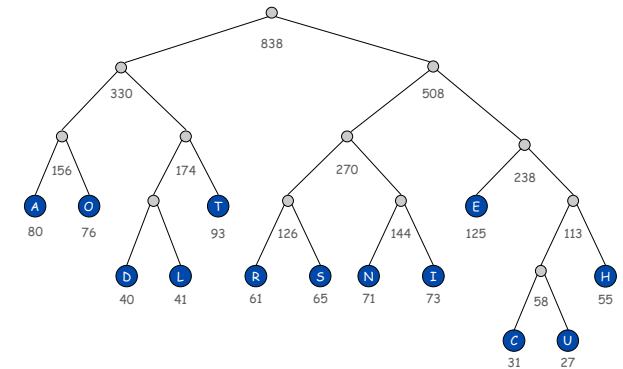
To compute Huffman code:

- Count frequencies p_s for each symbol s in message.
- Start with a forest of trees, each consisting of a single node corresponding to each symbol s with weight p_s .
- Repeat:
 - select two trees with min weight p_1 and p_2
 - merge into single tree with weight $p_1 + p_2$

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

Huffman Coding Example

Char	Freq	Huff
E	125	110
T	93	011
A	80	000
O	76	001
I	73	1011
N	71	1010
S	65	1001
R	61	1000
H	55	1111
L	41	0101
D	40	0100
C	31	11100
U	27	11101
Total	838	3.62



17

18

Huffman Tree Implementation

```
private class Node implements Comparable<Node> {
    char ch;
    int freq;
    Node left, right;

    // constructor
    Node(char ch, int freq, Node left, Node right) { ... }

    // print preorder traversal
    void preorder() {
        System.out.print(ch);
        if (isInternal()) {
            left.preorder();
            right.preorder();
        }
    }

    // compare by frequency
    int compareTo(Node b) { return freq - b.freq; }
    ...
}
```

19

Huffman Encoding Implementation

```
public HuffmanEncoder(String input) {
    private static int SYMBOLS = 128; // alphabet size (ASCII)
    private Node root;

    // tabulate frequencies
    int[] freq = new int[SYMBOLS];
    for (int i = 0; i < input.length(); i++)
        freq[input.charAt(i)]++;

    // initialize priority queue with singleton elements
    MinPQ<Node> pq = new MinPQ<Node>();
    for (int i = 0; i < SYMBOLS; i++)
        if (freq[i] > 0)
            pq.insert(new Node((char) i, freq[i], null, null));

    // repeatedly merge two smallest trees
    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();
}
```

20

Huffman Encoding

Theorem. [Huffman] Huffman coding is optimal prefix-free code.

Corollary. "Greed is good."

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
- Use heap for delete min and insert.
- $O(M + N \log N)$.

PQ implementation important if each symbol represents one English word or k-gram

M = file size
N = # distinct symbols

Difficulties.

- Have to transmit encoding (trie).
- Not optimal (unless block size grows to infinity)!

21

What Data Can be Compressed?

Theorem. Impossible to losslessly compress **all** files.

Pf.

- 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 even be encoded with ≤ 500 bits!

24

A Difficult File To Compress

One million pseudo-random characters (a - p)

```
fclkkaci fobjofmkgdcooicnfmcpcejfecabckjamolnhihkgobcjbngjiceelpfgeji ihppenefllhglfemdemahlpbi
ggmlmnefnhjelmgjncjoidlkhglhceinidmngnobkeglpnadanfbecoonbiehglmpnhkkamdfpacjmgjmaabpce
cplfbyamlidcekhlhfkmioljdnoaagiheiapaimlcnl1jni gypaanbmogkccoopmkmoifioeikefjidbadgdcphndpfj
aeapdjeofklpdeghi dbgcaiemajllhndnidei hbebi femacfadknhlbgincpmindogimeogeomlej f jgklkdgnhafoho
npjbmkapddhmpedncakejebmeknmeejnmenbmnfefdhhpmigbbjknjmobimamjjaaaffhlhiggaljbaijnebidpaeigd
gogchihodnlhahlhhoojdfacnhadhgkfahmeaebccacgeojgikcoapknlomfignanedmajinlompjoaifiaejbcjcdibp
kofcbmjiobbpdhfilfajkhfmpcngdneeinpnfaeladbbhi fechinkndpnlamackphekokigpddmmjnbngkllhibohdf
eagmc1llmdhafkldmimdbp1ggbbekcmhlkjocj1lcngeckfpakmpiaaanfddjllleinilaenbnikgfnjfcopbhgkhdg
mfpoehfmbkbiaignphogbke1phobonmfghpdmkfedkfcchceelkcofaldinljogafimaanelmfkocjkefkbmeqcgj
ifjcpjppnabdjoaafpbdafi fgcoibbcsofbbgigmnggefpmhbhghlbdjngeni dhgnfdlcmjdmoflhcogfjoldfjpaok
epndejmmbieaalkao fi fekdjkgedglgbi oacflfjla fbaeamgjl agbdgjl h efdca mhfmpfgohjphlmhegjechgpkklj
pndphfennanbmngpphncbkiaeknjhila fkegboila jdpocodpeodldjfcpi aloal feomjpbknhmpdmcpkgaeohfdm
cnegmibjka jcdcpjcpjgmihnhakihfgiia chfepfnlcooiepoapmdjniimfbolchkbkmbhkgconinkdchahcnhap
fdkiapikencegcjapkjkljgd1mgnrcpbakhjiidapblcgeeskkjaoihbnbigmhboengmediofgioofdphe1apijceajj
gcldcfodikal ehbcccbbcfakklmooobdmgdgkafbbkjndoi k f a k j c l b c h a m b c p a e f i n m e n p o o d a d o e c b g m f k k e a b i
laeogggchoekamaibhjibe f m o p p b h f b h f a p j n o d l o f e i h m j a h m e i p e j l f h l o e f g m j h j n l o m a p j a k h h j p n c o m i p p e a n b i k
k h e k p c f g b g k m l i p f i i k d k d c b o l o f h e l i p b k b j m j f o e m p c c n e a e b k l i b m c a d d l m j d c a j p m h a e e d b b f j a f c n d i a n l f c j
m m b f n p c d c c c o d e l d h m n b d j m e a j m b o c l k g g o j g h l o h l b h g j k h k m c l o h k g j a m f m c h k c h m i a d j g j h e f l c b k l f i f a c k b e c j
j o g p b k h l c m f n i p f l h m m i f p j m c o l d b e g h p c e k h g m m a h j p a b n o m n o k l d j c p p b c p c j o f n g m b d c p e e e i i c l m b m f j k h l
a n c k i d h m b e a n l a b n c n c p b h o a f a j j i c n f e e n p p o e k m l d d h o l n b d j a p b f c a j b l b o o i a e p f m m e o a e f d l m d c b a o d g e a h i m c
g p c a m m j l j o e b p f m g h o g f c k g m o m e c d j p m o d b c e m p i d f n l c g g p b f f o n c a j p n c o m a l g o i k e o l m i g l i i k j o l g o l f k d g i i j j
i o o l c k d h j b o f i o o l b a k a d j n e d l o d e a i j k l i i c n i o m a b l f d p j a f c f i n e e c h a f a a m h e i p e w e g j i b o o c m l a h j e k f i f
e f f m d h o a k l n i f d h c k m b o n b c h f h h c l e c j a m j i l d o n j d p i f n g b o j i a n p l j a h k i n d k d o a n l d e b a l m h j f o m i f m m c i k o l
j j h e b i d j p h d e p i b f g d o n j l f g i f i m i i p o g o c k p i d a m n k c p i p g l a f m l o a c j i b o g n b l e j n i k d o e f c d p f o m k i m f f g j
g i e l o c d e m b l i m f m b k f b h k e l k p f o e o k f o f o c h m i f l e e c b g l m f b n f n c j m e f n i h d c o e i f l i e m n o h f c d m b f e a b d m e a b
b a l g g f b a j d a m p l p h d g i m e h g l p i k b i p n k e e c k h i l c h h f a e a f b b f d m c j o j f p p o n g l k f d m h j p c i e o f c n j g k p i b c b i l f p
n j l e j k p p b h o h d g h l j l c o k h d o a h f m l g l b d k l i a j b m n k f c o k l h l e l h j h o i g i n a i m g c a b c f e a m j d n b f h o k j p h n l c b h c
j p g b a d a k o e c b j c a e b b a n h n f h n p k f b f p o h m n k l i g p g f k j a d o m d j j h l n f a i l f p c m n o l o l d j e k e o l h d k e b i f f e b a j j p c l j
h l l m e m e g c n k m k e o o g i l i j m m k o m l b k k a b e l m o d c o h d p p d a k b e l m l e j d n m b f m c j d e b e f n j i h n e j m m o g e e a f l d a b j o g f o
a e h l d c m k b n a f p c i e f h l o p i c i f a d b p p m f n g e c j h e f n k b j m l i d o d e h i c n f o n g e m d d e p c h k o d j a f e g p p l e d a k m b p
c m k c k h b f e i h p k a j g i n f n d o l f n l g n a d e f a m l f o o d i b h f k i a o f e e g p c j i l n d e p i e i h k p k g k p h n k g g j a o l n o l b j p o b j d
c e h g l e c k b h j l a f o c f i p e b c . . . .
```

25

A Difficult File To Compress

```
public class Rand {
    public static void main(String[] args) {
        for (int i = 0; i < 1000000; i++) {
            char c = 'a';
            c += (char) (Math.random() * 16);
            System.out.print(c);
        }
    }
}
```

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

```
% javac Rand.java
% java Rand > 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)

26

Intrinsic difficulty of compression.

- Short program generates large data file.
- Optimal compression algorithm has to discover program!
- Undecidable problem.

Q. So how do we know if our algorithm is doing well?

A. Want **lower bound** on # bits required by **any** compression scheme.

How compression algorithms work?

- Exploit biases of input messages.
- White patches occur in typical images.
- Word `Princeton` occurs more frequently than `Yale`.

Compression is all about probability.

- Formulate probabilistic model to predict symbols.
 - simple: character counts, repeated strings
 - complex: models of a human face
- Use model to encode message.
- Use same model to decode message.

Ex. Order 0 Markov model: each symbol s generated independently at random, with fixed probability $p(s)$.

Entropy

Entropy and Compression

Entropy. [Shannon 1948] $H(S) = - \sum_{s \in S} p(s) \log_2 p(s)$

- Information content of symbol s is proportional to $-\log_2 p(s)$.
- Weighted average of information content over all symbols.
- Interface between coding and model.

	p(a)	p(b)	H(S)
Model 1	1/2	1/2	1
Model 2	0.900	0.100	0.469
Model 3	0.990	0.010	0.0808
Model 4	1	0	0

	p(a)	p(b)	p(c)	p(d)	p(e)	p(f)	H(S)
Fair die	1/6	1/6	1/6	1/6	1/6	1/6	2.585

Theorem. [Shannon, 1948] If data source is an order 0 Markov model, **any** compression scheme must use $\geq H(S)$ bits per symbol on average.

- Cornerstone result of information theory.
- Ex: to transmit results of fair die, need ≥ 2.58 bits per roll.

Theorem. [Huffman, 1952] If data source is an order 0 Markov model, Huffman code uses $\leq H(S) + 1$ bits per symbol on average.

Q. Is there any hope of doing better than Huffman coding?

- A. Yes. Huffman wastes up to 1 bit per symbol.
- if $H(S)$ is close to 0, this matters
 - can do better with "arithmetic coding"
- A. Yes. Source may not be order 0 Markov model.

Entropy of the English Language

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 demonstrate. In a publication of New Scientist you could randomise all the letters, keeping the first two and last two the same, and readability would hardly be affected. My analysis did not come to much because the theory at the time was for shape and sentence recognition. Saberi's work suggests we may have some powerful parallel processors at work. The reason for this is surely that identifying content by parallel processing speeds up recognition. We only need the first and last two letters to spot changes in meaning."

A. Quite a bit.

31

Lossless Compression Ratio for Calgary Corpus

Year	Scheme	Bits / char
----	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
1988	SAKDC	2.47
1994	PPM	2.34
1995	Burrows-Wheeler	2.29
1997	BOA	1.99
1999	RK	1.89

Entropy	Bits/char
Char by char	4.5
8 chars at a time	2.4
Asymptotic	1.3

33

Entropy of the English Language

Q. How much information is in each character of English language?

Q. How can we measure it?

model = English text

A. [Shannon's 1951 experiment]

- Asked humans to predict next character given previous text.
- The number of guesses required for right answer:

# of guesses	1	2	3	4	5	≥ 6
Fraction	0.79	0.08	0.03	0.02	0.02	0.05

- Shannon's estimate = 0.6 - 1.3.

32

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.

34

LZW Algorithm

Lempel-Ziv-Welch. [variant of LZ78]



- Associate an integer with each **useful** string.
- When input matches string in ST, output associated integer.

Encoding.

- Find longest string s in ST that is a prefix of remaining part of string to compress.
- Output integer associated with s .
- Add $s \cdot x$ to dictionary, where x is next char in string to compress.

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

- String to be compressed: abaababb...
- $s = \text{abaab}$, $x = a$.
- Output integer associated with s ; insert abaaba into ST.

35

LZW Example

Input	Send
SEND	256
i	105
t	116
t	116
y	121
_	32
b	98
i	
t	258
t	
y	260
_	
b	262
i	
t	258
_	
b	
i	266
n	110
STOP	257

Dictionary

Index	Word	Index	Word
0		258	it
...		259	tt
32	_	260	ty
...		261	y_
97	a	262	_b
98	b	263	bi
...		264	itt
122	z	265	ty_
...		266	_bi
256	SEND	267	it_
257	STOP	268	_bin

36

LZW Implementation

Implementation.

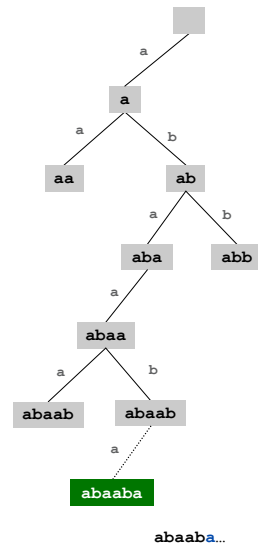
- Use trie to create symbol table on-the-fly.
- Note that prefix of every word is also in ST.

Encode.

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

Decode.

- Lookup index in array
- Output string
- Insert string + next letter.



37

LZW Encoder: Java Implementation

```
public class LZWEncoder {
    public static void main(String[] args) {
        String text = StdIn.readAll();
        StringST<Integer> st = new StringST<Integer>();
        int i;
        for (i = 0; i < 256; i++) {
            String s = Character.toString((char) i);
            st.put(s, i);
        }
        while (text.length() > 0) {
            String s = st.prefix(text);
            System.out.println(st.get(s));
            int length = s.length();
            if (length < text.length())
                st.put(text.substring(0, length + 1), i++);
            text = text.substring(length);
        }
    }
}
```

Annotations in the code:

- longest prefix match (pointing to `st.prefix(text)`)
- in real applications, encode integers in binary (pointing to `st.get(s)`)
- not the most efficient way (pointing to the `st.put` call)

38

LZW Decoder: Java Implementation

```
public class LZWDecoder {
    public static void main(String[] args) {
        ST<Integer, String> st = new ST<Integer, String>();
        int i;
        for (i = 0; i < 256; i++) {
            String s = Character.toString((char) i);
            st.put(i, s);
        }

        int code = StdIn.readInt();
        String prev = st.get(code);
        System.out.print(prev);

        while (!StdIn.isEmpty()) {
            code = StdIn.readInt();
            String s = st.get(code);
            if (i == code) s = prev + prev.charAt(0);
            System.out.print(s);
            st.put(i++, prev + s.charAt(0));
            prev = s;
        }
    }
}
```

in real applications,
integers will be encoded in binary

special case, e.g., for
"ababababab"

39

LZW Implementation Details

What to do when ST gets too large?

- Throw away and start over. GIF
- Throw away when not effective. Unix compress

40

LZW in the Real World

Lempel-Ziv and friends.

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

LZ77 not patented ⇒ 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

41

Summary

Lossless compression.

- Simple approaches. [RLE]
- 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, SVD, ...

Limits on compression. Shannon entropy.

Theoretical limits closely match what we can achieve in practice!

42